Springer Series in Reliability Engineering

Kjell Hausken Jun Zhuang *Editors*

Game Theoretic Analysis of Congestion, Safety and Security

Networks, Air Traffic and Emergency Departments



Springer Series in Reliability Engineering

Series editor

Hoang Pham, Piscataway, USA

More information about this series at http://www.springer.com/series/6917

Kjell Hausken · Jun Zhuang Editors

Game Theoretic Analysis of Congestion, Safety and Security

Networks, Air Traffic and Emergency Departments



Editors Kjell Hausken Faculty of Sciences University of Stavanger Stavanger Norway

Jun Zhuang Department of Industrial and Systems Engineering University at Buffalo, State University of New York Buffalo, NY USA

ISSN 1614-7839 ISSN 2196-999X (electronic) Springer Series in Reliability Engineering ISBN 978-3-319-13008-8 ISBN 978-3-319-13009-5 (eBook) DOI 10.1007/978-3-319-13009-5

Library of Congress Control Number: 2014956209

Springer Cham Heidelberg New York Dordrecht London

© Springer International Publishing Switzerland 2015

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

The publisher, the authors and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, express or implied, with respect to the material contained herein or for any errors or omissions that may have been made.

Printed on acid-free paper

Springer International Publishing AG Switzerland is part of Springer Science+Business Media (www.springer.com)

Preface

Game theoretic analysis of congestion, safety, and security is an interdisciplinary undertaking. Many researchers working on congestion have not extensively considered safety/security, and vice versa. However, significant interactions exist between the two research areas, which motivated this book. This book is intended to establish a new and enhanced current state of affairs within this topic, illustrate linkages between research approaches, and lay the foundation for subsequent research. Congestion (excessive crowding) is defined broadly to include all kinds of flows; e.g., road/sea/air traffic, people, data, information, water, electricity, and organisms. The book considers systems where congestion occurs, systems which may be in parallel, series, interlinked, or interdependent, with flows one way or both ways. Congestion models exist in abundance. The book makes ground by introducing game theory and safety/security. For the analysis to be game theoretic, at least two players must be present. For example, in Wang and Zhuang [1] one approver and a population of normal and adversary travelers are considered. Similarly, in Bier and Hausken [2], one defender and one attacker are considered, in addition to drivers who choose the more time-efficient of two arcs of different lengths. Multiple players can be adversaries with different concerns regarding system reliability; e.g., one or several terrorists, a government, various local or regional government agencies, companies, or others with stakes for or against system reliability. Governments, companies, and authorities may have tools to handle congestion, as well as ensure safety/security against various threats. The players may have a variety of individual concerns which may or may not be consistent with system safety or security. Much of the congestion literature is not game-theoretic, and does not extensively consider safety or security. Also, most game-theoretic analysis does not account for congestion. Volume 2 consists of nine chapters.

In "A Game Theory-Based Hybrid Medium Access Control Protocol for Congestion Control in Wireless Sensor Networks," Raja Periyasamy and Dananjayan Perumal consider a game-theoretical hybrid medium access control (GH-MAC) protocol in the face of wireless sensor networks congestion scenarios. GH-MAC is integrated with a game-based energy efficient time division multiple access (G-ETDMA) protocol for intra-cluster communication, as well as a game theory-based nanoMAC (G-nanoMAC) protocol for inter-cluster communication between the head nodes. The evaluation results show that GH-MAC would help reducing the energy consumption of the nodes in wireless sensor networks and increase the lifetime of the sensor networks.

In "Cooperative Games Among Densely Deployed WLAN Access Points," Josephina Antoniou, Vicky Papadopoulou-Lesta, Lavy Libman, and Andreas Pitsillides develop a method to mitigate the interference caused by individual wireless access points (AP) that share the same channels and cooperate by serving each other's clients. A small-size graphical game is used, where the underlying graph is a clique with heterogeneous edge weights. The equilibrium conditions are provided for APs to jointly serve each other's clients and achieve the maximum benefit for their clients.

In "Simulating a Multi-Stage Screening Network: A Queuing Theory and Game Theory Application," Xiaowen Wang, Cen Song and Jun Zhuang use simulations to study the game-theoretical screening strategies for the approver to balance congestion and security for a multi-stage security screening network facing strategic applicants. The Arena simulation software is used to build the screening system with three major components: arrival process, screening process, and departure process. A Matlab graphic user interface (GUI) is used to collect user inputs, then export data for Arena simulation through Excel, and finally export simulation from the results of the Arena to Matlab for analysis and visualization.

In "A Leader–Follower Game on Congestion Management in Power Systems," Mohammad Reza Salehizadeh, Ashkan Rahimi-Kian, and Kjell Hausken model a leader–follower game on congestion management. An operator first chooses multiple possible strategies which are announced to the generators (power generation companies). Each strategy consists of emission penalty factors for each power generation company, renewable power resources to reduce emission, and maximum prices generation companies can bid for selling electrical power. Second, each generator bids one price at which it is willing to sell its power, and a Nash-Supply Function equilibrium game is solved. Third, the operator performs congestion management and congestion-driven attributes and emission are obtained. Fourth, the operator's preferred strategy is selected using TOPSIS, a technique for order preference.

In "Network Throughput and Reliability: Preventing Hazards and Attacks Through Gaming—Part I: Modeling," Yupo Chan studies network throughput and reliability for preventing hazards and attacks by modeling a stochastic network, characterized by arcs and nodes that can fail unexpectedly. The author identifies not only tactics to prevent disruptions caused by natural/technological hazards and hostile tampering, but also strategies for defense budget allocations to balance maintaining critical infrastructure and defending against adversarial attacks.

In "Network Throughput and Reliability: Preventing Hazards and Attacks Through Gaming—Part II: A Research Agenda," following up on the modeling world in "Network Throughput and Reliability: Preventing Hazards and Attacks Through Gaming—Part I: Modeling," Yupo Chan provides a research agenda for pre-disaster mitigation/maintenance for critical civil infrastructure networks. Based on research to date, conditions under which a Pareto Nash equilibrium exists to prevent both hazards and attacks are provided. The author concludes with general principles of how to best defend networks of specific types against intelligent attacks, including both computational and behavioral considerations.

In "The Price of Airline Frequency Competition," Vikrant Vaze and Cynthia Barnhart present a game-theoretic model of airline frequency competition. Frequency competition affects congestion and each airline's market share and profit. Congestion affects and is affected by runway, taxiway, and airborne safety considerations. Myopic learning dynamics are analyzed and equilibria are determined for the two-player game and the N-player symmetric game. The price of airline frequency competition in the form of increased congestion and decreased profits is analyzed.

In "A Simulation Game Application for Improving the United States' Next Generation Air Transportation System NextGen," Ersin Ancel and Adrian Gheorghe develop a simulation game to evaluate the United States' next generation air transportation system NextGen. Such systems are characterized by congestion, energy shortage, delays, etc., which depend on technology, politics, environmental concerns, and interaction between players such as the government, the Federal Aviation Administration, military stakeholders, airlines, airport operators, and the general public.

In "A Congestion Game Framework for Emergency Department Overcrowding," Elizabeth Verheggen analyzes emergency department overcrowding as an El Farol Bar Game, illustrating the Tragedy of the Commons. Ambulance diversion, extensive wait times, and patient elopements cause overutilization and inefficient load balancing. Agent-based simulations reveal no statistically significant difference between two games and empirical observations, suggesting that a bar may be a good metaphor to understand congestion.

References

- 1. Wang X, Zhuang J (2011) Balancing congestion and security in the presence of strategic applicants with private information. Eur J Oper Res 212(1):100–111
- Bier V, Hausken K (2013) Defending and attacking networks subject to traffic congestion. Reliab Eng Sys Saf 112:214–224

Acknowledgments

The co-editor Dr. Jun Zhuang's co-editing effort for this book was partially supported by the United States Department of Homeland Security (DHS) through the National Center for Risk and Economic Analysis of Terrorism Events (CREATE) under award number 2010-ST-061-RE0001, and by the United States National Science Foundation (NSF) award numbers 1334930 and 1200899. However, any opinions, findings, and conclusions or recommendations in this document are those of the authors and do not necessarily reflect views of the DHS, CREATE, or NSF.

We thank the following scholars for reviewing chapters in this book:

Hamed Ahmadi Giuseppe Attanasi Tsz-Chiu Au Nick Bedard Geertje Bekebrede Ulrich Berger Steve Boyles Joel Bremson Marco Campenni Paola Cappanera Massimiliano Caramia David Carfi Edward Cartwright Damien Challet Andrew M. Colman Roberto Cominetti Xeni Dassiou Patrick deLamirande Wilfredo Yushimito Del Valle Loukas Dimitriou Pietro Dindo Anthony Downward

Ruud Egging Sjur Didrik Flåm Piero Fraternali Bernard Gendron David Gillen Richard J. Hamilton Fei He Hai-Jun Huang Miguel Jaller Eduard Jorswieck Max Klimm Maarten Kroesen Soundar Kumara Odd Larsen Elizabeth Lazzara Churlzu Lim Yingyan Lou Patrick Maillé Igor Stefan Mayer Sue McNeil Sebastiaan Meijer Federico Milano Michel Moreaux Stefano Moretti Renita Murimi Luis Gustavo Nardin Howard Ovens Kyoung Eun Park Vineet Madasseri Payyappalli Dananjayan Perumal Andreas Pitsillides Edward A. Ramoska Valentina Rotondi Walid Saad Javier Contreras Sanz Joe Scanlon Xiaojun Shan Cen Song Katerina Stanková Bruno Tuffin Vijay Venu Vadlamudi Sally Van Siclen Jonas Christoffer Villumsen Andreas Vlachos Xiaofang Wang

Qing Wang Jie Xu Lei Yang Elena Yudovina Jing Zhang Bo Zou

Contents

A Game Theory-Based Hybrid Medium Access Control Protocol for Congestion Control in Wireless Sensor Networks Raja Periyasamy and Dananjayan Perumal	1
Cooperative Games Among Densely Deployed WLAN Access Points Josephina Antoniou, Vicky Papadopoulou-Lesta, Lavy Libman and Andreas Pitsillides	27
Simulating a Multi-Stage Screening Network: A Queueing Theory and Game Theory Application	55
A Leader–Follower Game on Congestion Management in Power Systems	81
Network Throughput and Reliability: Preventing Hazards and Attacks Through Gaming—Part I: Modeling	113
Network Throughput and Reliability: Preventing Hazards and Attacks Through Gaming—Part 2: A Research Agenda Yupo Chan	141
The Price of Airline Frequency Competition	173

A Simulation Game Application for Improving the United States' Next Generation Air Transportation System NextGen	219
Ersin Ancel and Adrian Gheorghe	
A Congestion Game Framework for Emergency Department Overcrowding	255
Author Index	297
Subject Index	309

A Game Theory-Based Hybrid Medium Access Control Protocol for Congestion Control in Wireless Sensor Networks

Raja Periyasamy and Dananjayan Perumal

Abstract Network congestion occurs when the offered traffic load exceeds the available capacity of the network. In Wireless Sensor Networks (WSNs), congestion causes packet loss and degrades the overall channel quality, which leads to excessive energy consumption of the node to retrieve the packets. Therefore, congestion control (CC) is necessary to overcome the above-mentioned shortcomings to enhance the lifetime of the network. In this chapter, a game theory-based Hybrid Medium Access Control (GH-MAC) protocol is suggested to reduce the energy consumption of the nodes. GH-MAC is combined with Game-based energy-Efficient Time Division Multiple Access (G-ETDMA) protocol for intracluster communication between the cluster members to head node whereas Game theory-based nano-MAC (G-nanoMAC) protocol is used for intercluster communication between the head nodes. Performance of GH-MAC protocol is evaluated in terms of energy consumption, delay, and it is compared with conventional MAC schemes. The result thus obtained using GH-MAC protocol shows that the energy consumption is enormously reduced and thereby the lifetime of the sensor network is enhanced.

1 Introduction

The Wireless Sensor Network (WSN) is comprised of tiny embedded devices termed as "motes or nodes" that have inbuilt features for sensing, processing, and communicating information over wireless channels. A WSN consists of one or more sinks and a large number of sensor nodes scattered in a sensing area. A sensor node can sense the physical phenomenon, does the process of "raw" information,

D. Perumal (🖂)

R. Periyasamy

Department of ECE, Sri Manakula Vinayagar Engineering College, Pondicherry, India e-mail: rajashruthy@gmail.com

Department of ECE, Pondicherry Engineering College, Pondicherry, India e-mail: pdananjayan@pec.edu

[©] Springer International Publishing Switzerland 2015

K. Hausken and J. Zhuang (eds.), *Game Theoretic Analysis of Congestion, Safety and Security*, Springer Series in Reliability Engineering, DOI 10.1007/978-3-319-13009-5_1



Fig. 1 Upstream traffic model [1]

share the processed information with neighboring nodes, and reports information to the sink. The downstream traffic from the sink to the sensor node is usually a oneto-many multicast. The upstream traffic from sensor nodes to the sink is many-toone communication as shown in Fig. 1.

Due to the convergent nature of upstream traffic, congestion appears in the upstream direction. Congestion refers to a network state where a node or link carries much datum. It may deteriorate the quality of network service; frame or data packet loss and blocks the new connections were identified. In a congested network, response time slows down the network throughput. The fundamental reason for the occurrence of congestion is that the system resources cannot meet with the demand of the data transmission. The increase of system resources, such as bandwidth of the wireless channel, cache buffer size, and processing capability of the nodes, can be used to tackle congestion. Therefore, the approach that relies on increasing resources will inevitably cause the wastage of resources.

1.1 Types of Congestion

Two types of congestion could occur in WSNs as shown in Fig. 2. The first type named *node-level congestion* is common in conventional networks. It is caused by an overflow of buffer in the node and can result in packet loss, and increased queuing delay. Packet loss, in turn, can lead to retransmission and therefore



Fig. 2 a Node-level congestion and b link congestion

consumes additional energy. When wireless channels are shared by several nodes in WSNs using Carrier Sense Multiple Access (CSMA) protocols, collisions could occur, due to multiple active sensor nodes try to seize the channel at the same time. This can be referred as link-level congestion. Link-level congestion increases packet service time, and decreases both link utilization and overall throughput, and wastes energy at the sensor nodes. Both node-level and link-level congestions have direct impact on energy efficiency and quality of service (OoS) [1–3].

1.2 Congestion Control

Congestion must be efficiently controlled. The efficiency of CC protocol depends on how much it can achieve the following objectives:

- i. Energy efficiency requires to be improved in order to extend the lifetime of the system. Therefore, CC protocols need to avoid or reduce packet loss due to buffer overflow, and remain lesser control overhead that consumes less energy ii. Fairness needs to be guaranteed, so that each node can achieve fair throughput.

There are two general approaches to control congestion: network resource management and traffic control. The first approach tries to increase the network resources to mitigate congestion when it occurs. In wireless network, power control and multiple radio interfaces can be used to increase the bandwidth and weaken congestion. With this approach, it is necessary to guarantee precise and exact network resource adjustment in order to avoid over provided resources or underprovided resources. However, this is a hard task in wireless environment. Unlike the approaches based on network resource management, traffic control implies to control the congestion through adjusting the traffic rate at source nodes or intermediate nodes. When exact adjustment of network resource becomes difficult, this approach is helpful to save the network resources feasibly and efficiency. Most of the existing CC protocols belong to network resource management [4].

2 Energy Aware MAC Operation

The nodes are expected to operate on battery power for several years. Therefore, the energy conservation scheme plays a vital role in determining the lifetime of the network. Commercial standards such as IEEE 802.11 are not suitable for WSNs, because the nodes listen all the times and they cannot be used for multihop purposes [5]. The fairness of the WSN is defined as the number of unique packets received by a sink node from each source node. Most of the existing work guarantees simple fairness in every sensor node which obtains the same throughput to the sink. In fact, sensor nodes might be either outfitted with different sensors or geographically deployed in different place and they may have different importance or priority needed to gain different throughput. Hence the weighted fairness is required. The MAC protocol plays an important role in CC because the forwarding rate of packet depends on the MAC protocol. The protocols used in WSN are classified based on the location, data centric, hierarchical, and QoS, etc. These protocols are based on the minimum energy concept, clustering concept, and spanning tree and so on. Due to the limitation of battery power, we have to design the energy-efficient protocols for WSN.

2.1 Literature Survey

The MAC protocol, in addition to controlling medium access, can be designed to reduce the energy consumption of the radio in WSNs. Idle listening is often the largest source of energy waste, and duty cycling mechanism (i.e., sensors' radio gong to sleeping state) is considered as one of the necessary techniques to reduce energy consumption in WSN MAC protocols. An energy-efficient MAC protocol for WSN is called Sensor MAC (SMAC) [6]. SMAC is a MAC protocol specifically designed for wireless sensor networks (WSNs), which has better performance than IEEE 802.11 by setting duty cycles for each node which governs the nodes ON and OFF time. During communication process collision, overhearing, and control packet overhead will lead to energy wastage. MAC protocols have been developed to assist each node to decide when and how to access the channel. In SMAC, the duty cycle is low, but not adaptive, whereas in Timeout MAC (TMAC) the duty cycle is adaptable depending on the traffic. However, this protocol fails in heavy traffic. In Pattern MAC (PMAC) protocol, different bit patterns are used to sleep awake time schedule of nodes [7]. In PMAC, along with the pattern generated, the schedule for the duty cycle is also set. This has proven to be more energy efficient.

Some existing data dissemination schemes [8-10] can be configured or modified to be responsive to congestion. For example, directed diffusion can use in-network data reduction technique such as aggressive aggregation when congestion is detected. Other protocols, such as Pump Slowly Fetch Quickly (PSFQ) [3] can be adapted to avoid congestion. However, such approaches involve highly specialized parameter tuning, accurate timing configuration, and in-depth understanding of the protocol's internal operations. Congestion Detection and Avoidance (CODA) [11] is a congestion mitigation strategy, which uses both, buffer occupancy and channel load for measuring congestion levels in the network. It uses two strategies for handling both persistent and transient congestions. CODA performs a rate adjustment through traditional TCP-like Additive-Increase/Multiplicative-Decrease (AIMD) mechanism and thus often leads to the occurrence of packet loss. ARC [12] is a Linear Increase Multiplicative Decrease (LIMD) like algorithm. In an Analytical Rate Control (ARC), if an intermediate node overhears that the packets it sent previously are successfully forwarded again by its parent node, it will increase its rate by a constant α , otherwise it will increase its rate by a factor β (where $0 < \beta < 1$). The coarse rate adjustment could result in tardy control and introduce packet loss. Priority-based Congestion Control Protocol (PCCP) [13] introduces an efficient congestion detection technique addressing both node- and link-level congestion. However, it does not have any mechanism for handling prioritized heterogeneous traffic in the network.

Congestion Control and Fairness (CCF) [14] was proposed as a distributed and scalable algorithm that eliminates congestion within a sensor network and ensures the fair delivery of packets to the base station. The main spot light of CCF is its achieved fairness. However, its fairness is so simple as it results in a low throughput, especially, when some nodes do not have any packet to send. Priority-based Rate Control for Service Differentiation and Congestion Control (PRCSDCC) was proposed to distinguish high-priority real-time traffic from low-priority non-real-time traffic and service the input traffic based on its priority. However, PRCSDCC lacks flexibility due to static priority. Using Bit-Map-Assisted (BMA) MAC, each node is assigned a specific slot to transmit a one-bit control message if it has data to send; otherwise, its scheduled slot remains empty. After the contention period, the cluster head broadcasts its transmission schedule to the noncluster head nodes in the cluster and the system enters into the data transmission period. If none of the noncluster head nodes have data to send, the system proceeds directly to an idle period, which lasts until the next session. The nodes keep their radios off during the idle periods to save energy.

3 MAC Protocol

WSN is a novel communication paradigm involving the devices with low complexity that has limitations on processing capacity, memory and severe restrictions on power consumption [15–17]. The traffic in sensor network is often bursty and its energy wastage results from collisions, overhearing, control packet overhead and idle listening to the radio channel [18]. Thus, an effective MAC protocol is essential in determining the radio channel. The pertinent solution to save energy is to use a cluster based approach for the MAC scheme [19].

From the perspective of MAC layer the clustered network is divided into two distinct parts, i.e., intra- and intercluster domain. This chapter suggests a novel Hybrid MAC approach with Energy-efficient Time Division Multiple Access (ETDMA) and nanoMAC for intra- and intercluster domain respectively to reduce energy consumption. The main feature of the Hybrid MAC protocol is that it can adapt to either high or low level of contention in the network [20].

3.1 Overview of WSN

Clustering scheme organizes the nodes of the sensor network into two virtual domains, such as intracluster and intercluster domain as in Fig. 3. In the intracluster domain, the nodes sense the data and communicate with the cluster head directly



Fig. 3 Overview of WSN domain

within the cluster. Since the radio channel has high contention due to a large number of sensors in the intracluster domain, the ETDMA-based MAC protocol is utilized for achieving high energy efficiency. In the intercluster domain, the cluster head node communicates with the sink either directly (singlehop) or through another cluster head nodes (multihop). The number of nodes contending for the radio channel in intercluster is lesser compared to intracluster domain [21], and a CSMA-based MAC protocol (nanoMAC) is utilized for data transmission.

The frame structure of the Hybrid MAC protocol is shown in Fig. 4. In the intracluster domain, the cluster head assigns the time schedule based on TDMA to its nodes for data transmission [21]. The time slot is subdivided into mini-slots equal to the number of nodes in a cluster. This mini-slot carries a one-bit information of a node to determine whether it has the sensed data or not. If the node has no sensed data, its time slot is allocated to other nodes that have data to transmit. In intercluster domain, the cluster head nodes have data to transmit to sink or neighbor head node through a wireless channel by performing Carrier Sense (CS) before transmission. If a head node fails to get the channel, it goes to sleep mode and wakes up after a random time period and listens the channel again which leads to reduce the energy wastage and congestion. This feature contributes to increase the robustness of the Hybrid MAC protocol to synchronization and topology changes while enhancing its scalability to contention.



Fig. 4 Hybrid MAC frame structure

3.2 MAC Protocol for Intracluster Domain

In a conventional TDMA scheme, a node turns ON its radio during its assigned slot whether it has data to transmit or not, resulting in higher energy consumption. To reduce the energy consumption, ETDMA scheme is used, in which the node turns its radio OFF when it has no data to transmit. In addition, assigning dynamic time slot based on unpredictable traffic variations are difficult with conventional MAC schemes [22]. To efficiently assign the time schedule and minimize the energy consumption, the ETDMA MAC protocol is suggested for intracluster domain.

3.2.1 ETDMA MAC Protocol

The main objective of the ETDMA MAC protocol is to reduce the energy consumption due to avoiding idle listening and maintain low latency. In clustering approach, the data transmission of the noncluster head nodes are organized into rounds [23]. Each round consists of cluster set-up phase and steady-state phase as shown in Fig. 5.

(i) Setup Phase

During set-up phase, each node decides to become a cluster head based on its energy level. Elected cluster heads broadcast an advertisement message to all other nodes claiming to be the new cluster head by using nonpersistent CSMA. Each noncluster head node joins the cluster in which communications with the cluster head require a minimum amount of energy. Once the clusters are built, the system enters into the steady-state phase.



Fig. 5 Transmission periods of ETDMA MAC protocol

(ii) Steady-State Phase

The steady-state phase is divided into sessions. Each session consists of a contention period, data transmission period and an idle period as in Fig. 5. With N noncluster head nodes in the cluster, the contention period is exactly N slots. During each contention period, all nodes keep their radios ON. Using ETDMA MAC, each node is assigned a specific slot to transmit a one-bit control message if it has data to send or not.

After the contention period, the cluster head broadcasts its transmission schedule to the noncluster head nodes in the cluster and the system enters into the data transmission period. If the noncluster head nodes have no sensed data, the system proceeds directly to an idle period, which lasts until the next session. The nodes keep their radios OFF during the idle periods to save energy. When the session gets over, the next session begins with a contention period and the same procedure is repeated. The cluster head collects the data from all the source nodes and then forward the aggregated and compressed data to the sink directly or via a multihop path. After a predefined time, the system begins the next round and the whole process is repeated.

3.2.2 Energy Model of ETDMA MAC

This model describes the energy consumed by the sensor node in the intracluster domain. In ETDMA protocol, the sensor nodes keep their radio 'ON' during the whole contention period. After receiving the transmission schedule from cluster head, each source node sends its data packet to the cluster head over its scheduled time slot.

The energy consumed by each source node during a single session is given by

$$E_{\rm sn} = P_{\rm t}T_{\rm c} + (N-1)P_{\rm i}T_{\rm c} + P_{\rm r}T_{\rm ch} + P_{\rm t}T_{\rm d}$$
(1)

where

 $P_{\rm t}$, $P_{\rm r}$ and $P_{\rm i}$ are the power consumption during the transmission, reception, and idle mode, respectively

 $T_{\rm c}$ is the time required to transmit/receive a control packet

N is the number of noncluster head nodes within a cluster

 $T_{\rm ch}$ is the time required for cluster head to transmit/receive a control packet $T_{\rm d}$ is the time required to transmit/receive a data packet.

Each nonsource node stays idle during the contention period and keeps its radio OFF during the data transmission period. Thus, over a single session, the energy that it dissipates can be computed as A Game Theory-Based Hybrid Medium Access Control Protocol ...

$$E_{\rm in} = NP_{\rm i}T_{\rm c} + P_{\rm r}T_{\rm ch} \tag{2}$$

During the contention period of the *i*th session, cluster head node receives n_i control packets from noncluster head nodes and stays idle for $(N - n_i)$ contention slots. In the subsequent transmission period, the cluster head node receives n_i data packets from the noncluster head nodes. Hence, the energy expended in the cluster head node during a single session is given as

$$E_{\rm ch} = n_i (P_{\rm r} T_{\rm c} + P_{\rm r} T_{\rm d}) + (N - n_i) P_{\rm i} T_{\rm c} + P_{\rm t} T_{\rm ch}$$
(3)

where n_i is the number of source nodes in the *i*th session/frame. Therefore, the total system energy consumed in each cluster during the *i*th session is

$$E_{\rm si} = n_i E_{\rm sn} + (N - n_i) E_{\rm in} + E_{\rm ch} \tag{4}$$

Each round consists of k sessions, thus the total system energy dissipated during each round is computed as

$$E_{\text{round}} = \sum_{i=1}^{k} E_{\text{si}} \tag{5}$$

The average packet delay τ_d , is defined as the average time required for a packet to be received by the cluster head node and is given by

$$\tau_{\rm d} = \frac{NT_{\rm c} + T_{\rm ch} + n_i T_{\rm d}}{n_i} \tag{6}$$

The energy consumption and delay for the intracluster network are evaluated using the Eqs. (5) and (6), respectively. The nodes' battery energy is saved by providing the ETDMA scheduling based on the traffic load using the repeated game. The detailed description of repeated game theory is to be discussed in Sect. 5.

3.3 MAC Protocol for Intercluster Domain

In a conventional np-CSMA scheme [24, 25], a node with a frame is used to transmit the senses of the channel by using CS. If the channel is detected busy, the node waits for a random time interval for transmission to avoid collision. When two users sense the idle channel at the same time and transmit their frames, a collision occurs. This request for retransmission and results in high energy consumption of the sensor node. To minimize the energy consumption, nanoMAC protocol is suggested for intercluster domain.



Fig. 6 Transmission periods of nanoMAC protocol

3.3.1 NanoMAC Protocol

The nanoMAC protocol of CSMA/CA type is nonpersistent. With probability p, the protocol will act as nonpersistent and with probability (1 - p), the protocol will refrain from sending even before CS and schedule a new time to attempt for CS. Nodes contending for the channel do not listen constantly to the channel, contrary to the normal binary exponential backoff mechanism, but sleep during the random contention period.

When the backoff timer expires, the nodes wake up to sense the channel. This feature makes the CS time for nanoMAC a short, and saves the energy of sensor nodes to a greater extent. With one Request To Send (RTS) and Clear To Send (CTS) reservation, a maximum of 10 data frames can be transmitted using the frame train structure as in Fig. 6. The data frames are acknowledged by a single, common Acknowledgement (ACK) frame that has a separate ACK bit reserved for each frame. Only the corrupted frames are retransmitted and not the whole data packet.

3.3.2 Energy Model of NanoMAC

The transmission energy consumption model of the nanoMAC protocol is shown in Fig. 7. During data transmission, this model describes the energy consumption by taking average contention times, backoff times and frame collisions [26–30]. There are four different states: Arrive, Backoff, Attempt and Success. Arrive state is an entry point to the system for a node to transmit new data. Energy consumed for every arrival to one of these states. To reach the success state, all possible transitions starting from the arrival state and ending at the success state is calculated.

On the arrival of data, when a device finds the channel busy, it refrains from its transmission, and reaches the backoff state. When the channel is clear upon CS, the sensor node transmits RTS frame to the destination node and it waits for a CTS



frame and reaches the attempt state. On successful transmission of the RTS and reception of CTS, a transition to the success state is made. The success state represents a successful data exchange with the destination.

When the RTS frame collides, the device returns to the backoff state and no new data transmissions are made during this period. Backoff state represents the device's waiting period, trying to acquire the channel again. When the device detects the channel as vacant or idle, it transits to the attempt state by sending the RTS frame. When the channel is detected busy, it stays in the backoff state and repeats the process. The average energy consumption upon transmission from the point of packet arrival to the point of receiving an ACK frame is given by

$$E_{\text{TX}} = E_{\text{arrive}} + p_1 E(\mathbf{A}) + (1 - p_1) E(\mathbf{B})$$
 (7)

where

 E_{arrive} is the carrier sensing energy consumption when reaching the arrive state E(A) and E(B) are the energy consumption on each visit by the node to attempt state and backoff state and are given by

$$E(A) = p_2 E(S) + (1 - p_2) E(B)$$
(8)

and

$$E(B) = p_3 E(A) + (1 - p_3) E(B)$$
(9)

where

E(S) is the expected energy consumption upon reaching the success state from the attempt state

 $p_{\{1,2,3\}}$ are the different probabilities related to arriving at a certain state

The transmitter energy consumption can be simplified as

$$E_{\rm TX} = T_{\rm CS} M_{\rm RX} + p_{\rm b} \left(T_{\rm bb} + \frac{T_{\rm r}}{2} \right) M_{\rm slp} + p_{\rm b} E({\rm B}) + (1 - p_{\rm b})(1 - p_{\rm ers}) \left(T_{\rm bp} + \frac{T_{\rm r}}{2} \right) M_{\rm slp} + (1 - p_{\rm b}) p_{\rm ers} E({\rm A}) + (1 - p_{\rm b}) p_{\rm ers} \left(T_{\rm pr} + T_{\rm rts} \right) M_{\rm TX} + (1 - p_{\rm b})(1 - p_{\rm ers}) E({\rm B})$$
(10)

where

 $T_{\rm CS}$ is the time required for carrier sensing $M_{\rm RX}$ is the receiver power consumption $p_{\rm b}$ is the probability of finding channel busy during carrier sense $T_{\rm bb}$ is the incremented backoff time $T_{\rm r}/2$ is the average random delay $M_{\rm slp}$ is the sleep power consumption of the transceiver $T_{\rm bp}$ is the unincremented backoff time $p_{\rm ers}$ is the unincremented backoff time $T_{\rm pr}$ is the time required to transmit a preamble $T_{\rm rts}$ is the time required to transmit an RTS frame $M_{\rm TX}$ is the transmitter power consumption

The receiver energy consumption model of a packet for nanoMAC protocol is shown in Fig. 8. There are three different states: Idle, Reply, and Received. When the destination node receives an RTS packet, it transits to state Reply and forward the CTS packet to the source. When the destination node receives the valid data



packet from the source, it reaches the received state and sends an ACK frame to the source node. When the CTS packet transmitted by the receiver collides, it stays in idle state.

The average energy consumed by the receiver to receive data packet is given by

$$E_{\rm RX} = E(I) = \frac{(\mu + p_{\rm s}\theta)}{(p_{\rm s}p_{\rm senh})}$$
(11)

where

E(I) is the energy incurred in each visit of node to idle state

 $\boldsymbol{\mu}$ represents the energy model transitions from state idle

 $\boldsymbol{\theta}$ represents the energy model transitions from state reply

 $p_{\rm s}$ and $p_{\rm senh}$ are the probabilities of no collision during RTS or CTS transmission

The average packet delay τ_d , from the cluster head to the sink is calculated using Fig. 8 and is given by

$$\tau_{\rm d} = p_{\rm b} \left[T_{\rm bb} + \frac{T_{\rm r}}{2} + T_{E(\rm B)} \right] + (1 - p_{\rm b})(1 - p_{\rm ers}) \left[T_{\rm bp} + \frac{T_{\rm r}}{2} + T_{E(\rm B)} \right] + (1 - p_{\rm b})p_{\rm ers} \left[T_{E(\rm S)} \right]$$
(12)

where

 $T_{E(A)}$ is the time required to visit the attempt state $T_{E(B)}$ is the time required to visit the backoff state $T_{E(S)}$ is the time required to visit the success state

The energy consumption and delay are evaluated for intercluster communication by using the Eqs. (10), (11), and (12). The repeated game is applied in nanoMAC protocol to provide the effective sleeping time for CH which has reduced the energy consumption.

4 Hybrid MAC Protocol

From the perspective of the MAC layer, the clustered network is divided into two distinct parts, i.e., the intracluster domain and the intercluster domain. This chapter suggests a novel Hybrid MAC approach with the combination of ETDMA and nanoMAC, which is based on CSMA/CA for intra- and intercluster domain, respectively, to reduce collision and energy consumption. The main feature of the Hybrid MAC protocol is that it can adapt to either high or low level of contention in the network. The performance of the proposed Hybrid MAC protocol is evaluated in terms of the network lifetime and it is compared with conventional MAC schemes. The MAC layer provides efficient operation of a sensor network since it avoids congestion between data by not allowing two interfering nodes to transmit at

the same time. The MAC schemes for wireless data communication is categorized into a contention-based and schedule-based protocols [31]. Based on the sensing of the channel, in contention-based scheme, the sensor nodes keep their radio ON to transmit the message. The major disadvantage in using nonpersistent carrier sense multiple access (np-CSMA) protocol is the nodes compete to share the channel which leads to collision from nodes beyond one hop and it leads to hidden terminal problem. In the case of schedule-based MAC scheme, TDMA overcomes hidden terminal problem, but it requires efficiently scheduling of time to avoid idle listening of the channel. The major energy wastage in np-CSMA and TDMA MAC schemes are congestion, collisions, overhearing, control packet overhead, and idle listening to the channel [32, 33]. Therefore, an efficient MAC protocol has to consume less energy and this is achieved by using nanoMAC protocol, which has a sophisticated sleep algorithm and congestion avoidance technique in CSMA. However, the difficulty in decision making for data communication consumes more energy. The proposed Hybrid MAC is the combination of the ETDMA and nanoMAC protocol for intracluster and intercluster communication. The main objective of this chapter is that the game theory-based Hybrid MAC (GH-MAC) has to reduce the time duration to take the decision and forward the data effectively, which reduces energy consumption of the node.

4.1 System Model

The Hybrid MAC protocol is divided into two levels of communication process: the first level is intracluster communication, i.e., between the cluster member and the cluster head (CH), and the second level communication is between the CHs. The operations in ETDMA are divided into two rounds; one is cluster set-up phase while the other is steady-state phase. The cluster set-up phase checks the nodes based on the energy level, whether it can become CH node. The selected CHs broadcast an advertisement message to all nodes stating it is the new CH. Then the other non-cluster-head nodes, which require minimum energy to communicate with CH join together to form a cluster. With the formation of cluster, the system goes to steady-state phase [28, 34–37].

The categories of steady-state phase are contention period and frames. In the contention period, the nodes keep their radio ON, while the CH builds TDMA schedule and transmits it to all nodes within the cluster. A data slot is allotted to each node in a frame. The duration of each frame is fixed. The source node transmits its data to CH within its allocated time by turning ON its radio and all other times the radio is kept OFF. The different states of intercluster communication are Arrive, Backoff, Attempt, and Success states. In the arrive state, the node starts transmitting new data. To reach the success state, all possible transitions from the arrival state to the success state is calculated. On the arrival of data, when a device finds the channel busy, it refrains from its transmission, and reaches the backoff state. When the channel is clear upon CS, the source CH transmits the RTS frame to

the destination CH and it waits for a CTS frame and reaches the attempt state. On successful transmission of the RTS and reception of CTS, a transition to the success state is made. The success state represents a successful data exchange with the destination. When the RTS frame collides, the device returns to the backoff state and no new data transmissions are made during this failed period. Figure 9 shows the entire communication process of Hybrid MAC protocol.

5 Game Formulation

The game theory techniques have been widely applied to various engineering design problems in which the action (i.e., any activity) of one component has its impact on any other component. Game theory also addresses problems where multiple players with different objectives compete and interact with each other in the same system; such as mathematical abstraction which is useful for generalization of the problem. Therefore, game formulations are used, and a stable solution for the players is obtained through the concept of equilibrium. This chapter provides applications of game theory in wireless network and presents them in a layered perspective, emphasizing on which game theory could be effectively applied. Energy efficiency of MAC protocol in WSN is very perceptive to the number of nodes competing for the access channel to avoid congestion. Accurate analysis of congestion, collision probability, transmission probability, and so forth is very difficult for a MAC protocol by detecting channel in wireless medium [38]. Game theory provides the solution to solve this constraint of MAC protocol suggested for WSN by regulating the traffic load and providing an efficient sleeping time period to conserve the nodes energy, by applying repeated game for Hybrid MAC which is known as Game theory-based Hybrid MAC (GH-MAC) protocol.

5.1 Repeated Game

Repeated game theory is concerned with a class of dynamic games, in which a game is played for numerous times and the players can observe the outcome of the previous game before attending the next repetition [39]. To understand the concept of repeated games, let us start with an example, which is known as the Prisoner's Dilemma [40], in which two criminals are arrested and charged with a crime. The police do not have enough evidence to convict the suspects, unless at least one confesses. The criminals are in separate cells, thus they are not able to communicate during the process. If neither confesses, they will be convicted of a minor crime and sentenced for 1 month. The police offer both the criminals a deal. If one confesses and the other does not, the confess, they will be released and the other will be sentenced for 9 months. If both confess, they will be sentenced for 6 months. This game has a unique Nash equilibrium in which each player chooses to cooperate in a





single-shot setting. Now in the prisoner's dilemma, suppose one of the players adopts the following long-term strategy: (i) choose to cooperate as long as the other player chooses to cooperate, (ii) In any period the other player chooses to defect, and then choose to defect in every subsequent period. What should the other player do in response to this strategy? This kind of game is known as a repeated game with sequence of history-dependent game strategies.

5.2 Game Formulation for GH-MAC

Formally, the game is defined as $G = [N, A_i, U_i]$ and has the following three components [41]

- The set of players, $N = \{1, 2, 3, ..., x\}$
- The set of actions (strategy profile), A_i, available for a player 'i' to make a decision.
- The payoff (utility) function U_i resulting from the strategy profile.

The model considered in the analysis of the game consists of *N* homogeneous nodes in the sensor network. The players of the game are considered as nodes in WSN. The set of actions (strategies) available for the player '*i*' to make a decision, consists of all possible sleeping time ranging from the minimum level s_{min} to maximum level s_{max} . The nodes in the network play repeated game. All the nodes play the game simultaneously by choosing their individual strategies. This set of choice results in the payoff (utility) of the game.

Each round in WSN consists of data collection phase, aggregation phase, and transmission phase. The information available from previous round is used to work out strategies in future rounds. A source node is equally potential to many neighboring nodes within the interference range is concerned. The number of interfering nodes depends on the node density $\rho = N/A$, where A is the network area. Since the nodes are considered as homogeneous, the actions allowed by the nodes are same. All the nodes can transmit with any power level to make its successful transmission. If the nodes transmit an arbitrary high power level, it will increase the interference level of the other nodes, which leads to congestion. To overcome the effect of this high interference, the neighboring nodes in turn will transmit at higher power. This happens as a cascade effect and soon leads to a noncooperative situation. To control this noncooperative behavior, an equilibrium game strategy that imposes constraints on the nodes to act in cooperative manner, even in a noncooperative network is devised. The existence of some strategy sets $S_1, S_2, S_3, \dots, S_x$ for the nodes (1, 2, 3, ..., x) is assumed. In this game, if node 1 chooses its sleeping time strategy, $s_1 \in S_1$ and node 2 chooses its strategy $s_2 \in S_2$, and so on, then the set of strategies chosen by all 'x' nodes is given by

$$s = \{s_1, s_2, \dots, s_x\}$$
(13)

This vector of individual strategy is called a strategy profile. The set of all such strategy profiles are called the space strategy profile S'. At the end of an action, each node $i \in N$ receive a utility value as given by

$$u_i(\mathbf{s}) = u_i(s_i, s_{-i}) \tag{14}$$

where

 s_i is the strategy profile of the *i*th node

 s_{-i} is the strategy profile of all the nodes but for the *i*th node

The utility of each node depends not only on the strategy it picked, but also on the strategies of the other nodes. In the game, each node maximizes its own utility in a distributed fashion. The transmit power that optimizes individual utility depends on transmit powers of all the other nodes in the system. It is necessary to characterize a set of powers where the players are satisfied with the utility.

6 Simulation Results and Discussions

The analysis of the GH-MAC protocol is carried out using MATLAB 10. The parameters considered for the simulation is summarized in Table 1. The performances of the GH-MAC protocol are evaluated based on energy consumption and delay in terms of traffic load on the network.

Table 1 Simulation parameters Image: Comparison of the second			
	Parameters	Value	
	Number of nodes (<i>n</i>)	100	
	Area (A)	$100 \times 100 \ (m^2)$	
	Transmitting power (P_t)	462 (mW)	
	Receiving power (P_r)	346 (mW)	
	Power for idle listening (P_i)	330 (mW)	
	Data rate	2 (Mbps)	
	Data packet size	1,452 (bytes)	
	Control packet size	52 (bytes)	
	Control frame size for nanoMAC	18 (bytes)	
	Data frame size for nanoMAC	41 (bytes)	
	Data frame payload of nanoMAC	35 (bytes)	
	Device transmission distance	100 (m)	



6.1 Energy and Delay Analysis for Intracluster Network

Figure 10 shows the energy consumption with traffic load of intracluster domain for a single round. The energy consumption of the three schedules is based on MAC protocol such as Bit-Map-Assisted (BMA), ETDMA, and G-ETDMA. These are compared with 20 nodes in a cluster and five sessions/round (k). From the comparison, G-ETDMA provides better performance in terms of energy consumption than ETDMA and BMA for the entire traffic load. The energy consumption of G-ETDMA is 22 % less than ETDMA and 35 % less than BMA of traffic load 0.4. The reason for this improved performance in G-ETDMA is by avoiding idle listening and the CH node radio need not be kept ON for the entire time slot.

The average energy consumption of noncluster head nodes in a cluster for traffic load 0.4 with five sessions per round is shown in Fig. 11. G-ETDMA protocol performs better than ETDMA and BMA schemes when the number of noncluster







head nodes handled by the CH node is optimized for 33. As the number of noncluster head nodes in a cluster is larger, the contention period increases which result in greater energy consumption. If it is very less than the network, lifetime reduces due to the continuous transmission of sensed data by selective nodes, which leads to die out the node's battery immediately. Therefore, 33 noncluster head nodes are optimum in a cluster adapting with G-ETDMA scheme when compared with BMA and ETDMA protocol.

Figure 12 compares the three MAC techniques in terms of average packet delay. For higher traffic load, all the MAC schemes provide less delay. However, as the traffic load reaches a minimum, the average packet delay grows exponentially with BMA than G-ETDMA scheme. This is because in G-ETDMA protocol, the scheduling of nodes change dynamically according to the traffic variations in the network. This greatly reduces the energy consumption of nodes due to idle listening and thus maintains a good and lower delay performance.

6.2 Energy and Delay Analysis of Intercluster Network

Figure 13 shows the energy consumption analysis of single and multihop network, with and without game-based nanoMAC protocols. When the sink node is within the characteristic distance of 100 m, consumes less energy in a singlehop transmission. If the transmission distance is outside the characteristic distance, then the multihop communication is efficient and consumes less energy. From Fig. 13, the G-nanoMAC consumes 1.6 times less energy than multihop within the transmission distance of 100 m. When the hop distance is about 100 m (i.e., ten hops), the energy



consumption of singlehop increases approximately by the factor of 0.5 compared with multihop because of path loss.

G-nanoMAC outperforms beyond and up to the traffic load 0.5, the collisions may increase the energy consumption of the nodes. The comparison of normalized delay characteristics of nanoMAC and G-nanoMAC protocols are shown in Fig. 14.

The delay occurred at the reception of frame gradually increases with the traffic load due to the retransmission of entire frame when an error or collision occurring in this transmission period. In the G-nanoMAC protocol, a device sends 10 data frames of 41 bytes each, an ACK frame for the same transmission period and retransmits only the lost/collided frame under the consideration of traffic load. Thus, the delay offered in the network is 12 times lesser when compared to nanoMAC.



Fig. 14 Normalized delay versus traffic load





6.3 Energy and Delay Analysis for Hybrid MAC Protocol

Figure 15 shows the energy consumption with traffic load for a Hybrid MAC protocol and GH-MAC. The comparisons made for Hybrid MAC and GH-MAC with 100 nodes in a network area of $100 \times 100 \text{ m}^2$. GH-MAC is providing better performance in terms of energy than Hybrid-MAC up to the traffic load is 0.85. This is due to providing proper sleeping time to CH based on the traffic and beyond this traffic all the nodes try to transmit their data that make more collision and increase the energy consumption.

A comparison of delay characteristics of Hybrid MAC and GH-MAC protocols are shown in Fig. 16. Upon error or collision during this transmission period, the entire frame has to be retransmitted; hence the delay incurred in reception of frame gradually increases with the traffic load. GH-MAC outperforms when compared to Hybrid MAC due to providing the most effective sleeping time schedule using game theory to the nodes, which reduces the collisions of data packet transfer from source to sink.



7 Conclusion

The WSNs often experience congestion, so an advanced congestion control solution is required. The CC mechanism should differ from its sibling deployed in the Internet. A lot of research and solutions were published targeted to solve the congestion problem in resource restricted communications. The CC mechanisms made for WSNs constitute an effective algorithm to satisfy the traffic and controlling overflow. The key approach of GH-MAC is an active queue management mechanism to predict the overflow of intermediate device's buffers. In the GH-MAC protocol, energy and delay performance for offering traffic load has been discovered in the cluster based WSN. From these performances, it is evident that the G-ETDMA protocol for the intracluster communication achieves 25 % reduction in energy consumption compared to the BMA. It provides 15 % less packet transmission delay by incorporating proper dynamic scheduling schemes. G-nanoMAC protocols offered better performance in intercluster communication and its energy spent for data transmission is 20 % less than nanoMAC protocol. The delay performance for G-nanoMAC is considerably reduced by 12 % without any degradation in throughput compared with nanoMAC protocol. This efficient energy utilization in G-ETDMA and G-nanoMAC leads to energy reduction in GH-MAC for the entire communication process in WSN. This reduction in energy consumption and delay of the GH-MAC protocol can considerably extend the lifetime of the sensor network due to its defined congestion in WSN. The limitation of this proposed protocol requires additional hardware to make decisions using game theory. However, nodes consume less energy for making the decision quickly without using many calculations. In order to integrate this framework into secure network design, future work will be aimed at identifying and analyzing various network jammer types and realistically profiling the associated strategies and perceived payoffs.

References

- 1. Chen S, Yang N (2006) Congestion avoidance based on lightweight buffer management in sensor networks. IEEE Trans Parallel Distrib Syst 17(9):934–946
- Wang C, Member, IEEE, Li B, Senior Member, IEEE, Sohraby K, Senior Member, IEEE, Daneshmand M, Member, IEEE, Hu Y (2007) Upstream congestion control in wireless sensor networks through cross-layer optimization. IEEE J Sel Areas Commun 25(4):786–795
- Lin Q, Wang R-C, Sun L-J (2011) Novel congestion control approach in wireless multimedia sensor networks. J China Univ Posts Telecommun 18(2):1–8
- Nesa Sudha M, Valarmathi ML (2013) Collision control extended pattern medium access protocol in wireless sensor network. J Comput Electr Eng 39:1846–1853
- Yadav R, Varma S, Malaviya N (2008) Optimized medium access control for wireless sensor network. Int J Comput Sci Netw Secur 8(2):334–338
- Heidemann J, Ye W, Estrin D (2002) An energy efficient MAC protocol for wireless sensor networks. In: Proceedings of the IEEE info comm, Newark, pp 1567–1576
- 7. Zheng T, Radhakrishnan S, Saranga V (2005) PMAC: an adaptive energy-efficient MAC protocol for wireless sensor networks. In: Proceedings of the 19th IEEE international conference, parallel and distributed processing symposium, IPDPS' 05, April 2005
- Intanagonwiwat C, Govindan R, Estrin D (2000) Directed diffusion: a scalable and robust communication paradigm for sensor networks. In: Proceedings of 6th annual international conference on mobile computing and networking (MOBICOM'00), Aug 2000
- Wan CY, Campbell AT, Krishnamurthy L (2002) A reliable transport protocol for wireless sensor networks. In: Proceedings of the 1st ACM international workshop on wireless sensor networks and applications (WSNA'02), 28 Sep 2002, Atlanta, GA, USA. ACM, New York, NY, USA
- Wan CY, Eisenman SB, Campbell AT (2005) CODA: congestion detection and avoidance in sensor networks. In: Proceedings of the 1st international conference on embedded networked sensor systems, Nov 2005
- Woo A, Culler DE (2001) A transmission control scheme for media access in sensor networks. In: Proceedings of the 7th annual international conference on mobile computing and networking (MOBICOM'01), July 2001
- 12. Wang C, Li B, Sohraby K et al (2007) Upstream congestion control in wireless sensor networks through cross-layer optimization. IEEE J Sel Areas Commun 25(4):786–781
- Ee CT, Bajcsy R (2004) Congestion control and fairness for many-to-one routing in sensor networks. In: Proceedings of the 2nd international conference on embedded networked sensor systems (SenSys'04), Nov 2004
- Yaghmaee MH, Adjerohb DA (2009) Priority-based rate control for service differentiation and congestion control in wireless multimedia sensor networks. Comput Netw 53(11):17–24
- 15. Akyildiz F, Su W, Sankarasubramaniam Y, Cayirci E (2002) A survey on sensor networks. IEEE Commun Mag 40(8):102–114
- Balakrishnan H (2004) Opportunities and challenges in high rate wireless sensor networking. In: Proceedings of 29th annual IEEE international conference on local computer networks, Oct 2004, Florida, USA, pp 4–10
- 17. Demirkol I, Erosy C, Alagoz F (2006) MAC protocols for wireless sensor networks: a survey. IEEE Commun Mag 44(4):115–121
- Heinzelman WB, Chandrakasan AP, Balakrishnan H (2002) An application specific protocol architecture for wireless micro sensor networks. IEEE Trans Wireless Commun 1(4):660–670
- 19. Shakir M, Ahmed I, Peng M, Wang W (2007) Cluster organisation based design of hybrid MAC protocol in wireless sensor networks. In: Proceedings of 3rd international conference on networking and services, June 2007, Athens Greece, p 78
- El-Hoiydi A (2002) Spatial TDMA and CSMA with preamble sampling for low power adhoc wireless sensor networks. In: Proceedings of the 7th international symposium on computers and communications, July 2002, Italy, pp 685–692
- Yang H, Sikdar B (2007) Optimal cluster head selection in the LEACH architecture. In: Proceedings of IEEE international conference on performance, computing and communications, April 2007, New Orleans, LA, pp 93–100
- 22. Li J, Lazarou GY (2004) A bit map assisted energy efficient MAC scheme for wireless sensor networks. In: Proceedings of IEEE/ACM 3rd international symposium on information processing in sensor networks, April 2004, Berkeley, USA, pp 55–60
- Kleinrock L, Tobagi FA (2001) Packet switching in radio channels: Part I—Carrier sense multiple access modes and their throughput-delay characteristics. IEEE Trans Commun 23 (12):1400–1416
- 24. Ansari J, Riihijarvi J, Mahonen P, Haapola J (2007) Implementation and performance evaluation of nanoMAC: a low power MAC solution for high density wireless sensor networks. Int J Sens Netw 2(5/6):341–349
- 25. Bacci V, Chiti F, Morosi S, Haapola J, Shelby Z (2006) Performance evaluation of optimised medium access control schemes based on ultra wideband technology. In: Proceedings of 17th annual IEEE symposium on personal, indoor and mobile radio communications, Sept 2006, Helsinka, pp 1–6

- 26. Chen P, O'Dea B, Callaway E (2002) Energy efficient system design with optimum transmission range for wireless adhoc networks. In: Proceedings of international conference on communications, Aug 2002, vol 2, pp 945–952
- Shebli F, Dayoub I, Rouvaen JM (2007) Minimising energy consumption within wireless sensor networks. In: Proceedings of IEEE international conference on signal processing communications, Nov 2007, Dubai, pp 1–6
- Haapola J, Shelby Z, Pomalaza-Raez C, Mahonen P (2005) Multihop medium access control for WSNs: an energy analysis model. EURASIP J Wireless Commun Netw 4:523–540
- Vidhya J, Kalpana G, Dananjayan P (2009) Energy efficient hybrid MAC protocol for clusterbased wireless sensor network. Int J Comput Electr Eng 1(3):1793–8198
- Raja P, Dananjayan P (2012) Game theory based ETDMA for intra-cluster wireless sensor network. In: Proceedings of IEEE international conference on ICACCCT, pp 272–276
- Heinzelman WB, Chandrakasan AP, Balakrishnan H (2002) An application-specific protocol architecture for wireless micro sensor networks. IEEE Trans Wireless Commun 1(8):660–670
- 32. Mehta S, and Kwak KS (2010) An energy-efficient MAC protocol in wireless sensor networks: a game theoretic approach. EURASIP J Wireless Commun Netw 5, Article ID 926420
- Kasu SR, Bellana SK, Kumar C (2007) A binary countdown medium access control protocol scheme for wireless sensor networks. In: Proceedings of 10th international conference on information technology (ICIT 2007), pp 122–126
- Raja P, Dananjayan P (2013) Game theory based nanoMAC protocol for inter cluster wireless sensor network. Int J Comput Electr Eng (IJCEE) 5(6):606–610
- Demirkol I, Ersoy C, Alagoz F (2006) MAC protocols for wireless sensor networks. IEEE Commun Mag 44(4):115–121
- Raghunathan V, Schurgers C, Park S, Srivastava MB (2002) Energy-aware wireless micro sensor networks. IEEE Signal Process Mag 13(4):40–50, 123–127
- Stemm M, Katz RH (1997) Measuring and reducing energy consumption of network interfaces in hand-held devices. IEICE Trans Commun E80-B(8): 1125–1131
- Machado R, Tekinay S (2008) A survey of game-theoretic approaches in wireless sensor networks. Elsevier J Comput Netw 52:3047–3061
- Lin Q, Wang R, Sun L (2011) Novel congestion control approach in wireless multimedia sensor networks. J China Univ Posts Telecommun 18(2):1–8
- 40. Ratliff J. http://www.virtualperfection.com/gametheory
- Sengupta S, Chatterjee M, Kwait KA (2010) A game theoretic framework for power control in wireless sensor networks. IEEE Trans Comput 59(2):231–242

Cooperative Games Among Densely Deployed WLAN Access Points

Josephina Antoniou, Vicky Papadopoulou-Lesta, Lavy Libman and Andreas Pitsillides

Abstract The high popularity of Wi-Fi technology for wireless access has led to a common problem of densely deployed access points (APs) in residential or commercial buildings, competing to use the same or overlapping frequency channels and causing degradation to the user experience due to excessive interference. This degradation is partly caused by the restriction where each client device is allowed to be served only by one of a very limited set of APs (e.g., belonging to the same residential unit), even if it is within the range of (or even has a better signal quality to) many other APs. The current chapter proposes a cooperative strategy to mitigate the interference and enhance the quality of service in dense wireless deployments by having neighboring APs agree to take turns (e.g., in round-robin fashion) to serve each other's clients. We present and analyze a cooperative game-theoretic model of the incentives involved in such cooperation and identify the conditions under which cooperation would be beneficial for the participating APs.

Keywords Dense Wi-Fi access points • Unmanaged wireless deployment • Graph theory • Game theory • Graphical game • Cooperation

J. Antoniou (🖂)

V. Papadopoulou-Lesta School of Sciences, European University Cyprus, Nicosia, Cyprus e-mail: v.papadopoulou@euc.ac.cy

L. Libman School of Computer Science and Engineering, University of New South Wales, Sydney, Australia e-mail: lavy.libman@unsw.edu.au

A. Pitsillides Department of Computer Science, University of Cyprus, Nicosia, Cyprus e-mail: andreas.pitsillides@ucy.ac.cy

School of Sciences, University of Central Lancashire Cyprus, Pyla, Cyprus e-mail: jantoniou@uclancyprus.ac.cy

[©] Springer International Publishing Switzerland 2015 K. Hausken and J. Zhuang (eds.), *Game Theoretic Analysis of Congestion, Safety and Security*, Springer Series in Reliability Engineering, DOI 10.1007/978-3-319-13009-5_2

1 Introduction

1.1 Need for Cooperation in Unmanaged Wireless Environments

The low cost and easy deployment of wireless technologies, often results in the increase in the density of wireless networks in urban residential areas and consequently, unmanaged deployments of IEEE 802.11 (Wi-Fi) home networks that cause degraded experience for the home users. Very frequently in such deployments, an Access Point (AP) can be located within the range of dozens of other APs, competing for the limited number of channels offered by the IEEE 802.11 standard [1]. This is because such wireless networks are formed with no planning or central authority being present.

Wireless APs that operate in the same geographical region without any coordination typically cause degradation to their users' experience because, in the current standards, at any given time every terminal must be rigidly associated with one particular AP. This leads to a competition between collocated APs for the same communication resource (radio channel) and a reduced user quality of experience due to the resulting interference.

To mitigate the problem, it would be beneficial for individual APs that are in physical proximity to each other to form *cooperative groups*, where one member of the group would serve the terminals of all group members in addition to its own terminals, provided that the signal strength of the AP is sufficient to support the user activity. The rest of the APs in the group can be silent or even turned off, thereby reducing interference. The group members can take turns at regular intervals to serve all the terminals, as long as the coverage of the group does not change.

1.2 Cooperation and Improved User Experience

Clearly, it is important that such groups include only members whose signal strength and/or available bandwidth are sufficient to serve all group members without loss of user quality of experience. To maintain an incentive to cooperate, the group must be formed, so that the quality of experience for any individual terminal is not reduced; any reduction of signal quality at the physical layer (e.g., due to being served by an AP that is more distant than the 'native' one) should be compensated by the benefits from lesser interference at the MAC layer.

Since there is no centralized entity that can control the APs and force them to form cooperative groups, the creation of such groups must arise from a distributed process where each AP makes its own decisions independently and rationally for the benefit of itself and its terminals. We use game theory in order to model and investigate the feasibility of this cooperative operation among APs that selfishly maximize their own individual benefits under no central authority. More specifically, two nearby APs that decide to cooperate may enter into a *bilateral agreement*, so that they (i) transmit in nonoverlapping time slots and (ii) whenever one of the APs is transmitting, it serves the clients of the other AP as well as its own.

1.3 Cooperative Approach

The chapter considers and expands on earlier works where a game-theoretic approach toward reducing interference in dense wireless deployments was presented [3, 4, 13]. Specifically, in [4], a graphical strategic game model was presented and the basic case of two interacting APs was studied, showing sufficient conditions for Nash equilibria to result in cooperation.

The benefit from cooperation between two players (APs) is determined by the weight of an edge connecting them in an underlying graph (which is in turn derived from the signal quality of the communication links, which is dependent on the wireless network deployment, among other factors). In particular, we consider the case where the APs are located very densely. This is a very common case in buildings where many APs may be active at the same time.

In such dense placement of APs, interference appears between almost any pair of APs, which in turn degrades the signal quality of the communication link. Thus, the assumption of a desire for an agreement between two cooperative APs of mutual service of all their clients is realistic since dense placements of APs are already a reality, especially in big cities, with high population density and high density areas, as for example a shopping mall. So, the number of nearby APs (with high potential interference) is limited to a relative small number (the APs inside that center). Moreover, the high density of the placement of the APs implies that each AP is assigned a relatively small number of clients. Thus, it is very probable that APs can serve, during their ON transmission mode, not only their clients but also the clients of their neighbors that are in cooperation agreement with them. Finally, the presented case study assumes for its modeling that guarantees for the nonlocal clients exist, such that even though they agree to be served by one AP, they will never get higher priority than the particular AP's own clients; this can be achieved through suitable values of the *weights* on the corresponding edges of a related graph obtained according to the modeling assumed. The next subsection discusses further the graph representation.

The implementation of the game in practice will require some sort of an interaction (via a specially designed communication protocol) between the APs, but this only needs to be done once in a while (in the order of minutes) so the overhead involved is insignificant in the long run. The same can be said about any loss of throughput due to the reassociation overhead for clients to the new AP, it is not significant since it is incurred only infrequently.

1.4 Clique Graph Representation

The chapter presents a particular case study that models these placements through a corresponding directed, weighted *clique* graph; i.e., for every pair of nodes of the graph, there exists a bidirectional edge connecting them [7]. Note that no edge between two nodes implies a weight of zero (0). Using this modeling, the chapter investigates under which conditions *cooperation* between APs is forced. Cooperation between two nearby APs concerns a bidirectional agreement on their *ON-OFF* transmission modes as well as the service they provide: they (i) transmit in nonoverlapping time slots and (ii) when each one of the two APs is transmitting, it serves *besides* its clients, *the clients of the other AP*.

Investigation of the placements of the APs forming a clique graph is motivated by the fact that in a very dense placement of APs, interference appears between almost any pair of two APs. Investigating when cooperation is enabled between dense APs forming other graph topologies constitutes a future work.

1.4.1 Edge Weights of the Clique Graph

The weight of a (directed) edge connecting an AP to its neighbor represents the measurement of the signal power received by the neighbor when the AP is ON. SNR is a measure that compares the level of a desired signal to the level of background noise. It is defined as the ratio of signal power to the noise power and anything above 40 dB in wireless LANs is considered excellent, whereas anything above 25 dB is considered a very good signal. Anything within the signal range 15–25 dB is a low signal but can still be associated and is usually fast, however, lower than 10–15 dB is not desired. This signal is treated as interference when both the AP and its neighbor transmit at the same time over the same channel, but at the same time it indicates the strength of the (useful) signal that can be used to serve the neighbor's clients when the two APs are in a cooperation agreement.

While this chapter assumes that the weights in the graph correspond to the signal power, the weight can be any metric that represents the experienced quality of the user during cooperation. The model considered in this chapter focuses mainly on the quality of experience for the downlink traffic, which is typically the bulk of the communication in most applications, and ignores the uplink traffic which is usually minimal compared to the downlink.

Since SNR measures the success of transmitting packets, it is also a strong indication of the experienced quality of the user: high SNR results to minimal packet loss. Thus, high quality of service (directly related to the measure of quality of experience) is enabled as it concerns this parameter (success in transmitting packets). Furthermore, the SNR is measured on the downlink traffic, which is typically the bulk of the communication in most applications, and ignores the uplink traffic which is usually minimal compared to the downlink.

1.5 Security and Communication Concerns

The dense placements of APs are usually of relatively small size, compared to the number of clients in the area. Thus, it is reasonable to assume that APs can serve, during their ON transmission mode, not only their clients but also the clients of their neighbors that are in cooperation agreement with them.

It must be emphasized that cooperation between APs belonging to different service sets is not impossible even if they are required to use secure encryption/authentication (such as WPA in the 802.11 standard). At first sight, it may seem that such cooperation would require the APs to share their secret passwords, which is of course undesirable in an uncoordinated deployment. However, there are known techniques that allow a group of cooperating APs to establish a common "trust group," and share a secret key and/or authenticate their clients without a need to expose their passwords, as well as reduce the authentication delay when a client is handed off between neighboring APs [10].

Therefore, cooperation is still feasible even in the presence of encryption and authentication; i.e., when the traffic between each AP and its clients is encrypted and requires authentication with a secret password that precludes an AP from sharing others' clients. In other words, if AP1 and AP2 wish to cooperate, we need a mechanism that allows a client of AP1 to "prove" to AP2 that it knows the password for AP1, without actually disclosing the password. This is a relatively easy and long solved cryptographic problem; for example, in [10] a trust group or cloud is presented, establishing security key sharing for authentication of the members in the group, aiming to achieve a fast authentication in a 802.11 network.

Nevertheless, security restrictions and software/hardware incompatibilities may exist is such cooperative scenarios. Such issues do not pose a restriction to the modeling considered in this chapter, since they may be abstracted as non-cooperative APs that cause a constant additional interference that the cooperating APs (and hence the cooperating users) cannot avoid; (either through cooperations or not). In other words, these *non*-cooperative APs, result in a constant decrease of the quality of experience the AP can provide which is irrelevant to the cooperation agreements of the AP. So the modeling itself may ignore such issues.

Also, note that this modeling does not consider the uplink traffic sent *from* the clients of a given AP *to* the AP; the chapter investigates and tries to minimize interference only for communication *from* the AP to its clients. Although, this traffic is minimal, it is not zero. However, this chapter concentrates on the *most* of the communication between clients and AP, which is the one from the AP to the client. Incorporating the uplink communication too in a game-theoretic modeling consist a future work.

1.6 Chapter Contribution

The contribution in this chapter focuses on proving the necessary and sufficient conditions in order for all APs of the (clique) graph to join a cooperation group, where all APs maximize their users' quality of service. The main result proved in the chapter is that, if for *each* AP in the network, the signal strength received from any other AP (as represented by the weight of the respective graph edge) surpasses a certain threshold, then joining a cooperation group with all other APs is a best-response strategy (therefore leading to Nash equilibrium), which enables the given AP to provide a better quality of experience to its clients. More specifically, the agreement obtained in clique graphs consists of a (not necessarily equal) split of the total time period into nonoverlapping, nonzero time slots $T_1, T_2, ..., T_n$, such that during time slot T_i only the AP *i* is transmitting, serving all the clients in the clique.

Moreover, since the number of APs forming the clique is limited, it is necessary that these APs can satisfy the increased service demands obtained from the cooperation. It is shown that this translates to a constraint on the weights of the edges which becomes a necessary condition for a global agreement setting to be Nash equilibrium. Thus, we obtain a characterization of the necessary and sufficient conditions for agreement settings between the APs to be sustained in Nash equilibria.

2 Related Work

Game Theory has been extensively used in networking research as a theoretical decision-making framework, e.g., for mechanisms and schemes used for purposes of routing [25], congestion control [16], and resource sharing [20]. Coalitional game theory (or cooperative games) investigate the formation of cooperative groups, (coalitions) that allow the involving players to increase their gain by joining a coalition and receive (individually or by sharing between players of a common coalition) the benefits resulted by the cooperation. Coalitional Game Theory [19] has been successively utilized for enabling efficient communication of wireless networks in a number of works; for interference networks [14, 15], spectrum sensing [24], and other networking issues [19]. For a comprehensive work on coalitional games and their potential applications in communication networks see [21, 22].

Cooperative communication strategies, studied in [11], particularly focusing on multiple access networks using a coalitional game approach. The authors of [11] considered two different modes of operation showing an advantage of a dynamic setting compared to a static one. In the dynamic setting, users alternate at random between two modes of operation: an active state and a dormant state. Users within the same coalition share knowledge of their active states and a scheduler assigns the right to transmit within a coalition. Collisions occur when users belonging to

different coalitions transmit simultaneously. In this setup, the *grand coalition* formed by all users is sum rate optimal and is also stable in the sense that no user has incentive to leave the coalition. The second mode of operation is the static setting. Here, users are scheduled for transmission in a static round-robin fashion. However, members of the same coalition can share their right to transmit in a given time slot. In this case, the grand coalition although optimal can be unstable.

The benefits of cooperation between providers offering wireless access service to their respective subscribed customers through potentially multi-hop routers have been further investigated in [23]. In particular, the authors using coalitional game theory investigated the benefits of the providers when cooperating, i.e., pool their resources, such as spectrum and base stations and agree to serve each other's customers. They first considered the case where the providers share their aggregate revenues and modeled such cooperation using transferable payoff coalitional game theory [19]. Then, they considered the scenario that locations of the base stations and the channels that each provider can use have already been decided a priori and showed that the grand coalition is stable in this case, i.e., if all providers cooperate, there is always an operating point that maximizes the providers' aggregate payoff, while offering each such a share that removes any incentive to split from the coalition. They also examined cooperation when providers do not share their payoffs, but still share their resources, so as to enhance individual payoffs and showed that the grand coalition continues to be stable.

Ways of reducing interference in densely deployed wireless networks through cooperation were considered in [3] by making use of the *Prisoner's Dilemma/ Iterated Prisoner's Dilemma* game model and proposing a group strategy to motivate adjacent neighbors into cooperation. The Prisoner's Dilemma and Iterated Prisoner's Dilemma [12] has been widely used in research as models for motivating cooperation [2]. In particular, the publication of Axelrod's book in 1984 [6] was the main driver that boosted the concept of cooperation of entities with seemingly conflicting interests outside of game theory. The empirical results of the Iterated Prisoner's Dilemma (IPD) tournaments showed that group strategies performed extremely well and defeated well-known strategies in round-robin competitions in the 2004 and 2005 IPD tournaments [9].

3 A Graph-Theoretic, Cooperative Game

Here, we introduce the theoretical background used in [5] to develop their theoretical framework also described here next.

3.1 Background

3.1.1 Graph Theory

A weighted *digraph* $G(V, \vec{E}, \vec{W})$ consists of three types of elements, namely *nodes* constituting the set *V*, *edges* constituting the set \vec{E} , and finally the *weights* of the edges constituting the set \vec{W} . An edge *e* is specified by the *ordered* pair $e = (u, v) \in E$ iff there exists a (directional) connection (or link) *from* node *u* to node *v*, where $u, v \in V$ with positive weight w(u, v) > 0 in the graph *G*. w(u, v) is positive if and only if $(u, v) \in \vec{E}$.

For any node *v*, the set $Neigh_G$ denotes the set of (other) nodes of the graph connected to that node via a communication link, i.e., $Neigh_G = \{u \in V \mid (u, v) \in \vec{E}, (v, u) \in \vec{E}\}$. When clear from the context, we omit the graph subscript. A weighted digraph $G(V, \vec{E}, \vec{W})$ is a *clique* iff for every pair $u, v \in V$, there exists a pair $(u, v) \in \vec{E}, (v, u) \in \vec{E}$ and w(v, u) > 0.

A graph *G* is said to be of size *k* if |V| = k. Throughout the chapter, for an integer $n \ge 1$, denote $[n] = \{1, ..., n\}$. This chapter is concentrated on weighted, directed clique graphs. In the following, when we refer to a clique graph, we imply a weighted, directed clique graph.

3.1.2 Game Theory

A simple (one-shot) *strategic game* Γ is defined by a set of players N and a set of available strategies (behaviors) for each player $i \in N$, denoted as S_i . A (pure) profile σ of the game Γ specifies a particular strategy σ_i , one for each player $i \in N$. For a profile σ , each player i has an *utility* that depends on the specific player's current strategy as well as the other players' current strategies, and is denoted as $U_{\sigma}(i)$.

A profile σ of the game is said to be Nash equilibrium [17, 18] if, for every player, the utility in that profile is not smaller that the utility gained if that player unilaterally changes its strategy to any other $x_i \neq \sigma_i$, $x_i \in S_i$. Denote as σ^{x_i} the new profile obtained. Then, if σ is Nash equilibrium, it must be that $U_{\sigma}(i) \geq U_{\sigma^{x_i}}(i)$, for every $i \in N$.

Graphical Games are games in which the player strategies (and, consequently, their utilities) are related to, and defined in terms of, a graph G. The components of the graph (i.e., its nodes and edges) are said to capture the *locality* properties among the players, or, in other words, how the game is affected by the player's positioning [5] considers a graphical game.

3.2 The Mathematical Framework

Next, we present the mathematical framework and the corresponding graphical, cooperative game obtained as defined in [5].

3.2.1 The Graph

Consider a set of nodes $V = \{v_1, v_2, v_3, ..., v_n\}$ corresponding to the set of APs located in a geographical area. Each node in fact represents an entity comprising of one server (the AP) and a number of its clients (mobile stations). In an arbitrary case, we may assume that this number is approximately the same for all APs, so we may assume that each server serves one client. At a time *t*, say that server node *v* (an AP) is active or *ON*, if it is transmitting (to its own client). Otherwise, say that the node is inactive or *OFF*. We are concerned primarily with the decisions made by the server nodes (the APs); therefore, when we refer to a node, henceforth we refer to an AP.

Consider two nodes, u, v which are within range of each other's signals. In particular, node v (and its clients) may receive information from node u when it is OFF, while node u is ON, i.e., when node u broadcasts. If the quality of the signal received by node v (and its clients) from node u is above some lower bound value assumed by the node v, it is assumed that there exists a *directed* edge from node u to node v, denoted as (u, v). Moreover, the quality of the received signal is quantified by a weight w(u, v) associated with the directed edge (u, v). For purposes of normalization, the function $w(u, v) \in [0, 1]$ is defined. The weight value is analogous to the quality of the received signal, i.e., good quality reception corresponds to a value of w(u, v) close to 1. In [5] these values are assumed to be known; if necessary, they can be discovered through a distributed process initiated and executed locally at any AP.

Summing up, from this setting of APs and the communication of their clients, the following weighted graph, defined initially in [4], is derived:

Definition 1 (*The Weighted Digraph*) Consider a set of nodes (APs) $V = \{v_1, \dots, v_n\}$ located in a geographical area. Consider two nodes u and v, such that (the clients of) node v can receive transmissions from node u when (the server of) v is *OFF*, while node u is *ON*, with a quality above some lower bound value. Then, define a directed edge to exist from node u to node v, denoted as (u, v), of a positive weight $w(u, v) \in (0, 1)$ representing the quality of the received signal at node v.

Further assume that if w(u, v) > 0, then w(u, v) > 0 for all nodes $v, u \in V$.

An example weighted graph is illustrated in Fig. 1.





3.2.2 Time

As in [4], a basic unit of time period *T* is considered, e.g., 1 h, and we split the time period *T* into *x* smaller time slots T_1, T_2, \ldots, T_x , such that for each $T_k \in T$ there exists at least one node that in a group of nodes that may alternate between the *ON* and *OFF* node and remains *ON/OFF* for the whole time slot T_k . So, $\bigcup_{T_k} = T$ and $T_k \bigcap T_l = \emptyset$, $k \neq l$, i.e., the sum of the time slots is the time period *T* and no two time slots overlap. By $|T_k|$, we denote the time elapsed from the beginning of time slot T_k until the end of the time slot T_k .

Fix a time slot T_k . Then, denote:

 $ON(T_k) = \{ v \in V \mid \text{node } v \text{ is } ON \text{ in time slot } T_k \} \text{ and}$ $OFF(T_k) = \{ v \in V \mid \text{node } v \text{ is } OFF \text{ in time slot } T_k \}$

So, for node v, the time period T can be partitioned into two sets $ON_Time_T(v) = \{T_k \in T \mid v \in ON(T_k)\}$ and $OFF_Time_T(v) = \{T_k \in T \mid v \in OFF(T_k)\}$. Denote $Mode_T(v) = \{ON_T(v), OFF_T(v)\}$.

For the example graph of Fig. 1, $ON(t) = \{v, u_5, u_6\}$. For a sample time split of nodes v, u_1 , see Fig. 2.



Fig. 2 An example of a time split for nodes v, u_1 . In the shaded areas, the corresponding APs are ON

3.2.3 (Non-)Cooperative Neighbors

Given a weighted digraph *G*, for any node $v \in V$, Neigh(v) denotes the set of neighbors of node v in *G*, i.e., $Neigh(v) = \{u \in V \mid (u, v) \text{ and } (v, u) \in \vec{E}\}$. For the example graph of Fig. 1, $Neigh(v) = \{u_1, u_2, u_3, u_4, u_5\}$. Within a given time slot *T*, node *v* may be in *agreement* or in *cooperation* with some of its neighboring nodes. For any node *v*, being in agreement with node $u, u \neq v$ means that: (i) node *v* (or *u*) transmits *only* when *u* (respectively, *v*) does not transmit; and (ii) node *v* (*u*) serves the client(s) of the other node, in addition to its own client(s), during the time that it (*v* or *u*, respectively) is *ON*. The set of the neighbors of node *v* that are in agreement during time *T*, is denoted by $Coop_T(v) \subseteq Neigh(v)$. Thus, $ON_{_}$ $Time_T(v) \cap ON_{_}Time_T(u) = \emptyset$. On the other hand, the set of neighbors with which *v* is not in agreement with is denoted as $NCoop_T(v)$, where $NCoop_T(v) \subseteq$ Neigh(v) and it may be that $OFF_{_}Time_T(v) \cap OFF_{_}Time_T(u) \neq \emptyset$.

3.2.4 Experienced Quality

As in [4], for any node $v \in V$, the quality experienced by a node's clients can be quantified through measurements of the signal strength received at the clients of node v, at any time slot T_k during the time period T and it is denoted as $QoE_{T_k}(v)$. There are two cases: when the node is ON and when it is OFF.

• *ON Operation.* During any time slot *T_k*, if a node is *ON*, there are two possibilities: (i) none of its neighbors transmit at that time slot or (ii) some of them do. In the first case, the clients of the node experience the best quality (no interference), so we set the quality measure to be equal to a unity in this case, i.e.,

$$QoE_{T_k}(v) = 1. \tag{1}$$

In the second case, if the node is ON, some of its neighbors are ON during the time slot T_k . In this case, due to interference, the signal received is poorer and the experienced quality is degraded. This degradation is cumulative, i.e., increases with the total active time and signal strength of the competing nodes, captured by the weight w(u, v). Thus, in this case of the ON operation of node v the quality of its clients is denoted as:

$$QoE_{T_k}(v) = \max\{0, 1 - \sum_{\substack{u \in Neigh_T(v)\\ u \in ON(T_k)}} w(u, v)\}$$
(2)

• *OFF Operation.* When node v is *OFF*, the quality of the signal received at node v's clients, and hence the quality experienced, depends on the number of neighbors in cooperation with node v (i.e., in the set $Coop_{T_k}(v)$) that are *ON* at time T_k , and are serving the clients of node v as well as their own clients. If there exists only one such neighboring node u, the quality of the signal received at

node v is captured by the weight w(u, v) of edge (u, v). However, if there exist more than one neighboring nodes not in cooperation with node v that are ON at the same time, this results in interference received at node v, degrading the experienced quality at the node's clients.

As in [4], we assume that the service of v's clients is dominated by the *strongest* signal of its neighbors in the set $Coop_T(v_i)$, which is then degraded by the sum of the received signals from all the other neighbors that are ON at that same time (including other cooperating as well as the noncooperating neighbors). Thus, for the case of the *OFF* operation of node v, its experienced quality can be denoted as:

$$QoE_{T_{k}}(v) = \max\{0, (w(m, v) - \sum_{\substack{u \in Neigh_{T}(v), u \neq m \\ u \in ON(T_{k})}} w(u, v))\},$$
(3)

where $m = \arg \max_{\substack{u \in Coop_T(v) \\ u \in ON(T_k)}} \{w(u, v)\}.$

3.3 The Graphical Game

Using the above mathematical framework, we can now introduce resulting strategic game of [4]. This is a one-shot strategic game where the players are the nodes (APs). In any profile, given the decisions of the players (namely, whether to operate in *ON* or *OFF* mode) during each time slot $T_k \in T$, the utility of player $v \in V$ is equal to the quality of experience of its clients, i.e., $QoE_{T_k}(v)$. A formal definition of the game, defined initially in [4], follows:

Definition 2 Consider an one-shot strategic game Γ played on a weighted digraph $G(V, \vec{E}, \vec{W})$. The set of players of the game is *V*. A profile σ of the game is associated with the basic time period *T* of the scenario described. *T* is split into time slots T_1, T_2, \ldots, T_x , such that $\bigcup_{T_k} = T$ and $T_k \cap T_l = \emptyset, k \neq l, T_k$ corresponding to the time slot allowing alterations between *ON* and *OFF* operations of the nodes.

The strategy of any player (node) $v \in V$ in a profile σ is defined as follows:

$$\sigma_{v} = (ON_Time_{\sigma}(v), OFF_Time_{\sigma}(v), Coop_{\sigma}(v)),$$

where $ON_Time_{\sigma}(v) = \{T_{k} \in T \mid v \in ON_{\sigma}(T_{k})\}$
 $OFF_Time_{\sigma}(v) = \{T_{k} \in T \mid v \in OFF_{\sigma}(T_{k})\},$

and $ON_{\sigma}(T_k)$ ($OFF_{\sigma}(T_k)$) denotes the set of nodes in *V* that are in *ON* (*OFF*) operation during T_k according to σ . $Coop_{\sigma}(v) \subseteq Neigh(v)$ is the set of neighboring nodes of node *v*, with which node *v* has decided to cooperate with in σ . Cooperation

means that for each such cooperative neighbor u of node v, (i) $ON_Time_{\sigma}(v) \cap ON_Time_{\sigma}(u) = \emptyset$ and (ii) the two nodes are in agreement to serve each other's clients.

For player (node) v denote,

$$MaxCoop_{\sigma}(T_k, v) = \{m \in Neigh(v) | w(m, v) = \max_{\substack{u \in Coop_{\sigma}(v) \\ u \in ON_{\sigma}(T_k)}} w(u, v) \}$$

Then, the utility of player (node) v corresponding to the QoE for the ON/OFF operations of Eqs. (1)–(3), is given by:

$$U_{\boldsymbol{\sigma}}(v) = \sum_{T_k \in ON_Time_{\boldsymbol{\sigma}}(v)} ONU_{\boldsymbol{\sigma}}(v, T_k) \cdot |T_k| + \sum_{T_k \in OFF_Time_{\boldsymbol{\sigma}}(v)} OFFU_{\boldsymbol{\sigma}}(v, T_k) \cdot |T_k|$$
(4)

where the quality experienced for the ON operation of the node is determined by Eqs. (1) and (2):

$$ONU_{\sigma}(v, T_k) = \max\left\{0, 1 - \sum_{\substack{u \in Neigh(v)\\ u \in ON_{\sigma}(T_k)}} w(u, v)\right\}$$
(5)

and similarly

$$OFFU_{\sigma}(v, T_k) = \max\{0, (w(MaxCoop_{\sigma}(T_k, v), v) - \sum_{\substack{u \in Neigh(v)\\ u \in ON_{\sigma}(T_k)\\ u \neq MaxCoop_{\sigma}(T_k, v)}} w(h, i))\}$$
(6)

The linear model, which appears to simplify the relationship between interference and experienced quality provides, nevertheless, insight that allows for the problem at hand to be more generically understood, i.e., the fact that APs need to have strong links to each other in order to sustain cooperation. Future work will consider a more exact mathematical modeling of the relationship between interference and experienced quality.

Note that the time structure of the model is only a (theoretical) split (separation) of a continuum time period into a number of discrete time units. I.e., the model separates the whole time period into smaller continuous time units for purposes of ease of analysis, while overall the whole continuous time period is considered. Furthermore, while mixed strategies could be considered that would result in expected payoffs rather than deterministic payoffs, which are preferable for the specific case study.

4 Agreement Nash Equilibria for Clique Networks

Instances of the game Γ where the graph *G* is a clique of size *q* are mainly investigated in [5]. For simplicity, in this section, any time we refer to a graph we imply a clique graph. The work exploits the graph theoretical properties of clique topology in order to show when cooperation is feasible between the APs when they are placed close to each other.

4.1 Agreement Profiles

Next, a formal profile of cooperation between the players of the game, is defined. Then, for an arbitrary agreement profile, a corresponding simplified agreement profile (called *ordered agreement* profile) is considered, and equivalence regarding the utilities of the two players is shown (Sect. 4.2). Then, the corresponding simplified agreement profile is used to show stability (i.e., it is Nash equilibrium) of the original profile (Sects. 4.3 and 4.4). Specifically, they first define agreement profiles:

Definition 3 A profile σ is called **agreement** if all nodes of the graph *G* have agreed to cooperate with each other in σ .

They further define a subclass version of agreement profiles in which any AP has a single, continuing ON time slot. The profile is applicable for an arbitrary graph. The formal definition follows:

Definition 4 A profile σ is called **ordered agreement** if:

- it is an agreement profile; and
- assuming the time *T* is split into *n* time slots T_1, \ldots, T_n , where n = |V|, such that $T_i \cap T_j = \emptyset \ \forall i, j \in V$, any player $i \in V$ is *ON* at time slot T_i only and is *OFF* in all the remaining time slots $T_i, j \neq i$.

Applying ordered agreement profiles on clique graphs, observe:

Observation 1 Agreement profiles on Clique graphs Consider an agreement profile σ in the game Γ , over a clique graph G of size q. Then, consider any time slot $T_k \in ON_{\sigma}(T)$ and the corresponding (single) node k which is ON during T_k . Thus, by Eq. (5), $ONU_{\sigma}(k, T_k) = 1$. Also, the node k serves all other nodes of the clique during the time T_k . Thus, for any player $j \neq k$ (who is therefore OFF during time slot T_k), $OFFU_{\sigma}(j, T_k) = w(k, j)$.

4.2 Analysis of Agreement Profiles

Agreement profiles are equivalent to ordered agreement profiles as shown next. The equivalence relation is first defined:

Definition 5 A profile σ is said to be **equivalent** to another profile σ' if each player $i \in [q]$ has $U_{\sigma}(i) = U_{\sigma'}(i)$.

Equivalence between agreement and corresponding ordered agreement profiles is next shown, as presented in [5]:

Claim 1 Consider any agreement profile σ of a game Γ over a clique graph G of size q. Then, there exists an ordered agreement profile σ' , such that σ and σ' are equivalent.

Proof Consider any player *i* in the profile σ . Note that the time slots in which the player is ON may not be continuous. Let $T_{i_1} \cdots T_{i_x}$ be the time slots where the player *i* is ON in σ .

Note that, since σ is an agreement profile,

$$ONU_{\sigma}(i, ON_Time_{\sigma}(i)) = 1$$

and that for each $t_k \in ON_Time_{\sigma}(k), k \neq i$,

$$OFFU_{\mathbf{\sigma}}(i, t_k) = w(k, i).$$

Thus, summing up, from Eq. (4) and Observation 1, we get:

$$U_{\sigma}(i) = ONU_{\sigma}(i, ON_Time_{\sigma}(i)) \cdot |ON_Time_{\sigma}(i)| + \sum_{\substack{t_k \in ON_Time_{\sigma}(k)\\k \in [q] \setminus i}} OFFU_{\sigma}(i, t_k) \cdot |t_k|$$

$$= 1 \cdot |ON_Time_{\sigma}(i)| + \sum_{\substack{t_k \in ON_Time_{\sigma}(k)\\k \in [q] \setminus i}} w(k, i) \cdot |t_k|$$

$$= 1 \cdot |ON_Time_{\sigma}(i)| + \sum_{k \in [q] \setminus i} w(k, i) \cdot |ON_Time_{\sigma}(k)|$$

$$(7)$$

Now, construct a modified profile σ' , such that for each player $i \in [q]$, we set its ON mode to a *single* time slot T_i , which is equal to the total duration of its ON time slots in σ ; i.e., $|T_i| = |T_{i_1}| + |T_{i_2}| + \cdots + |T_{i_x}|$. Moreover, we start the ON time slot of the first player at time 0 and continue until time $|T_1|$, for the second player we start its (single) time slot at time $|T_1| + 1$ and continue until time $|T_1| + |T_2|$, and so on, so that player k is ON only during the time slot T_k and $|T_k| = |T_{k_1}| + |T_{k_2}| + \cdots + |T_{k_x}|$.

Note that, by construction, σ' is ordered, and it is an agreement profile (since the time slots allocated to different players are nonoverlapping). Moreover,

 $ON_Time_{\sigma'}(i) = ON_Time_{\sigma}(i)$ for each $i \in [q]$. It follows that for any player $i \in [q]$, $ONU_{\sigma'}(i, T_i) = 1$.

Also, note that for each one of the time slots $T_k \in T \setminus T_i$, only player k, which is ON and serves all other players. Thus, during the time slot T_k , the player i (which is in OFF mode) is served by the player k. So, $OFFU_{\sigma}(i, T_k) = w(k, i)$. It follows that,

$$\begin{split} U_{\sigma'}(i) &= ONU_{\sigma'}(i, ON_Time_{\sigma'}(i)) \cdot |ON_Time_{\sigma'}(i)| + \sum_{\substack{T_k \in ON_Time_{\sigma'}(k)\\k \in [q] \setminus i}} OFFU_{\sigma'}(i, T_k) \cdot |T_k| \\ &= 1 \cdot |ON_Time_{\sigma'}(i)| + \sum_{\substack{T_k \in ON_Time_{\sigma'}(k)\\k \in [q] \setminus i}} w(k, i) \cdot |T_k| \\ &= 1 \cdot |ON_Time_{\sigma'}(i)| + \sum_{\substack{k \in [q] \setminus i}} w(k, i) \cdot |ON_Time_{\sigma'}(k)|. \end{split}$$

Since $|ON_Time_{\sigma'}(i)| = |ON_Time_{\sigma}(i)|$ and for each player $k \in [q] \setminus i$, $|ON_Time_{\sigma'}(k)| = |ON_Time_{\sigma}(k)|$, the above equation becomes,

$$U_{\sigma'}(i) = 1 \cdot |ON_Time_{\sigma}(i)| + \sum_{k \in [q] \setminus i} w(k,i) \cdot |ON_Time_{\sigma}(k)|$$
(8)

By Eqs. (7) and (8), it follows that for $U_{\sigma'}(i) = U_{\sigma}(i)$, as claimed.

The following useful information is proved for agreement profiles [5]:

Claim 2 Assume an agreement profile σ for the game Γ over a clique graph G of size q. Then, for any player $i \in [q]$,

$$U_{\boldsymbol{\sigma}}(i) = |T_i| + \sum_{T_k \in T \setminus T_i} w(k, i) \cdot |T_k|.$$

Proof By Claim 1, assume without loss of generality that $\boldsymbol{\sigma}$ is an ordered agreement profile (equivalently, if σ' is the equivalent ordered agreement profile, then set $\boldsymbol{\sigma} = \sigma'$ without changing the utility values). Consequently, node *i* is the only node in *ON* operation during time slot T_i . Thus, by Eq. (5), $ONU_{\boldsymbol{\sigma}}(i, T_i) = 1$. Also, the node is *OFF* for the rest of the time and during any other time slot $T_k \in T \setminus T_i$ it is served by node *k* which is the only one node on *ON* operation during time slot T_k , for each $T_k \in T \setminus T_i$. Thus, by Eq. (6), $ONU_{\boldsymbol{\sigma}}(i, T_k) = w(k, i)$ for each $T_k \in T \setminus T_i$. We therefore obtain, by Eq. (4),

$$\begin{split} U_{\sigma}(i) &= \sum_{T_k \in ON_Time_{\sigma}(i)} ONU_{\sigma}(i, T_k) \cdot |T_k| + \sum_{T_k \in OFF_Time_{\sigma}(i)} OFFU_{\sigma}(i, T_k) \cdot |T_k| \\ &= ONU_{\sigma}(i, T_i) \cdot |T_i| + \sum_{T_k \in T \setminus T_i} OFFU_{\sigma}(i, T_k) \cdot |T_k|, \\ &= 1 \cdot |T_i| + \sum_{T_k \in T \setminus T_i} w(k, i) \cdot |T_k|, \end{split}$$

as claimed.

On the other hand, it is proved in [5]:

Claim 3 Assume an agreement profile σ of a game Γ over a clique graph G of size q. Then, no player can increase its utility by unilaterally decreasing its ON time and increasing its OFF time.

Proof By Claim 1, assume without loss of generality that σ is an ordered agreement profile. Consider any player $i \in [q]$. Then, by Claim 2, the node gets a utility of

$$U_{\sigma}(i) = |T_i| + \sum_{T_k \in T \setminus T_i} w(k, i) \cdot |T_k|.$$
(9)

If the player *i* increases its *OFF* period by some |t| > 0, it follows that its new *ON* time slot will be $|T'_i| = |T_i| - |t|$, thus $t \le |T_i|$. Let σ' be the resulting profile. Since σ is an agreement profile, there is no other player that is *ON* during time period *t*. So, its $OFFU_{\sigma'}(i)$ remains the same as in σ . Moreover, its $ONU\sigma'(i)$ must be decreased by *t*. Summing up, in the resulting profile its new utility becomes:

$$\begin{split} U_{\sigma'}(i) &= ONU_{\sigma'}(i, T'_i) \cdot |T'_i| + \sum_{T_k \in T \setminus T'_i} OFFU_{\sigma'}(i, T_k) \cdot |T_k| \\ &= 1 \cdot (|T_i| - |t|) + \sum_{T_k \in T \setminus T'_i} w(k, i) \cdot |T_k| \\ &= |T_i| - |t| + \sum_{T_k \in T \setminus T_i} w(k, i) \cdot |T_k| + 0 \cdot |t| \cdot \\ &= -|t| + \left(|T_i| + \sum_{T_k \in T \setminus T_i} w(k, i) \cdot |T_k| \right) \\ &= -|t| + U_{\sigma}(i). \end{split}$$

Since t > 0, it follows that $U_{\sigma}(i) > U_{\sigma'}(i)$, so that the player does not increase its utility by unilaterally increasing its *OFF* mode operation.

4.3 Necessary Conditions for Nash Equilibria

Using previous analysis on agreement profiles on clique graphs, the authors of [5] proceed to prove when such profiles are stable for the game, i.e., that when such profiles are Nash equilibria for the game. In particular, it is shown in this section that there exists a necessary condition for an agreement profile to be Nash equilibrium of the game. In the following sections, they provide sufficient conditions for stability of agreement profiles [5].

Proposition 1 Consider an agreement profile σ for a game Γ over a clique graph G of size q. Then, if σ is Nash equilibrium, then for any of the players $i, j \in [q]$, it holds that $w(i,j) \geq \frac{1}{2}$.

Proof Recall that by Claim 1, we may assume that σ is an agreement, continuous ON/OFF time slots profile. Also, by Claim 2, the utility of player *i* is:

$$U_{\sigma}(i) = |T_i| + \sum_{T_k \in T \setminus T_i} w(k, i) \cdot |T_k|.$$
⁽¹⁰⁾

Since σ is Nash equilibrium, no player that can increase its utility by unilaterally increase its *ON* time period. Assume now, by way of contradiction that there exists a player $i \in [q]$, such that $w(j, i) < \frac{1}{2}$ for some player $j \in [q]$. Assume now that the player *i* increases its *ON* time period to:

$$|ON_Time_{\sigma'}(i)| = |T'_i| = |T_i| + |t_i|,$$

where $t_j \in T_j$. Denote the resulting profile by σ' . Recall that node *j* was the only node in *ON* mode at the time slot T_j . Thus, during time t_j where the node *i* switches to be *ON*, the utility of node *i* is decreased due to the interference caused by node *j*. In particular, by Eq. (5), $ONU_{\sigma'}(i, T_i) = 1$ while $ONU_{\sigma'}(i, t_j) = (1 - w(j, i))$.

Moreover, observe that in σ' the player is *OFF* during time equal to:

$$OFF_Time_{\sigma'}(i) = T \setminus \{T_i \cup t_j\}$$

or

$$|OFF_Time_{\sigma'}(i)| = \sum_{T_k \in T \setminus T_i} |T_k| - |t_j|$$

For each $T_k \in T \setminus \{T_i \cup T_j\}$ where the node *i* is *OFF* and served by node *k*, by Eq. (6), its utility is $OFFU_{\sigma'}(i, T_k) = w(k, i)$. During the decreased time period $T_j \setminus t_j$, it also gets $OFFU_{\sigma'}(i, T_j \setminus t_j) = w(k, i)$.

In summary,

$$\begin{aligned} U_{\sigma'}(i) &= ONU_{\sigma'}(i, T'_i) \cdot |T'_i| + \sum_{T_k \in T'_i} OFFU_{\sigma'}(i, T_k) \cdot |T_k| \\ &= ONU_{\sigma'}(i, T_i) \cdot |T_i| + ONU_{\sigma'}(i, t_j) \cdot |t_j| + \sum_{T_k \in T \setminus \{T_i \cup T_j\}} OFFU(i, T_k) \cdot |T_k| \\ &+ OFFU(i, T_j \setminus t_j) \cdot (|T_j| - |t_j|) \\ &= 1 \cdot |T_i| + (1 - w(j, i)) \cdot |t_j| + \sum_{T_k \setminus T_i \cup T_j} w(k, i) \cdot |T_k| + w(j, i) \cdot (|T_j| - |t_j|) \\ &= |t_j| + U_{\sigma}(i) - 2w(j, i) \cdot |t_j|, \end{aligned}$$
(11)

by Eq. (10).

Since σ is Nash equilibrium, it must be that

$$U_{\sigma}(i) \ge U_{\sigma'}(i)$$

= $U_{\sigma}(i) - 2w(j,i) \cdot |t_j| + |t_j|,$

by Eq. (11). Thus, it must be that

$$2w(j,i)\cdot |t_j|\geq |t_j|.$$

Thus, it must be that $w(j,i) \ge \frac{1}{2}$, which is a contradiction since $w(j,i) < \frac{1}{2}$ by assumption. The proposition is therefore proved.

4.4 Sufficient Conditions for Nash Equilibria

In this section, we present the sufficient condition for an agreement profile to be Nash equilibrium of the game [5].

Theorem 1 Assume an agreement profile σ for the game Γ over a clique graph G of size q. If $w(i,j) \ge \frac{1}{2}$, for every pair $i, j \in [q]$, then σ is Nash equilibrium.

Proof Again, by Claim 1, we may assume without loss of generality that σ is an ordered agreement profile. By Claim 2, the utility of player *i* is:

$$U_{\sigma}(i) = |T_i| + \sum_{T_k \in T \setminus T_i} w(k, i) \cdot |T_k|.$$
(12)

Assume now that in contrary σ is *not* Nash equilibrium. Then, for at least one player $i \in [q]$, there exists an alternation of its ON/OFF time slots, so that he gets more (that is his utility is increased compared to σ). Note first that by Claim 3, the

node can not gain more by increasing its *OFF* time period. Accordingly, assume that the player unilaterally increases its *ON* time period as follows: $T'_i = T_i \bigcup_{\substack{t_k \in T_k, \\ T_k \in T \setminus T_i}} t_k$ where $t_i \ge 0$ for all k_i and $t_i \ge 0$ for at least one k_i . Thus

where $t_k \ge 0$ for all k, and $t_k > 0$ for at least one k. Thus,

$$|ON_Time_{\sigma'}(i)| = |T'_i| = |T_i| + \sum_{\substack{t_k \in T_k, \\ T_k \in T \setminus T_i}} |t_k|.$$

Recall that in σ , each node $k \in [q]$, is the only node on ON mode at the time slot T_k . Thus, during time $t_k \in T_k$ where the node *i* switches to ON operation in σ' , the node gets decreased utility due to the interference received by node *k*. In particular, by Eq. (5), $ONU_{\sigma'}(i, t_k) = 1 - w(k, i)$ for each such $t_k \in T_k$, while $ONU_{\sigma'}(i, T_i) = 1$.

Moreover, observe that in σ' , the player is *OFF* during time equal to:

$$OFF_Time_{\sigma'}(i) = T \setminus \{T_i \bigcup_{\substack{t_k \in T_k \\ T_k \in T \setminus T_i}} t_k\}$$

Thus,

$$|OFF_Time_{\sigma'}(i)| = \sum_{T_k \in T \setminus T_i} (|T_k| - |t_k|)$$

Thus, Eq. (4) becomes:

$$U_{\sigma'}(i) = \sum_{T_{k} \in ON_Time_{\sigma'}(v)} ONU_{\sigma'}(v, T_{k}) \cdot |T_{k}| + \sum_{T_{k} \in OFF_Time_{\sigma'}(v)} OFFU_{\sigma'}(v, T_{k}) \cdot |T_{k}|$$

$$= ONU_{\sigma'}(i, T_{i}) \cdot |T_{i}| + \sum_{\substack{T_{k} \in T_{k} \\ T_{k} \in T \setminus T_{i}}} ONU_{\sigma'}(i, t_{k}) \cdot |t_{k}| + \sum_{\substack{T_{k} \in T \setminus T_{i} \\ T_{k} \in T \setminus T_{i}}} OFFU_{\sigma'}(i, T_{k}) \cdot |T_{k}| - |t_{k}|)$$

$$= 1 \cdot |T_{i}| + \sum_{\substack{T_{k} \in T_{k} \\ T_{k} \in T \setminus T_{i}}} (1 - w(k, i)) \cdot |t_{k}| + \sum_{\substack{T_{k} \in T \setminus T_{i} \\ T_{k} \in T \setminus T_{i}}} w(k, i) \cdot (|T_{k}| - |t_{k}|)$$

$$= |T_{i}| + \sum_{\substack{T_{k} \in T_{k} \\ T_{k} \in T \setminus T_{i}}} |t_{k}| - 2 \cdot \sum_{\substack{T_{k} \in T_{k} \\ T_{k} \in T \setminus T_{i}}} w(k, i) \cdot |t_{k}| + \sum_{T_{k} \in T \setminus T_{i}} w(k, i) \cdot |T_{k}|$$
(13)

Since, by assumption, σ is not Nash equilibrium, it must be that $U_{\sigma'}(i) > U_{\sigma}(i)$. Thus, by Eqs. (12) and (13), combined with $U_{\sigma}(i) < U_{\sigma'}(i)$, it must be that

$$\begin{aligned} |T_i| + \sum_{T_k \in T \setminus T_i} w(k, i) \cdot |T_k| < |T_i| + \sum_{T_k \in T \setminus T_i \atop T_k \in T \setminus T_i} |t_k| - 2 \cdot \sum_{\substack{t_k \in T_k \\ T_k \setminus T_i}} w(k, i) \cdot |t_k| \\ + \sum_{T_k \in T \setminus T_i} w(k, i) \cdot |T_k| \end{aligned}$$

It follows that it must be that,

$$2 \cdot \sum_{\substack{t_k \in T_k \\ T_k \in T \setminus T_i}} w(k, i) \cdot |t_k| < \sum_{\substack{t_k \in T_k \\ T_k \in T \setminus T_i}} |t_k|$$
(14)

Since for all $i, j \in [q]$, $w(j, i) \ge \frac{1}{2}$,

$$2 \cdot \sum_{T_k \in T \setminus T_i} w(k,i) \cdot |t_k| \geq \sum_{T_k \in T \setminus T_i} |t_k|,$$

a contradiction to Eq. (14). It follows that σ is indeed Nash equilibrium.

4.5 A Characterization for Nash Equilibria

Proposition 1 implies that the condition $w(i,j) \ge \frac{1}{2}$, for any pair $i,j \in [q]$ is a necessary condition in order an agreement profile to be Nash equilibrium. Moreover, Theorem 1 implies that the condition is also a sufficient condition in order to have Nash equilibrium. Thus, based on [5],

Corollary 1 An agreement profile of the game Γ over a clique graph G of size q is Nash equilibrium if and only if $w(i,j) \ge \frac{1}{2}$, for all pairs $i,j \in [q]$.

5 Simulation Evaluation

The theoretical evaluation of the proposed model resulted in identifying necessary and sufficient conditions for cooperation to be possible. In this section, we demonstrate the potential benefits of cooperation by simulation. Specifically, we implement the scenario of densely deployed APs in a widely used network simulator (OPNET) and compare the performance, in terms of signal-to-noise ratio (SNR), between the default (noncooperative) case where all APs serve their own clients independently at the same time, and the cooperative case where all clients are served by one AP at a time.

The following scenarios are statically configured, i.e., the configuration of nodes does not change throughout the simulation (Fig. 3). A single content generating server in the background is assumed to be connected via a fast wired backhaul router to all the APs. Each scenario has an initial setup time, recorded in the results, which can be seen an oscillating initial behavior of the collected statistics. The authors do not consider this phase as part of the discussion of the collected results,



Fig. 3 The 9-AP configuration used by the simulation scenarios

as this is a simulator-generated behavior that is irrelevant to the effect of our model on the simulated behavior. At a time of 100 s after the start of the simulation, all the nodes start retrieving content from the server via FTP. All the wireless connections between the APs and the clients use standard Wi-Fi at a maximum 802.11 g rate (54 Mbps). All the active APs are configured to use the same resource (Channel 1), in order to study the interference resulting from simultaneous broadcasting of nearby APs due to lack of coordination.

Using this configuration, we study two separate scenarios:

- 1. the wireless clients are served by their own APs (default);
- 2. the wireless clients are served by the central AP (cooperative).

Furthermore, the simulations consider that there is no wireless noise, so a transmitted packet is always received perfectly by anyone in the range, unless a collision occurs.

We first show the results for the first (default) scenario. Figure 4 presents the SNR values recording in this case for each of the 9 players. Figure 5 presents the SNR of the second (cooperative) scenario. We observe that the SNR shows a sizeable increase (more for some users than others), clearly indicating the reduced interference in the environment, since the signal of the serving AP was not modified in any way.

To show that this has a positive effect on the user experience, we present in Figs. 6 and 7 the actual throughput results for the users in both scenarios. These figures



Fig. 4 The SNR values of the 9 users in a noncooperative scenario



Fig. 5 The SNR values of the 9 users in a cooperative scenario



Fig. 6 The throughput of the 9 users in a noncooperative scenario



Fig. 7 The throughput of the 9 users in a cooperative scenario

show that the increased SNR translates into a better throughput as well. Recall that the oscillating startup behavior is due to the simulation scenario setup phase, so the measurement of throughput is recorded after approximately 100 s when the application of FTP application is activated in the users.

6 Conclusions and Future Work

In this chapter, we considered a method to mitigate the interference caused by many individual wireless Access Points (AP) located in a dense area and transmitting using the same or overlapping channel via cooperation among the APs, such that they serve each other's clients at different times. We modeled the situation using a graphical game, particularly focusing on the case where the underlying graph is a clique with heterogeneous edge weights. We characterized the conditions under which agreements between all APs to jointly serve each other's clients are possible and achieve the maximum benefit for their clients.

Due to practical constraints, our results apply mainly in clique networks of small size. This is because Observation 1 relies on the assumption that in any time slot where an AP is ON, it can serve all the clients of other APs, which may be unrealistic for large cliques due to bandwidth limitations. While dense deployments resulting in large clique networks, where a client can be served by a large number of alternative APs with good signal strength, are arguably uncommon, an extension of our model to consider bandwidth-constrained equilibria that may arise in this case, as well as other more general graph topologies that fall short of a full clique, is left as a subject for future work.

The work presented in this chapter revealed that agreement behaviors of APs under some conditions are preferable for the APs. However, how such agreement behaviors are initiated and applied in practice, especially in a distributed manner, under no central authority, consists a next major step of our work, also considering security issues. Furthermore, investigating agreement profiles for deployments of APs that form other dense graph topologies are subject of future work. As a future work, one can target more complex simulation scenarios in terms of topologies and mix of applications as well as positioning of the users, and can experiment with a rotational scheme for serving users by the different APs in the neighborhood. This work is expected to lead to the design of a protocol leading to cooperation among neighboring APs.

Finally, let it be noted that various alternative physical layer technologies have been proposed in the recent literature to reduce interference between nearby APs, such as directional antennas and cooperative MIMO [8]. An extended model that would allow these ideas on cooperation among APs to be applied in combination with such technologies remains a topic for future work.

References

- 1. Akella A, Judd G, Seshan S, Steenkiste P (2005) Self-management in chaotic wireless deployments. In: Proceedings of in ACM MobiCom, pp 185–199
- 2. Antoniou J, Pitsillides A (2012) Game theory in communication networks: cooperative resolution of interactive networking scenarios. CRC Press, Hardcover, 152 pp
- 3. Antoniou J, Libman L, Pitsillides A (2011) A game-theory based approach to reducing interference in dense deployments of home wireless networks. In: 16th IEEE symposium on computers and communications (ISCC 2011)
- 4. Antoniou J, Lesta VP, Libman L, Pitsillides S (2012) Minimizing interference in unmanaged environments of densely deployed wireless access points using a graphical game model. In: 11th IEEE annual mediterranean ad hoc networking workshop (Med-Hoc-Net 2012), pp 75–82
- Antoniou J, Lesta VP, Libman L, Pitsillides A, Dehkordi H (2014) Cooperation among access points for enhanced quality of service in dense wireless environments. In: IEEE international symposium on a world of wireless, mobile and multimedia networks (WoWMoM 2014)
- 6. Axelrod RM (1984) The evolution of cooperation. BASIC Books, vol 4
- 7. Diestel R (2010) Graph theory. Springer, Heidelberg. ISBN 978-3-642-14278-9
- Gesbert D, Kountouris M, Heath Jt RW, Chae CB, Saizer T (2007) Shifting the MIMO paradigm: from single user to multiuser communications. IEEE Signal Process Mag 24 (5):36–46
- Grossman WM (2004) New tack wins prisoner's dilemma. http://www.wired.com/culture/ lifestyle/news/2004/10/65317. Accessed Oct 2004
- Hassan J, Sirisena H, Landfeldt B (2008) Trust-based fast authentication for multiowner wireless networks. IEEE Trans Mob Comput 7(2):247–261
- Karamchandani N, Minero P, Franceschetti M (2011). Cooperation in multi-access networks via coalitional game theory. In: Communication, control, and computing (Allerton), 49th annual allerton conference on, pp 329–336
- Kendall G, Yao X, Chong SY (2009) The iterated Prisoner's dilemma: 20 years on, ser. Advances In Natural Computation Book Series. World Scientific Publishing Co., 2009, vol 4
- Larsson E, Jorswieck E (2008) Competition versus collaboration on the MISO interference channel. IEEE J Sel Areas Commun 26(7):1059–1069
- Leshem A, Zehavi E (2007) Cooperative game theory and the gaussian interference channel. In: CoRR, vol abs/0708.0846
- 15. Li D, Xu Y, Liu J, Wang X, Wang X (2010) A coalitional game model for cooperative cognitive radio networks. In: Proceedings of the 6th international wireless communications and mobile computing conference, ser. IWCMC'10. ACM, pp 1006–1010
- Lopez L, Fernandez A, Cholvi V (2007) A game theoretic comparison of tcp and digital fountain based protocols. Comput Netw 51:3413–3426
- 17. Nash JF (1950) Equilibrium points in N-person games. In: Proceedings of National Academy of Sciences of the United States of America, vol 36, pp 48–49
- 18. Nash J (1951) Non-cooperative games. Ann Math 54(2):286-295
- 19. Peleg B, Sudholter P (2007) Introduction to the theory of cooperative games, 2nd edn. Springer, Heidelberg
- Rakshit S, Guha RK (2005) Fair bandwidth sharing in distributed systems: a game theoretic approach. IEEE Trans Comput 54(11):1384–1393
- Saad W (2010) Coalitional game theory for distributed cooperation in next generation wireless networks. Phd thesis, Department of Informatics, Faculty of Mathematics and Natural Sciences, University of Oslo
- 22. Saad W, Han Zhu, Debbah M, Hjorungnes A, Basar T (2009) Coalitional game theory for communication networks. Sig Process Mag IEEE 26(5):77–97

- 23. Singh C, Sarkar S, Aram A, Kumar A (2012) Cooperative profit sharing in coalition-based resource allocation in wireless networks. IEEE/ACM Trans Netw 20(1):69–83
- 24. Suris JE, DaSilva LA, Han Z, MacKenzie AB (2007) Cooperative game theory for distributed spectrum sharing. In: IEEE international conference on communications, pp 1006–1010
- 25. van de Nouweland A, Borm P, van Golstein Brouwers W (1996) A game theoretic approach to problems in telecommunication. Manage Sci 42(2):294–303

Simulating a Multi-Stage Screening Network: A Queueing Theory and Game Theory Application

Xiaowen Wang, Cen Song and Jun Zhuang

Abstract Simulation is widely used to study model for balancing congestion and security of a screening system. Security network is realistic and used in practice, but it is complex to analyze, especially when facing strategic applicants. To our best knowledge, no previous work has been done on a multi-stage security screening network using game theory and queueing theory. This research fills this gap by using simulation. For multi-stage screening, the method to determine the optimal screening probabilities in each stage is critical. Potential applicants may have access to information such as screening policy and other applicants' behaviors to adjust their application strategies. We use queueing theory and game theory to study the waiting time and the strategic interactions between the approver and the applicants. Arena simulation software is used to build the screening system with three major components: arrival process, screening process, and departure process. We use Matlab Graphic User Interface (GUI) to collect user inputs, then export data through Excel for Arena simulation, and finally export simulation from the results of the Arena to Matlab for analysis and visualization. This research provides some new insights to security screening problems.

X. Wang \cdot C. Song $(\boxtimes) \cdot$ J. Zhuang

X. Wang e-mail: xwang54@buffalo.edu

J. Zhuang e-mail: jzhuang@buffalo.edu

© Springer International Publishing Switzerland 2015 K. Hausken and J. Zhuang (eds.), *Game Theoretic Analysis of Congestion, Safety and Security*, Springer Series in Reliability Engineering, DOI 10.1007/978-3-319-13009-5_3

This research was partially supported by the United States National Science Foundation (NSF) under award numbers 1200899 and 1334930. This research was also partially supported by the United States Department of Homeland Security (DHS) through the National Center for Risk and Economic Analysis of Terrorism Events (CREATE) under award number 2010-ST-061-RE0001. However, any opinions, findings, and conclusions or recommendations in this document are those of the authors and do not necessarily reflect views of the NSF, DHS, or CREATE. Cen Song and Jun Zhuang are the corresponding authors.

Department of Industrial and System Engineering, University at Buffalo, Buffalo, USA e-mail: censong@buffalo.edu; songcen22@gmail.com

Keywords Security screening policy \cdot Two-stage queueing network \cdot Waiting time \cdot Game theory \cdot Imperfect screening

1 Introduction

1.1 Background

Nowadays, security screening are very important in many fields, including airport security screening [14], visa application [7], and customs inspection [15]. Screening process can not be perfect, and there exist *type I* and *type II* errors [5]. *Type I* error is the incorrect rejection of a true null hypothesis. In a screening system, if a good applicant is rejected, the approver is said to have a *type I* error. By contrast, *Type II* error is the failure to reject a false null hypothesis. In a screening system [21], if a bad applicant is approved, the approver is said to have a *type II* error. Because of these errors, it adds costs to the approver. Moreover, these errors could pose serious threats to the public, such as security and safety issues. Therefore, multi-stage screening process is introduced to reduce errors.

After the 9/11/2001 attack, the U.S. government requires 100 % scanning of all U.S. bound containers by radiation detection and nonintrusive inspection equipment at a foreign port before getting them loaded on vessels [3]. The Transportation Security Administration developed the Certified Cargo Screening Program for explosives on a passenger aircraft to get 100 % screening [25]. It results in a longer waiting time for the passengers to pass the security system. Moreover, reducing the probability of screening, although it would lead to less waiting time, may fail to catch some bad applicants. For each stage, low screening probability causes more errors while high screening probability causes congestion. In a visa application for a particular tourism destination, if a good applicant knew before he applied, that the waiting time is long, he would not apply. It is a loss for the destination country, because of the potential loss of economic contribution. In order to deter bad applicants and to attract good applicants to the most, a balance between congestion and intensity of screening should be achieved. In an airport security screening process, such a long security check time may result in passengers missing the flight. The flight schedule for those who missed will be rescheduled. It will result in seats unoccupied in the missed flight, which is such a waste. Meanwhile, the flight to which passengers are rescheduled to might not have enough seats for them, because company usually over sell tickets to maximize their benefits [6]. The passengers who got their itinerary changed are causing unbalance flows among the airline schedules. Due to butterfly effect [28], it will lead to many more problems in the future. According to the data of 1980-2012 annual passenger number for Newark Liberty International Airport [18], there are huge amount of people arriving at the airport. To avoid heavy congestion as well as to deter adversaries, a proper screening probability is required to the security screening.

Since 1970s, researchers have studied security screening with queueing models [9]. We study the multiple stages of screening system and find out how to predict applicants' behavior by applying queueing models and game theory. Based on the simulation results, we find the optimal strategy for the approvers.

Game theory is a study of strategic decision-making [17]. Strategic interdependence is present in a social situation when what is best for someone depends on others' choices. The best strategies for attackers and defenders are analyzed by balancing protection from terrorism and natural disasters and by considering resource allocation [31]. The optimal inspection policies for security agency balance the inspection probability and average delay time and consider the adversary strategic gaming behavior [8]. The optimal proportional inspection using game theory is analyzed to achieve the most cost-effective combination of deterrence and detection [2]. In the model, we assume that the decision makers are rational. Each player maximizes his payoff, given his beliefs about the strategies used by other players [24]. In this screening system, game players include applicants (good and bad) and approvers, whose actions impact each other. Knowing other players' potential best responses, players make their optimal decision. We apply game theory to construct the dynamic system of screening system in this paper.

Queueing theory is the mathematical study of waiting lines [23]. In queueing theory, a model is constructed so that queue lengths and waiting time can be predicted [23]. Single M/M/1, multiple M/M/c, single channel and multiple stages, and multiple channels and multiple stages models are used to estimate the truck delay in the seaport [29]. The M/M/N queuing models are developed to quantify truck congestion at primary scanning, and Monte-Carlo simulation is used to analyze the risk of containers missing vessels at secondary inspections [1]. An M/M/m queuing model is designed and applied into an airport security system to analyze the optimal number of security gates [20]. A queuing network and discrete event simulation are used to test the effects of baggage volume and alarm rate at the security screening checkpoint [4].

In general, there are three ways to study the phenomena (fact or occurrence): analytic modeling [22], simulation [12] and experiment [11]. As phenomena becomes more and more complex, analytic models may prove to be overly simplified, and some complex models cannot have analytical solutions. Meanwhile, sometimes experiments are not able to be performed or are too expensive to conduct. Simulation provides a way to meet our needs for cheaper, faster, and more practical data [16]. Compared to various simulation methods, computer simulation might be the most widely used one. Computer simulation is numerical evaluation using software designed to imitate the systems operation characteristics [10]. Different softwares can be chosen according to the different characteristics of the system. Pendergraft et al. [19] simulate an airport passenger and luggage screening security screening system in a discrete event way. We need to simulate screening system, which consists of queues and decision tree. It can be simulated by entity flows such as items or passengers constrained by conditions. In simulation, queueing models are often used for rough cut and condition setting [30]. We use queueing models to set conditions, making the screening system more accurate.

GUI in Matlab [13] provides point-and-click control of software application. Applied with user-defined functions, Matlab GUI can fulfill the following tasks without users knowing any command lines: data import/output, data analysis and plotting. In this paper, we apply GUI in Matlab for data analysis.

The rest of this paper is structured as follows: A description of the model is presented in Sect. 2. Designing Arena to simulate this screening system is discussed in Sect. 3. Matlab GUI design is discussed in Sect. 4. A numerical experiment and data analysis are provided in Sect. 5. The conclusion and discussion on some possible future work and application are provided in Sect. 6.

2 The Model

Figure 1 shows the flowchart of the screening system. Potential applicants classified as good and bad applicants decide whether to enter this system or not. Once they enter, they may go through several imperfect screening stages based on the screening probabilities at each stage. If they are screened, they would enter an M/M/1 service queue, which follows a first-in first-out rule. Based on the imperfect screening results, the suspected bad applicants are denied, while others are further to be determined to be screened or passed the system. Some applicants who are not screened at all are defined as good and they can pass the system.

2.1 Notation

Table 1 lists the notation that is used throughout this paper. We define the candidates as those who have the intention to enter the system but not certain enough to be classified as *potential applicants*. Those who exactly enter the system are defined as *applicants*. Applicants are divided into two groups: *good applicants* and *bad applicants*. There is probability P that the potential applicants are good. The potential applicants enter the system with a Poisson arrival rate of Λ , including



Fig. 1 Flowchart of screening system

	-
Λ	Potential applicants' total arrival rate
Р	Percentage of good potential applicants in potential applicants
μ	Service rate in screening point
$i = 1, 2, 3, \dots, N$	Stage number
Φ_1	Probability of screening in the first stage
Φ_{ig}	After passing the $(i - 1)$ th stage, probability of screening in <i>i</i> th stage
Φ_{ib}	After failed the $(i - 1)$ th stage, probability of screening in <i>i</i> th stage
α	Probability of reject good applicants
β	Probability of approving bad applicants
r_g	Reward of passing for good applicants
C _W	Waiting cost of good applicants
W	Total waiting time of good applicants
Wi	Waiting time in <i>i</i> th stage
r _b	Reward of passing for bad applicants
c _b	Cost of getting caught for bad applicants
R_g	Reward for approving good applicants for approver
R_b	Reward for denying bad applicants for approver
C_g	Cost for denying good applicants for approver
C_b	Cost for approving bad applicants for approver
N _{fb}	Simulation data: number of good denied
N _{fg}	Simulation data: number of bad approved
N _{rg}	Simulation data: number of good approved
N _{rb}	Simulation data: number of bad denied
U	Approver's utility
ug	Expected utility for good applicants before deciding entering
<i>u_b</i>	Expected utility for bad applicants before deciding entering
Pag	Calculated probability of passing for good applicants
Pab	Calculated probability of passing for bad applicants
P_{db}	Calculated probability of getting caught for bad applicants
<u>p1</u>	Represents Φ_1 in simulation
<u>p2</u>	Represents Φ_{2b} in simulation
<u>p</u> 3	Represents Φ_{2g} in simulation
pft	Represents P in simulation
tae	Represents α in simulation
tbe	Represents β in simulation
ти	Represents μ in simulation
lambda	Represents Λ in simulation
rewardg	Represents r_g in simulation
costw	Represents c_w in simulation
rewardb	Represents r_b in simulation
costb	Represents c_b in simulation

 Table 1
 The notation used in this chapter

good potential applicants' arrival rate of Λ_g and bad potential applicants' arrival rate of Λ_b . The one who takes charge of the system is defined as *approver*. The approver screens the selected applicants with a service rate of μ .

To simplify the model, we assume service rates at each stage are equal. The M/M/1 queueing model is applied to study screening process at each stage. The probabilities of screening at each stage are defined as $\Phi_i, \Phi_{ig}, \Phi_{ib}$ for i = 1, 2, ..., N. From the second stage, suspected good applicants and suspected bad applicants would have different screening probabilities. For suspected good applicants we use Φ_{ig} , while Φ_{ib} is for suspected bad applicants. When the applicants enter the system, the approver screens them according to the probability Φ_{ig} or Φ_{ib} . At stage 1, those who are not screened will pass immediately. At stage $i(1 \le i \le N)$, those who are not screened will be approved or denied immediately according to the last stage screening results. Those who are screened are facing three consequences after one of the three stages of screening: approved, denied or enter to next stage i + 1. At stage N, those who are not screened will be approved or denied immediately according to the last stage result. Those who are screened at the last stage will pass or fail according to their own attributes (good or bad).

We assume that *type I* and *type II* errors are the same at each stage. We define *type I* error probability as α , and *type II* error probability as β . While there is a probability α that good applicants would be screened as suspected bad applicants at each stage, there is a probability β of vice versa.

We define the good applicants as those who receive reward r_g when getting approved, and pay a cost c_w per unit waiting time. Similarly, we define the bad applicants as those who receive reward r_b when getting approved, and pay a penalty c_b while getting caught. As we assume that $r_b \gg c_w$, waiting cost is neglected at this case. On the other hand, we define the approvers as those who receive reward R_g when approving good applicants and reward R_b when denying bad applicants. The approver pays cost C_g when denying good applicants and cost C_b when approving bad applicants. We collect data after simulating for the number of good approved N_{rg} , the number of good denied N_{fb} , the number of bad approved N_{fg} and the number of bad denied N_{rb} . The approver's utility is defined as U. We define the calculated probability of passing for good applicants as P_{ag} and the waiting time in queue as w.

2.2 Payoffs of Applicants and Approver

Approver's expected utility is shown in Eq. (1), where he maximizes his utility payoff.

$$U = N_{rg}R_g + N_{rb}R_b - N_{fg}C_g - N_{fb}C_b \tag{1}$$

For applicants, we also use their utilities to scale their payoffs. Good applicants' utility is defined as u_g in Eq. (2), where he maximizes his utility payoff.

Simulating a Multi-Stage Screening Network ...

$$u_g = P_{ag}r_g - wc_w \tag{2}$$

$$P_{ag} = 1 - \left(\Phi_{1} \sum_{j=2}^{n} \left(\alpha^{j} (1-\alpha)^{n-j} \left(\prod_{i=2}^{n} \Phi_{ig} \left(\prod_{p=2}^{j} \frac{\Phi_{pb}}{\Phi_{pg}} \right) \right) \right) \right) \\ + \Phi_{1} (1-\alpha)^{n-1} \alpha \left(\prod_{i=2}^{n} \Phi_{ig} \right) + \Phi_{1} \alpha (1-\Phi_{2b}) \\ + \Phi_{1} \sum_{j=3}^{n} \left(\sum_{i=2}^{j-1} (1-\alpha)^{j-i-1} \alpha^{i} \left(\left(\prod_{k=2}^{j-1} \Phi_{kg} \left(\prod_{p=2}^{i} \frac{\Phi_{pb}}{\Phi_{pg}} \right) \right) (1-\Phi_{jb}) \right) \right) \\ + (1-\alpha)^{j-2} \alpha \left(\prod_{k=2}^{j-1} \Phi_{kg} \right) (1-\Phi_{jb}) \right)$$
(3)

To represent this series of problems, we use an M/M/1 queue as the queueing model, where there is a single server, unlimited waiting space, Poison arrival and exponential service time. Based on M/M/1 queue theory, at each screening point, we have waiting time W is shown in Eq. (4):

$$W = \frac{1}{\mu - \lambda} \tag{4}$$

For an N-stage screening, the total waiting time of good applicants is the summation of the screening waiting time at each stage, which are shown in Eq. (5).

$$w_{1} = \frac{1}{\mu - \Phi_{1}\Lambda}$$

$$w_{2} = \frac{1 - \alpha}{\mu - \Phi_{1}\Phi_{2g}\Lambda} + \frac{\alpha}{\mu - \Phi_{1}\Phi_{2b}\Lambda}$$

$$w_{3} = \frac{(1 - \alpha)^{2}}{\mu - \Phi_{1}\Phi_{2g}\Phi_{3g}\Lambda} + \frac{(1 - \alpha)\alpha}{\mu - \Phi_{1}\Phi_{2g}\Phi_{3b}\Lambda} + \frac{(1 - \alpha)\alpha}{\mu - \Phi_{1}\Phi_{2b}\Phi_{3g}\Lambda} + \frac{\alpha^{2}}{\mu - \Phi_{1}\Phi_{2b}\Phi_{3b}\Lambda}$$

$$\vdots$$

$$w_{n} = \frac{(1 - \alpha)^{n-1}}{\mu - \Phi_{1}\Phi_{2g}\Phi_{3g}\Phi_{4g}\dots\Phi_{ng}\Lambda} + \frac{(1 - \alpha)^{n-2}\alpha}{\mu - \Phi_{1}\Phi_{2b}\Phi_{3g}\Phi_{4g}\dots\Phi_{ng}\Lambda}$$

$$+ \frac{(1 - \alpha)^{n-2}\alpha}{\mu - \Phi_{1}\Phi_{2g}\Phi_{3b}\Phi_{4g}\dots\Phi_{ng}\Lambda} + \dots + \frac{(1 - \alpha)^{n-2}\alpha}{\mu - \Phi_{1}\Phi_{2g}\Phi_{3g}\Phi_{4g}\dots\Phi_{nb}\Lambda}$$

$$+ \frac{(1 - \alpha)^{n-3}\alpha^{2}}{\mu - \Phi_{1}\Phi_{2b}\Phi_{3b}\Phi_{4g}\dots\Phi_{ng}\Lambda} + \dots + \frac{\alpha^{n-1}}{\mu - \Phi_{1}\Phi_{2b}\Phi_{3b}\Phi_{4b}\dots\Phi_{nb}\Lambda}$$

$$w = \sum_{i=1}^{n} w_{i}$$
(5)


Fig. 2 Relationship between simulation and matlab GUI

Bad applicant's utility is defined as u_b is shown in Eq. (6), where he maximizes his utility payoff.

$$u_b = P_{ab}r_b - P_{db}c_b \tag{6}$$

We define the expected probability of passing for bad applicants as P_{ab} in Eq. (7), the expected probability of getting bad applicants caught as P_{db} in Eq. (8).

$$P_{ab} = 1 - \left(\Phi_{1}\sum_{j=2}^{n} \left((1-\beta)^{j}\beta^{n-j}\left(\prod_{i=2}^{n}\Phi_{ig}\left(\prod_{p=2}^{j}\frac{\Phi_{pb}}{\Phi_{pg}}\right)\right)\right) + \Phi_{1}\beta^{n-1}(1-\beta)\left(\prod_{i=2}^{n}\Phi_{ig}\right) + \Phi_{1}(1-\beta)(1-\Phi_{2b}) + \Phi_{1}\sum_{j=3}^{n} \left(\sum_{i=2}^{j-1}\beta^{j-i-1}(1-\beta)^{i}\left(\left(\prod_{k=2}^{j-1}\Phi_{kg}\left(\prod_{p=2}^{i}\frac{\Phi_{pb}}{\Phi_{pg}}\right)\right)(1-\Phi_{jb})\right)\right) + \beta^{j-2}(1-\beta)\left(\prod_{k=2}^{j-1}\Phi_{kg}\right)(1-\Phi_{jb})\right)$$

$$P_{db} = 1 - P_{ab}$$
(8)

By combining simulation and Matlab GUI, Fig. 2 shows the relationship between the tools we used.

3 Simulation

Arena is used to simulate the screening system. After each run, we get the numbers of applicants that are good but denied, bad but approved, good and approved and bad and denied. We define them as fake bad N_{fb} , fake good N_{fg} , real bad N_{rb} and real good N_{rg} . Figure 3 shows the whole structure of two-stage imperfect security screening process.



Fig. 3 Simulating two-stage imperfect security screening process

3.1 Simulating a Perfect One-Stage Screening System

Figure 3 part A excluding A2 shows the simulation of a one-stage imperfect screening process. First of all, we need to put a "Create" module named "Potential Arrival" to simulate potential applicants' arrival. It is Poison arrival as we described with the arrival rate of Λ , and "infinite" as max arrival. We define a variable "ar"

Name:		Entity Type:
Potential Arrival	•	Entity 1 👻
Time Between Arriva Type:	als Expression:	Units:
Expression	▼ POIS(0.001) ▼	Days 🔻
Entities per Arrival:	Max Arrivals:	First Creation:
1	Infinite	0.0

Fig. 4 Using "Create" module to simulate all applicants' arrivals

represent Λ , as Greek letters cannot be typed in Arena. "ar" has no value, and we assign a value to it in Sect. 5. Variable $\Lambda =$ "ar" per day holds the information that there is an average of "ar" persons arriving everyday. Figure 4 shows how to use "Create" module to simulate all applicants' arrivals. To set "Time Between Arrivals," there are several types to choose from: random[Expo], Schedule, Constant, and Expression. If you choose Expression, there are more Arena functions like WEIB and POIS you can choose from. In this case, we simulate a Poisson arrival process, then choose the "Type" of "Expression." The "Value Expression" box should be filled in with time value but not with rate value, "POIS(0.001)." The value of 0.001 is an approximate value, we may adjust it later in Sect. 5. Units should be defined correctly, so we put in "Days" here. In the last row, we define one entity at each arrival. "Max Arrival" is "Infinite" and first "Creation" is zero. The above simulates that every potential applicant follows a Poisson process, in an average of every 0.001 day to arrive.

Then, we divide potential applicants into two categories: good and bad. A "Decide" module named "type divide" can fulfill this task. We set it as "2-way" by chance, where chance is p, we use "pft" to present because p is a reserved variable in Arena. This module makes "pft" of potential applicants as good ones, followed by an "Assign" module named "good," which gives attribute "type" a value of 0. In some cases, entities in simulation need attributes to differentiate them from others. However, unlike many other program languages, these attributes can only be given values, but strings are prohibited. In this case, we give every entity an attribute named "type." Another "Assign" module named "bad" is added after the "False" output of "type divide" module, to assign attribute to bad applicants. If "type" equals zero, it means this entity represents good applicants. If "type" equals one, it means this entity represents bad applicants.

Figure 5 shows how to simulate two categories: good applicants and bad applicants. The variable "pft" can be found in "Variable" module, which is shown in Fig. 6. The "Initial Value" can be left blank for reading data from files in the

	Name:				Type:		
type devide	type d	evide		-	2-way by	Chance	•
	Percer	t True (0-100):					
0 9	. pft		- %				

Fig. 5 Simulating arrivals of good and bad applicants

Initi	al Values							vau P	8
	60								
4			_	_	_	_			
Vanab	e - Basic Name	Rows	Columns	Data Type	Clear Option	File Name	Initial Values	Report Statistics	ç.
1.	pft			Real	System		1 rows	Г	
2	p1			Real	System		1 rows	ir i	
							and the second se		
3	tae			Real	System		1 rows	Г	

Fig. 6 Setting probability of good applicants as a variable p

-



good	
Assignments:	
Attribute, type, 0	Add
	Edit
	Delet
04 0	ncel He

future. We assign a value, for example, "pft = 60", which means that there is a probability of 60 % that a potential applicant is good. The numbers by the modules stand for the numbers of entity going through this route. Figure 7 takes good

2

applicants as an example, and shows how to assign attribute and differentiate the two groups of potential applicants.

In Fig. 3, part A1 simulates potential arriving applicants and divides them into two types. The key part of the system is screening. The model we discuss here is a twostage screening model. We start from first-stage screening. It consists of two "Decide" modules and a "Process" module. One of the "Decide" module named "first stage" is to set as "2-way" by chance, whose chance is Φ_1 . This module makes Φ_1 of applicants step into the screening process. Meanwhile, the rest of applicants would get an immediate pass. The "Process" module named "screen" is set the action as "Seize Delay Release," having "1" resource as the approver with a service rate μ (Exponential Distribution). The "Seize" represents the process of getting a free resource. An entity might have to wait in a queue until the resource is free. "Delay" represents the process time, and "Release" corresponds to the action of an entity releasing a resource, so that the next entity in queue could use the resource. This module represents the screening process, using M/M/1 queue. At the end of screening, the approver has to decide whether to approve or deny according to the type of the applicants (good or bad). The other "Decide" module named "pass" is put here, setting as "2-way" by condition, where condition is "type ≤ 0 ". We use the variable "p1" refers to Φ_1 to decide the screening chance to the applicants at the first stage.

In Fig. 3, part A2 shows the simulation of the potential applicants' decisionmaking. To simulate congestion cost, we assume potential applicants quitting causes cost. We design the condition nodes to simulate this. The condition nodes have the condition $u_g/u_b > 0$, where u_b and u_g represent the expected utilities for good applicants and bad applicants, respectively. There are two ways to simulate the decision: Static and Dynamic. In this paper, we simulate the static decision. There is no update information for the later potential applicants before entering. Assuming the potential applicants know the information about the screening system and use the information to make decision on entering the system, we apply queueing theory to obtain waiting time and probability knowledge to obtain utility. For the applicants, if utility is greater or equal to zero, they will enter the system. Only potential good applicants will pay waiting cost, because we assume $r_b \gg c_w$ for bad applicants. We have w, P_{ag} , P_{ab} , and P_{db} for two-stage screening system in Eqs. (9)–(12), respectively.

$$w = \frac{1}{\mu - \Phi_1 \Lambda} + \frac{1 - \alpha}{\mu - \Phi_1 \Phi_{2g} \Lambda} + \frac{\alpha}{\mu - \Phi_1 \Phi_{2b} \Lambda}$$
(9)

$$P_{ag} = 1 - \left(\Phi_1 \Phi_{2g} \alpha (1 - \alpha) + \Phi_1 \Phi_{2b} \alpha^2 + \Phi_1 (1 - \Phi_{2b}) \alpha\right)$$
(10)

$$P_{ab} = 1 - \left(\Phi_1 \Phi_{2g} \beta (1-\beta) + \Phi_1 \Phi_{2b} (1-\beta)^2 + \Phi_1 (1-\Phi_{2b}) (1-\beta)\right)$$
(11)

$$P_{db} = \Phi_1 \Phi_{2g} \beta (1-\beta) + \Phi_1 \Phi_{2b} (1-\beta)^2 + \Phi_1 (1-\Phi_{2b})(1-\beta)$$
(12)

Name:			Type:	
good decide to	enter	~	2-way by Condition	
lf.				
Expression	~			
Value:				
rewardg * (1 ·	(p1 * p3 * tae * (1 · ta	e) + p1 * p2 * tae * tae •	+ p1 * (1 · p2) * tae)

Fig. 8 Simulating entering condition for potential good applicants

Meanwhile, if $(\mu - \Lambda \Phi_1)$ or $(\mu - \Lambda \Phi_1 \Phi_{2g})$ or $(\mu - \Lambda \Phi_1 \Phi_{2b})$ equals to zero, waiting time will be infinite. There is no point to entering the system, and it would eventually cause error in Arena. We add a new decision node to waive it out. Three "Decide" modules named as "good decide to enter", "bad decide to enter" and "waive out" are added before the screening process.

Figure 8 shows the simulation of the decision making for good applicants. For the flow of good potential applicants, they need to go through a "waive out," setting as "2-way by condition", where condition is $(\mu - \Lambda \Phi_1)$ or $(\mu - \Lambda \Phi_1 \Phi_{2g})$ or $(\mu - \Lambda \Phi_1 \Phi_{2b})$ equals to zero. The values for μ and Λ can be put in according to Sect. 5. If the specific condition is 'yes,' the entity will leave the system, which means the good applicants will quit applying. If the specific condition is "no", the entity will go through "good decide to enter," setting as "2-way by condition," where condition is $u_g \ge 0$. If the specific condition is 'yes,' the entity will continue to stay in system, which means the good applicant will finally decide to apply. If the specific condition is "no", the entity will leave the system, which means the good applicant will guit applying. The two conditions in "waive out" and "good decide to enter," are set as "if" Expression. We need to use Expression Builder in Tools menu to formulate the "Value" of the "Expression." The "Value" of "Expression" for "waive out" is mu - lambda * p1 == 0 ||mu - lambda * p1 * p3 == 0 ||mu - lambda * p1 * p3 == 0 ||mu - lambda * p1 * p3 == 0 ||mu - lambda * p1 * p3 == 0 ||mu - lambda * p1 * p3 == 0 ||mu - lambda * p1 * p3 == 0 ||mu - lambda * p1 * p3 == 0 ||mu - lambda * p1 * p3 == 0 ||mu - lambda * p1 * p3 == 0 ||mu - lambda * p1 * p3 == 0 ||mu - lambda * p3 == 0 ||mu - lambdalambda * p1 * p2 == 0. The "Value" of "Expression" for "good decide to enter" is written in the phrase of Arena as: rewardg * (1 - (p1 * p3 * tae * (1 - tae) + p1 * tae))p2 * tae * tae + p1 * (1 - p2) * tae))-costw * (1/(mu - lambda * p1) + (1 - lambda + p1))tae)/(mu - p1 * p3 * lambda) + tae/(mu - p1 * p2 * lambda)) > = 0.

Figure 9 shows the simulation of the decision-making for potential bad applicants. For the flow of bad potential applicants, they need to go through "bad decide to enter," setting as "2-way by condition," where condition is $u_b \ge 0$. If the specific condition is 'yes,' the entity will continue to stay in system, which means the bad applicants can apply. If the specific condition is 'no,' the entity will leave the system, which means the bad applicant would quit applying. The "Value" of "Expression" for "bad decide to enter" is written in the phrase of Arena as:

			Decide			?
Name:					Туре:	
bad decide t	o enter			~	2-way by Co	ondition
lf:						
Expression	~					
Value:	2					
rewardb * [1 · (p1 * p3	* tbe * (1 ·	be)+p1 * p2 *	(1 - tbe)	* (1 - tbe) +	e1 * (1

Fig. 9 Simulating entering condition for potential bad applicants



 $\begin{aligned} rewardb * (1 - (p1 * p3 * tbe * (1 - tbe) + p1 * p2 * (1 - tbe) * (1 - tbe) + p1 * \\ (1 - p2) * (1 - tbe))) - costb * (p1 * p3 * tbe * (1 - tbe) + p1 * p2 * (1 - tbe) * \\ (1 - tbe) + p1 * (1 - p2) * (1 - tbe)) > &= 0. \\ \end{aligned}$

Figure 10 simulates the screening process. In the "Logic" group, "Action" can be defined as Delay, Seize Delay, Seize Delay Release and Delay Release. We choose "Seize Delay Release". Then we add resource, which represents the server named as Approver, quantity = 1. "Delay Type" can be defined as Constant, Normal, Triangular, Uniform and Expression. As the server follows exponential distribution, we choose "Expression" and define the same unit as the arrival applicants: "Days." Allocation is Value Added. Expression is "EXPO(0.004)". The number 0.004 represents the service time in an average of every 0.004 day. We can adjust $\mu = \frac{1}{0.004}$ as in Sect. 5. Figure 11 simulates the approver that passes or fails applicants after screening. It depends on the attributes of applicants. The "Decide"

Name: Type: Pass Zway by Control	ndition
Pass Pass 2-way by Co	ndition
It: Named:	ls:
Attribute v type v	<= •
Value:	
0	

Fig. 11 Simulating the approver's decision on whether to pass or reject

module is defined as "2-way by Condition." If "Attribute: type ≤ 0 , which is type = 0," then good applicants can pass. Otherwise, it is bad applicants, who are to be rejected. All the entity flows end up in "Dispose" module, which represents the completion of the process.

3.2 Simulating an Imperfect Multi-stage Screening System

In Fig. 3, part A shows the simulation of the first stage imperfect screening process. In Fig. 3, part A3 simulates the *type I* and *type II* errors. Since there exist *type I* and *type II* errors, we use two "Decide" modules to model and simulate them. After a former "Decide" module named "pass," the entities will be divided into two categories: good applicants "(type = 0)" and bad applicants "(type = 1)". Following the entity flow of good applicants, we put a "Decide" module named "error adjustment g", setting as "2-way" by chance, where chance is $1 - \alpha$. It represents that the approver has *type I* error, leading to $100(1 - \alpha)$ % of good applicants, we put a "Decide" module named is $(1 - \beta)$. It will represent the approver have *type II* error, leading to $100(1 - \beta)$ % of bad applicants as bad and 100β % of bad applicants as good. We add two variables: "tae" to represent α and "tbe" to represent β . Good applicant group and bad applicant group have been updated to new *good applicants* group (good and bad), and new *bad applicant* group (good and bad).

To eliminate error, we apply multi-stage screening. In this simulation, we carry out two-stage screening to analyze. The two new groups of applicants step into second-stage screening. Both of them have the same scenario of modules as at the first-stage screening, expect that the screening probability is Φ_{2g} (for new *good ones*)/ Φ_{2b} (for new *bad ones*) instead of Φ_1 . In Fig. 3, part B shows the simulation of the second-stage imperfect screening. We use "p2" and "p3" to represent Φ_{2g} and Φ_{2b} , respectively. "Decide" Modules "second stage for good," "screen 2g," "pass 2g," "error adjustment gg" and "error adjustment bg" simulate second-stage screening process for good ones. Modules "second stage for bad," "screen 2b", "pass 2b," "error adjustment gb" and "error adjustment bb" simulate second-stage screening process for bad ones.

3.3 Designing Input and Output Functions in Arena

In this simulation, the inputs include service rate (μ) , percentage of potential good applicants in total potential applicants (pft), screening probabilities (p1, p2 and p3), *type I and type II* error (*tae* and *tbe*), rewards and costs for potential good applicants and bad applicants (*rewardg*, *costw*, *rewardb*, and *costb*). The outputs are numbers $(N_{fb}, N_{fg}, N_{rg} \text{ and } N_{rb})$. For the convenience of data analysis in Matlab, we add output for recording current screening probabilities, service and arrival rates.

Due to the large amount of input and output data, several "ReadWrite" modules and "Record" modules are added. "ReadWrite" module can be considered as a bridge between *Arena* and *Microsoft Excel*. For each Excel file read to Arena, it is called "Arena File Name". "ReadWrite" module is not read directly from Excel file name, but from "Arena File Name." In "File" module, there is a table of all the Excel File names that particular "Arena File Name." For input data, the type of module is set as "Read from File". The action to "Assignments" is put in each settings of the "ReadWrite" module. Data that read from this module assigns value to the "Assignments". Therefore, there are eleven "ReadWrite" modules for input. We name them as "ReadPft", "ReadMu", "ReadP1," "ReadP2," "ReadP3," "ReadCostb." Figure 12 shows how to set the modules for input. Figures 13 and 14 list the modules and files in the simulation. In Fig. 3, part C shows how to read the input.

Arena can save data in a .csv file, which can be opened in Excel. We need to collect four groups of data: N_{fb} , N_{fg} , N_{rg} and N_{rb} . To identify the four flows of entities, we add two "Decide" modules named "defergood" and "deferbad". After the two-stage screening processes, entities that the approver think as good applicants, go through "defergood." Its type is 2-way by "Condition," where condition is "type ≤ 0 " In other words, if its "type = 0," then it is a good applicant, called *realgood* here we count it to N_{rg} , on the other hand, if its "type = 1," then it is a bad applicant, called *fakegood*, here we count it to N_{fg} . Entities that approver thinks are bad applicants go through "deferbad." Its type is "2-way by Condition," where condition is "type ≥ 1 " In other words, if "type = 0," then it's a good applicant, called *fakebad*, we count it to N_{fb} , where as if its "type = 1," then it's a bad applicant, called *realgood*, we count it to N_{rb} . Figure 15 takes "defergood" as an example, which shows how to divide the entities.

	Name:			
PoodD#	ReadPft			
(eaurit	Туре:		Arena File Nar	me:
· · ·	Read from File	-	File 4	
	Recordset ID:		Record Numb	er:
1	Recordset 5	-	1	
adTae •	Assignments:			
	Variable, pft			Add
			1	Edit
	-			
			1	Delete
 ReadCo 	S			

Fig. 12 An example for "ReadWrite" module setting

	Name	Access Type	Operating System File Name	End of File Action	Initialize Option	Recordsets
1	File 1	Microsoft Excel 2007 (*.xlsx)	D:/thesis/documenting/server_parameters.xlsx	Dispose	Hold	3 rows
2)	File 2	Microsoft Excel 2007 (*.xlsx)	D:\thesis\documenting\screening_probabilities.xlsx	Dispose	Hold	3 rows
3	File 3	Microsoft Excel 2007 (*.xlsx)	D:\thesis\documenting\applicants_utility_parameters.xlsx	Dispose	Hold	4 rows
	File 4	Microsoft Excel 2007 (*.xlsx)	D:\thesis\documenting\probability_of_good.xisx	Dispose	Hold	1 rows

Fig. 13 List of "Files" in Arena

	Name	Туре	Arena File Name	Recordset ID	Record Number	Assignments
1 🕨	ReadPft	Read from File	File 4	Recordset 5	1	1 rows
2	ReadP1	Read from File	File 2	Recordset 9	1	1 rows
3	ReadMu	Read from File	File 1	Recordset 6	1	1 rows
4	ReadRewardg	Read from File	File 3	Recordset 1	1	1 rows
5	ReadP2	Read from File	File 2	Recordset	1	1 rows
6	ReadP3	Read from File	File 2	Recordset	1	1 rows
7	ReadTae	Read from File	File 1	Recordset 7	1	1 rows
8	ReadTbe	Read from File	File 1	Recordset 8	1	1 rows
9	ReadCostw	Read from File	File 3	Recordset 3	1	1 rows
10	ReadRewardb	Read from File	File 3	Recordset 2	1	1 rows
11	ReadCostb	Read from File	File 3	Recordset 4	1	1 rows





Fig. 15 Dividing entities for four groups of entities

Statistic - Advanced Process						
	Name	Туре	Expression	Report Label	Output File	
1 🕨	realgood	Output	defergood.N	realgood	D:\thesis\documenting\realgood.csv	
2	fakegood	Output	defergood.N	fakegood	D:\thesis\documenting\fakegood.csv	
3	realbad	Output	deferbad.Nu	realbad	D:\thesis\documenting\realbad.csv	
4	fakebad	Output	deferbad.Nu	fakebad	D:\thesis\documenting\fakebad.csv	

Fig. 16 List of outputs in "Statistic"

To collect the entity number that goes through "defergood" and "deferbad," we use "Statistic" module. In the "Statistic" table, we have four outputs: "realgood, fakegood, realbad, and fakebad." Their type is "Output", and expressions are "defergood.NumberOut True, defergood.NumberOut False, deferbad.NumberOut True, deferbad.NumberOut False." "Output File" will be in a .csv file with the path, as shown in Fig. 16.

3.4 Setting up Simulation

We collect data including N_{fb} , N_{fg} , N_{rg} and N_{rb} based on the difference in screening probabilities to run the simulation. There are p1, p2 and p3 screening probabilities throughout the imperfect two-stage screening system. Starting from 0 %, we take 5 %

each step to reach next level of screening probability. There are total $21 \times 21 \times 21 =$ 9, 261 sets of screening probabilities that needs to be run. We make use of replications in Arena to use the data from recording. In each replication, we change a set of screening probabilities. The replication length is the simulating experiment period. To make change to set of screening probabilities in every replication, we assign probabilities from particular row in the data column, which is row "c." We define "c" as equal to the current replication number, written as c = NREP.

Matlab is used to generate an excel of 9,261 rows \times 3 columns of probabilities. We can make another set of p1, p2 and p3 to do the simulation by generating a new set of inputs. This will be explained in Sect. 4. However, if the set of inputs changes, the run times should change as well. Figure 3 shows the simulation of a two-stage imperfect security screening system.

4 Designing a Graphic User Interface with Matlab

After coding, GUI is friendly for users to complete the tasks. We add the following functions to GUI: generating input for simulation, pulling data from simulation results, analyze data, and generate graphs.

First, we design the function modules to generate potential applicants arrival rate and service rate. Second, we design module to pull data from simulation and to find optimization and its condition. Third, we design an analyzing module for sensitive analyses. Last but not least, we design a record module for analyzed data. To realize these functions, we first create a new GUI, which goes to "HOME \rightarrow New \rightarrow Graphical User Interface" and naming it *screening*. Then, we divide the Blank GUI figure into 3 function areas by adding 3 "Button Group," named "Generating Inputs," "Optimization," and "Data Analysis."

4.1 Generating Input Parameters for Simulation

The input we need to generate includes the following: screening service rate μ , reward and cost for good applicants and bad applicants r_g , r_b , c_b , c_w , type I and type II error α , β , percentage of good applicants in total P, and screening probabilities Φ_1 , Φ_{2g} , Φ_{2b} . We use "static txt" to label inputs and "edit txt" for user to specify inputs. In total, 23 "static txt"s and 11 "edit txt"s are added to "Generating Inputs" Button Group in Fig. 17a.



Fig. 17 Generating inputs with approver's preference in GUI. a Generating inputs in GUI. b Defining the Approver's Preference

4.2 Calculating the Optimal Strategies Using Numerical Methods

In this section, we pull raw data and check their usability and then profile them. We add 4 "edit box"s to define approver's preference R_g , R_b , C_g and C_b for "Optimization" as shown in Fig. 17b.

Data reflecting N_{fb} , N_{fg} , N_{rg} and N_{rb} are saved in a .csv file. We find the raw data from Arena that have even rows writing 0 and odd rows writing data, which we need to profile.

After profiling the data, we can generate a matrix of utility according to different set of data. Then, we use "Max" function to find the optimization set of data from the matrix. Results of optimization will be shown on screen. In the meantime, a graph reflecting all the data in matrix is drawn to show how screening probability affects the approver's utility. We use green dots to show all the data points, and red diamond to point out the best strategy.

4.3 Designing Output Data Analysis

Sensitivity analysis is the study of how the uncertainty in the output of a mathematical model or system (numerical or otherwise) can be apportioned to different sources of uncertainty in the input. We want to see, once Φ_1 , Φ_{2g} or Φ_{2b} is fixed to 100 %, how the other screening probabilities affect the approver's utility and the

Fig. 18 Data visualization and analysis in Arena

Fixed P1 2-s.	. 💌	GO
Fixed P1 2-st	age screen	ning for U
Fixed P1 2-st	age screen	ning for Nrb
Fixed P2 2-st	age screen	ning for U
Fixed P2 2-st	age screen	ning for Nrb
Fixed P3 2-st	age screen	ning for U
Fixed P3 2-st	age screen	ning for Nrb
Compare per	fect and im	perfect screening
Compare 1-s	age and 2-	stage screening
Utility Range	n 1-stage i	mperfect screening

number of bad applicants getting caught. We add an Axe in GUI and put "Pop-up Menu" in "Data Analysis" function area. "Pop-up Menu" can give cases which we can call by adding a push button "GO" in its callback. Figure 18 shows the "Data Analysis" Part.

The layout of GUI is shown below in Fig. 19.



Fig. 19 Layout of the matlab GUI

(13)

5 Numerical Experiments

5.1 Data Sources for Input Parameters

We use the baseline according to the paper [27] to do a new numerical experiment. For good applicants, we have $r_g = 20$ and $c_w = 10$. For bad applicants, we have $r_b = 20$ and $c_b = 100$. For approver, we have $R_g = 5$, $R_b = 10$, $C_g = 3$, and $C_b = 20$. Using the data from Newark Liberty International Airport [18], we estimate the arrival rate and the service rate. The average number of arriving passengers is in Eq. (13).

$$\begin{split} N_{\text{arrive per year}} &= (34014027 + 33711372 + 33107041 + 33424110 + 35366359 \\ &+ 36367240 + 35764910 + 33078473 + 31893372 + 29428899 \\ &+ 29220775 + 31100491 + 34188701 + 33622686 + 32575874 \\ &+ 30945857 + 26626231 + 22255002)/(18) \\ &= 32038400 \end{split}$$

There are three terminals Terminal A, B, and C in this airport. Usually, five securities opening in a day in each terminal are expected. To simulate M/M/1 queue, the arrival rate for modeling can be defined in Eq. (14).

$$\Lambda = \frac{N_{\text{arrive per year}}}{\text{days of a year } \times \text{ servers}} = \frac{32038400}{365 \times 15} = 5851.76 \approx 5852$$
(14)

In Sect. 3.1, we have $1/5852 \approx 0.00017$, then the arrival setting is POIS (0.00017), Unit is "Days". We round it down to 0.0001 and round it up to 0.0002, to run twice.

According to experiences from airport security screening, the average screening time is 5 min per person. We estimate the service rate following the equation below:

$$\mu = \frac{\text{Minutes in a day}}{\text{Service Time per Person}} = \frac{60 \times 24}{5} = 288$$
(15)

Because $1/288 \approx 0.0034$, in Sect. 3.1, the service setting in screening is EXPO (0.0034), Unit will be "Days". We round it up to 0.004 and round it down to 0.003, to run twice.

We have two sets of Λ and μ to do the numerical experiment. They are $\Lambda = 10,000, \ \mu = 250$ with $\frac{\Lambda}{\mu} = 40$ and $\Lambda = 5000, \ \mu = 333$ with $\frac{\Lambda}{\mu} = 15$. We record entities going through the module to record data in a .csv file. We do 9,261 replications to set up different combinations of probabilities for each simulation and the replication length is 30 days.

5.2 Simulation Results: Optimal Strategies and Payoffs

We set $\frac{\lambda}{\mu} = 40$, where $\Lambda = 10,000$ and $\mu = 250$. We put arrival as POIS(0.001) and service as EXPO(0.04), where we have the best screening strategies for two-stage imperfect screening process with $\Phi_1 = 10\%$, $\Phi_{2b} = 10\%$, $\Phi_{2p} = 15\%$ and U = 9,505.

Red Diamond point shows the best strategy point. Figure 20 shows how Φ_{2b} and Φ_{2b} affect the approver's utility. Figure 20a shows an interesting jump when $\Phi_{2g} = \Phi_{2b}$. It seems that when the second stage has the equal probability of screening for potential good applicants and potential bad applicants, it's the worst case with smallest approver's utility. It neither attract potential applicants nor does any good to the approver. When $\Phi_{2g} > \Phi_{2b}$, the utility is greater than the case of $\Phi_{2b} > \Phi_{2g}$. In this case, approver should put much effort on screening probability in Φ_{2g} .

We set $\frac{\lambda}{\mu} = 15$, where $\Lambda = 5,000, \mu = 333$. We put arrival as POIS(0.002) and service as EXPO(0.03), where we have the best strategies for two-stage imperfect screening process with $\Phi_1 = 100 \%$, $\Phi_{2b} = 0 \%$, $\Phi_{2g} = 0 \%$, and U = 10,430. Figure 20b shows an interesting jump when $\Phi_{2g} = \Phi_{2b}$. It seems that when the second stage has the equal probability of screening for potential good applicants and potential bad applicants, it is the worst case with smallest approver's utility. It does not attract potential applicants or does any good to the approver. When $\Phi_{2g} > \Phi_{2b}$, the utility is greater than the one when $\Phi_{2g} > \Phi_{2b}$. Approver should put much effort on screening probability Φ_{2g} . We find that with two different ratio of λ and μ , the results are consistent. We can conclude that when $\Phi_{2g} = \Phi_{2b}$ is the worst case for approver; Since second-stage screening, Φ_{2g} is more important, on which approver should focus more.



Fig. 20 Illustration of approver's utility affected by screening probabilities $\Phi_{2g}(P_3)$ and $\Phi_{2b}(P_2)$. **a** $\lambda/\mu = 40$, **b** $\lambda/\mu = 15$



Fig. 21 Illustration of approver's utility affected by screening probabilities $\Phi_{2g}(P3)$, fixing $\Phi_{2b} = 10\%(P2)$ and $\Phi_1 = 10\%(P1)$ with $\frac{\lambda}{u} = 40$

In addition, we do another simulation to see how Φ_{2g} impacts utility when fixing $\Phi_1 = 10\%$, and $\Phi_{2b} = 10\%$. This is based on the optimal strategy for approver when $\frac{\lambda}{\mu} = 40$. Figure 21 shows that when $\Phi_{2g} \le 15\%$, approver's utility ≤ 0 ; when $\Phi_{2g} > 15\%$, approver's utility > 0. In this case, if the approver does not want to screen all applicants, they can screen 15 % of good applicants at the second stage, who are defined at the first stage. The results are based on this particular set of parameters. Changing input parameters may change the results accordingly.

6 Conclusion and Future Research Directions

In this research, we develop several modules in both Arena and Matlab to simulate and analyze an imperfect multi-stage screening system with screening errors. By setting conditions for decision modules and locating modules in different positions, we are able to control the applicants' flow in the system according to prespecified conditions. We use different settings and positions of decision nodes to control the arrival rate, classify good and bad applicants, and simulate type I and type II errors. We use replications to repeat simulation and get statistically significant results, with different sets of input to the same model. We can access and control the screening process based on the user/approver's preferences. By considering the potential applicants best responses and rates of arrival and service, the approver is able to find the best screening strategy to maximize her utility, based on simulation results.

To our best knowledge, this is the first research using simulation, queueing theory, and game theory to study a complex multi-stage screening system. In the future, we could extend this work to study more complex models with various distributions of approval and service processes. Besides, we can extend the simulated model as multi-stage screening imperfect model. On the other hand, the fast development of smart phones and social media makes it possible to use dynamic and real-time data to simulate and update the security screening process. For example, there is a recent smart phone-based app which enables passengers to post their waiting time in line for security screening at airports [26]. This enables both other passengers and the approver to get more accurate and timely information about the security screening. Future research could model and simulate dynamic systems, considering waiting time for the bad applicants.

References

- Bennett AC, Chin YZ (2008) 100 % container scanning: security policy implications for global supply chains. Master's thesis, Massachusetts Institute of Technology, Engineering Systems Division
- Bier VM, Haphuriwat N (2011) Analytical method to identify the number of containers to inspect at us ports to deter terrorist attacks. Ann Oper Res 187(1):137–158
- 3. Directorate-General Energy and Transport (2009) The impact of 100 transport. http://ec. europa.eu/transport/modes/maritime/studies/doc/2009_04_scanning_containers.pdf. Accessed July 2014
- 4. Dorton SL (2011) Analysis of airport security screening checkpoints using queuing networks and discrete event simulation: a theoretical and empirical approach. Embry-Riddle Aeronautical University, Daytona Beach
- Eckel N, Johnson W (1983) A model for screening and classifying potential accounting majors. J Acc Educ 1(2):57–65
- Fleming A (2014) Oversold flights, getting bumped and bumping. http://airtravel.about.com/ od/travelindustrynews/qt/bumped1.htm. Accessed July 2014
- 7. Gasson S, Shelfer KM (2007) IT-based knowledge management to support organizational learning: visa application screening at the INS. Inf Technol People 20(4):376–399
- Gaukler GM, Li C, Ding Y, Chirayath SS (2011) Detecting nuclear materials smuggling: performance evaluation of container inspection policies. Risk Anal 32(3):65–87
- 9. Gilliam RR (1979) Application of queuing theory to airport passenger security screening. Interfaces 9(4):117–123
- 10. Kelton WD, Sadowski RP, Swets NB (2010) Simulation with arena. McGraw-Hill, Higher Education, Boston
- 11. Kronik L, Shapira Y (1999) Surface photovoltage phenomena: theory, experiment, and applications. Surf Sci Rep 37(1):1–206
- Maddox MW, Gubbins KE (1997) A molecular simulation study of freezing/melting phenomena for lennard-jones methane in cylindrical nanoscale pores. J Chem Phys 107 (22):9659–9667
- 13. MathWorks (2014) MATLAB GUI. http://www.mathworks.com/discovery/matlab-gui.html. Accessed July 2014

- McCarley JS, Kramer AF, Wickens CD, Vidoni ED, Boot WR (2004) Visual skills in airportsecurity screening. Psychol Sci 15(5):302–306
- Merrick JRW, McLay LA (2010) Is screening cargo containers for smuggled nuclear threats worthwhile? Decis Anal 7(2):155–171
- Musa JD (2004) Software reliability engineering: more reliable software, faster and cheaper. Tata McGraw-Hill Education, New York
- 17. Myerson RB (1997) Game theory: analysis of conflict. Harvard University Press, Cambridge
- Newark Liberty International Airport (2014) Facts and information. http://www.panynj.gov/ airports/ewr-facts-info.html. Accessed July 2014
- Pendergraft DR, Robertson CV, Shrader S (2004) Simulation of an airport passenger security system. In: Proceedings of the 36th conference on Winter simulation, pp 874–878
- Regattieri A, Gamberini R, Lolli F, Manzini R (2010) Designing production and service systems using queuing theory: principles and application to an airport passenger security screening system. Int J Serv Oper Manage 6(2):206–225
- 21. Roxy P, Devore JL (2011) Statistics: the exploration and analysis of data. Oxford University Press, Oxford
- 22. Shustorovich E (1986) Chemisorption phenomena: analytic modeling based on perturbation theory and bond-order conservation. Surf Sci Rep 6(1):1–63
- 23. Sundarapandian V (2009) Queueing theory. probability, statistics and queueing theory. PHI Learning, New Delhi
- 24. Tadelis S (2013) Game theory: an introduction. Princeton University Press, Princeton
- Transportation Security Administration (2013a) Certified cargo screening program. http:// www.tsa.gov/certified-cargo-screening-program. Accessed July 2014
- 26. Transportation Security Administration (2013b) My TSA mobile application. http://www.tsa. gov/traveler-information/my-tsa-mobile-application. Accessed July 2014
- 27. Wang X, Zhuang J (2011) Balancing congestion and security in the presence of strategic applicants with private information. Eur J Oper Res 212(1):100–111
- Wolfram S (2002) Some historical notes. http://www.wolframscience.com/reference/notes/ 971c Accessed July 2014
- 29. Yoon D (2007) Analysis of truck delays at container terminal security inspection stations. Ph.D. thesis, New Jersey Institute of Technology, New York
- 30. Zeltyn Marmor YN, Mandelbaum et al SA (2011) Simulation-based models of emergency departments: operational, tactical, and strategic staffing. ACM Trans Model Comput Simul (TOMACS) 21(4):24
- Zhuang J, Bier VM (2007) Balancing terrorism and natural disasters-defensive strategy with endogenous attacker effort. Oper Res 55(5):976–991

A Leader–Follower Game on Congestion Management in Power Systems

Mohammad Reza Salehizadeh, Ashkan Rahimi-Kian and Kjell Hausken

Abstract Since the beginning of power system restructuring and creation of numerous temporal power markets, transmission congestion has become a serious challenge for independent system operators around the globe. On the other hand, in recent years, emission reduction has become a major concern for the electricity industry. As a widely accepted solution, attention has been drawn to renewable power resources promotion. However, penetration of these resources impacts on transmission congestion. In sum, these challenges reinforce the need for new approaches to facilitate interaction between the operator and energy market players defined as the generators (power generation companies) in order to provide proper operational signals for the operator. The main purpose of this chapter is to provide a combination of a leader-follower game theoretical mechanism and multiattribute decision-making for the operator to choose his best strategy by considering congestion-driven and environmental attributes. First the operator (as the leader) chooses K strategies arbitrarily. Each strategy is constituted by emission penalty factors for each generator, the amount of purchased power from renewable power resources, and a bid cap that provides a maximum bid for the price of electrical power for generators who intend to sell their power in the market. For each of the K strategies, the generators (as the followers) determine their optimum bids for selling power in the market. The interaction between generation companies is modeled as Nash-Supply Function equilibrium (SFE) game. Thereafter, for each of the K strategies, the operator performs congestion management and congestiondriven attributes and emission are obtained. The four different attributes are congestion cost, average locational marginal price (LMP) for different system buses,

Department of Electrical Engineering, College of Engineering, Marvdasht Branch, Islamic Azad University, Marvdasht, Iran e-mail: mohamadreza.salehizadeh@gmail.com

A. Rahimi-Kian Smart Networks Lab, School of ECE, College of Engineering, University of Tehran, Tehran, Iran

K. Hausken University of Stavanger, Stavanger, Norway

© Springer International Publishing Switzerland 2015 K. Hausken and J. Zhuang (eds.), *Game Theoretic Analysis of Congestion, Safety and Security*, Springer Series in Reliability Engineering, DOI 10.1007/978-3-319-13009-5_4

M.R. Salehizadeh (🖂)

variance of the LMPs, and the generators' emission. Finally, the operator's preferred strategy is selected using the Technique for Order Preference by Similarity to Ideal Solution (TOPSIS). The proposed procedure is applied to the IEEE reliability 24-bus test system and the results are analyzed.

Keywords Energy • Power systems • Independent system operator • Generators • Electricity market • Transmission congestion management • Leader–follower game

Nomenclature

$C(P_{it})$	Generator <i>i</i> 's cost function for power production when
	the operator chooses strategy $t, t = 1,, K$
$C_E(P_{it})$	Generator i 's emission cost function when the operator
	chooses strategy $t, t = 1, \dots, K$
ψ_{it}	Emission penalty factor imposed by the operator on the
	<i>i</i> th generator, $i = 1,, g$ when the operator chooses
	strategy $t, t = 1, \ldots, K$
r	Number of renewable power resources, $r \ge 0$
$P_{\text{ren},jt}$	Operator's amount of purchased renewable power from
	the resource located in <i>j</i> th bus in Mega Watt when the
	operator chooses strategy $t, P_{\text{ren}, jt} \ge 0, j = 1,, N$,
	$t = 1, \ldots, K$
$P_{\text{ren},t}$	Operator's amount of purchased renewable power when
	the operator chooses strategy $t, P_{\text{ren},t} \ge 0, P_{\text{ren},t} =$
	$\sum_{n=1}^{N} p_{n} + 1 = K$
	$\sum_{i=1}^{N} P_{\text{ren}, jt}, t = 1, \dots, \mathbf{K}$
$\beta_{\max t}$	Operator's market bid cap for limiting electricity price
/ max,t	on electrical power when the operator chooses strategy
	$t, t = 1, \dots, K$
a_i, b_i, c_i	Generator i's cost function coefficients
a_{Ei}, b_{Ei}, c_{Ei}	Generator i's emission cost function coefficients
P_{it}	Power produced by the <i>i</i> th generator, as determined by
	the operator to maximize social welfare within con-
	straints when the operator chooses strategy $t, t = 1$,
	, <i>K</i>
P_i^{\max}	Maximum power produced by the <i>i</i> th generator
P_i^{\min}	Minimum power produced by the <i>i</i> th generator
$P_{Th \ ki}^{\max}$	Thermal power flow limit of the transmission line
111,19	between buses k and j, $k, j = 1,, N$
$P_{St ki}^{\max}$	Stability power flow limit of the transmission line
5.,.9	between buses k and j, $k, j = 1,, N$
Ν	Number of transmission buses

$P_{line,kjt}$	Power flows across the transmission line between buses	
	k and $j, k, j = 1,, N$ when the operator chooses strat-	
	egy $t, t = 1, \ldots, K$	
LMP_{kt}	Locational marginal price of bus $k, k = 1,, N$ when	
	the operator chooses strategy $t, t = 1, \ldots, K$	
CC_t	The operator's decision matrix M 's x_{1t} element, i.e.,	
	congestion cost when the operator chooses strategy	
	$t, t = 1, \ldots, K$	
ave_LMP_t	The operator's decision matrix M 's x_{2t} element, i.e.,	
	average locational marginal price (LMP) for different	
	system buses when the operator chooses strategy	
	$t, t = 1, \ldots, K$	
var_LMP_t	The operator's decision matrix M 's x_{3t} element, i.e.,	
	variance of the LMPs when the operator chooses	
	strategy $t, t = 1, \dots, K$	
emission _t	The operator's decision matrix M 's x_{4t} element, i.e., the	
	g generators' emission when the operator chooses	
	strategy $t, t = 1, \dots, K$	
g_n	Number of generators connected to bus n	
D	Total electricity demand of all consumers	
d_n	Number of electricity demands connected to bus	
	$n, \sum_{k=1}^{N} d_n = D$	
P_{nD_k}	Active power consumption of the <i>k</i> th electricity demand	
R	connected to bus $n, k = 1, \dots, d_n, n = 1, \dots, N$	
SW_t	Social welfare when the operator chooses strategy	
	$t, t = 1, \dots, K$	
g	Number of generators	
β_{it}	The <i>i</i> th generator's market bid for trading his electrical	
	power in the market when the operator chooses strategy	
	$t, t = 1, \ldots, K$	
κ_t	The distance of the strategy $t, t = 1,, K$ from the	
	negative ideal strategy over the sum of distances of this	
	strategy from positive and negative ideal strategy	
λ_t	The estimated market clearing price for electrical power	
	when the operator chooses strategy $t, t = 1,, K$	
Y_i	The operator strategies' <i>i</i> th attribute, $i = 1,, 4$, i.e.,	
	the <i>i</i> th column of the operator's $K \times 4$ decision matrix	
	М	
Α	Operator's 4×4 comparison matrix of attributes	
М	Operator's $K \times 4$ decision matrix	
$S_{t}^{*} = \left[\psi_{1t}^{*}, \dots, \psi_{22t}^{*}, \right]$	The operator's preferred strategy	
$P_{max}^{*}, \beta_{max}^{*}$		
ICH, I MAX, I		

$S = \left\{ S_t = [\psi_{1t}, \dots, \psi_{3,t}, \right.$	The operator's set of K strategies	
$P_{\text{ren},t}, \beta_{\max,t}]$		
$\forall t = 1, \dots, K$		
PI_t	Performance index of the power system when the	
operator chooses strategy $t, t = 1,, K$		

Acronyms

Supply function equilibrium
Transmission congestion management
Technique for order preference by similarity to ideal solution
Locational marginal price
Positive ideal strategy
Negative ideal strategy

1 Introduction

1.1 Problem Statement

In most countries around the globe, electrical energy has been sold in a monopoly form for years. For example, in the United States, three sections of the power sector, i.e., generation, transmission, and distribution were traditionally bundled in a vertically integrated unit form, see Borenstein [5]. This form of electricity energy trading decreased the incentive for efficient power system operation, see Kirschen and Strbac [17]. From the late 1980s, deregulation and liberalization to introduce less restrictive regulation frameworks were introduced within the electricity industry, see Lai [23]. As a result, the vertically integrated unit form of the power sector became unbundled and restructured. This restructuring has increased the competition in the power market, and increased the need for a game-theoretic analysis of the conflicting and partly conflicting interests between the operator and generators, and between each generator. Since the early days of restructuring, competition of the market players defined as the generators (power generation companies), each maximizing his own profit, has posed various security and reliability challenges, and various other challenges, for the operator, consumers, and society at large.¹ The operator is a nonprofit organization maximizing social welfare and guaranteeing system reliability and security. In order to accomplish such

¹ In some references in the energy market literature "independent system operator" is used instead of "operator," and Generation Company or GenCo has been used instead of "generator." In order to make the chapter more readable, we use "operator" and "generator" throughout the chapter.

objectives, the operator performs congestion management for restricting power flows across transmission lines within their thermal and security limits. One of the major security challenges within the production of electrical power is transmission congestion which has been explored widely in the literature, see Kumar et al. [21], Zhang et al. [49], Kumar et al. [21]. The importance of transmission congestion management (TCM) has been highlighted recently by high and increasing penetration of renewable power resources at the transmission and distribution level, see Ahmadi and Lesani [2] and Kunz [22]. Renewable power is produced by resources which are naturally replenished on a human timescale. Examples are sunlight, wind, rain, tides, waves, and geothermal heat. The impact of renewable energies on TCM has been studied for the case of Germany in Kunz [22]. Besides the generators' competition and decisions, the operator's strategic choice, which impacts the generators, affect transmission congestion, as shown in Porter [33]. Although many approaches to TCM are possible and may generate various insights, this chapter provides a game-theoretic analysis to capture how the strategic interaction between one operator and many generators affects transmission congestion. The main purpose of this chapter is to provide a combination of a leader-follower game theoretical mechanism and multiattribute decision-making for the operator to choose his preferred strategy by considering congestion-driven and environmental attributes.

1.2 Literature Review

In general, "congestion" is defined as "an excessive accumulation" [30]. This phenomenon occurs frequently in different systems such as computer networks, see Telang et al. [42], urban traffic systems, see Sun et al. [41], wireless networks, see Dong et al. [11], social networks, see Wang et al. [47], Bier et al. [3], as well as power systems, see Bompard et al. [4]. In power systems, congestion is defined as a situation in which power flows across a transmission line exceeds its thermal or stability limits that could compromise system safety and cause system breakdown. As an example, see Kaplan [16] and Lin et al. [26], for more details about the Northeastern USA blackout in 2003. Due to the effects of congestion on system reliability and security, a proper strategy for "congestion management" needs to be developed. TCM refers to the preventive/corrective actions performed by the operator in order to mitigate the adverse impact of overload in the transmission, see Zou et al. [53]. Breakdown of congestion management has severe negative impact on safety and security. By maintaining transmission lines' power within their thermal and security limits, TCM improves the power system's safety and security. The importance of TCM is getting highlighted in the era of postrestructuring where competition between generators maximizing their own profit without considering power system reliability and security cause important challenge for the operator. In the pool form of a wholesale electricity market, various generators participate in the electricity market through offering their bids for trading electrical power as a commodity. Based on the provided bids, the operator clears the market and performs TCM. Providing proper methods for modeling competition between generators is important not only from an economic point of view, but also for ensuring that the power system operates acceptably with respect to reliability and security. Many attempts have been made to model players' competition in the electricity market, see Hobbs et al. [14], Hobbs [13], Son et al. [40], and de la Torre et al. [10]. However, only a few attempts have considered transmission congestion management while modeling competition in the electricity market, see Veit et al. [43], Krause et al. [19], Liu et al. [27], Sahraei-Ardakani and Rahimi-kian [37], Coneio et al. [9], and Lee [24]. In Veit et al. [43], the intense effect of intranetwork congestion on the bidding strategy of the players in the German network has been demonstrated. In Krause and Andersson [19], through an agent-based simulator, different TCM schemes assuming a perfect and oligopolistic structure have been evaluated. Also, in order to model the behavior of the players, a Q-learning approach has been deployed. In Liu et al. [27], through simulation, the effect of network constraints as well as congestion on Nash equilibrium has been demonstrated. For this purpose, a two-level optimization problem was modeled wherein in the first level the operator dispatch generation was included, and in the second level the generators adopted the Nash-SFE equilibrium, see Liu et al. [27].

In the original SFE concept, as introduced by Klemperer and Meyer [18], the goal was to determine a supply function, not in the sense of a specific mathematical expression, but to calculate the points of the function by solving a set of differential equations. The transfer of the original concept to real life applications, such as electricity markets and systems, narrowed down the calculations to a set of parameters of predefined functions, usually linear or quadratic. That makes these problems parametrical SFE models. A significant issue, regarding game models, is the existence and uniqueness of the solution. Since the first introduction of the SFE in game theory modeling, the concerns regarding existence and uniqueness of solution have not been overcome. A wide range of possible implementations and applications in real life problems have been considered.

It is still the case that, depending on the number of generators engaged in the game, the complexity of the overall problem, the similarities, and the uniformities of each individual problem, the SFE game may render multiple solutions (equilibria). It is not clear which of these multiple solutions is most qualified to represent the generators' strategic behavior. To date, only under very strong assumptions have SFE problems been solved when applied to real cases. Existence and uniqueness of a solution are very hard to prove, and often available only for very simple versions of the SFE model. Sensitivity analysis is often useful in order to examine the models' robustness and solvability. Further elaboration upon this is beyond the scope of this chapter, and sensitivity analysis is left for future research. Here, we assume that generators are able to determine a SFE.

In Sahraei-Ardakani and Rahimi-Kian [37], a SFE-based dynamic replicator model of generator's bids is developed. It is shown that when congestion occurs, some generators may increase their bids three times. In order to avoid the mixed strategy Nash-Cournot Equilibrium when congestion occurs, a leader–follower game theoretical approach was developed in Lee [24]. In Lee [24], it was assumed

that the generator at the receiving area of a congested line is the leader and the generator at the sending area has the follower position. Since transmission congestion may cause mixed strategy equilibrium in the Cournot model, Lee [25] proposed a leader–follower game. In contrast, our procedure is to provide a signal based on congestion-driven and environmental attributes enabling the operator to choose his best strategy.

In spite of drawing attention to the congestion problem in electricity market modeling in Veit et al. [43], Krause and Andersson [19], Liu et al. [27], Sahraei-Ardakani and Rahimi-Kian [37]. Conejo et al. [9], and Lee [24], there is a need to develop a procedure for giving the operator the possibility to select his strategy for transmission congestion management. In order to fulfill this need, we present a leader-follower game-theoretic approach for TCM in this present chapter. The method should be capable to consider both congestion-driven indices and emission. Game-theoretic approaches have been widely used for market modeling in the literature, see Ventosa et al. [44], Saguan et al. [36], and Zeng et al. [51]. As a practical study in this area, Lise et al. [29] has presented a game theoretical model for assessing the impact of competition on economic and environmental attributes of the electricity market. The model has been calibrated for eight Northwestern European countries, see Lise et al. [29]. However, there are few research works which have been devoted to game theoretical approaches to transmission congestion management, see for example Lee [24]. In this chapter, our emphasis on considering emission as an attribute besides congestion-driven indices is due to the current situation of the modern electricity industry. The electricity industry is considered as one of the largest producers of greenhouse gases. Carbon dioxide (CO₂), methane (CH₄), nitrous oxide (N₂O), and water vapor make up an important part of greenhouse gas emissions from the electricity industry, see (http://www.epa. gov/climatechange/ghgemissions/sources/electricity.html). In 2012, the electricity industry emitted 32 % of the greenhouse gases in the United States, see (http:// www.epa.gov/climatechange/ghgemissions/sources/electricity.html). The emission percentages from the other industries are 28 % transportation, 20 % other industries, 10 % commercial and residential, and 10 % agricultural, see http://www.epa.gov/ climatechange/ghgemissions/sources/electricity.html. Hence, the electricity industry is intended to be considered at the center of focus for emission reduction. This environmental concern besides the congestion problem in power systems, which might restrict cleaner generators to trade their energy, is considered in this chapter. The details of the proposed method are described in Sect. 2. Also, the relevant features of this chapter in comparison with other papers are depicted in Table 1. Most of the techniques in Table 1 account for system congestion, and thus these procedures can be easily adapted to include TCM signals. However, in none of them has a mechanism for the system operator been provided to evaluate his chosen strategy based on environmental and congestion-driven attributes.

References	Electricity market modeling technique	Considering congestion?	Providing a signal for TCM?
Hobbs et al. [15]	Cournot and supply function game and linear complementarity programming	No	No
Son et al. [40]	Hybrid coevolutionary programming	No	No
Veit et al. [43]	Agent-based	Yes	No
Krause and Andersson [19]	Q-learning	Yes	No
Liu et al. [27]	Two-level optimization and Nash-SFE	Yes	No
Sahraei-Ardakani and Rahimi-Kian [37]	Nash-SFE	Yes	No
Lee [24]	Leader-follower and Nash-Cournot	Yes	No
This chapter	Leader-follower and Nash-SFE	Yes	Yes

Table 1 Relevant features of the selected research works

1.3 Main Contributions and Structure of the Chapter

The main contribution of this chapter is to present a leader-follower game theoretical approach for TCM. Also, this chapter connects leader-follower game theory with multiattribute decision-making. The operator intends to determine his preferred strategy among the K arbitrarily chosen strategies by considering congestiondriven and environmental attributes. For this purpose, TOPSIS has been used as a multiattribute technique for determining the operator's preferred strategy. Although a few game theoretical approaches to transmission congestion management have been proposed, see Veit et al. [43], Krause and Andersson [19], Liu et al. [27], Sahraei-Ardakani and Rahimi-Kian [37] and Lee [24], none of them have discussed selecting the operator's preferred strategy in a manner similar to the one proposed in this chapter. Hence, the approach in this chapter is unique and difficult to compare with the other approaches. As a whole, the proposed approach solves practically a game between generators, given an input set of parameters together with K strategies announced by the operator, without the mutual participation of the latter in the game. Thereafter, the operator performs congestion management and uses TOPSIS to choose his preferred strategy. The rest of the chapter is organized as follows: Sect. 2 provides an overview of the proposed approach. The Nash-SFE model is presented in Sect. 3. An overview of TCM is provided in Sect. 4. Strategy selection using TOPSIS is described in Sect. 5. Section 6 is devoted to simulation results on the IEEE reliability 24-bus test system. The procedure could be applicable for any power system and energy market. However, the results could be different case by case. Finally, Sect. 7 concludes the chapter.

2 Overview of the Proposed Approach

Figure 1 provides the flowchart of the proposed approach. This proposed procedure presents a method for the operator to obtain the optimum strategy among K strategies by trading off between the four attributes congestion cost, average LMP, variance of the LMPs as well as the g generators' emission. The operator's th strategy, $t = 1, ..., t_{1}$ K, which has been considered throughout this chapter, has g + 2 components, that is, g emission penalty factors ψ_{it} imposed by the operator on the *i*th generator, i =1, ..., g, the amount $P_{\text{ren},t}$ of purchased renewable power by the operator, and the bid cap $\beta_{\max,t}$ imposed by the operator on all the g generators. The bid cap $\beta_{\max,t}$ restricts the generators' bidding to prevent too high electricity price. In this chapter, we suppose that the electricity market regulator gives the operator this option to choose his optimum strategy from a set of K considered strategies expressed as $S = \{S_t = [\psi_{1t}, \dots, \psi_{32t}, P_{\text{ren},t}, \beta_{\max,t}] \quad \forall t = 1, \dots, K\}.$ For this purpose, we develop what we refer to as a leader-follower game-theoretic procedure: Referring to Fig. 1, first, in block 1, the operator as the leader chooses a set of K strategies arbitrarily, using his insight about what may constitute good strategies, without yet knowing which strategy is actually optimal. For each of the operator's K strategies, in block 2 the generators' Nash-SFE game is analyzed. Generator i, i = 1, ..., g, determines his optimal market bid β_{it} for selling his electrical power. After observing the g^*K optimal bids, i.e., K bids from each of the g generators, in block 3 the operator performs TCM for each of the K strategies applying the four attributes congestion cost, average LMP, variance of the LMPs as well as emission. Finally, in block 4, the operator selects his preferred strategy through a multiattribute decision making procedure technique, i.e., TOPSIS. It is assumed that the operator is aware of the electricity demand through a load forecaster, i.e., predictor. The load forecaster is developed based on a data-driven model using historical load data and some other effective variables such as weather-related and temporal parameters. The ultimate goal of load forecasters is to predict the electricity load/demand accurately for some time steps into the future. See Muñoz et al. [31] for more details.

Summing up, this chapter connects leader–follower game theory with multiattribute decision making. In this regard, after running the above-mentioned leader– follower procedure for K different strategies for the operator, and applying bids from multiple generators for market electricity price, the operator performs TCM and conducts multiattribute decision-making, i.e., TOPSIS. The operator thus determines his preferred strategy by considering congestion-driven criteria of the system as well as emission. This procedure provides the operator sufficient insight for selecting his optimal strategy in real electricity market operation.



Fig. 1 Flowchart of the proposed approach

3 Initiation and Nash-SFE Model: Blocks 1 and 2

The operator chooses K strategies arbitrarily in block 1, applying his judgment. In block 2, the g generators know the operator's K strategies, but does not know which of the K strategies will eventually be chosen. Each operator must thus choose an optimal market bid β_{it} for selling his electrical power for each of the K strategies. Kirschen and Strbac [17] claim that "While the Cournot model provides interesting insights into the operation of a market with imperfect competition, its application to electricity markets produces unreasonably high forecasts for the market price." Although they do not prove this mathematically, the electricity market literature seems to have accepted this claim. In order to obtain what we believe is a more realistic model for the electricity market than the Cournot model, it is assumed that a generator's offered energy is related to market price through a supply function introduced in Kirschen and Strbac [17]. Therefore, a Nash-SFE model is considered instead of a Cournot model in block 2. The SFE model has been used in various references such as Kirschen and Strbac [17]. It is assumed that when the operator chooses strategy t, t = 1, ..., K, the *i*th generator's one and only control variable is the market bid β_{it} and the objective of the generators is to maximize their own profit. Afterwards, the generators inject their optimal bids to the congestion management which is handled by operator by the objective of maximizing social welfare. The formulations in this chapter are developed for linear bids. Without loss of generality, instead of linear bids, stepwise bids could be deployed. The generators' cost function for power production, derived from an input—output function, i.e., heat rate (Btu/h) as a function of power (MW), has despite its complexity frequently been modeled as a convex piecewise linear function, see Hobbs et al. [14]. The cost function has been frequently modeled as quadratic, see Visalakshi et al. [45], Abido [1], and Hobbs et al. [15]. We thus assume that for the operator's strategy t, t = 1, ..., K, generator *i*'s cost function for power production is quadratic in power production P_{ii} , as follows:

$$C(P_{it}) = \frac{1}{2}a_i P_{it}^2 + b_i P_{it} + c_i$$
(1)

where a_i , b_i , c_i are cost function coefficients and P_{it} is power produced by the *i*th generator, i = 1, ..., g when the operator chooses strategy t, t = 1, ..., K. In (1) we assume $a_i > 0$ to ensure convexity, and $c_i \ge 0$ to ensure $C(P_{it}) \ge 0$ when $P_{it} = 0$ since with no power production, some cost in the power plant still exists. The parameter b_i is often assumed to be positive, see Saber et al. [35], Visalakshi et al. [45], Abido [1], and Zhang et al. [50]. Power production P_{it} is a free choice variable determined by the operator to maximize social welfare in Sect. 4. The marginal cost, which is the first derivative of the cost function, is:

$$MC(P_{it}) = a_i P_{it} + b_i \tag{2}$$

which gives a price $MC(P_{it})$ linearly dependent on quantity P_{it} which is generator *i*'s power production when the operator chooses strategy *t*. The relationship between

price and quantity is called a "supply function". In this chapter, we consider a linear supply function. We assume that the slope of the inverse supply function $SF^{-1}(P_{it})$ equals the slope of the marginal cost function $MC(P_{it})$, i.e., a_i . Furthermore, the intercept of the inverse supply function $SF^{-1}(P_{it})$ with the vertical axis when $P_{it} = 0$ equals generator *i*'s market bid β_{it} which is his one and only control variable and expresses the price at which he is willing to sell his produced power at quantity P_{it} when the operator chooses strategy t, t = 1, ..., K. We thus express generator *i*'s inverse supply function when the operator chooses strategy t as

$$SF^{-1}(P_{it}) = a_i P_{it} + \beta_{it} \tag{3}$$

which is the price at which generator *i* sells his power production at quantity P_{it} . The market clears at the estimated clearing price λ_t at which all the *g* generators sell their produced power P_{it} , i = 1, ..., g, t = 1, ..., K. This implies

$$SF^{-1}(P_{it}) = \lambda_t = a_i P_{it} + \beta_{it} \Leftrightarrow P_{it} = \frac{\lambda_t - \beta_{it}}{a_i}$$
(4)

As described earlier, the control variable of the *i*th generator when the operator chooses strategy *t* is β_{it} , chosen within the range of $0 \le \beta_{it} \le \beta_{\max,t}$, i = 1, ..., g; t = 1, ..., K. Each generator maximizes his own profit by choosing his optimal value of β_{it} . In these formulas, it is assumed that $a_i > 0$. Without loss of generality, the formulation can be easily modified for the cases in which $a_i = 0$. The modified formulation is provided in the Appendix of this chapter. In the generator's game in block 2, decisions of each generator are coupled to other generators' decisions through the market clearing price λ_t for electrical power. On the other side, in the operator's problem nodal prices are defined. Thus, the market signal that is defined in the two problems may appear to be inconsistent. However, the inconsistency is resolved through a balance between the total amount of power generation by the generators and consumer demand. This constraint is formulated as

$$\sum_{i=1}^{g} P_{it} + P_{\operatorname{ren},t} = D \tag{5}$$

where g is the number of generators and D is total electricity demand of all consumers. In this chapter, it is assumed that the consumers are price takers and a single-sided market is considered. In the case of a double-sided market, consumers offer their bids to the market. In contrast to the inverse supply function of generators in (3) which is an increasing function with respect to the produced power (P_{it}), the inverse demand function is a decreasing function with respect to demand. In this chapter, inverse demand has not been modeled. In the case of a double-sided market, D is treated in the same manner as P_{it} in the formulations of this section. Since we consider a single-sided market where consumers are price takers and do not offer their bids into the market, the consumer demand for electricity is constant expressed as *D*. The demand *D* is communicated to the generators and operator in the sense that all parameters are common knowledge for all players. Equation (5) expresses that the sum of all generated power by all the *g* generators, and the operator's amount of purchased renewable power $P_{\text{ren},t}$, should meet the demanded and thus consumed power *D*.

By defining

$$d = D - P_{\operatorname{ren},t} \tag{6}$$

Equation (5) is rewritten as:

$$\sum_{i=1}^{g} P_{ii} = d \tag{7}$$

From (4) and (7), we have:

$$\sum_{i=1}^{s} \frac{\lambda_t - \beta_{it}}{a_i} = d \tag{8}$$

From (8), the value of the market clearing price λ_t is determined as follows:

$$\lambda_{t} = \lambda_{t}(\beta_{1t}, \beta_{2t}, \dots, \beta_{gt}) = \frac{d + \sum_{i=1}^{g} \frac{\beta_{ii}}{a_{i}}}{\sum_{i=1}^{g} \frac{1}{a_{i}}} \quad \text{when } a_{i} > 0 \quad \forall i = 1, \dots, g$$
(9)

Assuming that $a_i = 0$ when $i = 1, ..., g_0$, and $a_i > 0$ when $i = g_0 + 1, ..., g$, $1 \le g_0 \le g$, the Appendix shows that when $a_i = 0$ for at least one generator, λ_t becomes:

$$\lambda_t = \frac{\beta_{1t} + \dots + \beta_{g_0 t}}{g_0} \tag{10}$$

when $a_i = 0$ for at least one generator, $1 \le g_0 \le g$, $\forall i = 1, ..., g$

That is, regardless how many generators g_0 have $a_i = 0$, when $1 \le g_0 \le g$, (10) is used instead of (9). It is assumed that the generators are not aware of the transmission system's constraints. Hence, in the event of congestion, λ_t will not be equal to the locational marginal prices (LMPs). In the Nash-SFE model in this block 2, we have assumed that a generator is located in just one bus.

The profit function of generator i is as follows:

$$\Pi_{it} = \lambda_t (\beta_{1t}, \beta_{2t}, \dots, \beta_{gt}) P_{it} - C(P_{it}) - C_E(P_{it}) \psi_{it}$$
(11)

where ψ_{it} [\$/ton], i = 1, ..., g; t = 1, ..., K is the emission penalty factor imposed by the operator on the *i*th generator, for the area in which generator *i* is located, λ_t is the estimated market clearing price when the operator chooses strategy t, and $C_E(P_{it})$ is the emission cost function of the *i*th generator when the operator chooses strategy t, t = 1, ..., K. The formulation of the emission cost function is as follows:

$$C_E(P_{it}) = \frac{1}{2}a_{Ei}P_{it}^2 + b_{Ei}P_{it} + c_{Ei}$$
(12)

where $a_{Ei} > 0$ to ensure convexity, and $c_{Ei} \ge 0$ to ensure $C_E(P_{it}) \ge 0$ when $P_{it} = 0$ since in the event of no power production, the boiler still needs to be kept warm. Furthermore, the parameter b_{Ei} can be positive or negative, see Chaturvedi et al. [6], but is often assumed to be negative, see Visalakshi et al. [45], Abido [1], and Zhang et al. [50]. Negative b_i means that the convex increase of $C_E(P_{it})$ occurs for larger values of P_{it} . Equation (10) expresses that generator *i*'s profit Π_{it} depends on the operator's strategy *t*, his own power production P_{it} , cost functions, etc., but also on all the *g* generators' market bids $\beta_{1t}, \beta_{2t}, \ldots, \beta_{gt}$ which impact the estimated market clearing price λ_t . That is, the *g* generators' profits impact each other through the price mechanism affected by their market bids. The *g* generators' profits additionally impact each other indirectly through the operator who affects all generators and regulates congestion. The *i*th generator maximizes Π_{it} by choosing his optimal market bid $\beta_{it}, i = 1, \ldots, g, t = 1, \ldots, K$. Hence, we write:

$$\frac{\partial \prod_{it}}{\partial \beta_{it}} = 0 \tag{13}$$

Differentiating (13) with respect to β_{it} yields:

$$\frac{\partial \lambda_t}{\partial \beta_{it}} P_{it} + \frac{\partial P_{it}}{\partial \beta_{it}} \lambda_t - \frac{\partial C(P_{it})}{\partial P_{it}} \frac{\partial P_{it}}{\partial \beta_{it}} - \frac{\partial C_E(P_{it})}{\partial P_{it}} \frac{\partial P_{it}}{\partial \beta_{it}} \psi_{it} = 0$$
(14)

By constituting the values of each term, the optimal value for β_{it} is obtained:

$$\beta_{it}^{\text{opt}} = \frac{-M_{1t}a_iM_{6t} - M_{3t}a_iM_{6t} + a_iM_{5t}M_{7t} + b_iM_{7t}a_iM_{6t} + a_iM_{8t}\psi_{it}M_{7t}M_{5t} + \psi_{it}b_{Ei}a_iM_{6t}}{a_iM_{6t}M_{2t} + a_iM_{6t}M_{4t} - M_{7t} + M_{7t}a_iM_{6t} - \psi_{it}M_{7t}M_{8t} + \psi_{it}a_iM_{6t}M_{7t}M_{8t}}$$
(15)

where

$$M_{1t} = \frac{1}{a_i^2 \sum_{j=1}^g \frac{1}{a_j}} \frac{d + \sum_{\substack{j=1\\j\neq i}}^g \frac{\beta_{ji}}{a_j}}{\sum_{j=1}^g \frac{1}{a_j}}$$
(16)

$$M_{2t} = \frac{1}{a_i^2 \sum_{j=1}^{g} \frac{1}{a_j}} \left[\frac{1}{a_i \sum_{j=1}^{g} \frac{1}{a_j}} - 1 \right]$$
(17)

A Leader-Follower Game on Congestion Management in Power Systems

$$M_{3t} = \frac{1}{a_i} \frac{d + \sum_{j=1}^{g} \frac{\beta_{ji}}{a_j}}{\sum_{j=1}^{g} \frac{1}{a_j}} \left[\frac{1}{a_i \sum_{j=1}^{g} \frac{1}{a_j}} - 1 \right]$$
(18)

$$M_{4t} = \frac{1}{a_i^2} \frac{1}{\sum_{j=1}^g \frac{1}{a_j}} \left[\frac{1}{a_i \sum_{j=1}^g \frac{1}{a_j}} - 1 \right]$$
(19)

$$M_{5t} = d + \sum_{\substack{j=1\\j\neq i}}^{g} \frac{\beta_{jt}}{a_j} \tag{20}$$

$$M_{6t} = \sum_{j=1}^{g} \frac{1}{a_j}$$
(21)

$$M_{7t} = \frac{1}{a_i} \left[\frac{1}{a_i \sum_{j=1}^{g} \frac{1}{a_j}} - 1 \right]$$
(22)

$$M_{8t} = \frac{a_{Ei}}{a_i} \tag{23}$$

As it is observed from (15), by considering β_{it} as the only control variable of generator *i*, a unique optimum bid has been obtained for each of the generators. In order to consider generator *i*'s power limits and market bid cap, the following changes are taken into consideration in the model:

$$\begin{cases} \text{if } P_{it} > P_i^{\min} & \text{then } \beta_{it}^{\text{opt}} = \lambda_t - a_i P_i^{\min} \\ \text{if } P_{it} < P_i^{\max} & \text{then } \beta_{it}^{\text{opt}} = \lambda_t - a_i P_i^{\max} \\ \text{if } \beta_{it} < \beta_{\max,t} & \text{then } \beta_{it}^{\text{opt}} = \beta_{\max,t} \\ \text{if } \beta_{it} > b_i & \text{then } \beta_{it}^{\text{opt}} = b_i \end{cases}$$
(24)

The process above applies for all the operator's K strategies.

4 Operator's Transmission Congestion Management: Block 3

This section analyzes block 3 of Fig. 1. Transmission congestion management refers to the actions that the operator performs for keeping all constraints within their limits. By providing the bids β_{it} , i = 1, ..., g, t = 1, ..., K, to the operator, the operator

performs congestion management by determining the power P_{it} , i = 1, ..., g; t = 1, ..., K, produced by the generators. That is, the operator maximizes social welfare within constraints for each of his *K* strategies. The operator is exclusively concerned about maximizing social welfare, and ignores other objectives such as, e.g., efficiency. The objective function to be maximized when the operator chooses strategy t, t = 1, ..., K, is as follows:

Max
$$SW_t = -\sum_{i=1}^{g} \left(\frac{1}{2} a_i P_{it}^2 + \beta_i P_{it} + c_i \right)$$
 (25)

The constraints include active power balance in each bus, and thermal and stability limits of each transmission line which are formulated as follows in Conejo et al. [9]:

$$\sum_{j=1}^{g_n} P_{jt} - \sum_{k=1}^{d_n} P_{nD_k} + P_{\text{ren},nt} = \sum_{k=1}^N \sum_{j=1}^N P_{line,ijt} \quad \forall n = 1, \dots, N$$
(26)

$$0 \le \left| P_{line,kjt} \right| \le \min \left\{ P_{Th,kj}^{\max}, P_{St,kj}^{\max} \right\} \quad \forall k, j = 1, \dots, N$$
(27)

where $|P_{line,kjt}|$ is the absolute value of $P_{line,kjt}$. Equations (26) and (27) are referred to as the power system equations. Equation (26) shows active power balance in bus *n*. That is, the power production of all generators and renewable power resources connected to bus *n* minus the active power consumption of all electricity demands connected to this bus equals the sum of power flows across the lines connected to bus *n*. Equation (27) limits power flows across transmission lines by considering thermal and stability limits. In the event of congestion occurring in a transmission line, power across that line reaches its limit. The power generation limits of each generator are:

$$P_i^{\min} \le P_{it} \le P_i^{\max} \quad i = 1, 2, \dots, g \tag{28}$$

That is, the operator chooses the g generators' production levels P_{it} to maximize social welfare in (25) given the three constraints in (26)–(28). In the proposed TCM of (25)–(28), some simplifications have been made. For example, reactive power is not considered and it is assumed that the voltage magnitude for all transmission buses is equal to 1. If power flow in a transmission line exceeds its stability and/or thermal limits as expressed in (27), the safety and security of the power system will be compromised. For example, consider the system's performance index (PI_i), i.e.,

$$PI_{t} = \sum_{j=1}^{N} \sum_{i=1}^{N} \left(\frac{P_{\text{line,ijt}}}{P_{ST,ij}^{\text{max}}} \right)^{2h}, h \ge 1, \quad t = 1, \dots, K$$
(29)

If power flow of a transmission line exceeds it stability limit, i.e., $P_{line,ijt} > P_{ST,ij}^{max}$ for one or several lines so that the ratio in (29) exceeds 1, then PI_t becomes large which compromises security. The parameter *h* is chosen substantially above 1 for critical power systems. The security and thermal constraints in (27) are thus used to maintain safety and security of the power system. This incorporates safety and security concerns into transmission congestion management.

Once the $g P_{it}$'s have been determined, the four attributes are determined as follows:

• The congestion cost: By solving the above-mentioned optimization (25)–(28), a locational marginal price (LMP) is assigned to each transmission bus. In general, LMP is defined as the minimum cost of supplying an incremental load at a specific location, see Liu et al. [28]. The calculation process of LMP has been provided in Conejo et al. [8]. We calculate the congestion cost as:

$$CC_t = \sum_{j=1}^{N} \sum_{k=1}^{N} P_{line,kjt} \left(\text{LMP}_{kt} - \text{LMP}_{jt} \right) \quad \forall t = 1, \dots, K$$
(30)

• The average locational marginal price (LMP) for different system buses is:

$$ave_LMP_t = \frac{1}{N} \sum_{k=1}^{N} LMP_{kt} \quad \forall t = 1, \dots, K$$
 (31)

• The variance locational marginal price (LMP) for different system buses is:

$$var_LMP_t = \frac{1}{N} \sum_{k=1}^{N} (LMP_{kt} - ave_LMP_t)^2 \quad \forall t = 1, \dots, K$$
(32)

• The g generators' emission is:

$$emission_t = \sum_{i=1}^{g} C_E(P_{it}) \quad \forall t = 1, \dots, K$$
(33)

where *N* is the number of transmission buses and $P_{line,kjt}$ is the power flows between buses *k* and, k, j = 1, ..., N. Linking to (30)–(33), the four attributes are expressed as *K*-dimensional vectors:

$$Y_1^T = \begin{bmatrix} CC_1 & CC_2 & \dots & CC_k \end{bmatrix}$$
(34)

$$Y_2^T = [ave_LMP_1 \quad ave_LMP_2 \quad \dots \quad ave_LMP_K]$$
(35)
$$Y_3^T = \begin{bmatrix} var_LMP_1 & var_LMP_2 & \dots & var_LMP_K \end{bmatrix}$$
(36)

$$Y_4^T = [emission_1 \ emission_2 \ \dots \ emission_K]$$
(37)

In this chapter, we consider four attributes which we think are representative and realistic to illustrate the method. In principle, an arbitrary number of attributes could have been chosen. The generalization is straightforward. Summing up, the operator's TCM consists in choosing the g optimal power production levels P_{it} 's for the generators within the constraints in (26)–(28) to maximize social welfare in (25), i = 1, ..., g, t = 1, ..., K. The g generators' production levels P_{it} are used to determine the four attributes in (30)–(33).

5 Operator's Strategy Selection Using TOPSIS: Block 4

This section analyzes block 4 of Fig. 1, where the operator selects his preferred strategy by using a multiattribute decision making technique. In this section, we have used TOPSIS for this purpose. TOPSIS is a powerful decision making tool with several applications such as calculation of scenario degree of severity, see Salehizadeh et al. [38], recruiting new staffs for companies, see Chen [7], and bridge risk assessment, see Shih et al. [39]. The step-by-step calculation process, which has been used in Salehizadeh et al. [38], is utilized in this section for selecting a preferred strategy for the operator:

- (i) Constituting the decision matrix *M*. To this end, using (34)–(37), congestion cost Y₁^T = [CC₁ CC₂ ... CC_k], average of the LMPs Y₂^T = [ave_LMP₁ave_LMP₂...ave_LMP_k], variance of the LMPs Y₃^T = [var_LMP₁ var_LMP...var_LMP_k], and the *g* generators' emission Y₄^T = [emission₁ emission₂...emission_k] for all system buses for each of the *K* strategies are considered as the decision making attributes. Each column and row of the decision matrix represents an attribute and an operator's strategy, respectively. The operator's decision matrix *M* is shown in Table 2.
- (ii) Normalizing the decision matrix M. In order to transform different types of attributes into a common scale, the decision matrix needs to be normalized. If X_i is the *j*th column of the decision matrix, it is normalized by applying

Table 2 The operator's decision matrix M			Attrib	utes		
			Y_1	<i>Y</i> ₂	<i>Y</i> ₃	Y_4
	The operator's	S_1	<i>x</i> ₁₁	<i>x</i> ₁₂	<i>x</i> ₁₃	<i>x</i> ₁₄
	strategy	S_2	<i>x</i> ₂₁	<i>x</i> ₂₂	<i>x</i> ₂₃	<i>x</i> ₂₄
		S_K	$x_{K 1}$	x_{K2}	<i>x</i> _{<i>K</i> 3}	x_{K4}

 $X_j / ||X_j||_{\infty}$, where $||X_j||_{\infty} = \max(|x_{1j}|, \dots, |x_{Kj}|)$ is the infinity norm of X_j . The *tj*th element of the normalized decision matrix is represented by y_{ij} . Instead of the infinity norm, other types of norms such as the Euclidean norm could be applied for normalization.

(iii) Considering the comparison matrix *A* of attributes. The operator forms a comparison matrix *A* of the attributes with $a_{jk}, j = 1, ..., 4; k = 1, ..., 4$, elements which reflects the attributes' comparative degrees of importance:

$$A = \begin{bmatrix} a_{11} & a_{12} & a_{13} & a_{14} \\ a_{21} & a_{22} & a_{23} & a_{24} \\ a_{31} & a_{32} & a_{33} & a_{34} \\ a_{41} & a_{42} & a_{43} & a_{44} \end{bmatrix}$$
(38)

where $a_{tj} > 0$ is a strictly positive parameter. This matrix is a reciprocal 4×4 and the relationship among its elements is as, $a_{tj} = 1/a_{jt}$, where a_{tj} is the importance of the *t*th attribute with respect to the *j*th attribute and 4 is the number of considered attributes. The matrix of attributes is consistent if $a_{tj}a_{jk} = a_{tk}$ for all values of t, j, k = 1, ..., 4, see Saaty et al. [34]. The consistency of the matrix is examined, see Lee et al. [25]. In the case of inconsistency the operator is asked to revise it.

(iv) Applying the 'eigenvector technique' proposed by Saaty et al. [34] to the comparison matrix of attributes, the weight of each attribute, which reflects its pure degree of importance, is calculated as in Yoon et al. [48]: If λ_{max} is the maximum eigenvalue of matrix *A*, the normalized related eigenvector (W^T) shows the attributes' weights:

$$W^{T} = \begin{bmatrix} w_{1} & w_{2} & w_{3} & w_{4} \end{bmatrix}$$
(39)

(v) Obtaining the positive/negative ideal strategies. By multiplying the normalized decision matrix's columns by the attributes' weights which are obtained from (39), the V matrix is yielded. Its *tj*th element is as follows:

$$v_{tj} = w_j y_{tj} \tag{40}$$

The positive ideal strategy (PIS) and the negative ideal strategy (NIS) are obtained from the V matrix. The PIS is a vector whose *j*th element is computed as follows:

$$v_j^+ = \max_{t=1,\dots,K} \quad v_{tj} \tag{41}$$

The NIS is also a vector whose *j*th element is computed as follows:

$$v_j^- = \min_{t=1,\dots,K} \quad v_{tj} \tag{42}$$

(vi) The distance of the *t*th strategy to the PIS and NIS, respectively, is calculated as follows:

$$L_t^+ = \sqrt{\sum_{j=1}^4 (v_{tj} - v_j^+)^2} \quad t = 1, 2, \dots, K$$
(43)

$$L_t^- = \sqrt{\sum_{j=1}^4 (v_{tj} - v_j^-)^2} \quad t = 1, 2, \dots, K$$
(44)

(vii) The strategies are ranked based on the following index:

$$\kappa_t = \frac{L_t^-}{\left(L_t^+ + L_t^-\right)} \qquad t = 1, 2, \dots, K \tag{45}$$

Finally, the strategy with minimum κ_t is selected by the operator. Since the operator first chooses *K* strategies, generator *i* thereafter chooses one optimal bid β_{it} for each operator strategy *t*, *t* = 1,..., *K*, and the operator finally performs TCM and chooses his preferred strategy S_t^* , a unique solution is guaranteed. Convergence to a unique solution is guaranteed in the iterative algorithm as determined by minimum κ_t in (45).

Let us summarize the procedure leading up to the minimum κ_t which determines the operator's preferred strategy $S_t^* = \begin{bmatrix} \psi_{1t}^*, \dots, \psi_{32t}^*, P_{\text{ren},t}^*, \beta_{\text{max},t}^* \end{bmatrix}$ for one specified $t, t = 1, \dots, K$, where * means preferred strategy. In Sects. 3 and 4, the operator examines K strategies S_1, S_2, \dots, S_K for the four attributes congestion cost, average locational marginal price (LMP) for different system buses, variance of the LMPs, and the *g* generators' emission. Arbitrarily many attributes could have been considered, but here we illustrate the procedure with four attributes. The K strategies can be expressed as $S_1 = [\psi_{1,1}, \dots, \psi_{32,1}, P_{\text{ren},1}, \beta_{\text{max},1}], S_2 = [\psi_{1,2}, \dots, \psi_{32,2}, P_{\text{ren},2}, \beta_{\text{max},2}], \dots, S_K = [\psi_{1K}, \dots, \psi_{32K}, P_{\text{ren},K}, \beta_{\text{max},K}]$. The operator intends to select one of these strategies by considering various values for these four attributes. First, we consider all $S_t, t = 1, \dots, K$, strategies and the procedure of Sects. 3 and 4 is performed by considering each of the $S_t, t = 1, \dots, K$ strategies as the operator's strategy. Suppose that the generators' optimal bids from the procedure of Sect. 3 have been obtained as: $[\beta_{1t}, \dots, \beta_{32t}], t = 1, \dots, K$, and after performing congestion management in Sect. 4, the four attributes have been calculated as $[CC_t, ave_LMP_t, var_LMP_t, emission_t], t = 1, ..., K$. Next, through TOPSIS the following $K \times 4$ decision matrix, equivalent to Table 2, is considered:

$$m = \begin{bmatrix} Y_1 & Y_2 & Y_3 & Y_4 \end{bmatrix}$$

$$= \begin{bmatrix} CC_1 & ave_LMP_1 & var_LMP_1 & emission_1 \\ CC_2 & ave_LMP_2 & var_LMP_2 & emission_2 \\ \dots & \dots & \dots & \dots \\ CC_{K-1} & ave_LMP_{K-1} & var_LMP_{K-1} & emission_{K-1} \\ CC_K & ave_LMP_K & var_LMP_K & emission_K \end{bmatrix}$$
(46)

Using this decision matrix, the procedure in this Sect. 5 determines the minimum κ_t which determines the operator's preferred strategy $S_t^* = \left[\psi_{1t}^*, \dots, \psi_{32t}^*, P_{\text{ren},t}^*, \beta_{\text{max},t}^*\right]$ and generator *i*'s optimal bid β_{it} for one specific value of each of the four attributes.

6 Simulation Results

The IEEE 24-bus test system or the IEEE reliability test system (RTS) is one the IEEE standard test systems which is used for reliability studies, see Grigg et al. [12]. This system has been frequently used to illustrate the proposed methods developed in this area. For examples, see Salehizadeh et al. [46], Conejo et al. [9], and Orfanos et al. [32]. Hence, in this chapter, we have chosen this system for implementing the proposed leader-follower game. The method is performed on an IEEE 24-bus test system and the results are analyzed. A diagram of the tested power system is depicted in Fig. 2. In a power system, transmission buses (or nodes) are connected together through transmission lines which are depicted in this figure. It is assumed that a wind turbine is connected to transmission bus 1. The market is single-sided, i.e., the electricity consumers do not participate in the SFE game. The information related to load, lines, and generation limits is based on the case 24_ieee_rts in MATPOWER 4.1, see Zimmerman et al. [52]. In order to impose further stress on the transmission system, each line's capacity limits are reduced 0.6 times that of case 24_ieee_rts. Information about cost function and emission cost function coefficients are brought in Table 3. An example of the operator's strategy t is provided in Table 5. The emission cost function coefficients are selected similar to the data provided in [35]. In the real-world case studies, these coefficients are obtained based on type of fuel as well as physical functioning of different parts of the generator such as the boiler and the turbine.

In this study, 32 generators are considered and the synchronous condenser connected to bus 14 is omitted from the list of generators. The comparison matrix *A* of the attributes is specified as in Table 4. All simulations are done using MATLAB 7.13 (R2011b) and the MATPOWER package also used in Zimmerman et al. [52]. The role of the MATPOWER package is to calculate the optimal power



Fig. 2 Diagram of the tested power system

Generator no.	a	b	с	a _{Ej}	b_{Ej}	C _{Ej}
	Cost function coefficients		Emission cost function coefficients			
1	0.01	130	400.6849	0.0300	-2.3000	301
2	0.01	130	400.6849	0.0300	-3.4000	103
3	0.0283	16.0811	212.3076	0.0849	-4.3000	104
4	0.0283	16.0811	212.3076	0.0849	-5.6000	205
5	0.01	130	400.6849	0.0300	-4.8000	195
6	0.01	130	400.6849	0.0300	-3.3000	93
7	0.0283	16.0811	212.3076	0.0849	-2.1000	414
8	0.0283	16.0811	212.3076	0.0849	-1.7000	106
9	0.1053	43.6615	781.5210	0.3159	-2.4000	309
10	0.1053	43.6615	781.5210	0.3159	-1.8000	102
11	0.1053	43.6615	781.5210	0.3159	-3.3000	205
12	0.0143	48.5804	832.7575	0.0429	-4.2000	107
13	0.0143	48.5804	832.7575	0.0429	-3.7000	103
14	0.0143	48.5804	832.7575	0.0429	-2.9000	108
15	0.6568	56.564	86.3852	0.0300	-3.2000	204
16	0.6568	56.564	86.3852	0.0320	-2.4000	103
17	0.6568	56.564	86.3852	0.0543	-4.3000	110
18	0.6568	56.564	86.3852	0.0432	-2.1000	313
19	0.6568	56.564	86.3852	0.0646	-3.6000	111
20	0.0167	12.3883	382.2391	0.0501	-2.5000	106
21	0.0167	12.3883	382.2391	0.0501	-1.4000	104
22	0.0004	4.4231	395.3749	0.0012	-5.3000	293
23	0.0004	4.4231	395.3749	0.0012	-3.2000	94
24	0.01	0.001	0.0010	0.0300	-2.5000	189
25	0.01	0.001	0.0010	0.0300	-1.7000	85
26	0.01	0.001	0.0010	0.0300	-4.8000	383
27	0.01	0.001	0.0010	0.0300	-3.5000	89
28	0.01	0.001	0.0010	0.0300	-2.3000	485
29	0.01	0.001	0.0010	0.0300	-2.2000	101
30	0.0167	12.3883	382.2391	0.0501	-2.7000	304
31	0.0167	12.3883	382.2391	0.0501	-1.4000	305
32	0.0098	11.8495	665.1094	0.0294	-3.2000	401

Table 3 Cost function and emission cost function coefficients

flow and calculate the LMPs, generation schedules, and power flow of each transmission line. As the first step of our study, see Table 5 for an example of the operator's strategy *t*. From (1)–(23), the Nash-SFE equilibrium has been obtained and the optimal bids β_{it} are depicted in Fig. 3. In this simulation and by the assumption of this case study, the obtained power of the generators exceeds the power generation limits. For example, consider the first generator. His β_{it} is 130 but

	Congestion $cost (CC_t)$	Average of LMPs (<i>ave_</i> LMP _t)	Variance of LMPs (var_LMP _t)	The g generators' emission (<i>emission</i> _t)
Congestion cost (CC_t)	1	2	4	4/3
Average of LMPs (<i>ave</i> _LMP _t)	1/2	1	2	2/3
Variance of LMPs (<i>var_LMP_t</i>)	1/4	1/2	1	1/3
The g generators' emission (<i>emission</i> _t)	3/4	3/2	3	1

 Table 4
 Comparison matrix A of the attributes

```
Table 5An example of theoperator's strategy t
```

$\psi_{it}, i = 1, \dots, 32$	$\beta_{\max,t}$	$P_{\mathrm{ren},t}$
10	140	100

his bid has been obtained as 26.1099. The reason for this difference is the obtained power, i.e., $P_{it} = -1.0373 \times 10^4$, i = 1, whilst $P_i^{\text{max}} = 20$ and $P_i^{\text{min}} = 16$. The obtained power exceeds the boundary. Based on (24), $\beta_{it}^{\text{opt}} = \lambda_t - a_i P_i^{\text{min}} =$ $26.2699-0.0100 \times 16 = 26.1099$ which is smaller than $\beta_{it} = 130$. Applying the g generators' optimal bids, the operator performs transmission congestion through using (25)–(28). As a result of the TCM procedure, the obtained LMPs are depicted in Fig. 4. This figure reveals differences between the LMPs of all transmission buses. This observation reveals that congestion has occurred in the transmission



Fig. 3 The optimal bid β_{it} provided by generator *i*, *i* = 1 ,..., *g*, when the operator chooses strategy *t*, *t* = 1 ,..., *K*



Fig. 4 Locational marginal price LMP_{kt} for transmission bus k, k = 1, ..., N when the operator chooses strategy t, t = 1, ..., K

system. From (30), the congestion cost is calculated and obtained equal to 8.0305e + 003 [\$/h].

The effect of each of the operator's K strategies on the TCM result has been studied. From the comparison matrix A of attributes (Table 4), it is deduced that congestion cost and emission are the most important attributes. Hence, first, we have assessed the impact of each of the operator's K strategies on these attributes. Figure 5 shows the effect of purchased renewable power on congestion cost and the g generators' emission. As this figure shows, congestion cost decreases by increasing the amount of purchased renewable power from 0 to about 100 MW and increases afterwards. The rate of increase rises quickly from 285 to 300 MW. This observation shows that increasing the amount of purchased renewable power and its penetration to the power system significantly affects congestion cost and requires further consideration. Also, Fig. 5 shows the effect of purchased renewable power on the g generators' emission. By increasing the amount of purchased renewable power to 285 MW, the g generators' emission decreases. However, from 285 to 300 MW, a rapid increase in the g generators' emission occurs. Hence when the power system is under stress, increased renewable penetration may cause enhanced system emission.

The effect of the operator's market bid cap β_{max} (specifying the range within which the generators can bid electricity prices β_{it} in the market) on congestion cost and emission has been shown in Fig. 6. The parameters of the emission cost function obviously affect the value of emission. To analyze the effects of each parameter on the results in detail, sensitivity analysis needs to be applied which is suggested for future research. Figure 6 shows that by increasing the bid cap β_{max}



Fig. 5 The effect of purchased renewable power $P_{\text{ren},jt}$ on congestion cost CC_t and the *g* generators' *emission*_t, j = 1, ..., N, t = 1, ..., K



Fig. 6 The effect of the operator's market bid cap $\beta_{\max,t}$ on congestion cost CC_t and the *g* generators' *emission*_t, t = 1, ..., K



Fig. 7 The obtained normalized attributes CC_t , $emission_t$, ave_LMP_t , var_LMP_t : for K strategies, K = 61

from 0 to 8, the congestion cost decreases but the emission increases. This reveals that increasing the bid cap β_{max} in this interval causes pollutant generators to be more dispatched. However, this production schedule reduces total congestion. By increasing the bid cap β_{max} above 8, congestion cost and emission remain almost constant. Comparing Fig. 5 with Fig. 6 shows that the purchased renewable power provides a wider control factor for the operator than the market bid cap. Our simulation results show that the emission penalty factor is even more restricted than the bid cap from the viewpoint of control flexibility. Hence, we consider the purchased renewable power as a crucial control factor for choosing the operator's preferred strategy in the leader-follower game-theoretic TCM proposed in Fig. 1. For this purpose, first, the values of the bid cap and emission penalty factors are specified according to Table 5 and the purchased renewable power is set equal to 0. Next, this value is increased to 5 MW and K = 61 operator strategies are considered. For each of the K strategies, the four different attributes congestion cost, average of the LMPs, variance of the LMPs as well as the g generators' emission are obtained. The attributes have been normalized by dividing them by their maximum value. The normalized attributes are depicted in Fig. 7.

Applying the Saaty technique in Saaty et al. [34] (described in Sect. 5) on the matrix of attributes of Table 4 yields the pure weight of the attributes. The pure weights are depicted in the first row of Table 6.

Applying (38)–(45) yields the distance of each strategy to the PIS, NIS and κ_i indices. These three indices have been depicted in Fig. 8. The minimum of the κ_i

Congestion $cost (CC_t)$	Average of LMPs (<i>ave_LMP_t</i>)	Variance of LMPs (<i>var_LMP_t</i>)	The g generators' emission $(emission_t)$
0.4	0.2	0.1	0.3
953.59	26.76	0.168	4409.55

Table 6 The pure weights and the attributes of the operator's preferred strategy



Fig. 8 The distance L_t^+ to the positive ideal strategy PIS, the distance L_t^- to the negative ideal strategy NIS, and the κ_i index, t = 1, ..., K

index occurs for the operator's 22nd strategy, out of K = 61. The amount of purchased renewable power for the operator's 22nd strategy is equal to 105 MW. The attributes of this preferred strategy are shown in the second row of Table 6.

7 Conclusion

In this chapter, a leader–follower game on transmission congestion management (TCM) in power systems has been proposed. The leader is an independent system operator and the followers are generators (power generation companies). The operator purchases guaranteed renewable power for providing incentives for promotion of these environmental friendly resources. The leader, i.e., the independent system operator, determines a strategy consisting of three types of components: a bid cap for the maximum bid the generators can offer for selling its electricity in the market, the operator's amount of purchased renewable power, and the emission

penalty factors for the area in which generators exist. Each generator offers his bid to the market for selling his electrical power. The generators interact based on a Nash-supply function equilibrium (SFE) game. The operator performs TCM and calculates the four attributes congestion cost, average locational marginal price (LMP), variance of the LMPs as well as emission. In the TCM model, the upper and lower limits of power flow, accounting for thermal factors and stability, ensure that the power system functions safely and securely. The leader-follower game is analyzed for multiple operator strategies and finally, a multiattribute decisionmaking approach using the Technique for Order Preference by Similarity to Ideal Solution (TOPSIS) has been implemented by the operator for selecting his preferred strategy by considering both emission and congestion-reflective attributes. The approach solves practically a game between generators, given an input set of parameters together with K strategies announced by the operator, without the mutual participation of the latter in the game. Thereafter, the operator performs congestion management and uses TOPSIS to choose his preferred strategy. As a whole, the main purpose of this chapter is to provide a combination of a leaderfollower game theoretical mechanism and multiattribute decision-making for the operator to choose his best strategy by considering congestion-driven and environmental attributes. We implemented the developed method on an IEEE reliability 24-bus test system by considering 32 generators and a wind farm connected to bus 1. It was also assumed that the operator provided a comparison matrix of attributes. In this case study, the simulation analysis shows that among the operator's strategic choices, the amount of purchased renewable power provides a wider control flexibility to generate different strategies in the leader-follower game-theoretic congestion management. Simulation of the developed procedure provided the operator insight to select his preferred strategy accounting for TCM, environmental issue, and maximizing social welfare. Although a few game theoretical approaches to transmission congestion management have been proposed in the literature, see e.g., Veit et al. [43], Krause and Andersson [19], Liu et al. [27], Sahraei-Ardakani and Rahimi-Kian [37] and Lee [24], none of them have discussed selecting a preferred strategy by the operator in a manner consisting of the four blocks proposed in this chapter. By implementing the operator's obtained strategy, the operator obtains the preferred real market results. In this direction, developing a long-term learning mechanism for the operator could be considered for future research. Uncertainties due to prediction of renewable power resource outputs could affect the results of the developed leader-follower game. This point could be also investigated in future studies. This chapter provides a useful tool for policy makers, operators, generators, consumers, and associated actors.

Acknowledgments We thank anonymous referees for their helpful suggestions.

Appendix When $a_i = 0$ for at Least One Generator

Assume that $a_i = 0$ when $i = 1, ..., g_0$, and $a_i > 0$ when $i = g_0 + 1, ..., g$, $1 \le g_0 \le g$. Taking the limit for (9) gives:

$$\lambda_{t} = \lim_{(a_{1},...,a_{g_{0}}) \to (0,...,0)} \frac{d + \sum_{i=1}^{g} \frac{\beta_{it}}{a_{i}}}{\sum_{i=1}^{g} \frac{1}{a_{i}}}$$
(47)

Separating the terms related to the g_0 generators with $a_i = 0$, and the $g - g_0$ generators with $a_i > 0$, gives:

$$\lambda_{t} = \lim_{\substack{(a_{1},\dots,a_{g_{0}})\to(0,\dots,0)}} \frac{d + \sum_{i=g_{0}+1}^{g} \frac{\beta_{ii}}{a_{i}} + \sum_{i=1}^{g_{0}} \frac{\beta_{ii}}{a_{i}}}{\sum_{i=g_{0}+1}^{g} \frac{1}{a_{i}} + \sum_{i=1}^{g_{0}} \frac{1}{a_{i}}} = \lim_{\substack{(a_{1},\dots,a_{g_{0}})\to(0,\dots,0)}} \frac{d + \sum_{i=g_{0}+1}^{g} \frac{\beta_{ii}}{a_{i}}}{\sum_{i=g_{0}+1}^{g} \frac{1}{a_{i}} + \sum_{i=1}^{g_{0}} \frac{1}{a_{i}}}{\sum_{i=g_{0}+1}^{g} \frac{1}{a_{i}} + \sum_{i=1}^{g_{0}} \frac{1}{a_{i}}} = 0 + \lim_{\substack{(a_{1},\dots,a_{g_{0}})\to(0,\dots,0)}} \frac{\sum_{i=g_{0}+1}^{g_{0}} \frac{\beta_{ii}}{a_{i}}}{\sum_{i=g_{0}+1}^{g} \frac{1}{a_{i}} + \sum_{i=1}^{g_{0}} \frac{1}{a_{i}}}{\sum_{i=g_{0}+1}^{g} \frac{1}{a_{i}} + \sum_{i=1}^{g_{0}} \frac{1}{a_{i}}} = 0 + \lim_{\substack{(a_{1},\dots,a_{g_{0}})\to(0,\dots,0)}} \frac{\sum_{i=g_{0}+1}^{g_{0}} \frac{\beta_{ii}}{a_{i}}}{\sum_{i=g_{0}+1}^{g} \frac{1}{a_{i}} + \sum_{i=1}^{g_{0}} \frac{1}{a_{i}}}}{\sum_{i=g_{0}+1}^{g} \frac{1}{a_{i}} + \sum_{i=1}^{g_{0}} \frac{1}{a_{i}}} = 0 + \lim_{\substack{(a_{1},\dots,a_{g_{0}})\to(0,\dots,0)}} \frac{\sum_{i=g_{0}+1}^{g} \frac{\beta_{ii}}{a_{i}}}{\sum_{i=g_{0}+1}^{g} \frac{1}{a_{i}} + \sum_{i=1}^{g_{0}} \frac{1}{a_{i}}}}{\sum_{i=g_{0}+1}^{g} \frac{1}{a_{i}} + \sum_{i=1}^{g_{0}} \frac{1}{a_{i}}} + \cdots + \lim_{\substack{(a_{1},\dots,a_{g_{0}})\to(0,\dots,0)}} \frac{\sum_{i=g_{0}+1}^{g} \frac{\beta_{g_{0}}}{a_{i}}}{\sum_{i=g_{0}+1}^{g} \frac{1}{a_{i}}} + \sum_{i=1}^{g_{0}} \frac{1}{a_{i}}}} = \lim_{\substack{(a_{1},\dots,a_{g_{0}})\to(0,\dots,0)}} \frac{\beta_{g_{0}}}{\sum_{i=g_{0}+1}^{g} \frac{1}{a_{i}} + \sum_{i=1}^{g_{0}} \frac{1}{a_{i}}}}{g_{0}}}$$

$$= \frac{\beta_{1t}}{0+g_{0}} + \cdots + \frac{\beta_{g_{0}t}}{0+g_{0}}} = \frac{\beta_{1t} + \cdots + \beta_{g_{0}t}}{g_{0}}}$$

$$(48)$$

The calculated λ_t will be substituted in (10) and the rest of the calculation procedure will be the same as what we developed for the case of $a_i > 0; i = 1, ..., g$.

References

- Abido MA (2003) A niched Pareto genetic algorithm for multiobjective environmental/ economic dispatch. Int J Electr Power Energy Syst 25(2):97–105
- Ahmadi H, Lesani H (2014) Transmission congestion management through LMP difference minimization: a renewable energy placement case study. Arab J Sci Eng 39(3):1963–1969
- Bier VM, Hausken K (2013) Defending and attacking a network of two arcs subject to traffic congestion. Reliab Eng Syst Saf 112:214–224
- 4. Bompard E, Correia P, Gross G, Amelin M (2003) Congestion-management schemes: a comparative analysis under a unified framework. IEEE Trans Power Syst 18(1):346–352
- 5. Borenstein S, Bushnell J (2000) Electricity restructuring: deregulation or reregulation. Regulation 23:46–52
- Chaturvedi KT, Pandit M, Srivastava L (2008) Hybrid neuro-fuzzy system for power generation control with environmental constraints. Energy Convers Manage 49(11):2997–3005
- 7. Chen CT (2000) Extensions of the TOPSIS for group decision-making under fuzzy environment. Fuzzy Sets Syst 114(1):1–9

- Conejo AJ, Castillo E, Mínguez R, Milano F (2005) Locational marginal price sensitivities. IEEE Trans Power Syst 20(4):2026–2033
- Conejo AJ, Milano F, García-Bertrand R (2008) Congestion management ensuring voltage stability. In: Proceedings of IEEE power and energy society general meeting-conversion and delivery of electrical energy in the 21st century, 2008 IEEE
- de la Torre S, Contreras J, Conejo AJ (2004) Finding multiperiod Nash equilibria in poolbased electricity markets. IEEE Trans Power Syst 19(1):643–651
- Dong S, Yang Q, Fu F, Kwak KS (2013) Distributed link scheduling for congestion control in multihop wireless network. In: Proceedings of international conference on IEEE wireless communications and signal processing (WCSP), pp 1–5
- 12. Grigg C, Wong P, Albrecht P, Allan R, Bhavaraju M, Billinton R, Singh C (1999) The IEEE reliability test system-1996. A report prepared by the reliability test system task force of the application of probability methods subcommittee. IEEE Trans Power Syst 14(3):1010–1020
- Hobbs BF (2001) Linear complementarity models of Nash-Cournot competition in bilateral and POOLCO power markets. IEEE Trans Power Syst on 16(2):194–202
- Hobbs BF, Helman U, Pang JS (2001) Equilibrium market power modeling for large scale power systems. In: Proceedings of IEEE power engineering society summer meeting, vol 1, pp 558–563
- 15. Hobbs BF, Rothkopf MH, O'Neill RP, Hung-po C (ed) (2001) The next generation of electric power unit commitment models, vol 36. Springer, New York
- Kaplan SM (2009) Electric power transmission: background and policy issues. US Congressional Research Service, pp 4–5, 14 April 2009
- 17. Kirschen D, Strbac G (2004) Fundamentals of power system economics. John Wiley and Sons, New York
- Klemperer PD, Meyer MA (1989) Supply function equilibria in oligopoly under uncertainty. Econom J Econom Soci 57(6):1243–1277
- 19. Krause T, Andersson G (2006) Evaluating congestion management schemes in liberalized electricity markets using an agent-based simulator. In: Proceedings of IEEE power engineering society general meeting
- Kumar A, Sekhar C (2013) Comparison of Sen transformer and UPFC for congestion management in hybrid electricity markets. Int J Electr Power Energy Syst 47:295–304
- Kumar A, Srivastava SC, Singh SN (2005) Congestion management in competitive power market: a bibliographical survey. Electric Power Syst Res 76(1):153–164
- 22. Kunz F (2013) Improving congestion management: how to facilitate the integration of renewable generation in Germany. Energy J 34(4):55–78
- 23. Lai LL (ed) (2001) Power system restructuring and deregulation: trading, performance and information technology. John Wiley and Sons, Chichester
- 24. Lee KH (2014) Strategy equilibrium in Stackelberg model with transmission congestion in electricity market. J Electr Eng Technol 9(1):90–97
- Lee AH, Yang CN, Lin CY (2012) Evaluation of children's after-school programs in Taiwan: FAHP approach. Asia Pac Educ Rev 13(2):347–357
- 26. Lin S, Fletcher BA, Luo M, Chinery R, Hwang SA (2011) Health impact in New York city during the northeastern blackout of 2003. Public Health Rep 126(3):384
- Liu Y, Wu FF (2007) Impacts of network constraints on electricity market equilibrium. IEEE Trans Power Syst 22(1):126–135
- Liu Z, Tessema B, Papaefthymiou G, van der Sluis L (2011) Transmission expansion planning for congestion alleviation using constrained locational marginal price
- Lise W, Linderhof V, Kuik O, Kemfert C, Östling R, Heinzow T (2006) A game theoretic model of the Nnrthwestern European electricity market—market power and the environment. Energy Policy 34(15):2123–2136
- 30. Merriam-Webster dictionary. http://www.merriam-webster.com/dictionary/congestion
- Muñoz A, Sánchez-Úbeda EF, Cruz A, Marín J (2010) Short-term forecasting in power systems: a guided tour. In: Proceedings of handbook of power systems II. Springer, Heidelberg, pp 129–160

- Orfanos GA, Georgilakis P, Hatziargyriou ND (2013) Transmission expansion planning of systems with increasing wind power integration. IEEE Trans Power Syst 28(2):1355–1362
- 33. Porter K (2002) The implications of regional transmission organization design for renewable energy technologies. National Renewable Energy Laboratory, USA
- 34. Saaty TL, Vargas LG (2001) Models, methods, concepts and applications of the analytic hierarchy process, vol 1. Kluwer Academic Publishers, Boston
- Saber AY, Venayagamoorthy GK (2010) Intelligent unit commitment with vehicle-to-grid—a cost-emission optimization. J Power Sources 195(3):898–911
- 36. Saguan M, Keseric N, Dessante P, Glachant JM (2006) Market power in power markets: game theory vs. agent-based approach. In: Proceedings of IEEE/PES transmission and distribution conference and exposition, Latin America, 2006, TDC'06, pp 1–6
- Sahraei-Ardakani M, Rahimi-Kian A (2009) A dynamic replicator model of the players' bids in an oligopolistic electricity market. Electr Power Syst Res 79(5):781–788
- Salehizadeh MR, Rahimi-Kian A, Oloomi-Buygi M (2014) A multi-attribute congestiondriven approach for evaluation of power generation plans. Int Trans Electr Energy Syst. doi:10.1002/etep.1861
- Shih HS, Shyur HJ, Lee ES (2007) An extension of TOPSIS for group decision making. Math Comput Model 45(7):801–813
- 40. Son YS, Baldick R (2004) Hybrid coevolutionary programming for Nash equilibrium search in games with local optima. IEEE Trans Evol Comput 8(4):305–315
- Sun H, Wu J, Ma D, Long J (2014) Spatial distribution complexities of traffic congestion and bottlenecks in different network topologies. Appl Math Model 38(2):496–505
- 42. Telang NG, Jordan KE, Supekar NS (2013) US Patent No. 8,516,121. US Patent and Trademark Office, Washington, DC
- 43. Veit DJ, Weidlich A, Krafft JA (2009) An agent-based analysis of the German electricity market with transmission capacity constraints. Energy Policy 37(10):4132–4144
- Ventosa M, Baillo A, Ramos A, Rivier M (2005) Electricity market modeling trends. Energy Policy 33(7):897–913
- Visalakshi S, Baskar S (2011) Multiobjective decentralized congestion management using modified NSGA-II. Arab J Sci Eng 36(5):827–840
- 46. Salehizadeh MR, Rahimi-Kian A, Oloomi-Buygi M (2015) Security-based multi-objective congestion management for emission reduction in power system. Int J Electr Power Energy Syst 65:124–135
- 47. Wang X, Zhuang J (2011) Balancing congestion and security in the presence of strategic applicants with private information. Eur J Oper Res 212(1):100–111
- Yoon KP, Hwang CL (eds) (1995) Multiple attribute decision making: an introduction. Sage 102–104
- 49. Zhang YP, Jiao LW, Chen SS, Yan Z, Wen FS, Ni YX, Wu F (2003) A survey of transmission congestion management in electricity markets. Power Syst Technol 8:1–9
- 50. Zhang Y, Gong DW, Ding Z (2012) A bare-bones multi-objective particle swarm optimization algorithm for environmental/economic dispatch. Inf Sci 192:213–227
- 51. Zeng M, Luan F, Zhang J, Liu B, Zhang Z (2006) Improved ant colony algorithm (ACA) and game theory for economic efficiency evaluation of electrical power market. In: Proceedings of international conference on IEEE computational intelligence and security, 2006, vol 1, pp 849–854
- Zimmerman RD, Murillo-Sánchez CE, Gan D (1997) A MATLAB power system simulation package. http://www.pserc.conrnell.edu/matpower. Accessed 8 Jan 2006
- 53. Zou X, Luo X, Peng Z (2008) Congestion management ensuring voltage stability under multicontingency with preventive and corrective controls. In: Proceedings of IEEE power and energy society general meeting-conversion and delivery of electrical energy in the 21st century, 2008, pp 1–8

Network Throughput and Reliability: Preventing Hazards and Attacks Through Gaming—Part I: Modeling

Yupo Chan

Abstract By facilitating activity flows, infrastructure networks such as transportation, telecommunication and power grids are essential to the economy, quality-oflife and security of a society. However, as evidenced in recent historical events such as the collapse of the Minneapolis Highway Bridge and attacks on a country's defense-information infrastructure, a network is only useful if it is reliable, secure, and functioning properly. In other words, it has to be devoid of unexpected failures due to natural/technological disasters and outside attacks. A stochastic network, characterized by arcs (links) and nodes that can fail unexpectedly, is proposed to mimic such unpredicted interruptions. Through such a stochastic-network model, we identify tactics to prevent disruptions caused by natural/technological hazards and hostile tampering. Strategically, we can also advance public-policy options by determining an appropriate budget needed not only to maintain our infrastructure, but also to guard against adversarial attacks. The latter are accomplished by *imputing* the value of network security *scientifically* based on the cost of disrupting economic and noneconomic transactions, rather than using traditional cost accounting.

Keywords Network security • Multicriteria decision-making • Stochastic network • Cooperative game • Non-cooperative game • Nash equilibrium • Network reliability • Expected maximum-flow

1 Introduction

It is widely recognized that transportation, communication and other infrastructure networks are essential to economic prosperity, security, and quality of life. Unfortunately, due to random and deliberate incidents, such an infrastructure network is often compromised [36]. Such events might be triggered by breakdowns, adverse

Y. Chan (🖂)

University of Arkansas at Little Rock, Little Rock, USA e-mail: ychan@alum.MIT.edu

[©] Springer International Publishing Switzerland 2015

K. Hausken and J. Zhuang (eds.), *Game Theoretic Analysis of Congestion, Safety and Security*, Springer Series in Reliability Engineering, DOI 10.1007/978-3-319-13009-5_5

weather, earthquakes, floods, or acts of terrorism. Due to resource constraints, an infrastructure network has only limited capacity for repair. In addition, lengthy infrastructure repairs and their diverse failure modes distinguish them from simple equipment failure, where often replacing a component is all that is required to "fix the problem." While we will recognize differences between communication, transportation and power networks, a goal of our analysis is to obtain insights that will hold across these modal networks under a common mathematical paradigm.

An infrastructure network has to be reliable to perform its function. In this regard, network reliability has been defined differently in the transport versus telecommunication community. In transportation, it is mainly the probability that a trip on a particular path can be made successfully within an acceptable time (or cost) threshold. Travel-time (or travel-cost) reliability is therefore a measure of the stability of the sojourn duration (or sojourn expense), and is subject to fluctuations in flow and capacity. In communication, it is defined in terms of connectivity, or the probability that specific origin-destination (O-D) pairs in a network remain connected when arcs (links) are subject to complete failures. Due to current efforts in counter-terrorism and to mitigate against natural and technological disasters, increasing attention is given to the study of network reliability [7, 12, 25, 30, 37, 41, 49]. For the purpose of this chapter, reliability is broadly defined as the ability of a network to perform a required function, under adverse environmental and operational conditions, for a stated time period [23].

Related to network reliability is throughput, or the largest amount of traffic that can get through. For example, bridge capacities determine the traffic movement in a city striding across a river. In a mathematical model, we are concerned with the cut sets, or the bottlenecks, which determine the maximum flow that can be effected between any O-D pair. Chan et al. [10] provided linear-programming procedures to compute the expected (average) throughput in a stochastic network, a network where components can fail unexpectedly. When an arc fails, it becomes unavailable; and when it functions, it becomes available. Instead of elaborate Monte-Carlo simulations, they provided a tight analytical lower-bound (or an accurate conservative estimate) of the expected throughput. Even for realistic-size networks, they found that such a conservative estimate of the throughput can be obtained within minutes of calculation time, rather than days of simulation time. Parallel efforts include the work of Nojima [37] and Chen et al. [12] for capacity reliability, defined as finding the probability that a transportation system can accommodate a given demand level at an acceptable level of service, while taking the route choice behavior into account. This reliability definition can include connectivity reliability as a special case.

Compared with communication networks, transport-network reliability and throughput is a subject receiving relatively recent attention. In communication networks, the concept of two-terminal reliability dates back to as early as 1970, if not before [19]. In transportation, the concept of two-terminal reliability in a stochastic network becomes popular since the early 1990s [48]. Early work in both fields focuses on natural and technological hazards rather than threats or attacks. With threats or attacks, each researcher appears to hold a different perspective. Aiming for uniformity, we seek a consistent body-of-knowledge to improve

network reliability-and-throughput, including when the network is under attack. We wish to reduce the negative impacts of *both* natural and deliberate disruption on traffic flow as well as economic and non-economic transactions. After almost three decades of background research, we wish to report our progress to date. At the same time, we also need to further identify some common *strategies* and *tactics* that can avoid both. This long-term and comprehensive approach distinguishes our work from other network-design research. Meanwhile, a disclaimer is made in that we have been concentrating on connectivity reliability, rather than travel-cost reliability. Another disclaimer is that we would leave in the "Future Extension" section in the companion chapter to explicitly incorporate the impact of congestion on user or other types of network equilibriums.

The disruption of a transport network results in a net loss in local, state, and national economies because economic transactions are stopped or delayed. It also causes civil disorder and upsets people's quality-of-life. Here, we propose important technical advances for "prevention" rather than "cure." Hazards are prevented with anticipatory maintenance (or reconditioning) and upgrades (or improvement) on the most vulnerable components [34]. At the same time, attacks can be mitigated by limiting access to and placing extra security at strategic facilities, thus *hardening* the facility against exploitation. While we guard against hazards by network improvement, we guard against adversarial attacks by *gaming*. In general, it is much less costly to anticipate hazards and prevent attacks than to perform the required reparations and replacement afterwards. This is particularly the case when the emerging body of work on resilience interprets "hardening" not just in the sense of resistance to attack but also recovery.

Here is the layout of the rest of the chapter. Section 2 will review current literature on the subject, ending with a summary of the state-of-the-art. In Sect. 3, our research to date on this subject is then presented by first laying out the back-ground, followed by a summary of our research to date in Sect. 4. The latter is discussed in terms of our current results and the loose ends. The conclusion in Sect. 5 leads toward the following chapter, in that our current findings raise some new questions that have yet to be investigated. In a following, companion chapter, we put our research in perspective with a few other similar works, ending by posing some open questions in Sect. 1. These consist of ad hoc questions to which we have found partial answers. These answers are organized into five possible extensions in Sect. 2, including both modeling and computational suggestions. In Sect. 3, labeled 'Applications,' the implications of these research extensions are discussed in depth. The chapter ends with Sect. 4, emphasizing directions for further research.

2 Literature Review

In the reference section are listed dozens of literatures on the subject of network throughput and reliability. While there is in-depth treatment of the individual topics of throughput and reliability, there is but a dearth of information on examining both topics in a single model. It turns out that between throughput and reliability improvements, an *efficient frontier* can be formulated, representing the "win–win" solutions one can attain within an available improvement budget [33]. These *Pareto optima* on the efficient frontier are the non-dominated (or viable) solutions worthy of closer examination. Trades can now be made between the two performance measures-reliability and throughput-among these viable solutions, resulting in a desirable network design, such as a network improvement strategy. In this way, one can comprehensively assess the benefits of network improvement, when both reliability and throughput gains are considered.

Related to network reliability and throughput is *survivability*, which is defined here to mean the capacity to withstand a degradation or an assault. With a few exceptions, early research on survivability is mainly concerned with natural and technological degradation of components in the absence of adversarial tampering [1, 5]. In other words, only network *functionality* is analyzed, leaving out the important subject of network *security*. Instead of the typical independence assumption that underlines component failures, Kalyoncu and Sankur [27] considered dependency among damaging events. They consider secondary failures triggered by an initial failure. In order to inflict the worst damage, an initial attack may also be followed by a "complementary" attack. This concern paves the way for determining network security against a planned deliberate disruption. Considering the shattering damage it can inflict, such "double whammy" incidents are of particular interest to the attacker. Obviously, they are of equal concern to the defender.

As a forerunner to our work, Lyle et al. [33] introduced invulnerability as the measure of network security. This complements the traditional reliability measure for natural and technological disasters. A game was set up between the attacker and the defender to model exploitation and protection of network assets respectively. Following classic game theory [3], the goal of the defender-for instance a state department of transportation-is to ensure the network is invulnerable to attack. On the other side, the attacker (say the terrorist) wishes to inflict the worst damage. Most interestingly, a Nash equilibrium for this game was found, a condition in which all the participants play their best offensive or defensive strategies and cannot unilaterally improve, resulting in a stable solution. For a game where the defender hardens selected network components to guard against an attacker, the word posturing is used, as was done during the Cold War where a "first strike" was perceived to be unable to destroy the adversary-only elicit a counter offensive. This way, neither side has the motivation to upset the fragile stability. This points toward not only the need for a robust (reliable) network against hazards, but also a secure design to guard against an attack. Going forward, we will refer the defender as "player *Blue*" and the attacker as "player *Red*" for simplicity.

In Lyle et al. [33], a *valuation* of network security was found. In other words, the model allows the monetary value of network-security to be *imputed* through the use of opportunity cost. A critical component, for instance a collapsed bridge, is assessed at a loss that is equal to the opportunity cost associated with missing the facility. To the best of our knowledge, this is a precedent-setting occasion, during

which a monetary figure is *analytically* derived for the worth of security. Such valuation goes well beyond the replacement cost of the bridge and includes the implicit value of the bridge in facilitating economic and non-economic transactions, as manifested in traffic flow. However, the work of Lyle et al. only considers network reliability.

Related to our work is the literature on network interdiction. On a network path, Washburn and Wood [50] examined the maximal likelihood of defender Blue's interdiction of attacker Red's traffic and at the same time the minimal likelihood of Red's evasion. In short, Red exploits the worst scenario for Blue, while Blue minimizes Red's evasion. For a multicriteria game, Washburn and Wood [50] established that a pair of Pareto equilibrium strategies exists for the attacker and defender if both players have "ideal strategies." Unlike a mixed strategy that involves two or more defensive or offensive strategies, an ideal strategy is defined as the *unique* way to defend and attack a subject network. Hong [24] developed a strategic interdiction model of a non-cooperative game of network flow. A strategic model of network interdiction was presented, where two players have complete information and simultaneously choose their strategies. The Red adversary chooses a flow rather than a single path to attack. If there are multiple paths in a network, Red can attack them all at once, thus deploying a frontal aggregate advance en mass. In return, the Blue defender chooses a blockage strategy rather than intercepting a single arc, in that Blue can block multiple arcs simultaneously. By virtue of this, the model becomes more tractable and gives sharper results on equilibrium behavior. The network interdiction game is neither a zero-sum game nor even a strictly competitive game. Because of this, it is necessary to use an un-conventional solution technique to find all the equilibriums. Notice that Washburn and Wood [50] and Hong [24] mainly focused on logistical networks. Since we are interested inter-disciplinary applications, please be mindful that the physical and protocol layers of a communication network may introduce additional considerations in an interdiction problem for a communication network.

Replacing a network-centric approach with a policy-analytic approach, Fernadez and Puerto [17] presented a zero-sum multicriteria matrix-game, modeling the rivalry between Blue and Red where there are more than a single metric for performance (such as including both reliability and throughput in our study). A zerosum game was modeled, meaning that the Red's gain is exactly the Blue's loss. The relationship between single-criterion and multicriteria games was delineated, and they showed that linear programming can be used to solve both problems. To further apply such a non-cooperative game, Zhuang and Bier [52] identified equilibrium strategies where the opponents engage in a *rivalry* to outwit one another. They proposed a fully endogenous, *sequential* model of resource allocation for the defender in countering both terrorism and natural and technological disasters. Together with the work of Hausken et al. [22], it represents one of the very few works concerned with both hazards and attacks.

3 Background on Our Research

By seeing what has been accomplished to date on this subject, we wish to claim that our research is transformative rather than incremental. In this section, we wish to convey the complexity of this subject, suggesting a "steep learning curve" that we have climbed over the last three decades, and the challenges that still lie ahead. Instead of employing obscure mathematics, we choose to provide this complex background via simple, motivating examples. These selected examples and the accompanying graphics are an integral part of our presentation. The examples are strategically placed at the appropriate junction where clarity is needed to make an important point, leading toward a more formal discussion in Sect. 4. The reader is invited to spend the extra time wading through these examples.

To start with, some basic premises governing a stochastic network are in order [10]:

- The state of component *i* is independent of the state of component *j*, for all *i* and *j*, i ≠ *j*.
- A node or arc is either *up* or *down*.
- Without loss of generality, any node can be modeled as an arc; i.e., a node is split into a pair of nodes connected by an arc, resulting in *n* arcs altogether for the network.

The network is at state S_k ($k = 1, ..., 2^n$) with probability P_k , together with the corresponding throughput or max flow $F(\cdot)$. Correspondingly, the expected (or mean) throughput of the stochastic network is: $V(\mathbf{x}) = \sum_{k=1}^{2^n} F(\mathbf{x}, S_k) \cdot P_k$, where $\mathbf{x} = (x_1, x_2, ..., x_i, ...)$ is a vector of path flows x_i , or arc flows with entries x_{ij} .

An upper bound on the expected (or mean) throughput V_U could represent a peak-hour throughput while a lower bound V_L the off-peak throughput, where the former taxes the design capacity and the latter does not. Such bounds are scientifically derived from *Jensen's Inequality*, which in its simplest form states that the secant line of a convex function lies above the graph of the function, thus forming an upper bound of the affected part of the function. If **x** is a vector of random variables and *f* is a concave function, the algebraic inequality then applies in vector form for **x** as $E[f(\mathbf{x})] \leq f(E[\mathbf{x}])$. For max flow, Jensen's inequality suggests that the expected maximum flow must be no bigger than the maximum network-flow when arc capacities are set at their expected values.

Example 1: Lower Bound and Upper Bound on Expected (or Mean) Throughput. Consider the example network shown in Fig. 1 where the two attributes for each arc represent availability and capacity respectively. The arc-path formulation for the 3 flow paths is: $s \rightarrow 1 \rightarrow 3 \rightarrow t$, with flow x_1 ; $s \rightarrow 1 \rightarrow 4 \rightarrow t$, with flow x_2 ; and $s \rightarrow 2 \rightarrow 4 \rightarrow t$, with flow x_3 . Let us define r_{ij} as the availability on arc (i, j), the paths have these reliability expressions: $R(1) = r_{s1}r_{14}r_{4t} = 1 \cdot 1 \cdot 0.9 = 0.9$, and similarly, R(2) = 1, and R(3) = 0.05.



The corresponding Lower Bound (LB) on expected (or mean) max-flow is

$$\max V_{L}(\mathbf{x}) = 0.9x_{1} + x_{2} + 0.05x_{3}$$
subject to
$$\operatorname{arc}(s, 1) \qquad x_{1} + x_{2} \le 5$$

$$\operatorname{arc}(s, 2) \qquad x_{3} \le 6$$

$$\operatorname{arc}(1, 3) \qquad x_{1} \le 2$$

$$\operatorname{arc}(1, 4) \qquad x_{2} \le 4$$

$$\operatorname{arc}(2, 4) \qquad x_{3} \le 5$$

$$\operatorname{arc}(3, t) \qquad x_{1} \le 4$$

$$\operatorname{arc}(4, t) \qquad x_{2} + x_{3} \le 7$$

$$x_{1}, x_{2}, x_{3} \ge 0$$

$$(1)$$

Notice that in the two inequalities governing flow on path 1, the x_1 can be reduced to $x_1 \le 2$, the tighter of the two LBs. The same is true for x_3 .

Similarly, the Upper Bound (UB) can be defined as

$$\max V_{U}(x) = x_{1} + x_{2} + x_{3}$$

subject to

$$\operatorname{arc}(s, 1)x_{1} + x_{2} \leq 5$$

$$\operatorname{arc}(s, 2) \quad x_{3} \leq 0.6$$

$$\operatorname{arc}(1, 3) \quad x_{1} \leq 2$$

$$\operatorname{arc}(1, 4) \quad x_{2} \leq 4$$

$$\operatorname{arc}(2, 4) \quad x_{3} \leq 2.5$$

$$\operatorname{arc}(3, t) \quad x_{1} \leq 3.6$$

$$\operatorname{arc}(4, t) \quad x_{2} + x_{3} \leq 7$$

$$x_{1}, x_{2}, x_{3} \geq 0$$

(2)

Notice that the only difference between the LB & UB models is in the coefficients in the objective function and the right hand sides (RHS). Here the RHSs are $r_{ij}u_{ij}$, such as $r_{s2}u_{s2} = (0.1)(6) = 0.6$ for arc(s, 2).

Before we proceed any further, we wish to show a bit more of the computational complexity of the network problem we are solving, again through the graphical example in Fig. 2 below. While flow is *linear* for this network, the two-terminal reliability from origin s to destination t shows nonlinearity and analytical complexity in its functional form. The probability that flow from the source reaches the sink equals the sum of the probabilities of being in a state where the sink is connected to the source. For the simple series network the only state where the sink is connected to the source is when both components are functioning. The corresponding reliability is r_1r_2 . Similarly, the probability flow from the source reaches the sink for the parallel network is the sum of the probabilities for the states where at least one of the components is up. It is equal to $r_1 + r_2 - r_1r_2$. These two fundamental relations can be applied recursively to larger networks to reduce the network to a single-component network. This method of calculating network reliability is known as *network* reduction [11, 38]. The main difficulty with this method is that as the number of components n grows, the number of states (k) grows exponentially $k = 2^n$. The generally high component reliability within a communication network, however, leads to the conclusion that the most probable states are those where only very few of the network components are in a failed state. This allows the process of enumerating the most probable states to be simplified, rendering the computation feasible [15]. While the most-probable-states model is a practical way to bypass the curse of dimensionality for reliable networks, it becomes computationally inefficient for large networks with many failing components.

To fix ideas, the current statistical-reliability of this network can be quantified using the network-reduction (factoring) method [39]. The s-to-t reliability function, r_{st} , has seven nonlinear terms, consisting of products of availabilities r_{ij} for each of the relevant arcs (i, j):



120

$$r_{st} = r_{s1}r_{13}r_{3t} + r_{s1}r_{14}r_{4t} + r_{s2}r_{24}r_{4t} - r_{s1}r_{13}r_{3t}r_{14}r_{4t} - r_{s1}r_{14}r_{s2}r_{24}r_{4t} - r_{s1}r_{13}r_{3t}r_{52}r_{24}r_{4t} + r_{s1}r_{13}r_{3t}r_{14}r_{s2}r_{24}r_{4t}$$
(3)

In network-optimization modeling, convexity is a highly desirable analytic property, since it helps to determine whether a global (true) optimum is obtained. Unfortunately, little can be stated regarding the convexity of the mathematical expression for network reliability as shown. For real-world networks, Shier [44] stated that "network reliability problems are … among the most challenging of all computational problems."

Example 2: Design against Natural and Technological Hazards Let us start with a simple, generic infrastructure network. Each arc in Fig. 2 has a *capacity*, or the most traffic that it can handle. An up-or-down binomial distribution (a Bernoulli random variable) on the arc represents *arc availability*. For the work performed by the authors to date, component availability is a given number based on historical records. An arc availability of 1 means it is fully functioning, while a 0 means it is down. An availability of 0.8 means the arc is functional 80 % of the time, with a 20 % downtime. In this Figure, the normal (Roman) typeface represents arc availability, while the italics typeface represents capacity. Thus, arc(s, 2) has an availability of 0.8 and a capacity of 6.

The propensity for failure, as encoded by an arc's *availability*, makes this a *stochastic*, instead of a deterministic, network. Notice that capacity and availability can also be specified for nodes, in addition to arcs. Well-established network-reliability and throughput analysis can be used to provide a more robust network [32], facilitating better traffic flow from origin *s* to destination *t*. To make the network in Fig. 2 more robust, we made selected arcs more reliable and provided them with more capacity. For instance, the arc availability for (*s*, 1) is upgraded in Fig. 2 from 0.8 to 0.9—or by an increment of 0.1. The same arc is provided with an additional 2.5 units of capacity, improving it from 5 to 7.5.

In order to better endure potential natural and technological hazards, which may disable some but not all connections from node *s* to *t*, the capacity and availability upgrades are intended to maximize expected-flow (or the average throughput) and provide better O-D connectivities. While *availability* improvements strengthen the *two-terminal* reliability from *s* to *t* [14], they also improve the *expected throughput*. In other words, when components are functioning, more traffic can get through *on the average*.

When deliberate attacks are possible, traditional performance measures of reliability and throughput *must* be supplemented with *security* measures. Traditional methods for calculating network reliability assume component failures are independent, or failure in one component does not trigger failure in another component [6, 14, 44]. This is an acceptable assumption if the network components are only subject to rare, random failures due to natural degradation. Unfortunately, network components are now also subject to failures from deliberate tampering. We like to think of random, natural failures in network components as caused by "Mother Nature," who does not necessarily exploit the component where failure will cause the most disruption. However, malicious adversaries will choose to attack the critical components that contribute most to network performance. This leads to the conclusion that resilience against random failures alone, while important, is not necessarily a good measure to ensure network security in today's world, where acts of terror abound. Cox [16] examined probability models for failures in packetswitched data networks caused by attacks on most-loaded nodes. Such networks are inherently resilient to attacks if and only if they have enough spare capacity at each node (typically about 10 k % more than would be required in the absence of attacks, if the attacks are focused on the k most heavily loaded nodes).

Another problem with traditional reliability measures is that a strategy that makes a network component impervious to exploitation has no *measurable* (quantifiable) value (in the literature), although such a strategy has obvious *security* and *economic* value. Theoretically, one can expend an inordinate amount of resource and still fail to render a network "hardened" against an attack–when one hardens the wrong arc. Hardening is different from infrastructure upgrades that guard against a hurricane or earthquake. An upgrade carries with it a cost tag and an improvement in the "safety design factor," while hardening is an abstract concept that defies costing and valuation (to date). This problem motivated the search for a quantitative measure that an adversary uses to select the *most vulnerable* network component (the Achilles heel) to exploit. Accordingly, *hardening* such targets to thwart exploitation should have quantifiable value [43]. Haimes [20] suggested: "The fortification [hardening] of critical infrastructure systems, by reducing the states of their vulnerability and enhancing their resilience, could diminish their attractiveness as a targeted system."

Example 3: Illustrating Complexity of the Basic Metrics Following a prevailing practice with such a stochastic network, a binomial distribution of arc availability and capacity is assumed—an arc is either functioning or disabled with the corresponding probabilities. In Fig. 1, we show that even with a binomial distribution, the problem as captured by the combination of all the above examples is extremely difficult to solve. For example, arc(s, 1) functions 80 % of the time and fails 20 % of the time. When an arc fails, it diverts traffic to an alternate path. Thus, when arc (s, 1) fails, the flow can be diverted from path s-1-3-t to path s-2-4-t, assuming that all components are functional on the alternate path. To tackle the nonlinearity and analytical complexity of s-to-t reliability, numerical factoring and reduction algorithms are among the computationally feasible solutions as mentioned. To model travel time, one can also use a binomial distribution of arc costs. For example, travel time can be defined by the un-delayed arc cost and the delay cost respectively [8]. When an arc "fails," a time penalty of "delay cost" is imposed to account for congestion and queuing. These two costs, corresponding to the "operational" and "failed" operational modes, can be used to model a trafficdependent flow. Correspondingly, total network travel cost (including delays due to congestion) can be determined from the resulting flow pattern. This approach by passes explicitly computing the equally complicated travel-time equivalent of two-terminal reliability. $\hfill \Box$

In short, procedures such as factoring-and-reduction and considering most probable states make computing the highly nonlinear network-reliability measure feasible. Short of enumerative-simulation replications, upper-and-lower *bounds* can be used to approximate such system performances as expected throughput [10, 33, 42].

4 Our Research Results

Developed independently, our work addresses a similar problem by *explicitly modeling a network topology of nodes and arcs*, which is absent in the Zhuang and Bier's work. Instead of a conceptual framework, node-arc representation facilitates *physical implementation* of improvement and hardening decisions on a variety of physical networks. Through such a formulation, guidelines for network *designs* (including its topology) can be derived (at least in a meso-level of detail). We opt for an engineering design that is most robust and most secure. While network topology is considered in the work by Lou and Zhang [31], it concentrates only on operational strategies. This contrasts with our approach where both tactics and strategy. In our judgment, it is only by placing a monetary value on security that priorities can be established, and hence resource allocation can be made between various tactics and strategies.

To date, computational experiments with our model on three *realistic networks* suggest that a robust equilibrium-design exists that can thwart the worst possible damage. Notice the use of bounds leads toward an approximate solution, but the approximation does not affect the result that an equilibrium *exists*. Most important, the Pareto Nash equilibrium is stable irrespective of the unit costs associated with availability versus capacity improvement and how one wishes to trade between them. In other words, we have successfully identified circumstances for a "win–win–win" design—a design that is robust enough to guard against damage irrespective of how an adversary weighs the importance of throughput versus reliability. One may think of this as a shopper that is fortunate enough to locate a "cheaper" and "better quality" shirt. Naturally, it leads right away to a "win–win–win" purchase decision; the cheaper (the first 'win') and better quality (second 'win') shirt is the one to buy irrespective of how she trades price against quality (third 'win')—in other words, a "no brainer."

In trading between throughput versus reliability improvement, the trade can be *totally compensatory, totally non-compensatory,* or somewhere in between. Throughput and reliability are *totally compensatory* if the lack of one can be made up by more of the other. An example is the defender's belief that lack of "reliability" can be compensated by more "throughput." In other words, a "shot gun"

approach would likely hit the target even though the chance of each "bullet hitting the target precisely is small." The two criteria are *totally non-compensatory* if the decision-maker only cares about the more prominent of the two, to the exclusion of the remaining one. An example is "crisis management," a short-term practice that is driven only by an imminent threat (say the throughput will be totally severed by the attacker), in which the defender will rush to increase arc capacity, forgetting that such threats can be mitigated equally well by making the network more reliable, hence increasing the expected throughput. Between totally compensatory and noncompensatory, we can incorporate a whole host of preference structures (or styles) for decision-making for both Blue and Red [21].

4.1 Formal Mathematical Programming Models

Before we formulate the models, let us compile and define the following symbols in one place [42]:

- B_c the budget for capacity improvements
- B_r the budget for reliability improvements
- *B* total fixed budget for network improvement, including both capacity and reliability improvements
- c_{ij} the initial capacity of a directed arc from node *i* to node *j*
- d_{ij} the amount of improvement made to the arc capacity c_{ij}
- d vector of arc-capacity improvements for all arcs
- p_{ij} the initial availability of an arc(*i*, *j*)
- P_k probability of a network being at state k
- r_{ij} the amount of improvement to arc availability p_{ij}

$$\mathbf{r}$$
 a vector of r_{ij} , or $\mathbf{r} = (\leftarrow r_{ij} \rightarrow)^T$

- R_{st} the two-terminal reliability of a network, or *s*-*t* reliability
- R_{ij} the two-terminal network reliability after the removal of an arc (i j)
- S_k network is at state k
- *s* a source node (or the origin)
- t a sink or terminus node (or the destination)
- V expected (or mean) throughput of the stochastic network
- V_U the upper bound of expected (mean) throughput, or expected (mean) maximum-flow, or expected flow for short
- V_L the lower bound of expected (mean) throughput, or expected (mean) maximum-flow, or expected flow for short
- V_f the Expected Flow including the worst effect of arc loss on the network
- V_c the Expected-flow-damage Utility (a continuous variable that is 0–1 ranged)
- V_r the Reliability-damage Utility (a continuous variable that is 0–1 ranged)
- x_i the flow through a path *j*
- x_{ij} the flow through a directed arc(*i*, *j*)

- **x** vector of arc flows, where $\mathbf{x} = (\langle x_{ij} \rangle)^T$ for arc flow, or $\mathbf{x} = (\langle x_j \rangle)^T$ for path flow
- X_1 the feasible region (or polyhedron) of an expected-flow network-improvement model
- X_2 the feasible region (or polyhedron) of a reliability improvement model
- *X* combined feasible region (polyhedron) of an expected throughput and reliability-improvement model, or $X = X_1 \cap X_2$
- y_{ij} a game variable, a continuous variable to denote the percentage time an arc (i, j) is hardened
- z_{ij} dual game-variable, a continuous variable to denote the percentage time arc (i, j) is tampered with

Example 2 can now be formalized in the following mathematical programming model that maximizes expected flow:

$$\max V_{U}(d, r)$$
subject to
$$\sum_{k} x_{sk} - V_{U} = 0 \qquad (4)$$

$$\sum_{k} x_{jk} - \sum_{i} x_{ij} = 0 \quad \text{for all } j$$

$$V_{U} - \sum_{i} x_{it} = 0$$

$$\sum_{(i,j)} c_{ij} d_{ij} \leq B_{c} \qquad (5)$$

$$\sum_{(i,j)} p_{ij} r_{ij} \leq B_{r}$$

$$p_{ij} + r_{ij} \leq 1 \qquad \text{for all } (i,j)$$

In the maximization objective function, the *expected* flow $V_U(\mathbf{d}, \mathbf{r})$ —after accounting for arc failures in a stochastic network—is highlighted to be a function of arc capacity and availability improvements. Thus, the larger the arc capacities and the higher the availabilities, the expected flow will increase. In this model, x_{ij} , V_U , \mathbf{d} , \mathbf{r} are non–negative variables. All other symbols are fixed parameters, including the capacity- and reliability-improvement budgets B_c and B_r Eq. (4) is the set of typical network conservation-of-flow constraints. Equation (5) is the set of constraints governing capacity and availability improvements. For example, the first constraint in Equation set (5) suggests that the flow on $\operatorname{arc}(i, j)$ is limited by the expected capacity on this arc, i.e., the product of the expanded arc capacity $(c_{ij} + d_{ij})$ and the improved arc availability $(p_{ij} + r_{ij})$. According to Chan et al. [10], the objective-function value is an upper bound on the expected throughput. This model is re-written in shorthand as max $\{V_U | \mathbf{d}, \mathbf{r} \in X_1\}$, where X_1 is the polyhedron formed by network-improvement constraint sets denoted by Eqs. (4) and (5).

Meanwhile, the numerical result as shown in Example 4 below can be derived from the following model that maximizes flow-damage utility (with hardening game):

$$\max V_{f}(\mathbf{x}, y_{ij})$$

subject to
$$V_{f} \leq V_{U} - X_{ij}(1 - y_{ij}) \text{ for all } (i, j)$$
(6)
$$\sum_{(i,j)} y_{ij} = 1$$

$$\mathbf{x}, \mathbf{d}, \mathbf{r} \in X_{1}$$

Then there is the following group of constraints we label collectively as $\mathbf{r} \in X_2$, or the feasible region (polyhedron) of a reliability-improvement model:

$$p_{ij} + r_{ij} \leq 1 \qquad \text{for all} \quad (i,j)$$

$$\sum_{(i,j)} p_{ij} r_{ij} \leq B_r \qquad (7)$$

$$\mathbf{r}, y_{ij} \text{ nonnegative}$$

Notice the cost of hardening a target is not included in this model. Rather, the valuation of a hardening strategy is to be imputed as shown in Eq. (14), as will be discussed. Consistent with a mixed-strategy game, the defensive game-variables y_{ij} sum to 1 as enforced by the corresponding constraint in Eq. (6). An interesting result from this Model is that the shadow prices z_{ij} that correspond to each of the game constraints automatically sum to 1 as well. These shadow prices represent an optimal Red offensive strategy. Suppose we drop the last set of improvement constraints $\mathbf{x}, \mathbf{d}, \mathbf{r} \in X_1$, defining the feasible region (polyhedron) of the Expected-flow network-improvement model. The current Model boils down to a classic linear program (LP) of a two-person, zero-sum, non-cooperative game. Barring degeneracy, there is a solution to a properly formulated LP—in this case a stable strategy for both Blue and Red that is referred to as *Nash equilibrium*. Furthermore, the variables y_{ij} can assume fractional (not just binary) values, suggesting a mixed-strategy game.

The most comprehensive model is one that maximizes reliability *and* flow-damage utility:

$$\max V_{f}$$

$$\max V_{r}$$
subject to
$$V_{f} \leq V_{U} - x_{ij} (1 - y_{ij}) \qquad \text{for all } (i, j)$$

$$V_{r} \leq (R_{ij} - R_{st} + 1) (1 - y_{ij}) + y_{ij} \qquad \text{for all } (i, j) \qquad (8)$$

$$\sum_{(i,j)} y_{ij} \leq 1$$

$$B_{c} + B_{r} = B$$

$$\mathbf{x}, \mathbf{y}, \mathbf{d}, \mathbf{r} \in X$$

$$y_{ij} \quad \text{nonnegative}$$

where $X = X_1 \cap X_2$, or the intersection of the two polyhedral X_1 and X_2 . The model can be *decomposed* into the gaming part and the improvement part. The former is an LP, while the latter a nonlinear program in general. The gaming constraints are:

$$V_{f} \leq V_{U} - x_{ij}(1 - y_{ij}) \qquad \text{for all } (i, j)$$

$$V_{r} \leq (R_{ij} - R_{st} + 1)(1 - y_{ij}) + y_{ij} \qquad \text{for all } (i, j) \qquad (9)$$

$$\sum_{(i,j)} y_{ij} \leq 1$$

whereas the improvement constraints are

$$B_{c} + B_{r} \le B$$

x, **y**, **d**, **r** $\in X$ (10)

The two parts can possibly be solved separately and iteratively (until convergence is obtained), as suggested by Lai [29].

4.2 Summary Results

To assess the value of component failures, we propose to measure it by a *utility function* [13, 2], which will allow Blue-or-Red's decisions to be *valuated*. A monetary valuation will lead toward an informed resource-allocation decision. Here, a 0–1 ranged utility function is used to measure the value the attacker and defender places on a particular asset, with unity being the highest value and zero being the lowest. A preliminary multicriteria-optimization model was implemented by including both *flow-damage utility* V_f and *reliability-damage utility* V_r [42]. Reliability-damage utility, V_r , quantifies the compromise when a network has been damaged by an adversary. Reliability-damage utility is formally defined as the difference in network reliability before-and-after a hazard or an attack: Damage = Reliability (before) – Reliability (after), and correspondingly: Reliability-

damage utility = $V_r = 1$ – Damage. The latter formula follows the convention that one wishes to *maximize* a utility function (as contrasted with "damage," which is to be minimized).

Similarly, flow-damage utility, V_f is introduced. Unlike reliability-damage utility, V_r , throughput means "max flow." We are mindful of the "max-flow/min-cut theorem" that suggests the that the maximum flow is determined by the minimum cut-set of the arcs that will disconnect the origin from the destination [18]. Recall the $V(\cdot)$ stands for the expected (or mean) max-flow (or throughput) from source *s* to sink *t*. Accordingly, expected flow-damage utility is defined in two steps. First, we define flow damage, or F_{damage} , as

$$F_{\text{damage}} = \text{Normalized decrease in throughput} = [V(A) - V(A - (i, j))]_{\text{norm}}$$
 (11)

where A is the set of arcs in the network; (i,j) is the bottleneck arc taken out by Red. Thus the compromise in throughput is simply the difference between the integral network and the network after a critical arc has been taken out by the adversary. F_{damage} is normalized against V(A) so that the resulting utility function will be ranged between 0 and 1. Again, following the convention that a utility function is to be maximized, we define flow-damage utility as:

$$V_f = \left(1 - F_{\text{damage}}\right)_{\min} \tag{12}$$

The subscript "min," or the "minimum value for flow damage utility," suggests that a malicious attacker would inflict the *most* damage. Again, please remember that a normalized "damage" metric, V_{f^5} is 0–1 ranged. When flow damage, F_{damage} , is maximized, the complement $(1 - F_{\text{damage}})$ is minimized. In the case of natural and technological hazards, we *anticipate* the worst scenario in which a hurricane or earthquake happens to take out the most critical components in the network, even though Mother Nature is not necessarily malicious. In this regard, both damage utilities take into account the performance degradation of the entire network after a deliberate, adversarial attack. Through a utility function, we have now captured the consequence of critical component failures in a 0–1 valuation scale, whether it is the result of natural hazards or malicious attacks. When both improvement and hardening are considered, a *combined* multi-attribute utility function is further defined below, allowing the trades between reliability damage and flow damage—what appear to be incommensurate "apples" versus "oranges" comparisons—to be performed.

In a classic zero-sum non-cooperative game, Blue and Red each decides on a single arc to protect (harden) and to exploit (attack). This constitutes a *pure-strategy* game in which both sides commit to a unique course of action. Oftentimes, it is assumed that the two opponents are playing such a game repeatedly. In a scoring matrix of "payoffs," the percentage of times Blue hardens an arc(i, j) is designated for a row of the game matrix—a square matrix with a row, and a column for each target arc as shown below. Similarly, the percentage of times Red attacks an arc (i, j) arc(i, j), is designated for each column. Now we have a *mixed-strategy* game, in which an attack is launched between two or more arcs over time, and,

correspondingly, a defense is set up for these arcs. The damage-utility outcome for each hardening/exploitation intersection between the Blue and Red players can be identified by the cells in the following payoff matrix, which is constructed for the network in Fig. 2.

In the matrix shown in Table 1, the diagonal "payoff" elements suggest that if Blue has hardened $\operatorname{arc}(1,3)$, Red's attempt to attack is to no avail. Blue's hardening strategy maintains network integrity even after Red's attack—retaining the same flow-damage utility, *V*, without any compromise. On the other hand, the off-diagonal payoffs show performance degradation should Blue harden a component other than the one Red attacks. For example, if Red attacks the *bottleneck* $\operatorname{arc}(s, 2)$, networkflow performance—as measured in the normalized throughput utility *V*—is now degraded to $V - x_{s2}$ after the attack, where x_{s2} is the flow on $\operatorname{arc}(s, 2)$ that has been eliminated. A similar game matrix can be repeated for the reliability utility function.

From this matrix, the utility/value function (before 0-1 ranging) may be generated by calculating the expected value of each Red strategy. For example, the expected value of Red's attack on (*s*, 1), corresponding to column 1, is

$$Vy_{s1} + (V - x_{s1})(y_{s2} + y_{13} + y_{14} + y_{24} + y_{3t} + y_{4t})$$

= $Vy_{s1} + (V - x_{s1})(1 - y_{s1})$
= $Vy_{s1} + V - V y_{s1} - x_{s1} + x_{s1}y_{s1}$
= $V - x_{s1}(1 - y_{s1})$ (13)

Now let us range the column expected-value between 0 and 1. A rational Red will choose the column that renders Blue the lowest payoff possible.

Example 4: Flow-Damage-Utility Maximization with a Hardening Game Aside from preventing natural failure, a component is now hardened against adversarial attacks in the model output shown in Fig. 3. By examining arc(s, 1), for example, we clearly see that the resulting availability and capacity improvements are different from Fig. 1. Switching from the expected throughput of Fig. 2 to the flow-damage utility as a figure-of-merit, valuation is now placed on the potential damage. As a result, we observe that resources are placed in a more "spread out" fashion in order to anticipate "all moves" by the adversary. When a game-theoretic model is fully implemented *with hardening*, the most-frequently used arcs-(s, 1) and (4, t)—receive the most hardening, lest these become the natural target of exploitation.

In Fig. 3, the percentages denote the frequency with which an arc is *hardened* to make it impervious to attack. In other words, we are playing a mixed-strategy game, in which a player adopts two or more strategies, where hardening a particular arc constitutes a strategy. Thus, $\operatorname{arc}(s, 1)$ is hardened 47 % of the time, while $\operatorname{arc}(4, t)$ is hardened 48 % of the time. In the context of a mixed-strategy game, this means that —roughly speaking—we only harden $\operatorname{arcs}(s, 1)$ and (4, t) respectively only "half of the time," leaving them "undefended" the remaining "half of the time." Unlike a

	Player red							
	arc	(<i>s</i> , 1)	(<i>s</i> , 2)	(1, 3)	(1, 4)	(2, 4)	(3, t)	(4, t)
Player blue	(<i>s</i> , 1)	Λ	$V - x_{\rm s}$	$V - x_{13}$	$V - x_{14}$	$V - x_{24}$	$V - x_{3t}$	$V - x_{4t}$
	(<i>s</i> , 2)	$V - x_{s1}$	V	$V - x_{13}$	$V - x_{14}$	$V - x_{24}$	$V - x_{3t}$	$V - x_{4t}$
	(1, 3)	$V - x_{s1}$	$V - x_{s2}$	Λ	$V - x_{14}$	$V - x_{24}$	$V - x_{3t}$	$V - x_{4t}$
	(1, 4)	$V - x_{s1}$	$V - x_{s2}$	$V - x_{13}$	Λ	$V - x_{24}$	$V - x_{3t}$	$V - x_{4t}$
	(2, 4)	$V - x_{s1}$	$V - x_{s2}$	$V - x_{13}$	$V - x_{14}$	V	$V - x_{3t}$	$V - x_{4t}$
	(3, t)	$V - x_{s1}$	$V - x_{s2}$	$V - x_{13}$	$V - x_{14}$	$V - x_{24}$	V	$V - x_{4t}$
	(4, t)	$V - x_{s1}$	$V - x_{s2}$	$V - x_{13}$	$V - x_{14}$	$V - x_{24}$	$V - x_{3t}$	V

	Davioti matriv	T ayon many
,		•
	٩	Ś
1	C	2
	0	5



pure strategy game that commits to a single strategy, remember that a mixed strategy game introduces uncertainty in gaming. Within rounding errors, the hardening percentages among all arcs sum to unity: 47 + 3 + 3 + 48 = 100 % as expected. Here, the two separate problems—managing reliability-damage and flow-damage—are linked by the budget constraint $B_c + B_r \le B$, where B_c is the budget for capacity improvement, B_r for availability improvement, and B is the total available budget. While B is given, both B_c and B_r to be determined.

Blue has a limited budget *B* to prepare for both natural hazards and deliberate attacks. An efficient spending between availability and capacity investment is derivable from this model. An exchange ratio will account for the difference in unit costs of capacity versus reliability improvement. In such a model, we wish to account for the loss of economic and non-economic activities when a critical facility, like a bridge, is taken out by an adversary [51]. We propose that the monetary value of this loss be obtained by comparing the budgets to achieve two Nash equilibriums. The first equilibrium is obtained without hardening, leaving the network subject to an attack. The other is obtained with the desired hardening protection, making it invulnerable to attack. This is illustrated in Fig. 4, in which an efficient frontier for reliability-damage utility and flow-damage utility is shown for both scenarios, where the former V_r is shown as the horizontal axis and the latter V_f as the vertical axis. Here, an efficient frontier contains the Pareto optimal solutions, or the viable ways to protect the infrastructure for a given budget.



Fig. 4 Imputing the worth of security (Adapted from [30])

Here is an important illustration of another advantage of a *utility-theoretic* model. In our mathematical programming model, suppose the game constraints with hardening—as partially represented by the game matrix above as shown in Eq. (9)—are replaced by the game constraints without hardening. We start with an initial budget: *B* (no hardening), which is a feasible budget that effects the reliability and capacity improvements in the absence of hardening. The budget is now increased to *B* (hardening), until the resulting efficient frontier is upgraded to the frontier when hardening is performed. Graphically, this is shown when the "+" signs are now aligned with the " \Box " symbols, where the several "boxes" \Box of Fig. 4 represent the *efficient frontier* of the hardening decision. To the extent that one wishes to buy the best security, the efficient frontier represents the maximum reliability-damage utility and flow-damage utility that can be accomplished with a given budget, or the "best bang for the buck" in colloquial expression. When the "+" signs are now aligned with the " \Box " symbols, the improved network is now valuated the same as the one with hardening in place.

Let us take the difference between the final budget and the initial budget:

Value of hardening
$$= B(\text{hardening}) - B(\text{no hardening})$$
 (14)

The difference represents the monetary value of a hardening strategy, or the worth of security to ensure normal economic and non-economic transactions! The procedure uses the *tangible* cost of infrastructure improvement to *impute* the *intangible* value of hardening. The security insurance amounts to the economic and non-economic value of maintaining normal- transaction traffic, now conveniently valuated in dollars. Let us put it another way: Suppose a capital outlay is used to ensure security; we equate the potential budget shortfall to the opportunity cost of the security *compromise*.

This valuation procedure is consistent with the "asset management" procedure widely practiced today, in as much as asset management is a process of resource allocation and utilization [9, 28, 31, 34, 47]. Yet our procedure is simpler and more precise, since an accurate accounting is always brought to question in traditional asset valuation methodologies [35, 4]. Notice we are following the basic definition of asset management, which is "a systematic process of maintaining, upgrading, and operating physical assets cost-effectively" [28]. However, we emphasize that there is merit in looking at assets as a whole, including the disruption cost to commerce and life style, above and beyond the replacement cost of lost assets [9]. In other words, the disruption cost should be counted in addition to replacement costs used in existing asset management practices. Moreover, the concept of marginal cost regarding marginal survivability improvement has important implications that deserve room for further discussion. Notice that it cannot be argued that arcs with lower upgrading costs relative to survivability improvement should be more attractive for investment, since the marginal cost to guard against attack is not known until after the imputation valuation in Eq. (14). Furthermore, the valuation is

for the entire network, rather than a particular arc; hence marginal analysis is not suitable for deciding between upgrade versus hardening.

4.3 Some Loose Ends

Through the above background discussion, it is clear that the research on network throughput and reliability is highly complex, yet infantile in its accomplishments. Without a firm grasp of the state-of-the-art, there really is no way to lay out a research agenda, where we wish to bring some unity and maturity in the field. Having acquired the background body-of-knowledge, let us now propose several paradigms for innovative model building, followed by a research plan, and ending with the implications for both applications and future extensions. The goal is to advance a *unifying* body of knowledge on the subject.

As "elegant" as it might sound, the model is a highly nonlinear mixed-integerprogram (MIP), consisting of both continuous and discrete variables that cannot be solved easily. In our investigation to date, the following propositions have been proven, serving as the basis for a possible solution to this complex model. The propositions are taken directly from our published results [42]. To date, they offer a better understanding of the properties of a complex infrastructure network; what AWAITS FURTHER WORK IS AN EFFICIENT SOLUTION ALGORITHM. This is by no means a trivial task, considering the huge size of most infrastructure networks.

Proposition 1 The nonlinear model can be approximated by a linear model. Through this linearization, the approximate solution is asymptotically close to optimal. \Box

Proposition 2 There exists a unique range-equalized compromise-solution to the linearized version of the model. \Box

Proofs of these two propositions are documented in Schavland et al. [42] and will not be repeated here. Proposition 1 is based on the commonly recognized fact that communication networks are highly reliable (i.e., reliability > 0.95). For a highly reliable network, the final availabilities of network components are not very different from the initial availabilities after improvement. For a highly reliable network, it is also clear that it pays to invest in capacity improvement rather than reliability improvement. The payoff is simply better. Through this linearization, the approximate solutions to the LP game in Model (8) will be close to optimal. Through this procedure, Model (8) can now be decomposed into the improvement portion and the hardening portion according to Lai [29]. As most communication networks are highly reliable, reliability improvement is much less important in comparison to hardening. The hardening component is the dominant part of the decomposed model. Given the linearization, Model (8) can effectively be solved by the linear-approximation procedure for the hardening game. Aside from obvious computational convenience from linearizing the highly complex reliability utility-
function, therefore, Proposition 1 also potentially allows interesting analytical properties to be discerned—a task we wish to propose to our colleague researchers.

Since now that Model (8) can be turned into a multicriteria LP, Proposition 2 follows from Proposition 1. Now both the flow-damage utility $V_{\rm f}$ and the reliabilitydamage utility V_r are to be ranged to unity by way of *range equalization* [45]. With these numerical procedures, the outcome space showing the efficient frontier is examined. In this context, range equalization means that the solution is independent of the relative weight placed on the throughput versus reliability utilityfunctions. It has been shown in Schavland et al. [42] that neither is the equilibrium affected by the difference in unit costs of availability or capacity improvement. (This is a formalization of our earlier example of a "cheaper and better-quality shirt," which is the indisputable buy, irrespective of how the consumer weighs price against quality.) In this case, the robust solution is a Pareto Nash equilibrium, or a viable, stable solution considering both reliability and throughput. In other words, Proposition 2 suggests a robust equilibrium solution. Notice this result was obtained for the generic network model formulated above. We expect this result to hold for telecommunication, transportation and power networks that can be modeled under this mathematical paradigm.

Setting aside improvements to prevent hazards, if one wishes to simply assess a defensive strategy against adversarial tampering, the model degenerates into a linear program (LP), since the nonlinear constraints associated with network design can now be dropped. Performance evaluation can then be obtained directly from such an LP model, instead of from the nonlinear MIP model. Obviously, the LP model is much simpler. Since the gaming LP is a specialized form of the general model, we conjecture that the LP solution is also a robust defense-or-offensive strategy, or again a Pareto Nash equilibrium. We suggest that the model is now decomposable into the network-design part for natural hazards and the gaming part for malicious attacks. Each part can be applied by itself or used in conjunction with the other. For computational efficiency, it is conceivable that the two parts can be solved iteratively, with the results combined at a later stage, as suggested by Lai [29]. Instead of formal procedures involving Lagrangian or Karash-Kuhn-Tucker duality, degrees of individual optimality and feasibility-membership functions can be used directly to make tradeoffs between objectives and constraints, respectively, if necessary. This technique is judged to be robust and adept in planning for both hazards and attacks—the subject of this and the following chapter. Obviously, THE COMPUTATIONAL PROCEDURE THAT GUARANTEES CONVERGENCE HAS YET TO BE FOUND.

In participating in a game, we wish to validate that those with better knowledge about other stakeholders tend to do better than those who do not. Under this situation, we assess the role of *epistemic knowledge*, which is defined as knowledge supporting a belief, truth or hypothesis. *Epistemic utility* is a measure of the value of this knowledge in supporting a hypothesis, i.e., a measure of the importance of the information about the other game players relative to the belief that the other players are going to adopt a certain strategy. Epistemic utility theory has so far furnished us with a number of arguments for some of the central norms governing partial beliefs. Of course, some are stronger than others, and it seems that none is yet decisive; each relies on a premise that we might reasonably question [40]. For example, is it legitimate to employ the notion of expected utility when the belief function by the lights of which the expected utility is calculated is not a probability function?

Stirling [46] proposed a theory of *multi-agent coordinated decision-making* using epistemic utility. It represents an alternate theory of multi-agent decision-making, termed *satisficing games*, which is designed to address the limitations that are imposed by the assumption that players are motivated exclusively by self-interest. However, it does not abandon the principles of value and performance that are vital to rational behavior, which includes both common good and self-interest. Collaboration can be facilitated by replacing individual rationality with a less rigid model that permits decision makers to expand their spheres of interest beyond immediate, individual concerns.

This asymmetrical behavior may tell more about this competitor than the performance metric that she commands at the present moment. In some communication games this game information, which we may call *side information* needs to be considered together with the actual information being transferred. An example is the "Supply chain game" by Janakiram et al. [26], an empirical game played among up to 18 members in which there is intensive human interaction. In this light, the game players really participate in a *communications game*. Information value has so far been discussed in terms of the utility derived from the information being communicated. Accordingly, payoff functions associated with this value will tend to be associated with 'operational' utility, and therefore the decision process the information supports. There is, however, additional information deliberately or inadvertently generated in competitive communication environments, ranging from "blowing smoke" to "body language." This type of information is the game information; it is responsible for the style of play and performance of the game being considered. Players' actions will generate and affect this information. An example is the amount of asymmetry between players, in which one player in the business world may engage in a certain marketing strategy in excess of other players. Another example arises in information warfare, in which *posturing* can be an effective deterrent to cyber attacks on a computer network.

5 Conclusions

Capitalizing on decades of investigations and recent progress, we are confident a unifying paradigm can be formalized mathematically to model both natural hazard and deliberate tampering. This is to be accompanied by efficient solution algorithms for practical day-to-day engineering and other technical problems. Insights gathered to date suggest the following scientific premises:

• Natural hazard and deliberate attacks can be thwarted separately or in combination. In both cases, the network is improved for better reliability and throughput.

- The monetary worth of ensuring network security can be *imputed* from a mathematical model directly, independent of any conventional asset-management accounting-schemes.
- Through linearization of the highly nonlinear model, the analytical properties of such a model can be formally delineated. Based on these properties and the use of upper-and-lower performance bounds, the model can be solved efficiently, where an upper bound accompanies the lower bound as the "liberal" and "conservative" estimates.

Aside from the methodological "loose ends" mentioned above, the above findings raise some new questions that have yet to be investigated. Throughout our discussions, we emphasized the virtue of 'prevention' rather than 'cure.' Thus far, we have argued that it is better investing in enhancing an infrastructure ahead of time. However, such preparation is justified if only this investment means a reduction of cost for the decision maker. Otherwise, it does not make any sense investing in preparation for rare events or attacks instead of saving the money to react in a proper manner. It boils down to how can one compare the cost of improving an infrastructure in preparation for a disaster of a certain magnitude "x" and the cost of recovering such an infrastructure for the same disaster. In practice and especially for minor disasters, rehabilitating a highway infrastructure might be quick and easy. Its cost could even be inexpensive, depending on factors such as pre-established agreements and preparation. In real-world applications, there are whole hosts of issues that have yet been examined in this chapter. The following chapter is meant to go further into these issues by systematically organizing them, identifying possible resolutions, and end with how these resolutions-both modeling and computational advances-can make a difference in real-world applications.

Acknowledgments This chapter draws heavily from the published works of the author. The original publication sources are properly cited throughout. Certain sentences may be paraphrased from these sources, and certain figures are reproduced with or without editing. All the published research cited in this fashion was funded by the U.S. Department of Defense (DOD) and was performed while the author was a government employee. The DOD support is gratefully acknowledged. The author also wishes to acknowledge the assistance of Henry Shyllon who assembled some of the reference publications used in this chapter upon the author's suggestion. Mr. Shyllon also formatted the first draft of this chapter according to the publisher guidelines.

References

- 1. Abdel-Rahim A, Oman P, Johnson B, Tung L-W (2007) Survivability analysis of large-scale intelligent transportation system networks. Transp Res Rec 2022:9–20
- 2. Aguirregabiria V, Mira P (2010) Dynamic discrete choice structural models: a survey. J Econ 156:38–67
- Aliprantis CD, Chakrabarti SK (2000) Games and decision making. Oxford University Press, New York
- 4. Amekudzi A, McNeil S (eds) (2008) Infrastructure reporting and asset management. ASCE Press, Reston, VA

- Balakrishnan A, Mirchandani P, Natarajan HP (2009) Connectivity upgrade models for survivable network design. Oper Res 57(1):170–186
- Ball MO, Provan JS (1982) Bounds on the reliability polynomial for shellable independence systems. SIAM J Algebraic Discret Methods 3:166–181
- 7. Bell MGH (2000) A game theory approach to measuring the performance reliability of transport networks. Transp Res 34B:533–545
- 8. Bell MGH (2003) The use of game theory to measure the vulnerability of stochastic networks. IEEE Trans Reliab 52(1):63–68
- 9. Cambridge Systematics Inc (2009) An asset-management framework for the interstate highway system, NCHRP project 20-74. National Cooperative Highway Research Program, Transportation Research Board, Washington
- Chan Y, Yim E, Marsh A (1997) Exact and approximate improvement to the throughput of a stochastic network. IEEE Trans Reliab 46(4):473–486
- Chari MK, Provan JS (1996) Calculating k-connectedness reliability using Steiner bounds. Math Oper Res 21(4):905–921
- Chen A, Yang H, Tang WH (1999) A capacity related reliability for transportation networks. J Adv Transp 33(2):183–200
- Chen A, Kongsomsaksakul S, and Zhou Z (2007) Assessing network vulnerability using a combined travel demand model. Pre-prints of the transportation research board 2007 meeting, Washington, DC, Transportation Research Board
- 14. Colbourn CJ (1987) The combinatorics of network reliability. Oxford University Press, Oxford
- Colbourn CJ, Harms DD (1994) Evaluating performability: most probable states and bounds. Telecommun Syst 2:275–300
- 16. Cox LA Jr (2009) Making telecommunications networks resilient against terrorist attacks. In: Bier VM, Azaiez MN (eds) Game theoretic risk analysis of security threats. International Series in Operations Research and Management Science. Springer, New York
- Fernadez FR, Puerto J (1996) Vector linear programming in zero-sum multicriteria matrix games. J Optim Theory Appl 89(1):115–127
- Ford LR Jr, Fulkerson DR (1962) Flows in networks. Princeton University Press, Princeton, New Jersey
- Frank H, Frisch IT (1970) Analysis and design of survivable networks. IEEE Trans Commun COM-18:501–519
- 20. Haimes YY (2011) On the complex quantification of risk: systems-based perspective on terrorism. Risk Anal 31(8):1175–1186
- 21. Haimes YY, Steuer RE (2002) Research and practice in multiple criteria decision making. Lecture notes in economics and mathematical systems, vol 487. Springer, Berlin, 552 pp
- 22. Hausken K, Bier VM, Zhuang J (2009) Defending against terrorism, natural disaster, and all hazards. In: Bier VM, Azaiez MN (eds) Game theoretic risk analysis of security threats. International Series in Operations Research and Management Science. Springer, New York
- 23. Holyland A, Rausand M (1994) System reliability theory: models and statistical methods. Wiley, New York
- 24. Hong S (2013) Strategic network interdiction. Working paper, Korea Institute of Public Finance, Seoul
- Iida Y (1999) Basic concepts and future directions of road network reliability analysis. J Adv Transp 33(2):125–134
- 26. Janakiram M, Keskinocak P, Maku T, Xia S (2011) Supply chain game. OR/MS Today, pp 38-42
- Kalyoncu H, Sankur B (1992) Estimation of survivability of communication networks. Electron Lett 28(19):1790–1791
- 28. Krugler PE, Chang-Albitres CM, Pickett KW, Smith RE, Hicks IV, Feldman RM, Butenko S, Kang DH, Guikema SD (2007) Asset management literature review and potential: applications of simulation, optimization, and decision analysis techniques for right-of-way and

transportation planning and programming. Report FHWA/TX-07/0-5534-1, Texas Department of Transportation Research and Technology Implementation Office, Austin

- 29. Lai Y-S (1995) IMOST: interactive multiple objective system technique. J Oper Res Soc 46 (8):958–976
- Lam WHK, Xu G (1999) A traffic flow simulator for network reliability assessment. J Adv Transp 33(2):159–182
- Li Z, Sinha KC (2004) A methodology for multicriteria decision-making in highway asset management. Pre-Print #04-4756, Transportation Research Board 2004 Meeting, Washington, DC
- 32. Lou Y, Zhang L (2011) Defending transportation networks against random and targeted attacks. Transp Res Rec: J Transp Res Board 2234:31–40
- Lyle D, Chan Y, Head E (1999) Improving information-network performance: reliability versus invulnerability. IEEE Trans 31:909–919
- 34. McLaughlin BJ, Murrell SD, DesRoches S (2011) Anticipating climate change. Civil Engineering, April, pp 50–55
- 35. Midwest Regional University Transportation Center (2002) Synthesis of national efforts in transportation asset management. Project 01–01, Department of Civil and Environmental Engineering, University of Wisconsin, Madison, Wisconsin
- Natural Hazards Center (2009) A . NCHRP Project 24:20–59 (TCRP Project J-I0E, published as Natural Hazards Informer (No. 4, September), University of Colorado at Boulder. Boulder, Colorado)
- 37. Nojima N (1999) Performance-based prioritization for upgrading seismic reliability of transportation network. J Natl Disaster Sci 20(2):57-66
- Page LB, Perry JE (1989) Reliability of directed networks using the factoring theorem. IEEE Trans Reliab 38(5):556–562
- Perry JE, Page LB (1994) Reliability polynomials and link importance. IEEE Trans Reliab 43 (1):51–58
- 40. Pettigrew R (2010) Epistemic utility theory. In: Working Paper, Department of Philosophy, University of Bristol, UK
- 41. Sanchez-Silva M, Daniels M, Lleras G, Patino D (2004) A transport network reliability models for the efficient assignment of resources. Transp Res Part B 39:47–63
- 42. Schavland J, Chan Y, Raines R (2009) Information security: designing a stochastic-network for reliability and throughput. Naval Res Logist 56(7):625–641
- 43. Science Applications International Corporation (2009) Costing asset protection: an all-hazard guide for transportation agencies. NCHRP Project 17:20–59 (National Cooperative Highway Research Program Report 525, Transportation Research Board, Washington, D.C.)
- 44. Shier D (1991) Network reliability and algebraic structures. Oxford University Press, Oxford
- 45. Steuer R (1986) Multiple criteria optimization: theory, computation, and application. Wiley, Englewood Cliffs
- 46. Stirling WC (2002) Satisficing games: moving beyond individual rationality. In: Working Paper, Electrical and Computer Engineering Department, Brigham Young University, Provo, Utah
- Transportation Research Board (2010) Special issue on "Asset Management at Work". In: TR New 270
- 48. Wakabayashi H, Iida Y (1992) Upper and lower bounds of terminal reliability of road networks: an efficient method with Boolean algebra. J Natl Disaster Sci 14:29–44
- 49. Wakabayashi H (2008) Travel time reliability indices for highway users and administrators under uncertain environment. In: Paper ID 676, proceedings of the 10th international conference on application of advanced technologies in transportation, Athens, 28–30 May 2008, 19 pp
- 50. Washburn A, Wood K (1995) Two-person zero-sum game for network interdiction. Oper Res 43(2):243–251

- Zhang L, Levinson D (2008) Investing for reliability and security in transportation networks. Transp Res Rec: J Transp Res Board 2041:1–10 (Transportation Research Board of the National Academies, Washington, D.C., 2008)
- 52. Zhuang J, Bier VM (2007) Balancing terrorism and natural disasters—defensive strategy with endogenous attacker effort. Oper Res 55(5):976–991

Network Throughput and Reliability: Preventing Hazards and Attacks Through Gaming—Part 2: A Research Agenda

Yupo Chan

Abstract It is clear from the discussions in the previous chapter that, while there has been work done on improving network throughput and reliability, only a plethora of literature discusses a unified model to address both hazards and attacks. This is the motivation behind the previous chapter and the current one, where an attempt is made to bridge this gap. In spite of our initial modeling and computational experiments, no one has yet established the *formal* conditions under which a Pareto Nash equilibrium exists to prevent both hazards and attacks. Such an equilibrium includes the necessary *posturing* to discourage terrorist attacks, and to perform *preventive maintenance and upgrade* on our critical civil infrastructure ahead of a hazard. Meanwhile, general principles of how to best defend systems of specific types against intelligent attacks are emerging that can help system managers to allocate resources to best defend their systems. The research framework outlined in this chapter will gain insights into hazards and attack mitigation for a variety of infrastructure networks. In addition, there is much that can be done in the behavioral science perspective.

Keywords Network security • Multicriteria decision-making • Stochastic network • Shapley value • Information value theory • Posturing

1 Background

It is clear from the discussions in the previous chapter that, while there has been work done on network throughput and reliability, only a plethora of literature discusses a unified model to address both hazards and attacks. It turned out that the works of the authors are among the scanty few that have addressed a network design problem with a survivability emphasis, allowing critical physical components to be identified for upgrade, maintenance, and hardening in terms of arcs and nodes. In the

Y. Chan (🖂)

University of Arkansas at Little Rock, Little Rock, Arkansas, USA e-mail: ychan@alum.MIT.edu

[©] Springer International Publishing Switzerland 2015

K. Hausken and J. Zhuang (eds.), *Game Theoretic Analysis of Congestion, Safety and Security*, Springer Series in Reliability Engineering, DOI 10.1007/978-3-319-13009-5_6

previous chapter, the authors have summarized their research findings on the subject of network throughput and reliability, emphasizing strategies to prevent hazards and attacks through gaming. By doing so, they have identified some "loose ends," some of which they are working on. To respond to the open questions raised in the previous chapter, we wish to propose a research agenda on the subject. To accomplish this goal, we first summarize the closely related literature in Table 1. Here, the models in Table 1 are introduced using a common unifying taxonomy. This would in addition to clarify the target problems in question—facilitate a comparison between the different approaches to solve similar problems.

We will briefly summarize the major findings to date, which will put our recommended research plan in perspective. In chronological order, five sets of authors are identified in Table 1: [22, 35, 42, 43], [50]. We are disappointed that there are few other related works, if any, which address both hazards and attacks. Nevertheless, we judge these five selections would serve our goal the best, as all are focused on a similar subject, albeit with differences. The contributions of each will be discussed under "Stochasticity or risk assumption," "Arc/node or system component attributes," "Metrics for network or system performance," and "Network or system equilibrium and other properties." Notice that we are particularly interested in stochastic *network* representation of a physical infrastructure, under the premise that any policy-analytic findings will have to be translated back into physical implementation on facilities that are represented as arcs and nodes.

Once again, a unique feature of the previous chapter and the current chapter is the integration of hazard-and-attack preventions in a proactive manner, rather than in a reactionary manner. While hazard prevention has been a historical subject of interest, the tragic event of 9/11 has awakened our sensitivity toward malicious and asymmetric attacks. In our opinion, Lyle et al. [35] represents a pioneering piece of work on the subject, way ahead of the 9/11 incident. Their most significant contribution is that it generalizes a classic *network interdiction* problem, defined as the monitoring or halting of an adversary's activity on a network. Instead of following a leader-follower paradigm, a proactive stance is adopted to identify a Pareto Nash equilibrium, which will discourage an adversary action altogether. Another contribution is the utility-theoretic approach, which gives *value* to efforts that make a network component invulnerable to exploitation. Unlike interdiction, a component is hardened against exploitation, instead of being removed from exploitation. In lieu of a cost accounting approach, the worth of the resulting network security is *imputed* from unit costs of network improvement to guard against hazards.

In their work published in 2007, Zhuang and Bier responded to the 9/11 event by examining an appropriate investment level for homeland security. Their findings are interesting, in that when increased defensive investment causes the attacker (henceforth referred to as the Red player) to redouble his efforts, defensive investment against terrorism will not decrease the probability of a successful attack as much as the defender (henceforth referred to as the Blue player) might have expected. This will in general tend to reduce the effectiveness of protecting a large number of targets against intentional attacks, and therefore increase the relative desirability of protection from natural disaster (and of both hazards and attacks).

Table I Com	parison of alternate approaches			
Model	Stochasticity or risk assumption	Arc/node or system component attributes	Metrics for network or system performance	Network or system equilibrium and other properties
[35] [35]	For both natural/technological hazards <i>and</i> deliberate attacks, nodes and arcs (representing physical facilities) that fail in any combination throughout the network are specifically analyzed. This allows the most critical components to be iden- tified for hazard-and-attack prevention	Without loss of generality, a facility can be represented by an arc. An arc is characterized by its availability, while congestion effect is simply modeled by the arc capacity. As such, traffic congestion is not modeled explicitly	As a departure from many other literatures, a utility-theoretic approach is taken. Reliability damage utility is defined for an O-D pair, while throughput damage reliability is not mod- eled. The utility function now represents the "value/payoffs" in a game matrix	It was found that the nonlinear model can be approximated by a linear model. Nash equilib- rium is established computa- tionally. A deliberate attempt is made to quantify the damage of an attack. In lieu of a cost accounting approach, the worth of network security is imputed from unit costs of network improvement to guard against hazards
Zhuang and Bier [50]	Balancing protection from ter- rorism and natural disasters, policy related insights are gained from the model. While risk is explicitly modeled, there is no physical network topol- ogy other than the number of attack targets. The results of this paper are intended to pro- vide mainly qualitative insights. There is a disclaimer on using the model "in a numeric manner in support of specific decisions"	Protection from terrorism will tend to become less cost-effec- tive for the Blue player as the number of targets grows, due to the ability of the Red player to redirect his attack effort to less defended targets. The model describes the Red player choice by a continuous level of effort rather than a discrete choice (i.e., attack or not). There is a "cost function" for multiple targets versus a single target	In the single-target case, increased defensive investment can lead a Red player to either increase his level of effort or decrease his level of effort. When increased defensive investment causes the Red player to redouble his efforts, defensive investment against terrorism will not decrease the probability of a successful attack. This points toward a Blue player's optimal alloca- tion of resources between pro- tection from intentional attacks and from natural disasters	Insights were gained in the nature of equilibrium defensive strategies, emphasizing the importance of intelligence in counterterrorism—to anticipate not only the Red player's choice of targets, but also the likely Red player responses to defensive investments. This means either a reduction or increase in the effectiveness of protection from intentional attack, and can therefore affect the relative desirability of investing in protection from natural disasters
		-		(continued)

Table 1 Comparison of alternate approaches

Table 1 (cont	inued)			
Model	Stochasticity or risk assumption	Arc/node or system component attributes	Metrics for network or system performance	Network or system equilibrium and other properties
Smith and Lim [43]	In the network interdiction and fortification problem, a Red player is to compromise certain network elements before the Blue player's action inopera- ble. Stochasticity is treated as part of the general problem, including stochastic shortest path instances	The typical node/arc attributes are used, such as cost, capacity, and even availability. A node/ arc is interdicted by either increasing its cost or compro- mising its capacity. The only conceptual equivalence of "hardening" is to making a component inoperable, in order to deny Red player's use. There seems to be no consideration to make a component totally invulnerable	Most conventional perfor- mances apply, such as profit maximization, least cost, max flow, and reliability. However, only one objective is consid- ered at a time during each stage of the Stackelberg game. It appears that multiple criteria have yet to be considered in this publication or other inter- diction literature	The models describe a Stac- kelberg competition rather than a Stackelberg equilibrium. For an equilibrium to occur, the leader must know ex ante that the follower observes his action. The follower must have no means of committing to a future non-Stackelberg fol- lower action and the leader must know this
Schavland et al. [42]	Continuing the pioneering effort by Lyle et al., a network component (as represented by arcs) can fail in any combina- tion throughout the network for both natural/technological haz- ards <i>and</i> deliberate attacks	As assumed by Lyle et al., an arc is characterized by its availability, while congestion effect is modeled by the arc capacity. This assumption allows important insights to be gained on the relative cost- effectiveness of capacity versus availability improvement	A complete utility-theoretic model is constructed. In addi- tion to reliability damage util- ity, flow damage utility is defined for each O-D pair. Trades are then performed between reliability- and flow- damage utilities. How one trades between throughput and reliability is modeled ranging from a totally compensatory and noncompensatory prefer- ence structure	In addition to the Lyle et al. findings, there exists a unique range-equalized robust equilib- rium to the linearized model. Pareto Nash equilibrium and security-worth imputation can be obtained irrespective of the relative unit costs of availabil- ity versus capacity improve- ment, or the game players' preference structure

(continued)

Table 1 (continued)

,				
Model	Stochasticity or risk assumption	Arc/node or system component attributes	Metrics for network or system performance	Network or system equilibrium and other properties
Hausken et al. [22]	In the presence of risk, the Blue player chooses tradeoffs between investments in protec- tion against natural disaster only, protection against terror- ism only, and hazard-and- attack protection. A policy- analytic model is proposed, rather than a model on a phys- ical network	It is assumed that costs increase linearly with the investment level. A game player advan- taged with a sufficiently low normalized unit cost of invest- ment relative to that of its opponent prefers to move first, which deters the opponent entirely, causing maximum utility for the first mover and zero utility to the deterred sec- ond mover, who prefers to avoid this game	A quantitative value is assigned to a target asset, including the assessed value of its destruc- tion. When hazard-and-attack protection is sufficiently cheap, it jointly protects against both natural disaster and terrorism. As the cost increases, either pure natural disaster protection or pure terrorism protection joins in, dependent on which is more cost-effective	When protecting targets that have relatively low value to potential terrorists. Blue play- ers will more often wish to invest in protection from natu- ral disasters and/or hazard-and- attack protection, and less often wish to invest in protection against terrorism alone. Like- wise, protections that are effective against only a rela- tively few terrorist threats may be less desirable than protect- ing against a variety of threats

Smith and Lim [43] continued the line of work on network interdiction by a rather comprehensive synthesis book-chapter. First, they provide a literature review of interdiction research and the theoretical/methodological fundamentals behind solving interdiction and fortification problems. Then they discuss advances in interdiction algorithms and then provide an analysis of fortification optimization techniques. Finally, they conclude the chapter with a discussion of future research directions. It appears that multiple criteria have yet to be considered in this publication or other interdiction literature. The models describe a Stackelberg competition rather than a Stackelberg equilibrium. Both a two-stage and three-stage model are discussed, depending on whether Red takes the first move, followed by the reaction of Blue, or that Blue takes the first step in prevention, followed by Red's reaction, and in turn Blue's follow-up reaction.

Rather than following Zhuang and Bier's policy-oriented approach, Schavland et al. [42] chose to continue the work of Lyle et al. by completing the utility-theoretic approach on a physical network. Their work considers tradeoffs between arc availability versus capacity improvements in improving a network component —a real-world consideration for infrastructure network operators and managers. Most critically, a *Pareto Nash equilibrium* is identified irrespective of the preference structure of the defender (Blue player) or the attacker (Red player). As a result, a stable strategy against exploitation is clearly identified.

Hausken et al. [22] follow the line of thought of Zhuang and Bier. Again, they determine the resource allocation between preventing hazards versus attacks in terms of a *value metric* assigned to the asset being protected. Similar to the Zhuang and Bier article, no physical network is modeled. Rather, a policy-analytic model traces a dynamic game between the Red and Blue players in terms of "moves" similar to a chess game. Using a linear cost function, a player's move is determined by the assessed value of that particular move. It is interesting to note that it also uses a utility-theoretic approach to assess the value of an investment decision.

Instead of using a relative scale, Lyle et al. and Schavland et al. went one step further to quantify the "absolute" monetary value of hardening a network against attack by *imputing* from the model the cost of disrupting economic and noneconomic transactions when a network facility is destroyed. This is a marked departure from the conventional cost-counting procedure in which a monetary value is assigned based on the replacement cost of a destroyed facility. The monetary value, as derived based on this approach, is more defensible—in our opinion—in requesting a specific budget for attack mitigation in the public arena.

While all the five articles listed in Table 1 address both hazards and attacks, it is obvious that we are partial toward a physical network modeling approach rather than a policy-analytic approach. At the same time, we wish to learn from the policy-analytic experience. Our recent efforts since 2009 have been focusing on this approach, including the work of Van Hove et al. [47], Del Vecchio et al. [16], Chan and McCarthy [10, 11]. The Van Hove et al. work continues the network-centric tradition, while the Del Vecchio et al. work is more policy-oriented. While the network-centric bias allows us to identify critical network components to improve or to harden, it comes with a significant computational cost. To date, we are happy

to report that we have linearized a highly nonlinear reliability damage utility function V_r [42]. For a highly reliable network, which is the common case instead of the exception, the set of nonlinear constraints can also be linearized. Together with a proposed decomposed solution strategy to separate the hazard mitigation part of the model from the attack prevention part, significant progress has been made in the computational front.

Our venture into the policy science arena opened a new vista into the behavioral science world. Chan and McCarthy [10] discusses the multiparty dynamics of knowledge dissemination, acceptance and its translation into implementation. To understand the chemistry among participants, they review game theory and information value theory, both of which determine the effectiveness of collaborative research, considering the coexistence of cooperation and competition among participating enterprises. In Chan and McCarthy [11], an iterative experimental tool, consisting of a computational model component and a face-to-face gaming component, is proposed to empirically test the likely success of collaborative research, dissemination, and implementation in today's information-based economy.

Obviously, the best model is one that is comprehensive, rigorous, and computationally efficient. Accordingly, we propose five related research extensions to our findings to date. First, the model is extended to include a time dimension, considering the significant amount of time required to carry out infrastructure repairs or upgrades logistically speaking. This will also facilitate a dynamic game being played out, following a leader-follower sequence as implemented by Smith and Lim [43] and Hausken et al. [22]. Second, the single O-D model is extended to multiple O-D pairs for a more comprehensive network-wide monitoring strategy, extending the Smith and Lim [43] work from an interdiction application to the combined hazard-attack paradigm as defined in our research [34]. Third, a set of computational algorithms is outlined to solve the resulting models and beyond. Smith and Lim [43] observed that many interdiction studies have traditionally been limited to those problems for which linearization constraints can be applied to eliminate troublesome nonlinear terms. We echo this observation as outlined in Proposition 1 of the previous chapter. Smith and Lim further suggested examining bilinear programming and global optimization theory to obtain the most effective algorithms for this class of problems. This is corroborated by our proposed decomposed solution strategy by separating the hazard mitigation part of the model from the attack prevention part, as discussed in the previous Part-1 chapter. Fourth, the current twoperson zero-sum symmetric game is to be extended to an *n*-person nonzero-sum asymmetric cooperative and noncooperative game, which is more reflective of today's attacks by faceless terrorists, which often requires a coalition approach to defend. In this regard, we were motivated by the Del Vecchio et al. [16] work performed for the U.S. Department of Defense. Fifth, motivated by the work of Hausken et al. [22], important player behaviors are considered during gaming, which could affect the outcome of a game. We were also encouraged by our recent work in this area, as documented in Chan and McCarthy [10, 11]. Each research extension is discussed in detail below as a separate task.

2 Future Extensions

Having navigated through the "learning curve," we are confident that the subject is germane to serious investigation by interested parties, as mentioned. Obviously, there are several important questions that remain to be answered on this highly complex subject before significant progress can be made. For the general research community, we wish to organize these "loose ends" scientifically into five research extensions.

2.1 Extension 1: Comprehensive Real-Time Performance Monitoring

Ideally, one wishes to monitor the performance of a stochastic network continuously over time, for both peak and off-peak periods. This requires gathering performance statistics on a real-time basis, often by sensors installed locally or remotely. While real-time monitoring may be performed in communication networks or power grids on a routine basis, this is often not done on a majority of highways. Realistically, it is simply not possible to sample every highway network node and arc comprehensively. In the field, only selected arcs or nodes are monitored, and one commonly keeps only a sample log of traffic delay or arc failures over some monitoring time period. Such a record is destroyed after a certain time [36]. The challenge lies in *where* to sample and *what* statistics to keep.

To achieve this monitoring function, network representation is not a trivial exercise, particularly for infrastructure networks. As mentioned, infrastructure networks require a lengthy construction/reparation period. This means maintenance or upgrade of an arc would impair network performance for a long time, often introducing additional traffic congestion that again needs monitoring. Accordingly, let us consider a time-expanded replica of a sample network for one O-D pair. The network in Fig. 1 shows a monitoring strategy of six time increments, where the arcs are directed arcs going from left to right and from top to bottom. For ease of presentation, the arrows on these directed arcs are not shown. Each arc is characterized by three attributes: cost, capacity, and availability. Cost refers to the impedance in traversing the arc, representing either traversal time, monetary expenditure, or physical impedance. Capacity is the maximum throughput that is possible on the arc, where in the case of a communication network, it represents the bandwidth. Availability is, once again, the percentage time the arc is functioning.

Here the sample network, consisting of nodes O, a, b, c, d, e, and D, is replicated six times for the six time periods, with the unused arcs and nodes removed—i.e., the arcs that carry zero flow for a particular time increment are deleted. Should one wish to assess the throughput of such a network, it can be shown that there exists a temporally repeated flow that is maximal [19]. In other words, a maximal throughput can be obtained over time, similar to the static throughput discussed



Fig. 1 Reduced time-expanded, stochastic network (Adapted from [47])

earlier. Since the computational complexity is huge for such time-expanded stochastic networks, we are again comforted to know that reliability and throughput *bounds* can be estimated in lieu of obtaining exact values [18, 40, 47]. In our research, we also wish to estimate average delay and misrouting, particularly during peak hours. (Here, misrouting means "making the wrong turns" in highway networks or "dropped calls" in communication networks.)

There can be some very useful findings from this time-expanded network flow model (Fig. 1). First, *dynamic performance* measures can be obtained from such a model to accomplish real-time monitoring. One can also adjust the performance models to fit the monitoring strategy. For example, it is easy to adjust the time increment in the model to reflect the sensor sampling rate used in the field. Thus, more frequent sampling means a shorter time increment, as in monitoring hourly traffic to avoid congestion. On the other hand, a piezoelectric sensor needs to be monitored only over months to detect pavement deterioration in structural health monitoring, if not over years [28].

Figure 1 represents a sample time-expanded network. Each row of nodes of this network represents all the original network nodes at a given time increment. Network flow is directed from left to right, starting from a super (darkened) origin node and ending in a super (darkened) destination node. Arcs are directed downward one level for each time increment. The solid lines represent physical arcs of a network. For example, a unit of flow starts at an origin node O shown in the second column at time increment t - 4 and it ends up at the destination node D in time increment t - 3 by tracing through a subset of the physical arcs (O, a), (O, b) (O, c), (a, d), (b, e), (c, D), or (e, D). The bold stapled-dotted itinerary shows the longest traversal path from origin O in t - 6 to destination D in t - 0, spanning all six time increments.

The stapled arcs are simply artificial arcs linking the physical network to the super source O at the far left. Likewise, they link the physical network to the super sink at the far right. All together, this defines the seven rows of replicated nodes, spanning six time increments t - 6, t - 5, ..., t - 0. Many of the arcs in the time-expanded network—such as (O, a) at t - 0 and at t - 1—constitute "dead end" arcs. They are not part of any alternate paths between the super source and super sink. Being "dead-end" arcs, any flow on these arcs at the particular time increment does not contribute to the total source-to-sink flow. In other words, they do not go anywhere. For this reason, we remove these arcs from the network representation to simplify the model, as mentioned.

In monitoring the performance during peak and off-peak hours, one can expect increasing congestion as throughput increases during peak hours. The question is whether one can establish an *analytical* relationship between expected (or mean) travel time and traffic throughput in a *stochastic* network. This *analytical* result is worth exploring to trace the congestion effects accurately. The same can be said for the other performance metrics. In our time-expanded stochastic network flow model, it was found that misroutings (or the number and percentage of dropped calls in a communication network) can be discerned, measuring service degradation [47]. For future extensions, Losada et al. [32] provided a procedure to speed up recovery time following a potential disruption. Here, protection is not necessarily assumed to prevent facility failure altogether, but more aligned with cost-efficient levels of investments in protection resources. Hsieh and Lin [24] proposed a method to *update* resource allocation in an unreliable network when resource, demand, or the characteristic of the flow network changes, as during reparation or upgrade. Dai and Lin [15] proposed a family of service policies for dynamically allocating service capacities in a stochastic network for its day-to-day operation.

2.2 Extension 2: Multi-O-D Minimum-Cost Network Flow Model

With the exception of some heuristics, the state-of-the-art in modeling *stochastic* networks (as defined in this chapter) has been limited to one single O-D pair (i.e., limited to the two-terminal *s*-to-*t* case). Preliminary work to date suggests, however, that our models can be formally extended to *multi-O-D* flow. Although obtaining an exact value for the expected (average) network flow may be computationally infeasible, our preliminary work to date on multi-O-D models suggests that it is possible—once again—to obtain statistical lower and upper bounds on the expected flow through Jensen's inequality [38].

The lower bound (LB) on the expected flow represents the best case from the user perspective where total cost is less than the expected cost, or congestion is less than anticipated. In this LB sub-model, the objective function represents the total cost of operating the system and is a function of arc-path decision variables f_p^k .

The first set of constraints requires all O-D pair demands to be satisfied. The second set constrains maximum arc utilization to the expected capacity of the arc. The final constraint set simply establishes nonnegativity requirements for the decision variables.

$$\min Z_{\text{LB}} = \sum_{k} \sum_{p \in P^{k}} c_{p}^{k} f_{p}^{k}$$

subject to
$$\sum_{p \in P^{k}} f_{p}^{k} = b^{k} \quad \text{for all } k \in K$$

$$\sum_{K} \sum_{p \in P^{k}} b_{ij}^{kp} f_{p}^{k} \leq r_{ij} u_{ij} \quad \text{for all } (i,j) \in A$$

$$f_{p}^{k} \geq 0 \quad \text{for all } k \in K \quad p \in P^{k}$$

$$(1)$$

where

 $\begin{array}{ll} K & = \text{set of user O-D pairs (or commodities), with } k \in K \\ c_p^k & = \text{cost of the } k \text{th commodity flow on path } p \in P^k \text{, where } P^k \text{ is the set of paths for commodity } k, and the path cost equals to the sum of arc costs on the path } \\ f_p^k & = \text{amount of the } k \text{th commodity flow on path } p \in P^k \\ h_{ij}^{kp} & = 1 \text{ if arc } (i, j) \text{ lies on path } p \in P^k; 0 \text{ otherwise.} \end{array}$

Here the expected arc capacity, $E(u_{ij})$, is equal to $r_{ij} u_{ij}$.

The upper bound (UB) on the expected flow represents a worse case scenario from the user's perspective, where total cost is greater than expected. The first set of constraints in the following submodel requires all O-D pair demands b^k to be satisfied. However, each path variable f_p^k has an associated loss parameter R_p^k if at least one arc in the path is not totally available. This potential flow loss may lead to an influx of slack external flow at some or all O-D pair origin nodes, representing queuing. Consequently, system cost increases as slack external flow increases. The second set constraint set establishes nonnegativity requirements for the decision variables.

$$\min Z_{\text{UB}} = \sum_{k} \sum_{p \in P^{k}} c_{p}^{k} f_{p}^{k}$$

subject to
$$\sum_{p \in P^{k}} R_{p}^{k} f_{p}^{k} = b^{k} \text{ for all } k \in K$$

$$\sum_{K} \sum_{p \in P^{k}} b_{ij}^{kp} f_{p}^{k} \le u_{ij} \text{ for all } (i,j) \in A$$

$$f_{p}^{k} \ge 0 \text{ for all } k \in K \quad p \in P^{k}$$

$$(2)$$

Figure 2 represents the peak and off-peak traffic flow from a network flow model with three O-D pairs, of which the flow from node 1 to node 9 is shown. The lower bound (LB) on expected flow represents the off-peak scenario, when the total travel cost (in time units) is less than normal–or when congestion is better than



Fig. 2 Off-peak and peak hour flows

anticipated. The upper bound (UB) on expected flow represents a peak hour scenario, where the total cost is greater than normal—or congestion is worse than anticipated.

Most importantly, system cost (or network-wide congestion) increases can be traceable to local queue length increases in this model, which are to be explicitly modeled in a corresponding level of detail in a separate, more microscopic analysis. Accounting for queuing (at an aggregate level) and with tampering involved, Fig. 2 illustrates a "compromised" scenario, showing the flow pattern under the external threat. In the Figure, we show the flow path for the O-D pair 1–9 during off-peak hours (LB) and during peak hours (UB). For example, the compromise path is 1-2-5-6-9 during off-peak, and 1-4-5-6-9 during the peak period.

The numbers on each arc correspond to the off-peak and peak traffic volume, respectively. For this example, it can be verified that the off-peak LB throughput is smaller than the peak UB throughput. Thus, the total flow arriving at destination node 9, representing the throughput for O-D pair 1–9, is 8.82 + 5.18 = 14 units during off-peak and 9 + 6.654 = 15.654 during peak hours. Statistically speaking, however, this does not mean that *each* off-peak LB arc flow is smaller than the corresponding peak UB arc flow (when one compares, say, arc (1, 4) in the accompanying networks where it shows 4.75 for off-peak and 3.284 for peak). Due to the stochasticity and the presence of a compromise, the exact flow patterns during peak and off-peak periods are of particular interest to both the network users and the traffic monitor, and hence the raison d'être for this analysis.

Here are several possible extensions. Murray-Tuite and Fei [37] uses probabilities of target–attack method combinations that are degree-of-belief based and updated using Bayes' Theorem after evidence of the attack is obtained. The average capacity reduction for a particular target–attack method combination serves as an input to the traffic assignment—simulation package DYNASMART-P—to determine travel-time effects. Cappanera and Scaparra [7] identified the set of components to harden so as to minimize the length of the shortest path between a supply node and a demand node after a worst-case disruption of some unprotected components. An important extension to this multilevel model involves the use of more complex network-flow models in the lower-level user problem such as multiple origin-destination flow problems. In the same vein, Scaparra and Church [41] proposed a bi-level formulation. The top level problem decides on which facilities to fortify to minimize the worst-case efficiency reduction due to the loss of unprotected facilities. Worst-case scenario losses are to be avoided in the lowerlevel interdiction problem. Croxton et al. [14] presented several mixed-integerprogram (MIP) models, based upon variable disaggregation, for generic multi-O-D network flow with *nonconvex* piecewise linear costs. The challenge is to delineate the exact analytical properties of such a multi-O-D network before and after a compromise, which is the focus of Extension 2. It should be noted that electricity networks deal with reinforcing the entire network by assuming that only 1 or 2 failures can occur at once.

2.3 Extension 3: Solution Algorithms for Multi-O-D Flow

Due to the combinatorial nature of the failure states of a stochastic network, we have mentioned more than once that—for practical reasons—only bounds can be estimated for network performance, rather than a more precise measure. For mincost flow, bounds for the expected (average) minimum travel time (or cost), *Z*, of a "Multi-O-D Minimum-Cost Network-Flow" model have been identified below in Properties 1 through 4 [38]. They are captured in Eqs. (3) and (4). These properties may very well form the basis for efficient solution algorithms to be designed for this class of multi-O-D stochastic networks. Instead of diving into the details, let us simply summarize the result of our investigations to date. In general, the cost of the totally available (or perfectly reliable) network flow *Z*^{*} provides the absolute lower bound, or the best-possible least-congested flow. At the same time, the peak hour network flow *Z*^{*}_{UB} provides the worst possible upper bound for the traffic flow cost. These bounds envelop the off-peak network flow cost *Z*^{*}_{LB} and expected network flow cost *Z*_{EV} in between. Mathematically, the following inequality can be written:

$$Z^* \le Z^*_{LB} \le Z_{EV} \le Z^*_{UB}, \tag{3}$$

which shows that off-peak flow cost is always less than or equal to the expected value. As one improves the reliability, the network approaches the fully reliable configuration, and the min-cost bounds Z_{LB}^* , Z_{EV} , and Z_{UB}^* converge to the totally available min-cost Z^* . For an increasingly reliable network, the gap between the absolute lower bound Z^* and absolute upper bound Z_{UB}^* is likely to be tight. For the best case scenario, where availability is 100 % for all arcs, it is "snug tight," to the extent that the inequalities become equalities:

$$Z^* = Z^*_{LB} = Z_{EV} = Z^*_{UB}.$$
 (4)

Under this scenario, the travel times (costs) are invariant between peak and offpeak, otherwise the equalities in (4) could not hold.

These bounds will likely facilitate the approximate solutions, providing efficient solution algorithms for practical applications where precise solutions may not be necessary. Our experiments to date suggest that the lower bound of the expected maximum flow has been found to be much tighter than the upper bound [9]. It will be interesting to find out whether this result carries over to the expected minimum-cost Z. To the extent that "real traffic flow cost" is likely to be similar to the expected minimum flow cost $Z_{\rm EV}$, we are in fact speculating whether the gap between $Z_{\rm LB}^*$ and $Z_{\rm EV}$ or the gap between $Z_{\rm EV}$ and $Z_{\rm W}^*$ is likely to be smaller. If so, $Z_{\rm LB}^*$ or $Z_{\rm UB}^*$ can be used to estimate "real traffic flow cost." We would like to see how this conjecture can be confirmed, and if so, be explained theoretically. Jensen inequality, which delineates the inequality between $E(f(\mathbf{x}))$ and $f(E(\mathbf{x}))$, may be a good place to start.

Although approximation and bounding methods are available, their accuracy and scope are very much dependent on the properties (such as size and topology) of the network. Hui et al. [29] studied how the Cross-Entropy method can be used to estimate network reliability more efficiently, thereby refining on lower and upper bound methods. A new method that generates the sequential lower and upper bounds of all terminal reliability (ATR) based on greedy network factoring is proposed by Won and Karray [49]. They showed that their ATR-bound update method could find an acceptable ATR bound of a given network much faster than the exact method. Thus one might use the algorithm to find the *approximate* ATR bounds by terminating the algorithm when the lower and upper bounds are within a certain error threshold.

Aside from offering an algorithm, Szeto [46] and Szeto et al. [45] warned against the existence of Braess' Paradox [5] in transportation networks, in which the total travel cost or reliability could be worsened after network improvement, whether it be reliability improvement or others. This paradox arises mainly due to road users' route choice behavior; it does not apply to telecommunication or power networks where behavioral routing factors are absent. This constitutes a possible exception to the bounds of Eq. (3) used for approximating performance, and needs to be investigated in detail for transportation networks.

2.4 Extension 4: Multicriteria n-Person Cooperative Model with Hardening

The discussions so far concentrate on a *noncooperative* zero-sum game between Blue and Red, in which Blue faces up to Red squarely in the "battle field." Departing from this symmetric game, we would like to examine today's *asymmetric* game in which there are multiple Blue players teaming together as a coalition against a faceless Red player [16]. For example, local, state, and federal highway departments may be working together to thwart a sneak attack. In another example, the telecom network service providers (NSP) may be teaming up against terrorism. Gaming theory also includes *n*-person cooperative games, where the payoff of each player is based on the Shapley Value (SV). SV is the marginal contribution of a player to the coalition [48]. Through the SV, cooperative game theory suggests that player *i*'s reward should be the expected amount that player *i* adds to the *n*-player coalition [8]. In other words, the player receives its expected marginal contribution.

Let us say there are *multiple* telecom carriers participating in a subsidy-incentivized coalition against a Red player. These telecom carriers can receive subsidies (in addition to their regular revenue) as incentives to improve their collective security. This can take the form of subsidized rates, where government funds are available to encourage the cooperation between the carriers against terrorism. Thus the carriers can use the subsidy to improve security, without passing the cost to customers. This results in a more robust telecom network and a stronger united coalition of carriers. Most importantly, a coalition is formed solely based on revenue incentives, both commercially earned and subsidized.

This strategy can be modeled as an *n*-person, zero-sum, *cooperative* game [1, 46]. A multicriteria optimization model was developed by Del Vecchio et al. [16] to establish a strategic competition between an international NSP-coalition and any adversarial Red player. Del Vecchio et al.'s results, schematically shown in Fig. 3, show the tradeoffs between the coalition's total revenue f_1 (or the "common good") and revenue to a particular carrier f_2 (or the participant's "individual welfare"), forming a bicriteria optimization problem. In this illustration, we have normalized the range of both the flow damage utility and reliability damage utility between 0 and 1—a process called range equalization as mentioned in our last chapter. Thus the horizontal axis is labeled V_f and the vertical axis is V_r . Let us define the *ideal point* (1, 1) as the "best of all possible worlds" for flow damage utility and reliability damage utility. It is the goal that one wishes to achieve, but



hardly achievable in practice—resulting in a "shortfall." Suppose the "shortfall" from the ideal (1, 1) to any viable improvement strategy on the efficient frontier, or the "regret," is minimized.

The following is the formal representation of l_p -metric (p = 1, 2 and ∞) under reverse filtering [44].

$$l_p = \left\| \mathbf{v} - \mathbf{v}^i \right\|_p^{\pi} = \left[\sum_{k=1^q} \left(\pi_k \left| \nu_k - \nu_k^i \right| \right)^p \right]^{1/p}$$
(5)

where:

 v_k is the *k*th attribute, denoting flow-, or reliability- damage utility in our study; **v** is the vector of attributes, consisting of the two normalized entries V_c and V_r ; **v**^{*i*} is the ideal, or the point $(V_{c^*}^i, V_r^i) = (1, 1)$ in our case;

q is the number of criteria (or dimension of the efficient points being filtered), which is two in our study;

 π_k is the range equalization weight on criterion k to convert V_f to V_c , the normalized 0-1 ranged utility; and

p is the metric parameter (notice that, without ambiguity, p is the common usage on this subject; obviously, it is different from the previous usage of p for path designation).

The l_1 -metric (or the "diamond contour") corresponds to a *totally compensatory* tradeoff. Here, a shortfall in flow damage utility can be made up by a gain in reliability damage utility. The l_{∞} -metric (or the "square" contour) corresponds to a *totally noncompensatory* tradeoff. Under this scenario, exclusive attention is drawn to the shortfall, with no regard for the other metric. An intermediate metric, l_2 -metric (or the "circle" contour), is the "middle ground" between the two preference structures. While we chose to show the unit "balls" in Fig. 3, the only part of the contours of interest is the part at the "southwest," where the "contours" touch the efficient frontier.

Similar to the symmetric game, the revenue incentives—which result in the corresponding hardening strategies—induce a stable equilibrium, thereby discouraging adversarial attacks. To be differentiated from a symmetric game, however, there are three separate equilibriums in Fig. 3, should one examine where the diamond, circle, and rectangle contours touch the efficient frontier. Instead of a robust Nash equilibrium, three separate equilibriums result. This can be gleaned from Fig. 3, where two equilibriums are distinctly shown. The l_1 -metric solution is shown as a red "dot." The l_{∞} -metric solution is shown as a blue dot. The l_2 -metric solution lies somewhere in between on the efficient frontier and is not shown for clarity.

The important point is that there are three distinct dots, representing three separate equilibriums. Unlike the noncooperative symmetric game, it can be seen that the preference structure—as represented by the l_1 , l_2 , and l_{∞} -metric—does matter in this case, resulting in its own pair of performance: flow damage and reliability damage.

Two important research questions remain: (A) Under what conditions can a robust Pareto Nash equilibrium be obtained, if such an equilibrium exists. (B) How can the model be solved efficiently? Garcia et al. [20] proposed a decentralized

solution method for general network optimization, facilitated through fictitious game playing. In their algorithm, a network simulation is conducted at each iteration, converging toward the optimum (an equilibrium) when certain properties are discerned in the model. Due to its generality, their fictitious game-theoretic approach has profound implications on general mathematical-programming solution-algorithms. However, Garcia et al. modeled only a *game* between imaginary players with identical interests. Our cooperative game here involves forming a cooperative coalition against a faceless adversary. The game is therefore no longer among players with identical interests. The challenge is how to exploit the Garcia et al. paradigm in solving our cooperative model, in order to gain general computational insights that will go well beyond the infrastructure applications in the current discussion. Incorporating fictitious play is a good idea. However, researchers now strive to find a best response to the empirical distribution of the opponent's past play and do not usually consider some imaginary player with identical interest.

2.5 Extension 5: Shared Cognition

While we are following the casual modeling approach in most of our extensions, we are mindful of the "soft" factors that are not included in casual modeling. In addition to the list of soft factors already discussed, a related concept that we need to discuss is shared cognition, and what does it buy us in terms of team and organizational performance [6]? The concept of shared cognition can help us to explain what separates effective from ineffective teams by suggesting that in effective teams, members have similar or compatible knowledge, and that they use this knowledge to guide their (coordinated) behavior. Essentially, shared cognition could serve as an indicator of a team's "readiness" or "preparedness" to take on a particular task. In a more practical sense, shared cognition research can help establish an understanding of the elements of effective teamwork, which can in turn lead to better interventions for improving team performance. From this perspective, several fundamental questions regarding the nature of shared cognition emerge. These fall into four broad (and related) categories: (A) What is shared? (B) What does "share" mean? (C) How should "share" be measured? and (D) What outcomes do we expect shared cognition to effect? To answer these questions, we propose to design a set of games to be played by the stakeholder participants in a two-step iterative validation process, involving playing a game and measuring the outcome afterwards.

Meanwhile, Kovenock et al. [31] investigate individual behavior in a game of attack and defense of a weakest-arc network. Their experimental investigation unveiled interesting game properties that are not predicted by theoretical constructs such as Nash equilibrium. The authors claim that their results offer a more plausible explanation of some observed patterns of terrorist attacks. For example, it provides evidence that infrequent "periods of high terrorism" may simply be the result of asymmetric objectives and strategic interactions between the Red players and Blue players within a weakest-arc type of contest. Samuelson [39] advocated

evolutionary game theory. Unlike traditional game-theory models, which assume that all players are fully rational and have complete knowledge of details of the game, evolutionary models assume that people choose their strategies through a trial-and-error learning process in which they gradually discover that some strategies work better than others. In games that are repeated many times, low-payoff strategies tend to be weeded out, and an equilibrium may emerge.

Throughout this chapter, we discussed investments in protection against natural disaster only, protection against terrorism only, and hazard-and-attack protection. According to Hausken et al. [22], a game participant advantaged with a sufficiently low normalized unit cost of investment relative to that of its opponent prefers to move first, which deters the opponent entirely, causing maximum utility for the first mover and zero utility to the deterred second mover, who prefers to avoid this game. When hazard-and-attack protection is sufficiently cheap, it jointly protects against both the natural disaster and terrorism. As the cost increases, either pure natural disaster protection or pure terrorism protection joins in, depending on which is more cost-effective.

Two types of social behavior, cooperative and noncooperative, are considered. Initially, all players will act independently of each other, with the relatively simple goal of maximizing their own outcomes. This 'selfish' style of behavior could occur through mutual agreement among all players at the commencement of the game, or may simply occur as a result of players being unable to communicate their intentions to make binding commitments with each other. This latter point may be a result of the limited extent to which information is shared between players, even if initial commitments have been made.

Cooperative games allow players to communicate and make binding commitments with each other. The purpose of such player 'agreements' is to improve benefits otherwise obtainable through individual-based, noncooperative, games. These benefits should be clearly identifiable to the players, either prior to the start of the game or during the game. This could occur through the "lure" of improved utility, or alternatively the clear desire of a noncooperating player to maximize the outcome at the expense of the other players.

3 Applications

We envisage that the proposed research plan will potentially result in the following applications that are critical to ensure not only security in transportation and communication networks, but also in other critical civil infrastructures. We recall that the recommended research extensions are intended to facilitate

- 1. Real-Time Performance
- 2. Multi-O-D Minimum-Cost Network Flow
- 3. Solution Algorithms
- 4. n-person Cooperative Model with Hardening
- 5. Shared Cognition

In response, Table 2 serves to detail how each of the proposed research agenda will facilitate *tactical application* in real time, *strategic application* in monitoring network flow, *infrastructure management* by forming coalitions against attacks, *computational science advancement* by advancing solution algorithms, and *behavioral science advancement* by a better understanding of shared cognition. The eight cited articles—[7, 14, 15, 20, 29, 31, 32, 46]—illustrate not only that the respective authors have successfully carried out the extensions, but also used the results for "common good." There are two caveats. First, most of these advances are specific toward either hazard prevention or attack prevention, but not necessarily both. Second, many of these advances tend to be mode specific. Even though most illustrations here are on transportation networks [33], however, we have included communication networks in our models presented in these two chapters [2], and possible extensions to other multimodal networks. Most importantly, we are focusing on methodological extensions toward realism by concluding our chapter on the various practical research outputs.

Aside from model formulations and solution methodologies, it should be remembered that the author of these two chapters has analyzed successfully three *realistic* communication networks as provided by the defense intelligence community. Formal validation experiments can be performed as a follow-up to this current research effort. While empirical considerations, such as a data-gathering strategy, have been mentioned—as in our discussion under "Extension 1"—most of the empirical work is deferred to a follow-up effort. For the current state-of-the-art, we judge that this two-phase approach will yield the best payoff. Follow-on research will cover more practical issues, including not only validation, but also field implementation.

3.1 Research Output 1: Tactical Applications

Since effective incident management is often expensive; we wish to use analytical resources to ensure the best payoff in efficiency, reliability, safety, and security. Toward this, we strive to formulate an incident *prevention* strategy that has multiple tactical applications. Below are a couple of examples:

For special events such as the 2012 Summer Olympics in London, Bell et al. [4] sought paths with least average "costs," considering the worst-case attack probabilities. Bell's experiments to date confirm these useful tactical results, which were validated during the summer-2012 game:

- Routes in the set of desirable paths tend to share few common arcs. In other words, the recommended paths are "disjointed," fanning out to avoid a common passage, which forms a natural "Achilles' heel" for the Red player.
- As identified by the model, arcs where potential losses are largest are to be avoided in forming these paths, lest they become the targets for exploitation.

Extension	Model	Tactical applications	Strategic applications	Infrastructure management	Computational sci- ence advancements	Behavioral science advancements
1—Real- time perfor- mance monitoring	Dai and Lin [15]	They model general, multiclass queuing networks, parallel server systems, net- works of data switches, and queu- ing networks with alternate routes	Preemption of activities is allowed; otherwise, the feasi- ble set of resource allocations is reduced and the network is not sta- ble, even though throughput is optimal	Maximum pressure policies, defined through extreme resource allocations, guide complex sys- tems toward allocat- ing service capacities dynamically	Being able to address new net- work elements, maximum- pressure policies in a <i>sto-</i> <i>chastic</i> processing network are proven throughput optimal, and at the same time able to characterize the stability regions of the networks	
	Losada et al. [32]	They reduce the response time to more distant facili- ties following a worst-case attack	There is a value of the protection bud- get after which the marginal improve- ment in efficiency becomes negligible	They present a facil- ity protection model with facility recovery time, identifying the optimal allocation of protection resources in an uncapacitated median network	For the bi-level mixed-integer linear program, super- valid inequalities decomposition iter- ates faster than the other decomposi- tions, including benders decomposition	They provide a bi- level model where the upper level fashions the Blue player's protection decisions while the lower level fashions the Red player's optimal response
2—Multi-O- D minimum- cost network flow model	Cappanera and Scap- arra [7]	Red player wants to disrupt the shortest path, with each arc characterized by cost and delay	An attack will not destroy an arc; it will only delay the travel	Inclusion of protec- tion decisions into an optimization model can produce much sounder protection than simply using	The heuristic inter- diction solutions within a tree search scheme are effective at solving protection problems on	A 3-level model is proposed to mini- mize the network- flow cost, to maxi- mize the damage by the Red player, and (continued)

Table 2 Implications of research extensions

Table 2 (contin	ued)					
Extension	Model	Tactical applications	Strategic applications	Infrastructure management	Computational sci- ence advancements	Behavioral science advancements
				shortest path inter- diction, but with some offensive resources rendered ineffective meanwhile	networks with more than 200 nodes and almost 1,000 arcs, allowing future extension to multi- ple O-D flow models	to minimize the damage by the Blue player through hardening
	Croxton et al. [14]	Both travel cost and throughput are of interest		The computational procedure handles nonconvex piecewise linear costs that arise in telecommunica- tions, transportation, and logistics	They approximate the cost objective function with its lower convex enve- lope in the space of commodity flows, with the cost on each arc being rep- resented as a func- tion of the flow for each commodity— resulting in signifi- cant improvement in LP lower bounds	
3—Solution algorithms for multi-O- D flow	Hui et al. [29]	Network connectiv- ity is of main concern			With a better "sam- pling structure" and smart conditioning, the merge process (MP) and the Per- mutation Monte Carlo (PMC)	
	-	_				(continued)

Table 2 (contin	ued)					
Extension	Model	Tactical applications	Strategic applications	Infrastructure management	Computational sci- ence advancements	Behavioral science advancements
					schemes are supe- rior to the Crude Monte Carlo scheme. Meanwhile, the cross-entropy technique can be applied to further improve the MP and the PMC scheme	
	Szeto [46]	Travel cost reliabil- ity is of main concern	He assesses the effects of the num- ber of coalitions formed by Red players on total net- work expected cost and O–D travel cost reliability		It is not clear whether the occur- rence of stochastic Braess' paradox implies the occur- rence of reliability paradox	He presents cooper- ative game approaches to mea- sure network travel cost reliability, in which the Red players can team up cooperatively to make the worst-case scenario worse
4—Multicri- teria <i>n</i> -per- son coopera- tive model with hardening	Garcia et al. [20]	For the proposed fictitious game, a "player" may be created by aggregat- ing flow at the node of traffic entry into a computer network, fanning out to dis- tinct streams	Through finite sam- pling of the relative frequency distribu- tion over joint actions, the scheme takes into account correlated decisions	The game-theoretic model for decentral- ized optimization, i.e., distributed search and evalua- tion, takes advantage of the fact that many large-scale systems can be decomposed	They contribute substantially to the understanding of the joint strategy ficti- tious play algorithm, proving conver- gence to pure strat- egy Nash equilibriums in the	
						(continued)

	ned)					
Extension	Model	Tactical applications	Strategic applications	Infrastructure management	Computational sci- ence advancements	Behavioral science advancements
		according to desti- nation IP address. There are as many fictitious players as the total number of ingress nodes, with each player follow- ing a path consisting of a sequence of network arcs		into a set of distrib- uted parts with inde- pendent control authority, yet still acknowledge a sys- tem-wide metric for success	case of (i) perfect evaluation of cost and (ii) of noisy cost evaluation	
5 Shared cognition	Kovenock, et al. [31]	For each player, the probability of win- ning any given tar- get is determined by the resource levels that the respective players allocate to that target and the contest success function (CSF) that maps the two play- ers' resource alloca- tions into their respective probabili- ties of winning	The "auction" CSF, in which the player with the greater allocation to a target wins that target with certainty, and the "lottery" CSF, in which the probabil- ity of winning a tar- get equals the ratio of a player's resource allocation to the sum of the players' resource allocations to that target	Under both CSFs, both players' resource expendi- tures exceed their respective theoretical predictions. How- ever, behavior appears to conform to the comparative statics of Nash equi- librium for the parameters to conform increases the Red player's expen- increases and the Blue player's expen- diture decreases		An interesting find- ing is that the auc- tion CSF's theoretical predic- tion that the Red player uses a "gue- rilla warfare" strat- egy and the Blue player uses a "com- plete coverage" strategy is observed under both the auc- tion and lottery CSFs. This is inconsistent with Nash equilibrium behavior under the lottery CSF

We will show that there are much more to these results beyond their apparently intuitive arguments. Hu and Chan [25–27], for example, showed a computationally efficient, safe, and perhaps obscure way to route motorists around incidents (for the *prevention* of incidents). They also showed how to dispatch response vehicles to incidents in a counterintuitive way (for the "cure" of an incident that has occurred). In short, clever *operational* procedures can be formulated to *avoid* the worst consequences of incidents, hazards, and attacks. As shown in Table 2, Losada et al. [32] were able to reduce the response time to more distant facilities following a worst-case attack. This is of particular interest to local authorities who are the first responders, and at the same time have to live with the adverse consequences of incidents.

3.2 Research Output 2: Strategic Applications

A number of stakeholders are involved in a large-scale civil infrastructure. For this research, system operators and users are merely two examples, with public officials being the third [30]. While espousing different objectives, all stakeholders have a common interest to guard against incidents that may compromise the safety and security of the infrastructure network. On a state (provincial) level, networkimprovement budgets to avoid public-facility degradation and natural disasters have been routinely allocated, but an adequate budget to ensure security has yet to be considered. A prime reason is that the worth of security assurance has not been estimated scientifically and has been, at best, a "guestimate." We recommend imputing the dollar value of security assurance based on its potential disruption to economic and non-economic transactions, so as to justify such a budget request. In so doing, we have a very promising start in addressing this problem realistically. Preliminary results suggest that the cost of ensuring security is several times more expensive than routine infrastructure maintenance. Adequate maintenance and upgrade will prevent damage due to natural and technological hazards but are inadequate against malicious tampering. Such a critical distinction will facilitate public debate on the judicious allocation of public and private funds in making our civil infrastructures more secure. In light of the world in which we live today, there has never been a more cogent time to start this debate.

These days, it is equally important to model multiple parties working together in a coalition against *asymmetric* attacks, as in the case of an attack by today's "faceless" terrorists. As shown in Table 2, Szeto [46] tried to assess the effects of the number of coalitions formed by Red players on total network expected-cost and O–D travel-cost reliability. The post-9/11 world, the attack on a commuter train in Madrid (2004), recent Katrina hurricanes in Florida and Louisiana/Mississippi (2005), the collapse of the ill-fated I-35 W Bridge in Minneapolis (2007), the earthquake in Haiti (2010), the flood in Pakistan (2010), and the earthquake and

tsunami in Japan (2011) have sensitized our appreciation for security and safety. The importance of guarding against both hazards and attacks speaks amply for itself.

3.3 Research Output 3: Infrastructure Management

Let us concentrate on the improvement of a network to withstand natural hazards. rather than guarding against attacks. Unlike replacing a failed component in equipment, damage to public infrastructure networks is unique in that it takes significant time and effort to repair. Most importantly, it disrupts traffic flow and economic and noneconomic transactions during reparation. Likewise, a civil infrastructure also takes time and effort to upgrade and maintain. Anticipating and preparing for such lengthy periods of reparation and upgrade is of paramount importance. Accordingly, we wish to prevent unnecessary repairs and upgrades by monitoring an infrastructure network for unforeseen service disruptions continuously. We propose a *real-time* model to monitor *existing* network performance for any potential problems as well as to monitor the performance of networks being reconditioned, improved, or hardened. Cappanera and Scaparra [7] carried this one step further by proactively rendering Red's resources ineffective beyond regular interdiction. For computational efficiency, we propose that such a model focuses on the upper-and-lower bounds of performance measures. Where safety factors are typically included in any design, performing designs as if we can handle only a lower-than-actual amount of traffic represents an equally conservative approach in engineering practice. This translates to replacing the actual design capacity with a lower bound on the capacity—as if the lower figure is all the design can handle. One should note that such a conservative approach could potentially increase the solution cost by requiring higher than necessary investments. However, in the case of estimating expected maximum flow, we have established that the use of lower bounds achieves considerable savings in computation, since it is only a point estimate, instead of obtaining a full probability distribution on performance-as proposed by Monte Carlo simulation models [9].

As shown in Table 2, Hui et al. [29] suggested that with a better "sampling structure" and smart conditioning, the Merge Process (MP) and the Permutation Monte Carlo (PMC) schemes are superior to the Crude Monte Carlo scheme. Meanwhile, the Cross-Entropy technique can be applied to further improve the MP and the PMC scheme. Our computational experiments to date (in Extension 1 as described above) suggest that we are on track in monitoring network performance on a real-time basis. This finding—if verified—would allow us to transfer the know-how to many other infrastructure projects of current interests, including *online structural health monitoring*—representing again a potentially cross-cutting and transformative research result.

3.4 Research Output 4: Computational Science Advancements

For the *theoretical side* of our research, network flow is a classic problem rich in analytical properties. While steady progress is being made, the field has not seen as much excitement over the last couple of decades as in its founding days. We are convicted that the spinoff from this research will add significantly toward a resurgence of activities in this field. In the transportation literature, multiple O-D pair flow distribution has long been modeled in terms of the *Cournot-Nash spatial equilibrium* [10], which characterizes the non-cooperative competition among multiple 'players' identified by their commuting travel between O-D pairs. As shown in Table 2, Garcia et al. [20] carried this one step further by solving such an optimization model by a *fictitious-play* algorithm, based on an *imaginary* game between players with identical interests. The interaction of these players allows for exploration of the solution space *computationally* and, for some problems, ultimately arrives at the optimal solution.

Garcia et al. also contributed substantially to the understanding of the jointstrategy fictitious-play algorithm, proving convergence to pure strategy Nash equilibriums in the case of (i) perfect evaluation of cost and (ii) of noisy cost evaluation. These computational investigations recommend a *game-theoretic algorithm* that decomposes the problem into a set of distributed parts with independent control authorities yet still acknowledges a system-wide metric for success. Preliminary investigations also reaffirm a familiar result in the network literature, i.e., that if the underlying component-cost functions are convex, the algorithm converges almost surely to an optimal solution, suggesting that our network model and solution algorithm—being a close cousin of these rather sparse efforts—has a significant *computational* implication for science and engineering in general. This goes well beyond modeling and analyzing the immediate transportation or communication network problems.

3.5 Research Output 5: Behavioral Science Advancements

It is possible to have a mixture of social behaviors between the participants in a game. This could occur in cases where cooperative teaming arrangements occur to "counter" overzealous selfish behavior of an individual player. In this latter case, two distinct player groups may evolve: the single selfish player group and the cooperative player group, who would be opposed to the former group.

The notion of *Shapley's value* can now be introduced. Shapley's value (SV) suggests player *i*'s reward should be the expected amount that player *i* contributes to the coalition made up of *n* players. In other words, the SV is the expected marginal contribution. Suppose coalition players are assigned an average of such marginal contributions. If, however, one of the players receives more than the other

players, even though it contributed less, then the reward is not a SV, and therefore even though cooperative gaming has occurred, the basis on which the coalition has been formed is biased in favor of one of the players. The reason for this bias could be explored.

Coalition worth can also be examined, where the worth is derived from the individual reward that each player could receive as part of the coalition. The coalition worth could also be based on the negative utility for the selfish player not achieving the coalition bid. Ultimately, the division of the aggregate coalition worth among players would need to be considered. Again, SV may be considered for such a division.

As shown in Table 2, experiments by Kovenock et al. [31] uncovered that in a two-person game, a specific theoretical prediction that the Red player uses a "guerilla warfare" strategy and the Blue player uses a "complete coverage" strategy was observed under all conditions, which is inconsistent with Nash equilibrium behavior. Furthermore, the Red player uses a stochastic guerilla warfare strategy, by randomly attacking a subset of the targets and ignoring the remaining targets. This interesting and sometimes suboptimal behavior by the Red player calls for further empirical research. Obviously, other game-theoretic properties will be investigated also, including behavior modification, characterizing game information, myopic value of information, and side information from voluntary "table talk." In other words, all the features discussed in the above state-of-the-art review will be included. Upon a careful design of experiment, analyzing the resulting data will yield a fair amount of insight. From a scientific viewpoint, we could advance the emerging field in information value theory, including epistemic utility theory.

4 Conclusions and Recommendations

A major focus of this Part 2 chapter is how to add to the dearth of literature on preventing and re-mediating both hazards and attacks, particularly on a network level. In hazard-and-attack mitigation, there are generally two approaches. The first is a dynamic game involving a leader and a follower. In transportation networks, for example, road users are critical players in the system, besides the Blue (defender) and the Red (attacker) players in the security-focus literatures. In today's security-minded world, two games are being proposed, where the game among road users is a lower level problem that should be embedded in the higher level Blue–Red game. The model then iterates between these two levels to trace the evolution of the game. In spite of efficient computational schemes such as *Super-Valid-Inequalities Decomposition* [32], formalization of the theory behind such dynamic Stackelberg game is still forthcoming. The second is a static game with or without a time dimension; here a general equilibrium condition is sought as a way to discourage attacks. In spite of significant gains in modeling and computation, no one has yet established the *formal* conditions under which a Pareto Nash equilibrium exists

[13], a robust condition under which network degradation is prevented, and tampering is neither possible nor desired.

Garcia et al. [20] has proposed a decentralized solution method for general network optimization, facilitated through fictitious game playing. Their fictitious game-theoretic approach has profound implications on general mathematical-programming solution algorithms, proving convergence to pure-strategy Nash equilibriums in the case of (i) perfect evaluation of cost, and (ii) of noisy cost evaluation. If the underlying component cost-functions are convex, the algorithm converges almost surely to an optimal solution, suggesting that our network model and solution algorithm—being a close cousin of these rather sparse efforts—has a significant *computational* implication for science and engineering in general, going well beyond the applicational context of these two chapters. The following question remains: Can this finding be generalized to mixed-strategy games and cooperative games, and if so, how?

In the real world, arriving at such an equilibrium includes the necessary posturing to discourage terrorist attacks, and to perform preventive maintenance and upgrade on our critical civil infrastructure ahead of a hazard. It has profound strategic implications on the design and upkeep of large-scale public facilities, ranging from highways to power grids [17]. Only selected literatures speak on policy vs. tactical schemes to best allocate resources between hazards & attacks for specific systems [42, 50] with the remainder addressing either hazard [15] or attacks [31, 32]. According to Guikema [21], general principles of how to best defend systems of specific types against intelligent attacks are emerging that can help system managers allocate resources to best defend their systems. There is still work to be done to improve these models. The current game-theoretic models for intelligent threats are based on a number of assumptions for which the implications and accuracy have not yet been fully explored. The most recent literature as summarized in Table 2 provides "food for thought" in tackling these unsolved problems, although most of them address either hazard or attack prevention, but not both. Beyond gaming, equilibration to include behavioral factors such as posturing and side information is still in an infantile stage. Much work remains as to how to include each into the general models as outlined in Table 1, where five models for preventing hazards and attacks are assembled.

We claim that the above research framework will gain insights into a variety of infrastructure networks. While most of our discussions have been centered upon communication and transportation networks, Holmgren et al. [23] has fully endorsed our game-theoretic approach to defend electric power networks. The uncertainty regarding the outcome of an attack could be represented using stochastic variables. They pointed out, however, that the objectives of the Blue player and the Red player need not be the complete opposites of each other. Beside the antagonistic threats, electric power networks are subject to technical component failures and weather-related disruptions. We refer to them as "hazards" throughout this chapter. While we provided a multidisciplinary paradigm in our two chapters, the readers should always be mindful that the law of physics on electric current is a

major constraint that cannot be overlooked in electricity networks—an item that could be addressed in future work. Component failures due to natural phenomena are not necessarily independent. In electricity transmission, for example, the failure of one or more lines may cause the tripping of additional transmission assets (also known as "cascading failures"), which may ultimately lead to voltage instability and blackouts. It is clear that mode specific issues need closer attention, including the definition difference on reliability between the telecommunication community and the transportation community, as pointed out at the outset of the previous chapter.

Of equal importance is the different treatment of network congestion effects, by virtue of the physical difference in a telecommunication network and a transportation network. For example, Bell [3] specifies a fixed arc cost for the functioning and failed states with possible extension to traffic-dependent arc costs. Arc costs are set equal to undelayed arc travel times plus delay, where delay was calculated using the usual arc-specific vehicle service rate and half this value—yielding a cost for an arc as functions of arc flow. This "method of successive averages" can be applied in situations where arc costs are traffic-dependent. The method offers great flexibility in the way the game is specified, and can (in principle) allow for arc congestion and the introduction of multiple commodities, e.g., multiple origins and destinations. While no proof of convergence is known, Bell has yet to encounter an example where convergence to an equilibrium trip cost has not occurred rapidly.

Meanwhile, a whole host of behavioral issues remains unsolved. The *price of anarchy* is a concept in game theory that measures how the efficiency of a system degrades due to selfish behavior of its agents—a metric worth considering in future work. For example, consider many commuters trying to go from home to work. Let efficiency in this case mean the average time for a commuter to reach the destination. In the "centralized" solution, a central authority can tell each commuter which path to take in order to minimize the average travel time for all. In the "decentralized" version, each commuter chooses its own path. The *price of anarchy* measures the ratio between the average travel time in these two cases. Viewing the problem from this angle, there is much that can be done in the behavioral science perspective. In fact, this might very well be the most exciting development in this field.

Acknowledgments This chapter draws heavily from the published works of the author. The original publication sources are properly cited throughout. Certain sentences may be paraphrased from these sources, and certain figures are reproduced with or without editing. All the published research cited in this fashion was funded by the U.S. Department of Defense (DOD) and was performed while the author was a government employee. The DOD support is gratefully acknowledged. The author also wishes to acknowledge the assistance of Henry Shyllon who assembled some of the reference publications used in this chapter upon the author's suggestion. Mr. Shyllon also formatted the first draft of this chapter according to the publisher guidelines.

References

- Aliprantis CD, Chakrabarti SK (2000) Games and decision making. Oxford University Press, New York
- 2. Antoniou J, Pitsillides A (2012) Game theory in communication networks: cooperative resolution of interactive networking scenarios. CRC Press, Boca Raton
- 3. Bell MGH (2003) The use of game theory to measure the vulnerability of stochastic networks. IEEE Trans Reliab 52(1):63–68
- Bell MGH, Kanturska U, Schmocker J-D, Fonzone A (2008) Attacker-defender models and road network vulnerability. Roy Soc Lond Philos Trans A Math Phys Eng Sci 366:1893–1906
- 5. Braess D (1968) Uber ein paradoxen der verkehrsplanung. Unternehmensforschu 12:258-268
- 6. Cannon-Bowers JA, Salas E (2001) Reflection on shared cognition. J Organ Behav 22:195–202
- 7. Cappanera P, Scaparra MP (2011) Optimal allocation of protective resources in shortest-path networks. Transp Sci 45:64–80
- 8. Castrillo DP, Wettstein D (2000) Bidding for the surplus: a non-cooperative approach to the Shapley value. Working Paper 461.00, Universitat Autonoma de Barcelona, Spain
- 9. Chan Y (2005) Location, transport and land-use: modelling spatial-temporal information. Springer, Berlin, 930 pp (with Web-based software)
- Chan Y, McCarthy J (in press) Game-theoretic paradigms in collaborative research: part 1 theoretical background. Int J Soc Syst Sci
- Chan Y, McCarthy J (in press) Game-theoretic paradigms in collaborative research: part 2 experimental design. Int J Soc Syst Sci
- Chan Y, Yim E, Marsh A (1997) Exact and approximate improvement to the throughput of a stochastic network. IEEE Trans Reliab 46(4):473–486
- 13. Cominetti R, Correa JR, Stier-Moses NE (2009) The impact of oligopolistic competition in networks. Oper Res 57:1421–1437
- Croxton KL, Gendron B, Magnanti TL (2007) Variable disaggregation in network flow problems with piecewise linear costs. Oper Res 55(1):146–157
- Dai JG, Lin W (2007) Maximum pressure policies in stochastic processing networks. Oper Res 55(4):662–673
- 16. Del Vecchio JR, Chan Y, Bruso K (2014) A game-theoretic model for secure international communications. Working Paper, Department of Systems Engineering, University of Arkansas at Little rock, Little Rock, Arkansas
- Downward A, Zakeri G, Philpott AB (2010) On cournot equilibria in electricity transmission networks. Oper Res 58:1194–1209
- 18. Evans JR (1976) Maximum flow in probabilistic graphs-the discrete case. Networks 6:161-183
- 19. Ford LR Jr, Fulkerson DR (1962) Flows in networks. Princeton University Press, Princeton
- Garcia A, Patek SD, Sinha K (2007) A decentralized approach to discrete optimization via simulation: application to network flow. Oper Res 55(4):717–732
- 21. Guikema SD (2009) Game theory models of intelligent actors in reliability analysis: an overview of the state of the art. In: Bier VM, Azaiez MN (eds) Game theoretic risk analysis of security threats. International series in operations research and management science. Springer, New York
- 22. Hausken K, Bier, VM, Zhuang J (2009) Defending against terrorism, natural disaster, and all hazards. In: Bier VM, Azaiez MN (eds) Game theoretic risk analysis of security threats. International series in operations research and management science. Springer, New York
- Holmgren ÅJ, Jenelius E, Westin J (2007) Evaluating strategies for defending electric power networks against antagonistic attacks. IEEE Trans Power Syst 22:76–84
- Hsieh C-C, Lin M-H (2006) Simple algorithms for updating multi-resource allocations in an unreliable flow network. Comput Ind Eng 50:120–129
- Hu J, Chan Y (2005) A multi-criteria routing model for incident management. In: Proceedings, 2005 IEEE international conference on systems, man and cybernetics, Waikoloa, Hawaii, 10– 12 Oct 2005, pp 832–839
- 26. Hu J, Chan Y (2008) Dynamic routing to minimize travel time and incident risks. Paper ID 403, proceedings of the 10th international conference on application of advanced technologies in transportation, Athens, Greece, 28–30 May, 14 pp
- Hu J, Chan Y (2013) Stochastic incident-management of asymmetrical network-workloads. Transp Res Part C 27:140–158
- Huang GL, Song F, Wang XD (2010) Quantitative modeling of couple piezo-elastodynamic behavior of piezoelectric actuators bonded to an elastic medium for structural health monitoring: a review. Sensors 10:3681–3702
- 29. Hui K-P, Bean N, Kraetzl M, Kroese DP (2005) The cross-entropy method for network reliability estimation. Ann Oper Res 134:101–118
- 30. Jha MK, Okonkwo FO (2008) Analysis of system and user travel-time reliability in urban areas. In: Proceedings of the 10th international conference on application of advanced technologies in transportation, Paper ID 873, Athens, 28–30 May, 10 pp
- 31. Kovenock D, Roberson B, Sheremeta RM (2010) The attack and defense of weakest-link networks. CESIFO working paper no. 3211, Ifo Institute, Center for Economic Studies, Munich Society for the Promotion of Economic Research
- Losada C, Scaparra MP, O'Hanley JR (2012) Optimizing system resilience: a facility protection model with recovery time. Eur J Oper Res 217:519–530
- Lou Y, Zhang L (2011) Defending transportation networks against random and targeted attacks. Transp Res Rec J Transp Res Board 2234:31–40
- 34. Lownes NE, Wang Q, Ibrahim S, Ammar RA, Rajasekaran S, Sharma D (2011) A many-tomany game theoretic approach to measuring transportation network vulnerability. Transp Res Rec J Transp Res Board 2263:1–8
- 35. Lyle D, Chan Y, Head E (1999) Improving information-network performance: reliability versus invulnerability. IIE Trans 31:909–919
- Missouri Agency Records Retention Schedules (2014) https://www.sos.mo.gov/records/ recmgmt/retention/agency.asp. Accessed 20 May 2014
- Murray-Tuite PM, Fei X (2010) A methodology for assessing transportation network terrorism risk with attacker and defender interactions. Comput-Aided Civil Infrastruct Eng 25:396–410
- Robinson AR, Chan Y, Dietz DC (2006) Detecting a security disturbance in multi-commodity stochastic networks. Telecommun Syst 31(1):11–27
- 39. Samuelson L (1997) Evolutionary games and equilibrium selection. MIT Press, Cambridge
- 40. Sancho NGF (1988) On the maximum expected flow in a network. J Oper Res Soc 39:481-485
- Scaparra MP, Church RL (2008) A bilevel mixed-integer program for critical infrastructure protection planning. Comput Oper Res 35:1905–1923
- 42. Schavland J, Chan Y, Raines R (2009) Information security: designing a stochastic-network for reliability and throughput. Naval Res Logist 56(7):625–641
- 43. Smith JC, Lim C (2008) Algorithms for network interdiction and fortification games. In: Chinchuluun A, Pardalos PM, Migdalas A, Pitsoulis L (eds) Pareto optimality, game theory and equilibria. Springer, New York
- 44. Steuer R (1986) Multiple criteria optimization: theory, computation, and application. Wiley, Englewood Cliffs, New Jersy
- 45. Szeto WY, O'Brien L, O'Mahony M (2009) Measuring network reliability by considering paradoxes: multiple network demon approach. Transportation research record, No. 2090, Transportation Research Board, Washington, DC, pp 42–50
- Szeto WY (2011) Cooperative game approaches to measuring network reliability considering paradoxes. Transp Res Part C 19:229–241
- 47. Van Hove JC, Chan Y, Caromi R, Abbosh A (2013) Performance in stochastic communication networks: monitoring and prediction. Working Paper, Systems Engineering Department, University of Arkansas at Little Rock, Little Rock, Arkansas

- 48. van den Nouweland A, van Golstein Brouwers W, Groot Bruinderink R, Tijs S (1996) A game theoretical approach to problems in telecommunications. Manage Sci 42:294–303
- 49. Won J, Karray F (2011) A greedy algorithm for faster feasibility evaluation of all-terminalreliable networks. IEEE Trans Syst Man Cybern B Cybern 41(6):1600–1611
- 50. Zhuang J, Bier VM (2007) Balancing terrorism and natural disasters—defensive strategy with endogenous attacker effort. Oper Res 55(5):976–991

The Price of Airline Frequency Competition

Vikrant Vaze and Cynthia Barnhart

Abstract Competition based on service frequency influences capacity decisions in airline markets and has important implications for airline profitability and airport congestion. The market share of a competing airline is a function of its frequency share. This relationship is pivotal for understanding the impacts of frequency competition on airport congestion and on the airline business in general. Additionally, airport congestion is closely related to several aspects of runway, taxiway, and airborne safety. Based on the most popular form of the relationship between market share and frequency share, we propose a game-theoretic model of frequency competition. We characterize the conditions for Nash equilibrium's existence and uniqueness for the two-player case. We analyze myopic learning dynamics for the non-equilibrium situations and prove their convergence to Nash equilibrium under mild conditions. For the N-player symmetric game, we characterize all the pure strategy equilibria and identify the worst-case equilibrium, i.e., the equilibrium with maximum total cost. We provide a measure of the congestion level, based on the concept of price of anarchy and investigate its dependence on game parameters.

Keywords Airline competition • Price of anarchy • Airport congestion • Airline probability • Degree of inefficiency • Best response

1 Introduction

Since deregulation of the U.S. domestic airline business in 1978, airlines have used fare and service frequency as the two most important instruments of competition. Passengers have greatly benefited from fare competition, which has resulted in a

C. Barnhart

V. Vaze (🖂)

Thayer School of Engineering, Dartmouth College, Hanover, NH, USA e-mail: Vikrant.s.vaze@dartmouth.edu

Department of Civil and Environmental Engineering, Massachusetts Institute of Technology, Cambridge, MA, USA

[©] Springer International Publishing Switzerland 2015

K. Hausken and J. Zhuang (eds.), *Game Theoretic Analysis of Congestion, Safety and Security*, Springer Series in Reliability Engineering, DOI 10.1007/978-3-319-13009-5_7

substantial decrease in real (inflation adjusted) air fares over the years. The frequency competition has also resulted in the availability of more options for air travel. The benefits of increased competition to the airlines themselves are not as obvious. Throughout the post-deregulation period, airline profits have been highly volatile. Several major U.S. carriers have incurred substantial losses over the last decade with some of them filing for Chapter 11 bankruptcy and some others narrowly escaping bankruptcy.

Provision of excess seating capacity is one of the reasons often cited for the poor economic health of airlines. Due to the so-called S-curve relationship between market share and frequency share, an airline is expected to attract disproportionately more passengers by increasing its frequency share in a market [8]. To increase their profits, airlines engage in frequency competition by providing more flights per day on competitive routes. The increased frequencies by all competitors on these routes result in fewer passengers and fewer seats per flight. The average aircraft sizes in domestic U.S. markets have been falling continuously over the last couple of decades (until the recent economic crisis) in spite of increasing passenger demand [10]. Similarly, the average load factors, i.e., the ratio of the number of passengers to the number of seats, on some of the most competitive and high-demand markets have been found to be lower than the industry average.

Apart from the chronic worries about the industry's financial health, worsening congestion and delays at the major U.S. airports have become another cause of serious concern. Increases in passenger demand, coupled with decreases in average aircraft size have led to a great increase in the number of flights being operated, especially between major airports, leading to congestion. The Total Delay Impact Study commissioned by the Federal Aviation Administration (FAA) of the United States has estimated that in calendar year 2007, delays cost over \$8 billion to airlines and another \$17 billion to passengers [3]. Note that 2007 remains the year of historically high delays. Though various structural changes have happened in the U.S. airline industry after 2007, no comprehensive study analyzing delays in the subsequent years has been conducted.

Congestion is closely linked to safety and security considerations. Airborne safety considerations, such as minimum spacing requirements during various phases of a flight, can ultimately determine the level of congestion and delays. For instance, wake vortex avoidance considerations set the minimum longitudinal separation requirements during a landing aircraft's runway approach. This minimum separation is a primary factor that decides the maximum landing rate, which in turn determines the airport congestion levels during periods of bad weather and/or high volume. Conversely, congestion metrics, such as an increase in the number of flights, have been shown to have a direct connection to aviation safety, with increases in congestion leading to potential safety risks. This is true in the context of both ground maneuvers and airborne operations. For instance, airborne queuing in the airspace surrounding a congested airport is considered to be a major safety hazard. Through technological improvements and strict safety regulations, first world countries have managed to improve airborne safety considerably, despite a steady increase in congestion levels. Midair collisions are on a decline, with no

midair collisions involving scheduled passenger flights being reported in a 20-year period between 1988 and 2008 [4]. However, ground maneuvers continue to be a cause of concern.

Approximately 30 % of the all aviation accidents between 1995 and 2008 involving commercial transport aircraft were runway related, which is a category of on-ground accidents [38]. Transport Canada [34] found that an increase in flight traffic volume is typically associated with a disproportionately large increase in runway incursion potential. Using a single-runway model, they showed that a 20 % increase in flight volume results in a 140 % increase in the runway incursion potential. Furthermore, actual data on traffic volumes and runway incursions in the U.S. between 1988 and 1998 consistently indicates that traffic increases of the order of 2-6 % are associated with runway incursion rate increases of the order of 30-70 % [34]. Barnett et al. [5] found that the risk of runway collisions varies proportionately to the square of the number of flight operations. Thus, understanding the factors affecting airport congestion can provide valuable insights into the nature and causes of safety challenges.

Furthermore, security issues also have an impact on congestion levels. Adverse events such as the September 11 terrorist attacks reduce passenger demand and airport congestion. U.S. airport congestion went down considerably in the months after the attacks. Additionally, delays introduced due to enhanced security procedures also contribute to total delays to flights and passengers. Finally, terror scares, such as the one in August 2006 at London Heathrow causing 75 % flight cancelations, can lead to a large number of stranded and delayed passengers.

In summary, frequency competition affects airlines' capacity allocation decisions, which have a strong impact on airline profitability and airport congestion which in turn affects aviation safety. In this chapter, we propose a game-theoretic framework, which is consistent with the most prevalent model of frequency competition. We analyze and prove the existence, uniqueness, and convergence properties of equilibrium for the airline frequency competition game and prove the dependence of airport congestion and airline profits on the parameters of the frequency competition game. In this chapter, unless stated otherwise, we use the term frequency to refer to the number of flights scheduled over a specific time period (e.g., a day) by a specific airline carrier over a specific nonstop segment. Section 2 provides background on airline schedule planning and reviews the literature on frequency competition. Section 3 presents the N-player game model. Best response curves are characterized in Sect. 4. In Sect. 5, we focus on the two-player game. We provide the conditions for existence and uniqueness of a Nash equilibrium and discuss realistic ranges of parameter values. We then provide two different myopic learning models for the two-player game and provide proof of their convergence to the Nash equilibrium. In Sect. 6, we identify all possible equilibria in an N-player game with identical players and find the worst-case equilibrium. We then evaluate the price of anarchy and establish the dependence of airline profitability and airport congestion on airline frequency competition. We conclude with a summary of main results and provide some promising directions for future research in Sect. 7.

2 Frequency Planning Under Competition

The airline planning process involves decisions ranging from long-term strategic decisions such as fleet planning and route planning, to medium-term decisions about schedule development [7]. Fleet planning is the process of determining the composition of a fleet of aircraft, and involves decisions about acquiring new aircraft and retiring existing aircraft in the fleet. Given a fleet, the second step in the airline planning process involves the choice of routes to be flown, and is known as the route planning process. A route is a combination of origin and destination airports (occasionally with intermediate stops) between which flights are to be operated. Route planning decisions take into account the expected profitability of a route based on demand, fare, and revenue projections as well as the overall structure of the airline's network. Given a set of selected routes, the next step in the planning process is airline schedule development, which in itself is a combination of decisions about frequency, departure times, and aircraft sizes for each route, and aircraft rotations over the network.

Frequency planning is the part of the airline schedule development process that involves decisions about the number of flights to be operated on each route. By providing more frequency on a route, an airline can attract more passengers. Given an estimate of total demand on a route, the market share of each airline depends on its own frequency as well as on competitor frequency. The S-curve or sigmoidal relationship between market share and frequency share is a widely accepted notion in the airline industry [8, 25]. However, it is difficult to trace the origins and evolution of this S-shaped relationship in the airline literature [16]. Empirical evidence of the relationship was documented in some early post-deregulation studies and regression analysis was used to estimate the model parameters [31–33]. Over the years, there have been several references to the S-curve including Kahn [22] and Baseler [6]. In this chapter, we use a more general model that is compatible with the linear, as well as the S-curve assumptions. The mathematical expression for the S-curve relationship [8, 31] is given by:

$$\mathcal{M}_i = \frac{F_i^{\alpha}}{\sum_{j=1}^n F_j^{\alpha}} \tag{1}$$

for parameter α such that $\alpha \geq 1$, where \mathcal{M}_i = market share of airline *i*, F_i = frequency share of airline *i*, and *n* = number of competing airlines. Note that this S-curve model structure found in the existing literature assumes that the market share of an airline is dependent only on its frequency share.

The structure of airline business has evolved over the last few decades. Some recent empirical and econometric studies have focused on investigating the extent to which the S-curve phenomenon remains valid and relevant. The resulting conclusions are quite mixed. In one of the recent studies, Wei and Hansen concluded that by increasing service frequency, an airline can get a disproportionately high share of the market and hence there is an incentive for operating more frequent flights with smaller aircraft [36]. Their analysis was based on a nested Logit model of passenger demand in nonstop duopoly markets. Button and Drexler found limited evidence for the existence of the S-curve phenomenon in the 1990s [16]. But for data from the early years of this century, they found that the relationship between market share and frequency share was well approximated by a 45° straight line, which can be characterized by setting $\alpha = 1$ in Eq. (1). They, however, have warned the industry analysts that observed lack of empirical support for the existence of S-curve in the recent years does not necessarily mean that it does not affect airline behavior. Furthermore, the general model in Eq. (1) continues to be consistent with their conclusions with the only modification being that α is set to 1 rather than to a value greater than 1. Binggeli and Pompeo found that the S-curve still exists in markets dominated by legacy carriers but not so much in markets where low cost carriers (LCCs) compete with each other, with the latter being better approximated by a straight line relationship [9]. They suggest a rethinking of the S-Curve-based planning methods that have been "hard-wired" in the heads of many network planners for decades.

A recent study by Vaze and Barnhart [35] used the S-curve model of market share for modeling airline frequency decisions under competition and found a good fit to empirical frequency data. Furthermore, some of the prior studies involving passenger choice models belonging to the Logit family have expressed utility as an affine function of the logarithm of flight frequency [20, 26, 37]. All other variable values being equal, this formulation simplifies to the S-curve where the frequency exponent (α) in the S-curve relationship is equal to the utility coefficient of the natural logarithm of frequency.

In summary, recent evidence confirms that market share is an increasing function of frequency share and hence competition considerations affect frequency decisions in an important way. However, evidence is mixed about the exact shape of the relationship, in particular the exact value of the parameter α for different types of markets.

Many of these studies go on to discuss the financial implications of the S-curve. Button and Drexler [16] associate it with provision of "excess capacity" and an "ever-expanding number of flights," while O'Connor [25] associates it with "an inherent tendency to overschedule". Kahn goes even further and raises the question of whether it is possible at all to have a financially strong and yet highly competitive airline industry at the same time [22].

Despite continuing interest in frequency competition based on the S-curve phenomenon, there is only a very limited amount of literature on game-theoretic properties of such competition. Hansen [20] was the first notable attempt to analyze frequency competition over airline networks in a game-theoretic setting. He analyzed airline frequency competition in a hub-and-spoke environment and validated the results against actual data. Dobson and Lederer [18] modeled airline competition on both schedule and fare as a strategic form game, while Adler [1] modeled competition on fare, frequency, and aircraft sizes as an extensive form game. Vaze and Barnhart [35] also modeled frequency competition at a single airport as a strategic form game and validated the results against actual data.

studies adopted a successive optimizations approach to solve for a Nash equilibrium. Only Hansen [20], and Vaze and Barnhart [35] mention some of the issues regarding convergence through discussion of different possible cases. But none of the existing studies provides any convergence conditions. In this chapter, we describe and prove the convergence of two simple dynamics to a Nash equilibrium under mild conditions.

Not all the prior studies use a successive optimizations approach. In fact, Wei and Hansen [37] solve for equilibrium through explicit enumeration of the entire strategy space. Brander and Zhang [11], Aguirregabiria and Ho [2], and Norman and Strandenes [24] model airline competition as a dynamic game and estimate the model parameters using empirical data. None of the studies mentioned so far provide a rigorous treatment of the question of existence or uniqueness of a pure strategy equilibrium. Brueckner and Flores-Fillol [13], Brueckner [12] obtain closed form expressions for equilibrium decisions analytically, but they focus exclusively on symmetric equilibria while ignoring the possibility of any asymmetric equilibria. In this chapter, we provide a rigorous treatment of the existence and uniqueness issues while accounting for the symmetric as well as asymmetric equilibria.

Most of the previous studies involving game-theoretic analysis of frequency competition, such as Adler [1], Pels et al. [26], Hansen [20], Wei and Hansen [37], Dobson and Lederer [18], Hong and Harker [21], model market share using Logit-or nested Logit-type models. In these studies, the passenger utility function is typically an affine function of some transformation of frequency, e.g., logarithmic, inverse, or polynomial. It is interesting and important to note that a logarithmic transformation, in fact, corresponds directly to an S-curve, while the rest of the relationships can be substantially different from the S-shaped relationship, depending on the exact values of utility parameters.

All of these studies involve finding a Nash equilibrium or some refinement of it. But there is not sufficient justification of the predictive power of the equilibrium concept. Hansen [20] provides some discussion of the shapes of best response curves and stability of equilibrium points. But none of the studies has focused on any learning dynamics through which less than perfectly rational players may eventually reach the equilibrium state. Vaze and Barnhart [35] do provide some computational evidence of the convergence properties but not a rigorous mathematical justification.

None of the aforementioned studies rigorously characterize the efficiency loss resulting from S-curve-based airline frequency competition. A number of studies (such as Saraydar et al. [29, 30], Goodman and Mandayam [19]) characterize the efficiency loss in communication networks due to the selection of a higher level of transmitter power, compared to that at the system optimal, by individual players. The concept of signal-to-interference ratio considered in these studies is somewhat analogous to the notion of frequency shares in the airline context. While the specifics of the game settings in these studies and the associated forms of the individual players' utility expressions are quite different from those in an airline setting, our

results in Sect. 6 of this chapter are generally consistent with the conclusions drawn by these communication network research studies.

In this chapter, we use the most popular characterization of the S-curve model, as given by Eq. (1). The $\alpha = 1$ case is well suited for modeling markets dominated by LCCs, whereas markets dominated by legacy carriers can be suitably modeled using higher values of α . Thus, despite the mixed recent evidence about the exact shape of the market share-frequency share relationship, the model specified by Eq. (1) captures airline scheduling decisions well. We analyze a strategic form game among airlines with frequency of service being the only decision variable. We will only consider pure strategies of the players, i.e., we will assume that the frequency decisions made by the airlines are deterministic. We use the Nash equilibrium solution concept under the pure strategy assumption. The research contributions of this chapter are threefold. First, we make use of the S-curve relationship between market share and frequency share and analyze its impact on the existence and uniqueness of pure strategy Nash equilibria. Second, we provide reasonable learning dynamics and provide theoretical proof for their convergence to the unique Nash equilibrium for the two-player game. Third, we provide a measure of inefficiency, similar to the price of anarchy, of a system of competing profitmaximizing airlines in comparison to a system with centralized control. This measure can be used as a proxy to understand the effects of frequency competition on airline profitability and airport congestion.

Note that the main focus of this chapter is airline competition based on frequency which is a strategic decision made by the airlines several months ahead of operations. Fare competition is yet another important aspect of airline competition. Fare decisions are made much closer to the day of flight operation. They are affected by the demand segmentation and revenue management practices of the airline as well as those of the competing airlines. Obviously, frequency and fare decisions are interrelated because they together determine the profits. Frequency decisions are necessarily made before the actual fare values are realized. Thus, frequency decisions are made based on likely or expected values of fares which are based on the airline's knowledge and past experience. In this chapter, we model airlines' frequency decisions based on these likely or expected values of fares. The actual mechanics of airline fare competition are beyond the scope of this chapter.

3 Model

Let M be the total market size, i.e., the number of passengers wishing to travel from a particular origin to a particular destination on a nonstop flight. While the total market size can itself be affected by frequency changes through effects such as demand stimulation, we will focus mainly on the distribution of this total market demand across different competitors and assume the total market size to be a constant. In general, an airline passenger may have more than one flight in his/her itinerary. Conversely, two passengers on the same flight may have different origins and/or destinations. But for our analysis, we will ignore these network effects and assume the origin and destination pair of airports to be isolated from the rest of the network. Let $I = \{1, 2, ..., n\}$ be the set of airlines competing in a particular nonstop market. Although most of the major airlines today follow the practices of differential pricing and revenue management, we will assume that the air fare charged by each airline remains constant across all passengers. Let p_i be the fare charged by each airline *i*. Further, we will assume that the type and seating capacity of aircraft to be operated on this nonstop route are known. Finally, we will not account for the effect of flight departure times on market share, except for the extent to which it is indirectly captured by the flight frequencies. Note that we will assume airline frequency values to be a continuous variable in this analysis. Let S_i be the seating capacities for airline *i* and C_i be the operating cost per flight for airline *i*. Let α be the parameter in the S-curve relationship. A typical value suggested by the literature is around 1.5. To keep our analysis general, we make the following assumption.

Assumption 3.1 $1 < \alpha < 2$

Our results are applicable even in the case of a linear relationship between market share and frequency share by taking the limit as $\alpha \rightarrow 1^+$.

Let x_i be the frequency (i.e., the number of flights per day) of airline *i*. As per the S-curve relationship between market share and frequency share, the *i*th airline's share of the market (\mathcal{M}_i) is given by

$$\mathcal{M}_i = x_i^{\alpha} \bigg/ \sum_{j=1}^n x_j^{\alpha}$$
⁽²⁾

This is obtained by multiplying the numerator and denominator of the right-hand side (RHS) of Eq. (1) by $\left(\sum_{j=1}^{n} x_{j}\right)^{\alpha}$. The number of passengers (\wp_{i}) traveling on airline *i* cannot exceed the product of its market share (\mathcal{M}_{i}) and the total market size (M). Additionally, the number of passengers (\wp_{i}) cannot exceed the total number of available seats ($S_{i}x_{i}$). Therefore, the number of passengers (\wp_{i}) traveling on airline *i* is given by

$$\wp_i = \min\left(Mx_i^{\alpha} \middle/ \sum_{j=1}^n x_j^{\alpha}, S_i x_i\right)$$
(3)

Airline i's profit (also referred to as it's *payoff* in a game-theoretic context) is given by

$$\Pi_i = p_i * \min\left(M\frac{x_i^{\alpha}}{\sum_{j=1}^n x_j^{\alpha}}, S_i x_i\right) - C_i x_i \tag{4}$$

We will also make the following assumption.

Assumption 3.2 $C_i < p_i S_i$ for every i

In other words, the total operating cost of a flight is lower than the total revenue generated when the flight is completely filled. This assumption is reasonable because if it is violated for some airline *i*, then there is a trivial optimal solution $x_i = 0$ for that airline.

From here onward, our game-theoretic analysis proceeds as follows. In the next section (Sect. 4), we characterize the shapes of best response correspondences, that is, sets of optimal responses of a player as a function of the frequencies of the other player(s). This analysis, which focuses on the general frequency competition game model as described in this section, facilitates the subsequent analysis of Nash equilibria in Sects. 5 and 6. In our Nash equilibrium analysis, we first focus on the two-player case (in Sect. 5) and later extend the analysis to the symmetric N-player case (in Sect. 6). We restrict our N-player game analysis to only the symmetric player case primarily for tractability reasons. As shown in our analysis in Sect. 5, even in the two-player case, there can be as many as six Nash equilibria depending on the combination of parameter values. There can be a multitude of equilibria in frequency competition games with more players. In real-life airline markets, the parameters of airline frequency competition, such as, fares, seating capacities, and operating costs of competing airlines are often not too different from each other. Therefore, focusing on the symmetric player case is not that unrealistic. Furthermore, as shown in Sect. 6, a thorough analysis of the symmetric player case presents several valuable insights. In Sect. 6, we analyze both symmetric and asymmetric equilibria for the symmetric N-player case. In Sects. 5 and 6, we present major results as a sequence of theorems and proofs. In order to keep the discussion crisp and concise, we defer proofs of some of the theorems to the Appendix.

4 Best Response Curves

Let us define the effective competitor frequency as

$$y_i = \left(\sum_{j \in I, j \neq i} x_j^{\alpha}\right)^{1/\alpha} \tag{5}$$

and let us rewrite profit Π_i as follows:

$$\Pi_i = \min(\Pi'_i, \Pi''_i), \text{ where } \Pi'_i = \left(\frac{Mp_i x_i^{\alpha}}{x_i^{\alpha} + y_i^{\alpha}}\right) - C_i x_i \text{ and } \Pi''_i = p_i S_i \mathbf{x}_i - C_i \mathbf{x}_i \quad (6)$$

Note that the definition of effective competitor frequency is merely a convenient way of portraying the best response curves and does not necessarily have any obvious practical implication outside of its usage in describing the best response curves.



Fig. 1 Typical shapes of profit functions for three cases

 $\Pi_i = \Pi'_i$ if the seating capacity constraint is not binding. We call Π'_i as the *uncapacitated profit*. $\Pi_i = \Pi''_i$ if all seats are filled. We call Π''_i as the *full-load profit*. Π'_i is a twice continuously differentiable function of x_i . Π'_i has a single point of zero curvature at $x_i = y_i((\alpha - 1)/(\alpha + 1))^{1/\alpha}$ and the function is strictly convex for all lower values of x_i and strictly concave for all higher values of x_i . For a given combination of parameters α , M, p_i , C_i , S_i and a given effective competitor frequency y_i , the global maximum of Π_i falls under exactly one of the following three cases. These three cases are also illustrated in Fig. 1a–c, respectively.

Case A: $\prod_i' \leq 0 \forall x_i > 0$. Under this case, either a local maximum with x > 0 does not exist for \prod_i' or it exists but value of the function \prod_i' at that point in negative. In this case, a global maximum of \prod_i' is at $x_i = 0$. This describes a situation where the effective competitor frequency is so large that airline *i* cannot earn a positive profit at any frequency. Therefore, the best response of airline *i* is to have a zero frequency, i.e., not to operate any flights in that market.

Case B: Local maximum of Π'_i exists at some x > 0 and the value of the function Π'_i at that local maximum is positive and less than or equal to Π''_i . In this case, the unique global maximum of Π_i exists at the local maximum of Π'_i . In this case, the optimum frequency is positive and at this frequency, airline *i* earns the maximum profit that it could have earned had the aircraft seating capacity been infinite. Note that under this case, either $\Pi'_i \leq \Pi''_i \forall x_i > 0$ is true, or Π'_i and Π''_i curves intersect each other at two values of $x_i > 0$, even though Fig. 1b only illustrates the latter of these two possibilities.

Case C: A local maximum of Π'_i exists at some x > 0 and the value of the function Π'_i at this local maximum is greater than Π''_i . In this case, Π'_i and Π''_i intersect at two distinct points (apart from $x_i = 0$). The unique global maximum of Π_i exists at the point of intersection with highest x_i value. This describes the case where optimum frequency is positive and greater than the optimum frequency under the assumption of infinite aircraft seating capacity. At this frequency, airline *i* earns lower profit than the maximum profit it could have earned had the aircraft seating capacity been infinite.

Figure 2 shows a typical best response curve. $\Pi'_i(0) = 0$ and for very low positive values of x_i , $\partial \Pi'_i / \partial x_i$ is negative. Therefore, at the first stationary point



(the one with lower x_i value), the Π'_i function value will be negative. Moreover, as $y_i \to \infty$, $\Pi'_i(x_i) < 0 \forall x_i$. For a given combination of parameters α , M, p_i , C_i and S_i , there exists a threshold value of effective competitor frequency y_i such that, for any y_i value above this threshold, $\Pi'_i(x_i) < 0 \forall x_i > 0$ and therefore the best response of airline *i* is $x_i = 0$. Let us denote this threshold by y_{th} and the corresponding x_i value as x_{th} . At $x_i = x_{\text{th}}$ and $y_i = y_{\text{th}}$ (Point 3 in Fig. 2), $\Pi'_i = 0$, $\partial \Pi'_i / \partial x_i = 0$, $\partial^2 \Pi'_i / \partial x_i^2 \leq 0$. Upon simplification, we get

$$x_{\text{th}} = (\alpha - 1) \left(\frac{Mp_i}{\alpha C_i} \right) \text{ and } y_{\text{th}} = (\alpha - 1)^{\frac{\alpha - 1}{\alpha}} \left(\frac{Mp_i}{\alpha C_i} \right)$$
 (7)

Of course, at $y_i = y_{th}$, $x_i = 0$ is also optimal (Point 4 in Fig. 2). It turns out that it is the only y_i value at which there is more than one best response (optimal frequency) possible. This situation is unlikely to be observed in real-world examples, because the parameters of the model are all real numbers with continuous distributions. So the probability of observing this exact idiosyncratic case is zero. If we arbitrarily assume that in the event of two optimal frequencies, an airline chooses the greater of the two values (which is consistent with airline's desire to retain market share), then the best response correspondence reduces to a function, which we will refer to as the best response function. The existence of two different maximum values at $y_i = y_{th}$ means that the best response correspondence is not always convex-valued. Therefore, in the case of a general asymmetric game, a pure strategy Nash equilibrium may or may not exist for this game.

For y_i values slightly below y_{th} , the global maximum of Π_i corresponds to the stationary point of Π'_i in the concave part as described in case B above. Therefore, for y_i values slightly below y_{th} , at the stationary point of Π'_i in the concave part,

 $\Pi'_i < \Pi''_i$. However, as $y_i \to 0$, $argmax(\Pi'_i(x_i)) \to 0$. Therefore, $argmax(\Pi_i(x_i))$ exists at a point of intersection of Π'_i and Π''_i curves, as explained in case C above. For y_i values slightly above 0, at the stationary point of Π'_i in the concave part, $\Pi'_i > \Pi''_i$. By continuity, therefore, for some y_i such that $0 \le y_i \le y_{\text{th}}$, there exists x_i such that, $\Pi'_i = \Pi''_i, \partial\Pi'_i/\partial x_i = 0$ and $\partial^2\Pi'_i/\partial x_i^2 \le 0$. It turns out that there is only one such y_i value that satisfies these conditions. Let us denote this y_i value by y_{cr} , because this is the critical value of effective competitor frequency such that case B prevails for higher y_i values (as long as $y_i \le y_{\text{th}}$) and case C prevails for all lower y_i values. The value of y_{cr} and the corresponding x_i value, x_{cr} , is given by (Point 2 in Fig. 2),

$$x_{\rm cr} = \frac{M}{S_i} \left(1 - \frac{C_i}{\alpha p_i S_i} \right) \text{ and } y_{\rm cr} = \frac{M}{S_i} \left(1 - \frac{C_i}{\alpha p_i S_i} \right) \left/ \left(\frac{\alpha p_i S_i}{C_i} - 1 \right)^{1/\alpha} \right)$$
(8)

For $y_i = 0$, it is easy to see that Π_i is maximized when $x = M/S_i$ (Point 1 in Fig. 2). We will denote the range of y_i values with $y_i > y_{\text{th}}$ as region A, $y_{\text{cr}} \le y_i \le y_{\text{th}}$ as region B and $y_i < y_{\text{cr}}$ as region C.

In region C, Π_i is maximized for a unique x_i value such that $\Pi'_i = \Pi''_i$ and $\partial \Pi'_i / \partial x_i \leq 0$. The equality condition translates into,

$$(M/S_i)x_i^{\alpha-1} - x_i^{\alpha} = y_i^{\alpha} \tag{9}$$

The left-hand side (LHS) of Eq. (9) is strictly concave because $1 < \alpha < 2$. Further, the LHS is maximized at $x_i = \frac{\alpha - 1}{\alpha} \frac{M}{S_i}$, which corresponds to $y_i = (\alpha - 1)^{(\alpha - 1)/\alpha} \frac{M}{\alpha S_i}$. So for every y_i value, there are two corresponding x_i values satisfying this equation that correspond to the two points of intersection of the Π'_i and Π''_i curves. The one corresponding to the higher x_i value corresponds to $\alpha x_i > (\alpha - 1)(M/S_i)$. Therefore,

$$\frac{\partial x_i}{\partial y_i} = \alpha \left(\frac{y_i^{\alpha - 1}}{x_i^{\alpha - 2}} \right) \middle/ \left((\alpha - 1) \left(\frac{M}{S_i} \right) - \alpha x_i \right) < 0$$
⁽¹⁰⁾

and

$$\frac{\partial^2 x_i}{\partial y_i^2} = \left(\frac{\partial x_i}{\partial y_i}\right)^3 \frac{\alpha - 1}{\alpha} \frac{x_i^{\alpha - 3}}{y_i^{\alpha - 1}} (\alpha x_i + (2 - \alpha)\frac{M}{S_i}) + \frac{\partial x_i}{\partial y_i} \frac{\alpha - 1}{y_i} < 0$$
(11)

So, the best response curve is a strictly decreasing and concave function for all $0 \le y_i < y_{cr}$.

In region B, Π_i is maximized for a unique x_i value such that $\partial \Pi'_i / \partial x_i = 0$ and $\partial^2 \Pi'_i / \partial x_i^2 < 0$. Upon simplification, we get

The Price of Airline Frequency Competition

$$\frac{\partial x_i}{\partial y_i} = \frac{x_i}{y_i} \left(x_i^{\alpha} - y_i^{\alpha} \right) \left/ \left(\left(1 + \frac{1}{\alpha} \right) x_i^{\alpha} - \left(1 - \frac{1}{\alpha} \right) y_i^{\alpha} \right) \right.$$
(12)

and

$$\left(1+\frac{1}{\alpha}\right)x_i^{\alpha} - \left(1-\frac{1}{\alpha}\right)y_i^{\alpha} > 0$$
(13)

Therefore, the best response curve $x_i(y_i)$ in region B has zero slope at $x_i = y_i$, is strictly increasing for $x_i > y_i$ and strictly decreasing for $x_i < y_i$. For $x_i = y_i$, we get $x_i = y_i = \frac{\alpha M p_i}{4C_i}$.

Figure 2 describes a typical best response curve as a function of effective competitor frequency. Now we provide some intuition behind the shape of the best response curve.

In region C, the effective competitor frequency is so small that airline *i* attracts a large market share even with a small frequency. Therefore, the optimal frequency ignoring seating capacity constraints is so low that the number of seats is exceeded by the number of passengers wishing to travel with airline *i*. As a result, the optimal frequency and the maximum profit that can be earned by airline *i* are decided by the aircraft seating capacity constraint. In this region, the optimal number of flights scheduled by airline *i* is just sufficient to carry all the passengers that wish to travel on airline *i*. In this region, airline *i* has 100 % load factor at the optimal frequency. With increasing effective competitor frequency, the market share attracted by airline *i* reduces and hence fewer flights are required to carry those passengers. Therefore, the best response curve is strictly decreasing in this region. Once the effective competitor frequency decision. Thus, in the presence of no or very little amount of competition, airlines will try to provide flights such that they can just about satisfy the passenger demand subject to the seating capacity constraint.

In region B, the effective competitor frequency is sufficiently large so the number of passengers attracted by airline *i* does not exceed the seating capacity. Therefore, the aircraft seating capacity constraint is not binding in this region. The optimal frequency is equal to the frequency at which the marginal revenue equals marginal cost, which is a constant C_i . As the effective competitor frequency increases, the market share of airline *i* at the optimal frequency decreases and the load factor of airline *i* at the optimal frequency also decreases. At a large value, y_{th} , of effective competitor frequency, the load factor of airline *i* at its optimal frequency reduces to a value C_i/p_iS_i and the optimal profit drops to zero. Thus, in the presence of considerable competition, airlines optimize their schedule by providing sufficient frequency to attract market share and the seating capacity ceases to have the constraining effect.

For all values of effective competitor frequency above y_{th} , i.e., in region A, there is no positive frequency for which the airline *i* can make positive profit. Therefore, the optimal frequency of airline *i* in region C is zero. Thus, in the presence of too much competition, airlines find it best to stay out of the market.

5 Two-Player Game

Let *x* and *y* be the frequency of carrier 1 and 2, respectively. The effective competitor frequency for carrier 1 is *y* and that for carrier 2 is *x*. For any Pure Strategy Nash Equilibrium (PSNE), the competitor frequency for each carrier can belong to any one of the three regions, A, B, and C. So potentially there are nine different combinations possible. We define the *type* of a PSNE as the combination of regions to which the competitor frequency belongs at equilibrium. We will denote each type by a pair of capital letters denoting the regions. For example, if carrier 1's effective competitor frequency, i.e., *y*, belongs to region B and carrier 2's effective competitor frequency, i.e., *x*, belongs to region C, then that PSNE is said to be of type BC. Accordingly, there are nine different types of PSNE possible for this game, namely AA, AB, AC, BA, BB, BC, CA, CB, and CC.

Frequency competition among carriers is the primary focus of this research. However, it is important to recognize that frequency planning is just one part of the entire airline planning process. Frequency planning decisions are not taken in isolation, the route planning phase precedes the frequency planning phase. Once the set of routes to be operated is decided, the airline proceeds to the decision of the operating frequency on that route. This implicitly means that once a route is deemed profitable in the route planning phase, frequency planning is the phase that decides the number of flights per day, which is supposed to be a positive number. However, in AA, AB, BA, AC, or CA type equilibria, the equilibrium frequency of at least one of the carriers is zero, which is inconsistent with the actual airline planning process. Moreover, for ease of modeling, we have made a simplifying assumption that the seating capacity is constant. In reality, seating capacities are chosen considering the estimated demand in a market. If the demand for an airline in a market exceeds available seats on a regular basis, the airline would be inclined to use larger aircraft. Sustained presence of close to 100 % load factors is a rarity. However type AC, BC, CA, CB, and CC type equilibria involve one or both carriers having 100 % load factors. Zero frequency and 100 % load factors make all types of equilibria, apart from type BB equilibrium, suspect in terms of their portrayal of reality.

We will now investigate each of these possible types of pure strategy Nash equilibria of this game and obtain the existence and uniqueness conditions for each of them.

5.1 Existence and Uniqueness

Theorem 5.1 A type AA equilibrium cannot exist.

Proof If $x^* = 0$ then $\Pi_2 = p_2 * \min(M, S_2 y) - C_2 y$, which is maximized at $y = M/S_2$ because $C_2 < S_2 p_2$. So $y^* > 0$ whenever $x^* = 0$. So this type of equilibrium cannot exist.

Type AA equilibrium, if it exists, is characterized by both competing airlines having zero frequency indicating that it is not profitable for either airline to operate any flights in this market. However, Theorem 5.1 shows that this can never happen under our Assumption 3.2. This indicates that if one of the two airlines has zero frequency in a market, then the other airline will always find it profitable to operate flights at a small frequency in that market.

Theorem 5.2 A type AB (and type BA) equilibrium cannot exist.

Proof Type AB equilibrium exists if and only if $x^* = 0$, $y^* > 0$ and $\wp_2 < S_2 y^*$. As shown above, if $x^* = 0$ then, Π_2 is maximized at $y = M/S_2$ as long as $C_2 < S_2 p_2$. So $\wp_2 = M = S_2 y^*$ whenever $x^* = 0$. So this type of equilibrium cannot exist. By symmetry, type BA equilibrium cannot exist either.

An equilibrium of type AB or BA, if it exists, characterizes a situation where one of the competing airlines stays out of the market (i.e., offers zero frequency) and the second competing airline offers a frequency such that there is excess seating capacity available on its flights. Theorem 5.2 shows that this can never happen because whenever one airline offers zero frequency, the other airline will be able to capture the entire passenger demand in that market with a frequency value equal to the least frequency necessary to provide sufficient seating capacity for all passengers in that market.

Theorem 5.3 A type AC equilibrium exists if and only if

$$\frac{C_1}{S_1 p_1} > \frac{S_2}{S_1} \frac{1}{\alpha} (\alpha - 1)^{\frac{\alpha - 1}{\alpha}}$$
(14)

and if it exists then it is a unique type AC equilibrium.

Proof This type of equilibrium requires $x^* = 0$ and $y^* = M/S_2$. So if an equilibrium of this type exists then it must be the unique type AC equilibrium. For this equilibrium to exist, the only condition we need to check is that $\frac{M}{S_2} = y > y_{\text{th}} = (\alpha - 1)^{\frac{\alpha-1}{\alpha}} \frac{MP_1}{\alpha C_1}$ because for $y^* = \frac{M}{S_2}$, $x^* = 0$ is true if and only if $\Pi_1 < 0$, for all $x \ge 0$. So type AC equilibrium will exist if and only if

$$\frac{C_1}{S_1 p_1} > \frac{S_2}{S_1} \frac{1}{\alpha} (\alpha - 1)^{\frac{\alpha - 1}{\alpha}}$$
(15)

By symmetry, a type CA equilibrium exists if and only if

$$\frac{C_2}{S_2 p_2} > \frac{S_1}{S_2} \frac{1}{\alpha} (\alpha - 1)^{\frac{\alpha - 1}{\alpha}}$$
(16)

and if it exists, then it is the unique type CA equilibrium.

An equilibrium of type AC or CA, if it exists, characterizes a situation where one of the two competing airlines stays out of the market (i.e., offers zero frequency) and the other offers the least frequency value necessary to provide seating capacity for all passengers in that market. Theorem 5.3 indicates that such an equilibrium exists if and only if there is no incentive for the first airline to offer nonzero frequency. In other words, the only condition necessary and sufficient for the existence of such an equilibrium is that the second airline has an equilibrium frequency that leads to negative profits for the first airline at all nonzero frequency values it can offer.

Theorem 5.4 A type BB equilibrium exists if and only if

$$k \le \left(\frac{1}{\alpha - 1}\right)^{\frac{1}{\alpha}}, \frac{1}{k} \le \left(\frac{1}{\alpha - 1}\right)^{\frac{1}{\alpha}}, \text{ and } \frac{C_1}{S_1 p_1} \le \alpha \frac{k^{\alpha}}{1 + k^{\alpha}}, \frac{C_2}{S_2 p_2} \le \alpha \frac{1}{1 + k^{\alpha}}$$
 (17)

where $k = \frac{C_1p_2}{C_2p_1}$, and if it exists then it is a unique type BB equilibrium. Proof In type BB equilibrium, $x^* > 0, y^* > 0, \wp_1 \le S_1x$, and $\wp_2 \le S_2y$. Therefore, $\Pi_1(x^*, y^*) = \Pi'_1(x^*, y^*)$ and $\Pi_2(x^*, y^*) = \Pi'_2(x^*, y^*)$. So Π_1 and Π_2 are both twice continuously differentiable at (x^*, y^*) . So type BB equilibrium exists if and only if there exist x and y such that $\frac{\partial \Pi'_1}{\partial x} = 0, \frac{\partial \Pi'_2}{\partial y} = 0, \frac{\partial^2 \Pi'_1}{\partial x^2} \le 0, \frac{\partial^2 \Pi'_2}{\partial y^2} \le 0,$ $\Pi'_1 \ge 0, \Pi'_2 \ge 0, M \frac{x^2}{x^2 + y^2} \le S_1x$, and $M \frac{y^2}{x^2 + y^2} \le S_2y$. Solving the two First-Order Conditions (FOCs) simultaneously, we get

$$x = \frac{\alpha M p_1}{C_1} \frac{k^{\alpha}}{(1+k^{\alpha})^2} \text{ and } y = \frac{\alpha M p_1}{C_1} \frac{k^{\alpha+1}}{(1+k^{\alpha})^2}$$
(18)

So if this equilibrium exists, then it must be the unique type BB equilibrium.

The second-order conditions (SOCs) can be simplified to $k \leq \left(\frac{\alpha+1}{\alpha-1}\right)^{\frac{1}{\alpha}}$ and $\frac{1}{k} \leq \left(\frac{\alpha+1}{\alpha-1}\right)^{\frac{1}{\alpha}}$. Also the $\Pi'_1 \geq 0$ and $\Pi'_2 \geq 0$ conditions translate into $k \leq \left(\frac{1}{\alpha-1}\right)^{\frac{1}{\alpha}}$ and $\frac{1}{k} \leq \left(\frac{1}{\alpha-1}\right)^{\frac{1}{\alpha}}$, which make the SOCs redundant. Finally, the last two conditions translate into $\frac{C_1}{S_1p_1} \leq \frac{\alpha k^2}{1+k^{\alpha}}$ and $\frac{C_2}{S_2p_2} \leq \frac{\alpha}{1+k^{\alpha}}$. Therefore, type BB equilibrium exists if and only if the following conditions are satisfied

$$k \le \left(\frac{1}{\alpha - 1}\right)^{\frac{1}{\alpha}}, \frac{1}{k} \le \left(\frac{1}{\alpha - 1}\right)^{\frac{1}{\alpha}}, \frac{C_1}{S_1 p_1} \le \alpha \frac{k^{\alpha}}{1 + k^{\alpha}} \text{ and } \frac{C_2}{S_2 p_2} \le \alpha \frac{1}{1 + k^{\alpha}}.$$
 (19)

A type BB equilibrium characterizes a situation with both competing airlines offering nonzero frequencies and both of them having excess capacities at equilibrium. Theorem 5.4 shows that this equilibrium exists if and only if, (1) the maximum profit corresponding to a nonzero frequency value is nonnegative, and

(2) the nonzero frequency value with maximum profit corresponds to sufficient seating capacity, for each of the competing airline carriers. This equilibrium, if it exists, corresponds to frequency values determined solely based on the S-curve-based frequency competition between the competing carriers and is not affected by the seating capacity constraints.

Theorem 5.5 A type BC equilibrium exists if and only if the following three conditions are true

$$\frac{C_1}{p_1 S_1} \frac{S_1}{S_2} \le (\alpha - 1)^{\frac{\alpha - 1}{\alpha}}$$
(20)

$$\frac{k^{\alpha}}{1+k^{\alpha}} < \frac{1}{\alpha}, \frac{1}{1+k^{\alpha}} < \frac{1}{\alpha} \frac{C_2}{p_2 S_2}$$
(21)

$$\frac{1}{\alpha} \frac{C_1}{p_1 S_1} \ge \frac{1}{1 + \left(\frac{S_1}{S_2}\right)^{\frac{\alpha}{\alpha - 1}}}$$
(22)

where $k = \frac{C_1p_2}{C_2p_1}$, and if it exists then it is a unique type BC equilibrium. *Proof* This proof involves slightly lengthier manipulations. So we have deferred it to the Appendix. Please refer to the detailed proof of Theorem A.1. By symmetry, a type CB equilibrium exists if and only if

$$\frac{C_2}{p_2 S_2} \frac{S_2}{S_1} \le (\alpha - 1)^{\frac{\alpha - 1}{\alpha}}, \frac{1}{1 + k^{\alpha}} < \frac{1}{\alpha}, \frac{k^{\alpha}}{1 + k^{\alpha}} < \frac{1}{\alpha} \frac{C_1}{p_1 S_1} \text{ and } \frac{1}{\alpha} \frac{C_2}{p_2 S_2} \ge \frac{1}{1 + \left(\frac{S_2}{S_1}\right)^{\frac{\alpha}{\alpha - 1}}}$$
(23)

and if it exists, then it is a unique type CB equilibrium.

A type BC or CB equilibrium, if it exists, characterizes two competing airlines, each offering a nonzero frequency, one with excess seating capacity and the other without excess seating capacity. Theorem 5.5 indicates that such an equilibrium exists if and only if, (1) the maximum profit corresponding to a nonzero frequency value for the first airline is nonnegative, (2) the nonzero frequency value with maximum profit for the first airline corresponds to sufficient seating capacity, and (3) the second airline is not able to increase its profit by providing a higher frequency with excess capacity.

Theorem 5.6 A type CC equilibrium exists if and only if

$$\left((S_2/S_1)^{\frac{\alpha}{\alpha-1}} \right) / \left(1 + (S_2/S_1)^{\frac{\alpha}{\alpha-1}} \right) < \frac{1}{\alpha} \frac{C_1}{S_1 p_1}$$
(24)

and

$$1 / \left(1 + (S_2/S_1)^{\frac{\alpha}{\alpha - 1}} \right) < \frac{1}{\alpha} \frac{C_2}{S_2 p_2}$$
(25)

and if it exists then it is a unique type CC equilibrium.

Proof For type CC equilibrium, x > 0, y > 0, $\varphi_1 = S_1 x$ and $\varphi_2 = S_2 y$. Existence of local maxima of Π_1 at $x = x^*$ requires that $\frac{\partial \Pi'_1}{\partial x} < 0$. Similarly, existence of local maxima of Π_2 at $y = y^*$ requires that $\frac{\partial \Pi'_2}{\partial y} < 0$. So, for a type CC equilibrium to exist at (x, y), the necessary and sufficient conditions to be satisfied are $\frac{x^a}{x^a + y^a}M = S_1 x$, $\frac{y^a}{x^a + y^a}M = S_2 y$, $\frac{\partial \Pi'_1}{\partial x} \le 0$ and $\frac{\partial \Pi'_2}{\partial y} \le 0$. Solving the two equalities simultaneously we get

$$x = \frac{M}{S_1 \left(1 + \left(\frac{S_2}{S_1}\right)^{\frac{\alpha}{\alpha-1}}\right)}$$
(26)

and

$$y = \frac{M}{S_2 \left(1 + \left(\frac{S_1}{S_2}\right)^{\frac{2}{\alpha-1}}\right)}$$
(27)

Therefore, if a type CC equilibrium exists, then it must be a unique type CC equilibrium. The two first-order inequality conditions translate into

$$\frac{\left(\frac{S_2}{S_1}\right)^{\frac{\alpha}{\alpha-1}}}{1+\left(\frac{S_2}{S_1}\right)^{\frac{\alpha}{\alpha-1}}} < \frac{1}{\alpha} \frac{C_1}{S_1 p_1}$$
(28)

and

$$\frac{1}{1 + \left(\frac{S_2}{S_1}\right)^{\frac{\alpha}{\alpha-1}}} < \frac{1}{\alpha} \frac{C_2}{S_2 p_2} \tag{29}$$

These two inequalities together are necessary and sufficient conditions for a type CC equilibrium to exist. $\hfill \Box$

Type CC equilibrium, if it exists, characterizes two competing airlines each providing nonzero frequency which is the least frequency that provides sufficient seating capacity for all passengers intending to fly with that airline. Theorem 5.6

indicates that such an equilibrium exists if and only if neither of the two competing airlines is able to increase its profit by providing a higher frequency with excess capacity.

Thus, depending on operating cost, fare, and seating capacity values, the twoplayer game can admit multiple (up to six) equilibria. But there can never be more than one equilibrium belonging to each of the six types, namely, AC, CA, BB, BC, CB, and CC. Also, no equilibria belonging to any of the remaining three types, namely, AA, AB, and BA can exist. Furthermore, all the necessary and sufficient conditions for the existence and uniqueness of each type of equilibrium can be expressed in terms of only five unitless parameters namely, $r_1 = \frac{C_1}{p_1 S_1}$, $r_2 = \frac{C_2}{p_2 S_2}$, $k = \frac{C_1 p_2}{r_1}$, $l = \frac{S_1}{S_2}$, and α , out of which *l* can be expressed as a function of the rest as $l = k \frac{r_2}{r_1}$. So there are only four independent parameters, which completely describe a two-player frequency game. Interestingly, the total passenger demand *M* plays no part in any of the conditions.

5.2 Realistic Ranges of Parameter Values

Up to six different pure strategy Nash equilibria may exist for a two-player game depending on game parameters. Apart from α , the flight operating costs, seating capacities, and fares are the only determinants of these parameters. In order to identify realistic ranges of these parameters, we looked at all the domestic segments in the U.S. with exactly two carriers providing nonstop service. We obtained the average operating cost per flight leg for each segment from the Form 41 financial database, the average flight seating capacities per segment from the T100 segments database, and the average fares for nonstop and connecting passengers on each segment from the DB1B database. All data was obtained from the Bureau of Transportation Statistics (BTS) website for the first quarter of 2007 [15]. There are 157 U.S. domestic segments with exactly two carriers providing nonstop service. This amounted to 314 combinations of carriers and segments. Many of these markets cannot be classified as pure duopoly situations because passenger demand on many of these origin-destination pairs is served not only by the nonstop itineraries, but also by connecting itineraries offered by several carriers, often including the two carriers providing the nonstop service. Moreover, one or both endpoints for many of these nonstop segments are important hubs of one or both of these nonstop carriers, which means that connecting passengers traveling on these segments also play an important role in the profitability of these segments. Therefore, modeling these nonstop markets as pure duopoly cases can be a gross approximation. Our aim is not to capture all these effects in our frequency competition model but rather to identify realistic relative values of flight operating costs, seating capacities, and fares. Despite these complications, these 157 segments are the real-world situations that come closest to the simplified frequency competition model that we have considered. Therefore, data from these markets were used to narrow down our



Fig. 3 Histograms of k, $\frac{s_1}{s_2}$, and $\frac{c}{ps}$

modeling focus. Figure 3a–c show the histograms of k, $\frac{S_1}{S_2}$, and $\frac{C}{pS}$ respectively. All k values were found to lie in the range 0.4–2.5, all $\frac{S_1}{S_2}$ were in the range 0.5–2, and all $\frac{C}{pS}$ values were found to lie in the range 0.18–0.8. We will restrict our further analysis to these ranges of values only. In particular, for later analysis, we will need only one of these assumptions, which is as follows.

Assumption 5.1 $0.4 \le k \le 2.5$

For $\alpha = 1.5$, the conditions for type BB equilibrium were satisfied in 144 out of these 157 markets, i.e., almost 92 % of the time. Conditions for type AC (or CA) equilibrium were satisfied in 71 markets, of which 8 were such that the conditions for both type AC and type CA equilibrium were satisfied together. Conditions for type BC (or CB) equilibrium were satisfied in only 1 out of 157 markets and conditions for type CC equilibrium were never satisfied. In all the markets, the conditions for the existence of at least one PSNE were satisfied. Out of 157 markets, almost 55 % (86 markets) were such that type BB was the unique PSNE.

We have already proved that AA, AB, and BA type equilibria do not exist. Further, as discussed above, AC, CA, BC, CB, and CC type equilibria are suspect in terms of portrayal of reality. Therefore, type BB equilibrium appears to be the most reasonable type of equilibrium. Indeed, the data analysis suggested that the existence conditions for type BB equilibrium were satisfied in most of the markets. So for the purpose of analyzing learning dynamics, we will only consider the type BB equilibrium.

Now, we propose two alternative dynamics for the nonequilibrium situations.

5.3 Myopic Best Response Dynamic

Consider an adjustment process where the two players take turns to adjust their own frequency decision so that each time it is the best response to the frequency chosen by the competitor in the previous period. If x_i and y_i are the frequency decisions by the two carriers in period *i*, then x_i is the best response to y_{i-1} and y_{i-1} is the best





response to x_{i-2} , etc. We will prove the convergence of this dynamic for two representative values of α , namely $\alpha = 1$ and $\alpha = 1.5$. We chose these two values because they correspond to two disparate beliefs about the market share-frequency share relationship. There is nothing specific about these two values that makes the algorithm converge. In fact, given any value in between, we would probably be able to construct a proof of convergence. But due to space constraints, we will restrict our attention to these two specific values of α .

Let us define $X = x^{\alpha}$ and $Y = y^{\alpha}$. We will often use the X - Y coordinate system in this section. Without any loss of generality, we assume that $k = \frac{C_1 p_2}{C_2 p_1} \leq 1$. We will denote the best response functions as $x_{BR}(y)$ and $y_{BR}(x)$ in the x - y coordinate system and as $X_{BR}(Y)$ and $Y_{BR}(X)$ in the X - Y coordinate system. Consider a twodimensional interval I (as shown in Fig. 4) given by $x_{Ib} \leq x \leq x_{ub}$ and $y_{Ib} \leq y \leq y_{ub}$ where $y_{ub} = (\alpha M p_2)/(4C_2)$, $x_{ub} = x_{BR}(y_{ub})$, $y_{Ib} = y_{BR}(x_{ub})$, and $x_{Ib} = x_{BR}(y_{Ib})$.

Theorem 5.7 As long as the competitor frequency for each carrier remains in region *B*, regardless of the starting point: (a) the myopic best response algorithm will reach some point in interval *I* in a finite number of iterations, (b) once inside interval *I*, it will never leave the interval.

Proof Please refer to the detailed proof of Theorem A.2 in the Appendix. \Box

Next we prove that the absolute value of slope of each of the best response curves inside interval *I* is less than 1 in the X - Y coordinates. We will prove this for two representative values of α namely, $\alpha = 1.5$ and $\alpha = 1$.

Theorem 5.8 For $\alpha = 1.5$, the absolute value of slope of each of the best response curves inside interval I is less than 1 in the X - Y coordinates.

Proof Please refer to the detailed proof of Theorem A.3 in the Appendix. \Box

Theorem 5.9 For $\alpha = 1$, the absolute value of slope of each of the best response curves inside interval I is less than 1 in the X - Y coordinates.

Proof Please refer to the detailed proof of Theorem A.4 in the Appendix. \Box

In order to prove the next theorem, we assume that the absolute value of slope of each of the best response curves is less than 1 in interval *I*.

Theorem 5.10 If the absolute value of slope of each of the best response curves is less than 1 in interval I, then as long as the competitor frequency for each carrier remains in region B, regardless of the starting point, the myopic best response algorithm converges to the unique type BB equilibrium.

Proof Please refer to the detailed proof of Theorem A.5 in the Appendix. \Box

Theorem 5.11 Regardless of the starting point, the myopic best response algorithm converges to the unique type BB equilibrium as long as the following conditions are satisfied:

$$\alpha M p_1/(4C_1) \le x_{\text{th}}, \ \alpha M p_2/(4C_2) \le y_{\text{th}}, \ x_{\text{cr}} \le x_{\text{BR}}(y_{\text{th}}), \ y_{\text{cr}} \le y_{\text{BR}}(x_{\text{th}}), \ x_{\text{cr}} \le x_{\text{BR}}(y_{\text{cr}}) \text{ and } y_{\text{cr}} \le y_{\text{BR}}(x_{\text{cr}}).$$
(30)

Proof First we develop sufficient conditions under which the competitor frequency for each carrier remains in region B for all iterations $i \ge 2$, regardless of the starting point.

As proved in the description of the best response curve in Sect. 4, the shape of the best response curve $y_{BR}(x)$ is such that at x = 0, $y = \frac{M}{S_2}$. Initially it is strictly decreasing followed by a point of nondifferentiability (at x_{cr}) beyond which it is strictly increasing until a local maximum is reached at $x = \frac{\alpha M p_2}{4C_2}$. Beyond the local maximum, it is strictly decreasing again up to a point of discontinuity (at x_{th}), beyond which it takes a constant value 0. For $x \leq x_{th}$, the only candidates for global minima of the best response curve $y_{BR}(x)$ are x_{cr} and x_{th} . The only candidates for global maxima are x = 0 and $x = \frac{\alpha M p_2}{4C_2}$. If the y-coordinate at each of these four important points lies in the range $y_{cr} \le y \le y_{th}$, then $y_{cr} \le y_{BR}(x) \le y_{th}$, for all $x \leq x_{\text{th}}$. Similarly, if $x_{\text{BR}}(y)$ at y = 0, $y = y_{\text{cr}}$, $y = \frac{\alpha M p_1}{4C_1}$, and $y = y_{\text{th}}$ are all in the range $x_{cr} \le x \le x_{th}$, then $x_{cr} \le x_{BR}(y) \le x_{th}$ for all $y \le y_{th}$. So for any starting point x_0 such that $x_0 \le x_{th}$, the algorithm will remain in the region B of both carriers for all subsequent iterations. The only remaining case is when $x > x_{th}$ or $y > y_{\text{th}}$. This does not pose any problem because for all $x > x_{\text{th}}$, $x_{\text{BR}}(y_{\text{BR}}(x)) =$ $x_{\rm BR}(0)$ and $x_{\rm cr} \le x_{\rm BR}(0) = \frac{M}{S_2} \le x_{\rm th}$. So if the aforementioned conditions are satisfied, then regardless of the starting point, the algorithm will remain in the region B of both carriers for all iterations *i* such that $i \ge 2$ (as per Theorem 5.7).

For all the aforementioned conditions to be satisfied, it is sufficient to ensure that the upper bound conditions on the points of local maxima are satisfied and the lower bound conditions on the points of local minima are satisfied. Let us first look at the upper bounds on the points of local maxima. There are four such conditions per carrier, namely $\frac{M}{S_1} \le x_{\text{th}}, \frac{M}{S_2} \le y_{\text{th}}, \frac{\alpha M p_1}{4C_1} \le x_{\text{th}}, \text{and } \frac{\alpha M p_2}{4C_2} \le y_{\text{th}}, \frac{M}{S_1} \le x_{\text{th}}$ simplifies to $\frac{C_2}{S_2 p_2} \le \left(\frac{S_1}{S_2}\right) \left(\frac{1}{\alpha}\right) (\alpha - 1)^{\frac{\alpha-1}{\alpha}}$, which is the exact negation of the condition for existence of type CA equilibrium. Because we have assumed that the unique PSNE in this game is a type BB equilibrium, a type CA equilibrium cannot exist. Hence this condition is automatically satisfied. By symmetry, due to the nonexistence of a type AC equilibrium, the condition $\frac{M}{S_2} \le y_{\text{th}}$ is automatically satisfied. The remaining six conditions are as follows: $\frac{\alpha M p_1}{4C_1} \le x_{\text{th}}, \frac{\alpha M p_2}{4C_2} \le y_{\text{th}}, x_{\text{cr}} \le x_{\text{BR}}(y_{\text{th}}), y_{\text{cr}} \le y_{\text{BR}}(x_{\text{th}}),$ $x_{\text{cr}} \le x_{\text{BR}}(y_{\text{cr}})$, and $y_{\text{cr}} \le y_{\text{BR}}(x_{\text{cr}})$. If each of these conditions is satisfied then, using Theorems 5.8, 5.9, and 5.10, the myopic best response algorithm converges to the unique type BB equilibrium, regardless of the starting point.

5.4 Alternative Dynamic

This dynamic is applicable only in the part of region B where the payoff function is strictly concave for both players' payoff functions, i.e., we will consider the region where $\frac{\alpha-1}{\alpha+1} + \epsilon \leq \left(\frac{x}{y}\right)^{\alpha} \leq \frac{\alpha+1}{\alpha-1} - \epsilon$, where ϵ is any sufficiently small positive number. This requirement is not very restrictive. This condition is always satisfied at the type BB equilibrium, since it is the second-order optimality condition. Moreover, the $\frac{x}{y}$ values satisfying this condition cover a large region surrounding the type BB equilibrium. For example, for $\alpha \to 1^+$, this condition is always satisfied for *all* values of $\frac{x}{y}$, while for $\alpha = 1.5$, the condition translates approximately to $0.342 \le \frac{x}{y} \le 2.924$, which is a large range. Given this restriction, in order to provide a complete specification of the player utilities, we will define the player *i* payoff outside this region by means of a quadratic function of a single variable x_i . The coefficients are such that $u_i(x_i)$ and its first- and second-order derivatives with respect to x_i are continuous. Note that the choice of a quadratic function is for convenience. Quadratic form is the simplest that can be used to define strictly concave functions. Any other strictly concave functional form would work equally well.

Multiplying the payoff function by a positive real number is an order preserving transformation, which does not affect the properties of the game. We will multiply the payoff of player *i* by $1/p_i$. So $u_i = \prod_i/p_i$. This dynamic was proposed by Rosen [27]. Under this dynamic, each player changes strategy such that the player's own payoff would increase if all the other players held to their current strategies. The rate of change of each player's strategy with time is equal to the gradient of the player's payoff with respect to the player's own strategy, subject to constraints. For the

frequency competition game, where each player's strategy space is 1-dimensional, the rate of change of each player's strategy simply equals the derivative of the player's payoff with respect to the frequency decision, subject to the upper and lower bound on allowable frequency values. Therefore, the rate of adjustment of each player's strategy is given by

$$\frac{\mathrm{d}x_i}{\mathrm{d}t} = \frac{\mathrm{d}u_i(x)}{\mathrm{d}x_i} + b_{\min} - b_{\max} \tag{31}$$

The only purpose of the b_{\min} and b_{\max} terms is to ensure that the frequency values stay within the allowable range, $x_{\min} \le x \le x_{\max} \cdot b_{\min}$ will be equal to 0 for all $x > x_{\min}$ and will take an appropriate positive value at $x = x_{\min}$ to ensure that the lower bound is respected. Similarly, b_{\max} will be equal to 0 for all $x < x_{\max}$ and will take an appropriate positive value at $x = x_{\min}$ to ensure that the upper bound is respected. Similarly, b_{\max} to ensure that the upper bound is respected. As long as the competitor frequencies remain in region B for each carrier, the utilities are given by:

$$u_1(x, y) = M \frac{x^{\alpha}}{x^{\alpha} + y^{\alpha}} - \frac{C_1}{p_1} x$$
(32)

and

$$u_2(x, y) = M \frac{y^{\alpha}}{x^{\alpha} + y^{\alpha}} - \frac{C_2}{p_2} y$$
(33)

The vector of payoff functions u(x, y) is given by:

$$u(x, y) = [u_1(x, y), u_2(x, y)]$$
(34)

The vector of first-order derivatives of each player's payoff with respect to the player's own frequency is given by:

$$\nabla u(x,y) = \left[\frac{\partial u_1(x,y)}{\partial x}, \frac{\partial u_2(x,y)}{\partial y}\right]$$
(35)

The Jacobian of ∇u is given by:

$$U(x,y) = \begin{pmatrix} \frac{\partial^2 u_1(x,y)}{\partial x^2} & \frac{\partial^2 u_1(x,y)}{\partial x \partial y} \\ \frac{\partial^2 u_2(x,y)}{\partial y \partial x} & \frac{\partial^2 u_2(x,y)}{\partial y^2} \end{pmatrix}$$
(36)

Under this dynamic, the frequencies of the competing carriers will converge to the unique type BB equilibrium frequencies if we can prove that $(U(x, y) + U^T(x, y))$ is negative definite (as per [27]). The first-order derivatives are given by

The Price of Airline Frequency Competition

$$\frac{\partial u_1(x,y)}{\partial x} = \frac{M\alpha x^{\alpha-1}y^{\alpha}}{\left(x^{\alpha} + y^{\alpha}\right)^2} - \frac{C_1}{p_1}$$
(37)

and

$$\frac{\partial u_2(x,y)}{\partial y} = \frac{M\alpha y^{\alpha-1} x^{\alpha}}{\left(x^{\alpha} + y^{\alpha}\right)^2} - \frac{C_2}{p_2}$$
(38)

and the second-order derivatives are given by

$$[U(x,y)]_{11} = \frac{\partial^2 u_1(x,y)}{\partial x^2} = \frac{M\alpha x^{\alpha-2} y^{\alpha}}{(x^{\alpha} + y^{\alpha})^3} ((\alpha - 1)y^{\alpha} - (\alpha + 1)x^{\alpha}) < 0$$
(39)

$$[U(x,y)]_{22} = \frac{\partial^2 u_2(x,y)}{\partial y^2} = \frac{M\alpha y^{\alpha-2} x^{\alpha}}{(x^{\alpha} + y^{\alpha})^3} ((\alpha - 1)x^{\alpha} - (\alpha + 1)y^{\alpha}) < 0$$
(40)

$$[U(x,y)]_{12} = \frac{\partial^2 u_1(x,y)}{\partial x \partial y} = \frac{M \alpha^2 x^{\alpha-1} y^{\alpha-1}}{(x^{\alpha} + y^{\alpha})^3} (x^{\alpha} - y^{\alpha})$$
(41)

$$[U(x,y)]_{21} = \frac{\partial^2 u_2(x,y)}{\partial y \partial x} = \frac{M \alpha^2 x^{\alpha-1} y^{\alpha-1}}{\left(x^{\alpha} + y^{\alpha}\right)^3} \left(y^{\alpha} - x^{\alpha}\right)$$
(42)

Therefore,

$$\left[U(x,y) + U^{T}(x,y)\right]_{11} = 2\left[U(x,y)\right]_{11}$$
(43)

$$\left[U(x,y) + U^{T}(x,y)\right]_{22} = 2\left[U(x,y)\right]_{22}$$
(44)

and

$$\left[U(x,y) + U^{T}(x,y)\right]_{12} = \left[U(x,y) + U^{T}(x,y)\right]_{21} = 0.$$
(45)

Therefore, $(U(x, y) + U^T(x, y))$ is a diagonal matrix with both diagonal elements strictly negative. Therefore, $(U(x, y) + U^T(x, y))$ is negative definite. This is sufficient to prove that the payoff functions are diagonally strictly concave [27]. Therefore, under the alternative dynamic mentioned above, the frequencies of the competing carriers will converge to the unique type BB equilibrium frequencies.

In this section, we showed that two different myopic dynamics converge to the type BB equilibrium in case of the two-player game. This finding has important practical implications. In the airline industry, year after year most airlines operate flights on similar sets of segments. The set of competitors and the general properties of the markets remain stable in most cases over long periods of time. Therefore, airlines have opportunities to adapt their decisions primarily by fine tuning the

frequency values. We capture these adjustments by modeling the dynamics of the game. Our results show that such adjustments bring the airline decisions closer to the equilibrium decisions with each additional step. This in turn means that if there aren't any significant changes in market properties for a few consecutive planning cycles (e.g., for a period of several months to a few years), then the frequencies are expected to be very close to the equilibrium frequencies in such markets.

6 Impact on Airport Congestion and Airline Profitability

As explained earlier, one of the main objectives of this chapter is to establish the connection between frequency competition and airport congestion, and to characterize the effect of the parameters of frequency competition on airport congestion and airline profitability. In order to achieve this, we will use the concept of Price of Anarchy. Koutsoupias and Papadimitriou [23] initially suggested the idea of using the ratio of the total cost of the worst-case Nash equilibrium to the total cost of the system optimal solution as a measure of efficiency degradation due to selfish behavior of independent agents, although they did not use the term Price of Anarchy. They provided upper and lower bounds on this worst-case ratio in the context of efficiency degradation due to network congestion caused by selfish routing by internet users. Other researchers later extended these results to more complex games. There exists a large body of research characterizing the degree of inefficiency in the context of decision-making by a very large number of selfinterested agents, each with infinitesimally small size [28] (also known as *atomistic* agents). In this section, we will use the same basic concept, but apply it to the case of congestion caused by self-interested nonatomistic agents such as airlines. Note that there is a large body of research that analyzes the competition between airlines treating them as self-interested nonatomistic decision makers (e.g., [14, 17]). However, none of these existing studies provide insights into the price of anarchy for such games.

So far, our analysis has been restricted only to two-player games. The impact of frequency competition on airport congestion is likely to be closely related to the number of competing players. So we now extend our analysis to an N-player symmetric game where N is any integer greater than 1. As shown in our analysis in Sect. 5, even in the two-player case, the number of Nash equilibria can vary between 0 and 6 depending on the combination of parameter values. The number of equilibria in frequency competition games with more players can be very high. In real-life airline markets, the parameters of airline frequency competition, such as, fares, seating capacities, and operating costs of competing airlines are often not too different from each other. Therefore focusing on the symmetric player case is not that unrealistic. Furthermore, as shown below, a thorough analysis of the symmetric player case and then quantify the price of anarchy for airline frequency competition game in Sect. 6.2.

6.1 N-Player Symmetric Game

Now we extend the analysis to the N-player symmetric case, where $N \ge 2$. By symmetry, we mean that the operating cost C_i , the seating capacity S_i and the fare p_i is the same for all carriers. For the analysis presented in this subsection, it is sufficient to have $\frac{C_i}{p_i}$ constant for all carriers. However, for computing the price of anarchy in the next subsection, we need the remaining assumptions. We will simplify the notation and denote the operating cost for each carrier as C, seating capacity as S and fare as p. Under symmetry, the necessary and sufficient conditions for the existence of a type BB equilibrium for a two-player game reduce to a single condition, which is as follows.

Assumption 6.1 $\frac{\alpha pS}{C} > 2$

We will assume that this condition holds throughout the following analysis. Note that we will use the phrase "excess seating capacity" or simply "excess capacity" to describe airlines whose effective competitor frequency is in region B.

Theorem 6.1 In an N-player symmetric game, a symmetric equilibrium with excess seating capacity exists at $x_i = \frac{\alpha M p}{C} \frac{N-1}{N^2}$ for all *i* if and only if $N \leq \frac{\alpha}{\alpha-1}$ and if it exists, then it is the unique symmetric equilibrium.

Proof Please refer to the detailed proof of Theorem A.6 in the Appendix. \Box

Theorem 6.2 In a symmetric N-player game, there exists no asymmetric equilibrium where all players have a nonzero frequency and excess seating capacity.

Proof Let us assume the contrary. For a symmetric N-player game, let there exist an asymmetric equilibrium such that all players have a nonzero frequency and excess seating capacity. Let us define $\beta = \sum_{j=1}^{N} x_j^{\alpha}$ and $\omega_i = \frac{x_i^{\alpha}}{\sum_{j=1}^{N} x_j^{\alpha}}$. So

$$x_i = (\omega_i \beta)^{1/\alpha} \tag{46}$$

Substituting in the FOC, we get

$$\frac{C}{\alpha M p} \beta^{1/\alpha} = \omega_i^{\frac{\alpha-1}{\alpha}} - \omega_i^{\frac{2\alpha-1}{\alpha}}$$
(47)

Let us define a function $h(\omega_i) = \omega_i^{\frac{\alpha-1}{\alpha}} - \omega_i^{\frac{2\alpha-1}{\alpha}}$. The value of $h(\omega_i)$ is the same across all the players at equilibrium. For all $\omega_i > 0$, $h(\omega_i)$ is a strictly concave function. So it can take the same value at at most two different values of ω_i . So all ω_i s can take at most two different values. Let $\omega_i = v_1$ for *m* (such that $1 \le m < N$) players, and $\omega_i = v_2$ for the remaining N-m players. Let $v_1 > v_2$, without any loss of generality. $h(\omega_i)$ is maximized at $\omega_i = \frac{\alpha-1}{2\alpha-1}$. So $v_2 < \frac{\alpha-1}{2\alpha-1} < v_1$.

At equilibrium, each player's profit must be nonnegative. The profit for each player *i* such that $\omega_i = v_2$ is given by $Mp\omega_i - Cx_i$. But $x_i = \frac{\alpha Mp}{C}\omega_i(1 - \omega_i)$. So the

condition on nonnegativity of profit simplifies to $v_2 \ge \frac{\alpha-1}{\alpha}$. Therefore, $\frac{\alpha-1}{2\alpha-1} > v_2 \ge \frac{\alpha-1}{\alpha}$, which can be true only if $\alpha < 1$. This leads to a contradiction. So we have proved that for a symmetric N-player game, there exists no asymmetric equilibrium such that all players have a nonzero frequency and excess seating capacity.

Theorem 6.3 In a symmetric N-player game, there exists some n_{\min} such that for any integer n with $\max(2, n_{\min}) \le n \le \min(N - 1, \frac{\alpha}{\alpha - 1})$ there exist exactly $\binom{N}{n}$ asymmetric equilibria such that exactly n players have nonzero frequency and all players with nonzero frequency have excess seating capacity. There exists at least one such integer for $N \ge \frac{\alpha}{\alpha - 1}$. The frequency of each player with nonzero frequency equals $\frac{\alpha Mp}{n-1} \frac{n-1}{n^2}$.

Proof Please refer to the detailed proof of Theorem A.7 in the Appendix. \Box

From here onward, we will denote each such equilibrium as an *n*-symmetric equilibrium of an N-player game.

Theorem 6.4 Among all equilibria with exactly n players ($n \le N$) having nonzero frequency, the total frequency is maximum for the symmetric equilibrium.

Proof Please refer to the detailed proof of Theorem A.8 in the Appendix. \Box

Theorem 6.5 There exists no equilibrium with exactly n players with nonzero frequency such that $n > \frac{\alpha}{\alpha-1}$.

Proof As per Theorem 6.1, if $n > \frac{\alpha}{\alpha-1}$, there exists no equilibrium with all *n* players having excess capacity. We have also proved that the number of players without excess capacity can be at most one. So consider some equilibrium with one player without excess capacity. Let the market share of that player be *l* and let the equilibrium frequency of each of the remaining players be x_2 . Because of $n > \frac{\alpha}{\alpha-1}$, we get $\alpha > \frac{n}{n-1}$. For nonnegative profit at equilibrium we require, $\frac{Mp}{C} \frac{1-l}{n-1} \ge x_2$. From the FOC, we get $x_2 = \frac{\alpha Mp}{C} \frac{1-l}{n-1} \left(1 - \frac{1-l}{n-1}\right)$. Combining the two we get $1 \ge \alpha \left(1 - \frac{1-l}{n-1}\right)$. Thus,

$$\frac{n}{n-1} < \alpha \le \frac{n-1}{n+l-2} \le \frac{n-1}{n-1.5}$$
(48)

For this to be true, we need n < 2, which is impossible. Therefore, we have proved that there exists no equilibrium with exactly *n* players with nonzero frequency such that $n > \frac{\alpha}{\alpha-1}$.

In this subsection, we proved that for an N-player symmetric game, if $N \le \frac{\alpha}{\alpha-1}$, then there exists a fully symmetric equilibrium where the equilibrium frequency of each carrier at equilibrium is $\frac{\alpha Mp}{C} \frac{N-1}{N^2}$ and there exists no asymmetric equilibrium

with all *N* players having a nonzero frequency. On the other hand, if $N > \frac{\alpha}{\alpha-1}$, then there exists no equilibrium with all players having nonzero frequency. In either case, there exist exactly $\binom{N}{n}$ n-symmetric equilibria for each integer n < N such that $\max(2, n_{\min}) \le n \le \min(N-1, \frac{\alpha}{\alpha-1})$ for some $n_{\min} \ge 0$. Additionally, there may be asymmetric equilibria such that each asymmetric equilibrium has exactly one player with 100 % load factor, n - 1 more players with nonzero frequency and excess seating capacity and N - n players with zero frequency. We also proved that there always exists at least one equilibrium for an N-player symmetric game. The aforementioned types of equilibria are exhaustive, that is, there exist no other types of equilibria. As before, we realize that all the equilibria except for those where all players have a nonzero frequency and excess capacity are suspect in terms of their portrayal of reality. So the fully symmetric equilibrium appears to be the most realistic one. In addition, the fully symmetric equilibrium is also the worst-case equilibrium in the sense that it is the equilibrium which has the maximum total frequency, as will be apparent in the next section.

We proved that for some n' < N, if there exists no symmetric equilibrium for any $n \ge n'$, then there exists no asymmetric equilibrium for any $n \ge n'$ either. We also proved that for any given n, the total frequency at each asymmetric equilibrium having n nonzero frequency players is, at most, equal to the total frequency at the corresponding n-symmetric equilibrium. These results will help us obtain the price of anarchy in the next subsection.

6.2 Price of Anarchy

In any equilibrium, the total revenue earned by all carriers remains equal to Mp. The total flight operating cost to all carriers is given by $\sum_{i=1}^{N} (Cx_i) = C \sum_{i=1}^{N} x_i$. On the other hand, if there were a central controller trying to minimize total operating cost, the minimum number of flights for carrying all the passengers would be equal to $\frac{M}{S}$ and the total operating cost would be $\frac{MC}{S}$. Similar to the notion introduced by Koutsoupias and Papadimitriou [23], let us define the price of anarchy as the ratio of total operating cost at Nash equilibrium to the total operating cost under the optimal frequency. The denominator is a constant and the numerator is proportional to the total number of flights. Please note that the way we have defined price of anarchy, the denominator does not refer to a situation that is likely to be favored by the passengers. In particular, passenger welfare is not accounted for in our expression in the denominator. In fact, similar to the rest of the analysis, the optimization is performed from the perspective of the airlines alone. The denominator represents a situation where congestion is minimized (while ensuring that all passengers are carried), the total profitability of the airlines is maximized, and the total cost of carrying the passengers is also minimized. Thus, our measure of price of anarchy compares the worst-case Nash equilibrium to this congestion minimizing

(total profitability maximizing, and total passenger carrying cost minimizing) solution.

A large proportion of airport delays are caused by congestion. Congestion related delay at an airport is an (often nonlinearly) increasing function of the total number of flights. Therefore, the greater the total number of flights, the greater is the delay. Total profit earned by all the airlines in a market is also a decreasing function of the total frequency. Also, because the total number of passengers remains constant, the average load factor in a market is inversely proportional to the total frequency. Lower load factors mean more wastage of seating capacity. Thus total frequency is a good measure of airline profitability, total operating cost, airport congestion, and load factors. Higher total frequency across all carriers in a market means lower profitability, more cost, more congestion, and lower average load factor, assuming constant aircraft size. The greater the price of anarchy, the greater is the inefficiency introduced by the competitive behavior of players at equilibrium.

Theorem 6.6 In a symmetric N-player game, the price of anarchy is given by $\frac{\alpha pS}{C} \frac{n-1}{n}$, where n is the largest integer not exceeding $\min(N, \frac{\alpha}{\alpha-1})$.

Proof As per Theorem 6.3, a symmetric N-player game has $\sum_{\max(2,n_{\min})}^{\min(N,\frac{\pi}{\alpha-1})} \binom{N}{n}$ equilibria (for some $n_{\min} \ge 0$), such that each equilibrium has a set of exactly n players each with frequency $\frac{\alpha Mp n-1}{C n^2}$ and excess capacity, whereas remaining N - n players have zero frequency. Also, for any $n < \min(N, \frac{\pi}{\alpha-1})$, there may exist equilibria with exactly n players having nonzero frequency and one of them not having any excess capacity at equilibrium. However, the frequency under any equilibrium with exactly n players having nonzero frequency is at most equal to the corresponding n-symmetric equilibrium. In any equilibrium having n players with nonzero frequency and excess capacity, the total flight operating cost is given by $\alpha Mp \frac{n-1}{n}$, which is an increasing function of n. The total cost under minimum cost scheduling is $\frac{\alpha p S}{C} \frac{n-1}{n}$, which is an increasing function of n. Also, no equilibrium exists for $n > \frac{\alpha}{\alpha-1}$. Therefore, the price of anarchy is given by $\frac{\alpha p S}{\alpha-1}$, where n is the greatest integer less than or equal to min $(N, \frac{\alpha}{\alpha-1})$.

This expression has several important implications. The greater the α value, the greater is the price of anarchy. This means that as the market share-frequency share relationship becomes more and more curved, and goes away from the straight line, the price of anarchy is greater. So the S-curve phenomenon has a direct impact on airline profitability and airport congestion. Also, the greater the fare compared to the operating cost per seat (i.e., the greater is the value of $\frac{pS}{C}$), the greater is the price of anarchy. In other words, for short-haul (low *C*), high-fare (high *p*) markets the price of anarchy is greater. Finally, as the number of competitors increases the price of anarchy increases (up to a threshold value beyond which it remains constant).

The equilibrium results from this simple model help substantiate some of the claims mentioned earlier. The price of anarchy increases because of the S-shaped (rather than linear) market share-frequency share relationship. Therefore, similar to the suggestions by Button and Drexler [16] and O'Connor [25], the S-curve relationship tends to encourage airlines to provide excess capacity and schedule greater numbers of flights. The total profitability of all the carriers in a market under the worst-case equilibrium provides a lower bound on airline profitability under competition. This lower bound is an increasing function of the price of anarchy, which in turn increases with the number of competitors. Therefore, similar to Kahn's [22] argument, this raises the question of whether the objectives of a financially strong and highly competitive airline industry are inherently conflicting. In addition, these results also establish the link between airport congestion and airline competition. Airport congestion under the worst-case equilibrium is directly proportional to the price of anarchy. So the greater the number of competitors and the greater the curvature of the market share-frequency share relationship, the greater are airport congestion and delays.

7 Summary

In this chapter, we modeled airline frequency competition based on the S-curve relationship which has been well documented in the airline literature. Regardless of the exact value of the α parameter, it is usually agreed that market share is an increasing (linear or S-shaped) function of frequency share. Our model is general enough to accommodate both a linear and an S-shaped market share-frequency share relationship. We characterize the best response curves for each player in a multiplayer game. Due to the complicated shape of best response curves, we proved that there exist anywhere between 0 and 6 different equilibria depending on the exact parameter values, for a two-player game. All the existence and uniqueness conditions can be completely described by three unitless parameters in addition to α of the two-player game. Only one out of the six possible equilibria seemed reasonable in portraying reality. This equilibrium corresponds to both players having nonzero frequency and less than 100 % load factors. In order to narrow down the modeling effort, realistic ranges of parameter values were identified based on realworld data that come closest to the simplified models analyzed in this chapter. We proposed two different myopic learning algorithms for the two-player game and proved that under mild conditions, either of them converges to Nash equilibrium. For the N-player (for any integer $N \ge 2$) game with identical players, we characterized the entire set of possible equilibria and proved that at least one equilibrium always exists for any such game. The worst-case equilibrium was identified. The price of anarchy was found to be an increasing function of the number of competing airlines, the ratio of fare to operating cost per seat, and the curvature of the S-curve relationship.

In this chapter, we presented two central results related to the impact of competition on airport congestion and airline profitability. First, there are simple myopic learning rules under which less than perfectly rational players would converge to an equilibrium. This substantiates the predictive power of the Nash equilibrium concept. Second, the S-curve relationship between market share and frequency share has direct and adverse implications to airline profitability and airport congestion, as speculated in multiple previous studies. Thus, airline frequency competition has a profound effect on airline profitability and airport congestion, which in turn has important implications for aviation safety.

Note that, in this chapter, we have assumed flight frequency to be the main decision variable under consideration when analyzing airline competition, owing to the direct connection of frequency decisions with congestion and delays, which is the main focus of this chapter. However, other decisions including fares, aircraft sizes and network structures affect, and are affected by, frequency decisions. Network structures and fleet planning decisions are typically made before making frequency decisions, while fleet assignment and fare decisions are typically made after making frequency decisions. Thus, the decisions about network structures (e.g., hub locations) and fleet planning serve as inputs to the frequency planning process. For example, the existence of a hub at the origin and/or destination of a flight leg clearly affects the leg passenger demand; and available fleet types affect the seating capacities and operating costs, which are input parameters to our models. On the other hand, fares and fleet assignment decisions, which are made at a later stage of airline planning, need to account for the frequency decisions. A promising future research direction is to analyze the impacts of airline competition with an augmented set of decision variables that account for the interrelationship between frequency decisions and airline decisions at other stages of the planning process.

Finally, note that this analysis does not explicitly quantify the net effect of frequency competition on passenger welfare. While frequency competition has resulted in the availability of more options to travel, it has also resulted in higher delays and disruptions to passengers. Explicit modeling and evaluation of costs and benefits of airline frequency competition to the passengers is beyond the scope of this chapter and can serve as a useful next step.

Acknowledgments We thank Professor Asuman Ozdaglar from the Department of Electrical Engineering and Computer Science at the Massachusetts Institute of Technology for her valuable comments and guidance during this research and while writing this chapter.

A.1 Appendix

Theorem A.1 A type BC equilibrium exists if and only if the following three conditions are true

$$\frac{C_1}{p_1 S_1} \frac{S_1}{S_2} \le (\alpha - 1)^{\frac{\alpha - 1}{\alpha}}$$
(49)

$$\frac{k^{\alpha}}{1+k^{\alpha}} < \frac{1}{\alpha}, \frac{1}{1+k^{\alpha}} < \frac{1}{\alpha} \frac{C_2}{p_2 S_2}$$
(50)

$$\frac{1}{\alpha} \frac{C_1}{p_1 S_1} \ge \frac{1}{1 + \left(\frac{S_1}{S_2}\right)^{\frac{\alpha}{\alpha - 1}}}$$
(51)

where $k = \frac{C_{1P2}}{C_{2P1}}$, and if it exists then it is a unique type BC equilibrium (Same as Theorem 5.5 in the main text).

Proof In type BC equilibrium, $x^* > 0$, $y^* > 0$, $\wp_1 \le S_1 x$, $\wp_2 = S_2 y$, and $\Pi_1 = \Pi'_1$. So Π_1 is twice continuously differentiable at (x^*, y^*) . For local maxima of Π_2 at (x^*, y^*) , we need $\Pi'_2 = \Pi''_2$ and $\frac{\partial \Pi'_2}{\partial y} < 0$. A type BC equilibrium then exists if and only if there exists (x, y) such that $\frac{\partial \Pi'_1}{\partial x} = 0$, $\Pi'_2 = \Pi''_2$, $\frac{\partial^2 \Pi'_1}{\partial x^2} \le 0$, $\frac{\partial \Pi'_2}{\partial y} < 0$, $\Pi'_1 \ge 0$, and $M \frac{x^2}{x^2 + y^2} \le S_1 x$. The first two conditions translate into

$$\frac{x^{\alpha-1}y^{\alpha}}{\left(x^{\alpha}+y^{\alpha}\right)^{2}} = \frac{C_{1}}{\alpha M p_{1}}$$
(52)

and

$$\frac{y^{\alpha}}{x^{\alpha} + y^{\alpha}} = \frac{S_2}{M}y \tag{53}$$

Solving Eqs. (52) and (53) simultaneously, we get

$$x = \left(\frac{MC_1}{\alpha p_1 S_2^{-2}}\right)^{\frac{1}{\alpha - 1}} y^{\frac{\alpha - 2}{\alpha - 1}}$$
(54)

$$\left(\frac{yS_2}{M}\right)^{\frac{1}{\alpha-1}} - \left(\frac{yS_2}{M}\right)^{\frac{\alpha}{\alpha-1}} = \left(\frac{C_1}{\alpha p_1 S_2}\right)^{\frac{\alpha}{\alpha-1}}$$
(55)

The nonnegativity condition on airline 1's profit implies that $Mp_1 \frac{x^{\alpha}}{x^{\alpha}+y^{\alpha}} \ge C_1 x$. Substituting Eqs. (54) and (55), we get

$$\frac{yS_2}{M} \le \frac{1}{\alpha} \tag{56}$$

The LHS of Eq. (55) is a strictly increasing function of *y* for $\frac{yS_2}{M} < \frac{1}{\alpha}$. Therefore, there exists some *y* that satisfies Eq. (55) and inequality (56) if and only if

$$\left(\frac{1}{\alpha}\right)^{\frac{1}{\alpha-1}} - \left(\frac{1}{\alpha}\right)^{\frac{\alpha}{\alpha-1}} - \left(\frac{C_1}{\alpha p_1 S_2}\right)^{\frac{\alpha}{\alpha-1}} \ge 0$$
(57)

i.e., if and only if

$$\frac{C_1}{p_1 S_1} \frac{S_1}{S_2} \le (\alpha - 1)^{\frac{\alpha - 1}{\alpha}}$$
(58)

and if it exists, then it is unique. Therefore, if a type BC equilibrium exists, then it must be a unique type BC equilibrium. Simplifying the second-order condition (SOC) and substituting Eqs. (54) and (55), we get

$$\frac{yS_2}{M} \le \frac{\alpha + 1}{2\alpha} \tag{59}$$

Therefore, inequality (56) makes SOC redundant. First-order condition (FOC) on $\Pi'_2(y)$ simplifies to

$$\frac{x}{y} < \frac{C_2 p_1}{C_1 p_2}$$
 (60)

Substituting Eqs. (54) and (55), we get

$$\frac{yS_2}{M} > \frac{k^{\alpha}}{1+k^{\alpha}} \tag{61}$$

Therefore, there exists a y that satisfies Eq. (55), inequality (56), and inequality (61) if and only if the following three inequalities are satisfied.

$$\frac{k^{\alpha}}{1+k^{\alpha}} < \frac{1}{\alpha} \tag{62}$$

$$\left(\frac{C_1}{\alpha p_1 S_2}\right)^{\frac{\alpha}{\alpha-1}} > \left(\frac{k^{\alpha}}{1+k^{\alpha}}\right)^{\frac{1}{\alpha-1}} - \left(\frac{k^{\alpha}}{1+k^{\alpha}}\right)^{\frac{\alpha}{\alpha-1}} \left(\text{ which is equivalent to } \frac{k^{\alpha}}{1+k^{\alpha}} < \frac{1}{\alpha}\right)$$
(63)

$$\frac{1}{1+k^{\alpha}} < \frac{1}{\alpha} \frac{C_2}{p_2 S_2} \tag{64}$$

Finally, the last condition, i.e., the condition that the seating capacity cannot be exceeded by the number of passengers for airline 1, simplifies to
The Price of Airline Frequency Competition

$$\frac{x^{\alpha-1}}{y^{\alpha-1}} \le \frac{S_1}{S_2} \tag{65}$$

Substituting Eq. (54), we get

$$\frac{yS_2}{M} \ge \frac{C_1}{\alpha p_1 S_1} \tag{66}$$

Combining with inequality (56), we get

$$\frac{1}{\alpha} \ge \frac{yS_2}{M} \ge \frac{C_1}{\alpha p_1 S_1} \tag{67}$$

Therefore, there exists some y that satisfies Eq. (55), inequality (56), and inequality (66) if and only if

$$\left(\frac{C_1}{\alpha p_1 S_2}\right)^{\frac{\alpha}{\alpha-1}} \ge \left(\frac{C_1}{\alpha p_1 S_1}\right)^{\frac{1}{\alpha-1}} - \left(\frac{C_1}{\alpha p_1 S_1}\right)^{\frac{\alpha}{\alpha-1}} \Leftrightarrow \frac{1}{\alpha} \frac{C_1}{p_1 S_1} \ge \frac{1}{1 + \left(\frac{S_1}{S_2}\right)^{\frac{\alpha}{\alpha-1}}} \tag{68}$$

Therefore, type BC equilibrium exists if and only if conditions (45), (50), (51), and (55) are satisfied. \Box

Theorem A.2 As long as the competitor frequency for each carrier remains in region *B*, regardless of the starting point: (a) the myopic best response algorithm will reach some point in interval *I* in a finite number of iterations, (b) once inside interval *I*, it will never leave the interval (Same as Theorem 5.7 in the main text).

Proof Let us denote the frequency decisions of the two carriers after the *i*th iteration by x_i and y_i respectively. At the beginning of the algorithm the frequency values are arbitrarily chosen to be x^0 and y^0 . If $i \ge 0$ is odd, then $x^i = x_{BR}(y^{i-1})$ and $y^i = y^{i-1}$. If $i \ge 0$ is even, then $y^i = y_{BR}(x^{i-1})$ and $x^i = x^{i-1}$. Therefore, for all $i \ge 2$, x_i is a best response to some y and y^i is a best response to some x. Best response curve $x_{BR}(y)$ in region B has a unique maximum at $y = \frac{\alpha M p_1}{4C_1}$ with $x_{BR}\left(\frac{\alpha M p_1}{4C_1}\right) = \frac{\alpha M p_1}{4C_1}$. By symmetry, the best response curve $y_{BR}(x)$ in region B has a unique maximum at $x = \frac{\alpha M p_2}{4C_2}$ with $y_{BR}\left(\frac{\alpha M p_2}{4C_2}\right) = \frac{\alpha M p_2}{4C_2}$. $k \le 1$ implies that $\frac{\alpha M p_2}{4C_2} \le \frac{\alpha M p_1}{4C_1}$. Therefore, $y^i \le \frac{\alpha M p_2}{4C_2} = y_{ub}$ for all $i \ge 2$. $\frac{\partial x_{BR}}{\partial y} > 0$ for $y < \frac{\alpha M p_1}{4C_1}$. Therefore, for all odd $i \ge 3$, $x^i = x_{BR}(y^{i-1})$ is $x_{BR}(y_{ub}) = x_{ub}$. So for all $i \ge 3$, $y^i \le y_{ub}$, and $x^i \le x_{ub}$.

Let us now prove that the type BB equilibrium point (x_{eq}, y_{eq}) is contained inside interval *I*. y_{eq} is a best response to x_{eq} . Therefore, $y_{eq} \leq \frac{\alpha M p_2}{4C_2} = y_{ub}$. For $k \leq 1$,

$$x_{\rm eq} = \frac{\alpha M p_2}{4C_2} \frac{4k^{\alpha - 1}}{(1 + k^{\alpha})^2} \ge \frac{\alpha M p_2}{4C_2}$$
(69)

and

$$y_{\rm eq} = \frac{\alpha M p_1}{4C_1} \frac{4k^{\alpha+1}}{(1+k^{\alpha})^2} \le \frac{\alpha M p_1}{4C_1}$$
(70)

 $\frac{\partial x_{BR}}{\partial y} \ge 0$ for all $y_{eq} \le y \le y_{ub}$. As a result,

$$x_{\rm eq} = x_{\rm BR}(y_{\rm eq}) \le x_{\rm BR}(y_{\rm ub}) = x_{\rm ub}$$

$$\tag{71}$$

 $\frac{\partial y_{BR}}{\partial x} \leq 0$ for all $x_{eq} \leq x \leq x_{ub}$. As a result,

$$y_{\rm eq} = y_{\rm BR}(x_{\rm eq}) \ge y_{\rm BR}(x_{\rm ub}) = y_{\rm lb}$$
(72)

 $\frac{\partial x_{\text{BR}}}{\partial y} \ge 0$ for all $y_{\text{lb}} \le y \le y_{\text{eq}}$. As a result,

$$x_{\rm eq} = x_{\rm BR}(y_{\rm eq}) \ge x_{\rm BR}(y_{\rm lb}) = x_{\rm lb}$$

$$\tag{73}$$

Thus, we have proved that $x_{lb} \le x_{eq} \le x_{ub}$, and $y_{lb} \le y_{eq} \le y_{ub}$, that is, the type BB equilibrium is contained inside interval *I*.

Because of existence of a unique type BB equilibrium, the best response curves intersect each other at exactly one point denoted by (x_{eq}, y_{eq}) . Further, for all $x < x_{eq}$ and for all $y < y_{eq}$, the y_{BR} curve is above the x_{BR} curve and x_{BR} curve is to the right of y_{BR} curve. Also, for all $y < y_{eq}, x_{BR}(y) < x_{eq}$. Therefore, for all $x^i < x_{eq}$, if *i* is odd then $x^{i+1} = x^i$, $y^i < y^{i+1} \le y_{ub}$ and if *i* is even then $x^i < x^{i+1} < x_{eq}$, $y^{i+1} = y^i$. So in each iteration, either x^i or y^i keeps strictly increasing until $y^i \ge y_{eq}$. (Note that the two curves, $x_{BR}(y)$ and $y_{BR}(x)$, never get asymptotically close to each other. So step sizes in each iteration can be easily shown to be lower bounded by a positive number). In the very next iteration, $x^{i+1} = x_{BR}(y^i) \ge x_{eq}$ and $y^{i+1} = y^i \ge y_{eq}$. Thus $x_{Ib} \le x_{eq} \le x^{i+1} \le x_{ub}$ and $y_{Ib} \le y_{eq} \le y^{i+1} \le y_{ub}$. Thus we have proved part (a) of the theorem.

We have already proved that at the end of any iteration $i \ge 2$, $x^i \le x_{ub}$ and $y^i \le y_{ub}$. So for all *i* such that $x_{lb} \le x^i \le x_{ub}$ and $y_{lb} \le y^i \le y_{ub}$, all that remains to be proved is that $x_{lb} \le x^{i+1}$ and $y_{lb} \le y^{i+1}$. We first consider the case where *i* is even. $y^{i+1} = y^i$. Because $\frac{\partial x_{BR}}{\partial y} > 0$ for $y < \frac{\alpha M p_1}{4C_1}$ and $y_{ub} \le \frac{\alpha M p_1}{4C_1}$, therefore $\frac{\partial x_{BR}}{\partial y} \ge 0$ for all *y* such that $y_{lb} \le y \le y_{ub}$. Therefore,

$$y_{lb} \le y^i \le y_{ub} \Rightarrow x_{lb} = x_{BR}(y_{lb}) \le x_{BR}(y^i) = x^{i+1} \le x_{BR}(y_{ub}) = x_{ub}$$
 (74)

Therefore, $x_{lb} \leq x^{i+1} \leq x_{ub}$ and $y_{lb} \leq y^{i+1} \leq y_{ub}$.

Now consider the case where *i* is odd, that is, $x^{i+1} = x^i$. For all x^i such that $x_{eq} \le x^i \le x_{ub}, \frac{\partial y_{BR}}{\partial x} \le 0$. Therefore, $y_{lb} = y_{BR}(x_{ub}) \le y_{BR}(x^i) = y^{i+1}$. On the other hand, for all $x^i < x_{eq}$, $y^i < \frac{\alpha M p_1}{4C_1}$, $y^{i+1} = y_{BR}(x^i) > y^i \ge y_{lb}$. Therefore, if $x_{lb} \le x^i \le x_{ub}$, then $y_{lb} \le y^{i+1}$. Thus we have proved that $x_{lb} \le x^{i+1} \le x_{ub}$ and $y_{lb} \le y^{i+1} \le y_{ub}$, if *i* is odd. Therefore, for any *i* such that (x^i, y^i) is in interval *I*, (x^{i+1}, y^{i+1}) is also in interval *I*. We have proved part (b) of the theorem.

Theorem A.3 For $\alpha = 1.5$, the absolute value of slope of each of the best response curves inside interval I is less than 1 in the X - Y coordinates (Same as Theorem 5.8 in the main text).

Proof We will first prove that at $x = x_{ub}$, $\left|\frac{\partial Y_{BR}(X)}{\partial X}\right| < 1$.

$$\frac{\partial Y_{\rm BR}(X)}{\partial X} = -\alpha \frac{Y}{X} \frac{1 - \frac{Y}{X}}{(\alpha + 1)\frac{Y}{X} - (\alpha - 1)}$$
(75)

The denominator of the right-hand side (RHS) is always positive, due to the SOCs. At $x = x_{ub}$, $x \ge y_{BR}(x)$, and hence $\frac{\partial Y_{BR}(X)}{\partial X} \le 0$. For $\alpha = 1.5$, solving for the point where $\frac{\partial Y_{BR}(X)}{\partial X} = -1$ leads to a unique solution denoted by (x_{-1}, y_{-1}) , where

$$y_{-1} = \frac{9}{32} \frac{Mp_2}{C_2} \text{ and } x_{-1} = 3^{2/3} \frac{9}{32} \frac{Mp_2}{C_2}$$
 (76)

Because $x_{\rm ub} = x_{\rm BR} \left(\frac{\alpha M p_2}{4 C_2} \right)$, we get

$$\frac{4}{k} = \left(\frac{4C_2 x_{\rm ub}}{1.5Mp_2}\right)^{2.5} + 2\left(\frac{4C_2 x_{\rm ub}}{1.5Mp_2}\right) + \left(\frac{4C_2 x_{\rm ub}}{1.5Mp_2}\right)^{-0.5}$$
(77)

Define

$$f(x) = \left(\frac{4C_2 x_{\rm ub}}{1.5Mp_2}\right)^{2.5} + 2\left(\frac{4C_2 x_{\rm ub}}{1.5Mp_2}\right) + \left(\frac{4C_2 x_{\rm ub}}{1.5Mp_2}\right)^{-0.5}$$
(78)

f(x) is a strictly increasing function of x for $x \ge \frac{1.5Mp_2}{4C_2}$. $f(x_{ub}) = \frac{4}{k}$, $f(x_{-1}) \approx 6.96$. $f(x_{ub}) < f(x_{-1})$ if and only if k > 0.575 (approximately), which is always satisfied because one of the necessary conditions for the existence of type BB equilibrium requires that $k \ge (\alpha - 1)^{\frac{1}{\alpha}} = 0.5^{\frac{2}{3}} > 0.575$. Therefore, $x_{ub} < x_{-1}$. Thus, we have proved that at $x = x_{ub}$, $-1 < \frac{\partial Y_{BR}(X)}{\partial X} < 0$. Also, for $x \ge \frac{\alpha Mp_2}{4C_2}$, $\frac{\partial y_{BR}}{\partial x} \le 0$, therefore $y_{-1} = y_{BR}(x_{-1}) \le y_{BR}(x_{ub}) = y_{lb}$. Next, we will obtain the coordinates of the point (which turns out to be unique) such that $\frac{\partial x_{BR}(Y)}{\partial Y} = 1$, and prove that the y-coordinate at this point is less than y_{lb} . The condition,

$$\frac{\partial \alpha_{\text{BR}}(Y)}{\partial Y} = 1.5 \frac{X}{Y} \frac{\frac{X}{Y} - 1}{(1.5 + 1)\frac{X}{Y} - (1.5 - 1)} = 1 \text{ can be simplified to obtain}$$

$$x \approx 0.2029 \frac{1.5Mp_1}{C_1}$$
, and $y \approx 0.1091 \frac{1.5Mp_1}{C_1}$ (79)

Because $k \ge (\alpha - 1)^{\frac{1}{\alpha}} = 0.5^{\frac{2}{3}} > 0.5819$, we get $\frac{C_1 p_2}{C_2 p_1} > \frac{1.1091*1.5}{9}$. So we get $0.1091 \frac{1.5 M p_1}{C_1} < y_{-1} \le y_{\text{lb}}$. So the y-coordinate of the point at which $\frac{\partial x_{\text{BR}}(Y)}{\partial Y} = 1$ is less than y_{lb} . Because $\frac{\partial x_{\text{BR}}(Y)}{\partial Y} \ge 0$ throughout interval I, $0 \le \frac{\partial x_{\text{BR}}(Y)}{\partial Y} < 1$ for the $x_{\text{BR}}(Y)$ curve at $y = y_{\text{lb}}$.

Now let us obtain the coordinates of the point (which turns out to be unique) such that $\frac{\partial Y_{BR}(X)}{\partial X} = 1$ and prove that the x-coordinate of this point is less than x_{lb} . Solving for $\frac{\partial Y_{BR}(X)}{\partial X} = 1$ we get $y \approx 0.2029 \frac{1.5Mp_2}{C_2}$, and $x \approx 0.1091 \frac{1.5Mp_2}{C_2}$. In order to prove that $0.1091 \frac{1.5Mp_2}{C_2} < x_{lb} = x_{BR}(y_{lb})$ it is sufficient to prove that the y-coordinate of the point on the lower part of $x_{BR}(y)$ curve at which $x = 0.1091 \frac{1.5Mp_2}{C_2}$ is less than $y_{-1} = \frac{9}{32} \frac{Mp_2}{C_2}$. This is easy to prove because for $y < \frac{\alpha Mp_1}{4C_1}$, the $x_{BR}(y)$ curve lies below y = x line. Therefore, the y-coordinate corresponding to $x = 0.1091 \frac{1.5Mp_2}{C_2}$ is less than $0.1091 \frac{1.5Mp_2}{C_2}$ which is less than $\frac{9}{32} \frac{Mp_2}{C_2}$. Therefore, at $x = x_{lb}$, $\frac{\partial Y_{BR}(X)}{\partial X} < 1$. So far we have proved that $-1 < \frac{\partial Y_{BR}(X)}{\partial X} \le 0$ at $x = x_{ub}$ and $\frac{\partial Y_{BR}(X)}{\partial X} \le 1$ at $x = x_{lb}$. Therefore, $-1 < \frac{\partial Y_{BR}(X)}{\partial X} < 1$ for all x such that $x_{lb} \le x \le x_{ub}$. Also we have proved that $0 \le \frac{\partial x_{BR}(Y)}{\partial Y} < 1$ at $y = y_{lb}$ and $0 \le \frac{\partial x_{BR}(Y)}{\partial Y} < 1$ at $y = y_{ub}$. Therefore, $-1 < \frac{\partial x_{BR}(Y)}{\partial Y} < 1$ for all y such that $y_{lb} \le y \le y_{ub}$. Therefore for $\alpha = 1.5$, the absolute value of slope of each of the best response curves inside interval I is less than 1 in the X - Y coordinates.

Theorem A.4 For $\alpha = 1$, the absolute value of slope of each of the best response curves inside interval I is less than 1 in the X - Y coordinates (Same as Theorem 5.9 in the main text).

Proof For $\alpha = 1$, the X - Y coordinate system is the same as the x - y coordinate system. We will first prove that at $x = x_{ub}$, $\left|\frac{\partial Y_{BR}(X)}{\partial X}\right| < 1$. For $\alpha = 1$, $\frac{\partial Y_{BR}(X)}{\partial X} = -\frac{1}{2}\left(1 - \frac{Y}{X}\right) > -\frac{1}{2}$. We know that at $x = x_{ub}$, $\frac{\partial Y_{BR}(X)}{\partial X} \leq 0$. Therefore at $x = x_{ub}$, $\left|\frac{\partial Y_{BR}(X)}{\partial X}\right| < 1$. Next, we will obtain the coordinates of the point (which turns out to be unique) such that $\frac{\partial x_{BR}(Y)}{\partial Y} = 1$ and prove that the y-coordinate at this point is less than y_{lb} . Solving for $\frac{\partial x_{BR}(Y)}{\partial Y} = \frac{1}{2}\left(\frac{X}{Y} - 1\right) = 1$, we get

$$x = \frac{3Mp_1}{16C_1}$$
 and $y = \frac{Mp_1}{16C_1}$ (80)

/--- \

For $x \ge \frac{Mp_2}{4C_2}$, we have $\frac{\partial Y_{BR}(X)}{\partial X} \le 0$ and for $y \le \frac{Mp_1}{4C_1}$, we have $\frac{\partial x_{BR}(Y)}{\partial Y} \ge 0$. Also $y_{ub} = \frac{Mp_2}{4C_2} \le \frac{Mp_1}{4C_1}$. So $x_{ub} = x_{BR}(y_{ub}) \le x_{BR}\left(\frac{Mp_1}{4C_1}\right) = \frac{Mp_1}{4C_1}$. So we get $y_{lb} = y_{BR}(x_{ub}) \ge y_{BR}\left(\frac{Mp_1}{4C_1}\right)$. As per the FOCs,

$$\frac{\left(\frac{Mp_1}{4C_1}\right)}{\left(y_{\rm BR}\left(\frac{Mp_1}{4C_1}\right) + \frac{Mp_1}{4C_1}\right)^2} = \frac{C_2}{Mp_2} \Leftrightarrow y_{\rm BR}\left(\frac{Mp_1}{4C_1}\right) = \frac{Mp_1}{4C_1}\left(2\sqrt{k} - 1\right) \tag{81}$$

$$y_{\rm lb} \ge y_{\rm BR} \left(\frac{Mp_1}{4C_1}\right) = \frac{Mp_1}{4C_1} \left(2\sqrt{k} - 1\right) > \frac{Mp_1}{16C_1}$$
 (82)

because $k \ge 0.4$. Therefore, the y-coordinate of the point where $\frac{\partial x_{BR}(Y)}{\partial Y} = 1$ is less than y_{lb} .

Now, let us obtain the coordinates of the point (which turns out to be unique) such that $\frac{\partial Y_{BR}(X)}{\partial X} = 1$ and prove that the x-coordinate of this point is less than x_{lb} . Solving for $\frac{\partial Y_{BR}(X)}{\partial X} = 1$, we get

$$x = \frac{Mp_2}{16C_2}$$
 and $y = \frac{3Mp_2}{16C_2}$ (83)

Because
$$\frac{Mp_1}{4C_1} > y_{\text{lb}} > \frac{Mp_1}{16C_1}$$
, and $\frac{\partial x_{\text{BR}}(Y)}{\partial Y} > 0$ for $y < \frac{Mp_1}{4C_1}$, we get

$$x_{\rm lb} = x_{\rm BR}(y_{\rm lb}) > x_{\rm BR}\left(\frac{Mp_1}{16C_1}\right) = \frac{3Mp_1}{16C_1} > \frac{Mp_2}{16C_2}$$
(84)

The last inequality in (84) holds because $k \leq 1$. Therefore, the x-coordinate at the point where $\frac{\partial Y_{BR}(X)}{\partial X} = 1$ is less than x_{lb} . Thus we have proved that $-1 < \frac{\partial Y_{BR}(X)}{\partial X} \leq 0$ at $x = x_{ub}$ and $\frac{\partial Y_{BR}(X)}{\partial X} < 1$ at $x = x_{lb}$. Therefore, $-1 < \frac{\partial Y_{BR}(X)}{\partial X} < 1$ for all x such that $x_{lb} \leq x \leq x_{ub}$. Also, we have proved that $0 \leq \frac{\partial x_{BR}(Y)}{\partial Y} < 1$ at $y = y_{lb}$ and $0 \leq \frac{\partial x_{BR}(Y)}{\partial Y}$ at $y = y_{ub}$. Therefore, $-1 < \frac{\partial x_{BR}(Y)}{\partial Y} < 1$ for all y such that $y_{lb} \leq y \leq y_{ub}$. Therefore, for $\alpha = 1$, the absolute value of slope of each of the best response curves inside interval I is less than 1 in the X - Y coordinates.

Theorem A.5 If the absolute value of slope of each of the best response curves is less than 1 in interval I, then as long as the competitor frequency for each carrier remains in region B, regardless of the starting point, the myopic best response algorithm converges to the unique type BB equilibrium (Same as Theorem 5.10 in the main text).

Proof We have assumed that the absolute value of slope of each of the best response curves is less than 1 in interval *I*. Also we have proved that as long as the competitor frequency for each carrier remains in region B, regardless of the starting

point the myopic best response algorithm will reach some point in interval *I* in a finite number of iterations and once inside interval *I*, it will never leave the interval. Let (X_{eq}, Y_{eq}) be the type BB equilibrium point in the X - Y coordinate system. We define a sequence L(i) as follows:

$$L(i) = \begin{cases} |X^i - X_{eq}| & \text{if } i \text{ is odd} \\ |Y^i - Y_{eq}| & \text{if } i \text{ is even} \end{cases}$$
(85)

Let us consider any iteration *i* after the algorithm has reached inside the interval *I*. We will prove that once inside interval *I*, L(i) is strictly decreasing. Let us first consider the case where *i* is odd. $L(i) = |X^i - X_{eq}|$. In the $(i + 1)^{th}$ iteration, *X* value remains unchanged. Only the *Y* value changes from Y^i to Y^{i+1} .

$$L(i+1) = |Y^{i+1} - Y_{eq}| = |Y_{BR}(X^i) - Y_{BR}(X_{eq})| = \left| \int_{X_{eq}}^{X^i} \left(\frac{\partial Y_{BR}(X)}{\partial X} \right) dX \right|$$

$$\leq \left| \int_{X_{eq}}^{X^i} \left| \frac{\partial Y_{BR}(X)}{\partial X} \right| dX \right| < \left| \int_{X_{eq}}^{X^i} 1.dX \right| = |X^i - X_{eq}| = L(i)$$
(86)

We have proved that once inside interval I, L(i) is strictly decreasing for odd values of *i*. By symmetry, the same is true for even values of *i*. Moreover, L(i) = 0 if and only if $X = X_{eq}$ and $Y = Y_{eq}$. Therefore, L(i) is a decreasing sequence which is bounded below. So it converges to the unique type BB equilibrium point. \Box

Theorem A.6 In an N-player symmetric game, a symmetric equilibrium with excess seating capacity exists at $x_i = \frac{\alpha M p}{C} \frac{N-1}{N^2}$ for all *i* if and only if $N \le \frac{\alpha}{\alpha-1}$ and if it exists, then it is the unique symmetric equilibrium (Same as Theorem 6.1 in the main text).

Proof The utility of each carrier i is given by

$$u_i(x_i, y_i) = M \frac{x_i^{\alpha}}{x_i^{\alpha} + y_i^{\alpha}} - \frac{C}{p} x_i$$
(87)

where $y_i = \left(\sum_{j=1, j \neq i}^N x_i^{\alpha}\right)^{1/\alpha}$ is the effective competitor frequency for player *i*. From the FOCs, we get

$$x_i = \frac{\alpha M p_i}{C_i} \frac{x_i^{\alpha} y_i^{\alpha}}{\left(x_i^{\alpha} + y_i^{\alpha}\right)^2}$$
(88)

In the symmetric game, $\frac{C_i}{p_i}$ is the same for every player *i*. Let it be denoted as $\frac{C}{p}$. In general, this symmetric game may have both symmetric and asymmetric equilibria. In a symmetric equilibrium, $x_1 = x_2 = \ldots = x_N$. Assume excess seating capacity

for each carrier. Substituting in the FOCs we get $y_i = (N-1)^{1/\alpha} x_i$. Therefore, $x_i = \frac{\alpha Mp}{C} \frac{N-1}{N^2}$ for all *i* is the unique solution. Therefore, we have proved that if an equilibrium exists at this point, then it must be the unique symmetric equilibrium of this game. In order to prove that this point is an equilibrium point, we need to prove that the SOC is satisfied, the profit at this point is nonnegative and seating capacity is at least as much as the demand for each carrier. The SOC is satisfied if and only if

$$\frac{\partial^2 U_i}{\partial x_i^2} = \frac{M \alpha x_i^{\alpha-2} y_i^{\alpha}}{\left(x_i^{\alpha} + y_i^{\alpha}\right)^3} \left((\alpha - 1) y_i^{\alpha} - (\alpha + 1) x_i^{\alpha} \right) \le 0 \Leftrightarrow N \le \frac{2\alpha}{\alpha - 1}$$
(89)

The condition of nonnegativity of profit is satisfied if and only if

$$\frac{\alpha Mp}{C} \frac{N-1}{N^2} C \le \frac{Mp}{N} \Leftrightarrow N \le \frac{\alpha}{\alpha - 1}$$
(90)

The condition of excess seating capacity is satisfied if and only if

$$\frac{\alpha Mp}{C} \frac{N-1}{N^2} S \ge \frac{M}{N} \Leftrightarrow N \ge \frac{\alpha \frac{pS}{C}}{\alpha \frac{pS}{C} - 1}$$
(91)

which is always true for $\alpha \frac{pS}{C} > 2$. Thus the symmetric equilibrium exists if and only if $N \leq \frac{\alpha}{\alpha-1}$.

Theorem A.7 In a symmetric N-player game, there exists some n_{\min} such that for any integer n with $\max(2, n_{\min}) \le n \le \min(N - 1, \frac{\alpha}{\alpha - 1})$ there exist exactly $\binom{N}{n}$ asymmetric equilibria such that exactly n players have nonzero frequency and all players with nonzero frequency have excess seating capacity. There exists at least one such integer for $N \ge \frac{\alpha}{\alpha - 1}$. The frequency of each player with nonzero frequency equals $\frac{\alpha M p}{n} \frac{n-1}{n^2}$ (Same as Theorem 6.3 in the main text).

Proof Let us denote this game as G. Consider any equilibrium having exactly n players with nonzero frequency. Let us rearrange the player indices such that players i = 1 to i = n have nonzero frequencies. Let us consider a new game which involves only the first n players. We will denote this new game as G'. An equilibrium of G where only the first n players have a nonzero frequency is also an equilibrium for the game G' where all players have nonzero frequency. As we have already proved, the equilibrium frequencies of each of the first n players must be equal to $\frac{\alpha Mp}{C} \frac{n-1}{n^2}$. This ensures that any of the first n players will not benefit from unilateral deviations from this equilibrium profile. In order to ensure that none of the remaining N - n players has an incentive to deviate, we must ensure that the effective competitor frequency for any player j such that j > n must be at least equal to y_{th} . This condition is satisfied if and only if

$$n^{\frac{1}{\alpha}} \frac{n-1}{n^2} \frac{\alpha Mp}{C} \ge (\alpha - 1)^{\frac{\alpha - 1}{\alpha}} \frac{Mp}{\alpha C}$$
(92)

i.e.,

$$n^{\frac{1-\alpha}{\alpha}} - n^{\frac{1-2\alpha}{\alpha}} \ge \frac{(\alpha-1)^{\frac{\alpha-1}{\alpha}}}{\alpha^2} \tag{93}$$

LHS of inequality (93) is an increasing function of n for $n \leq \frac{\alpha}{\alpha-1}$. Also the RHS is a decreasing function of α (this can be verified by differentiating the log of RHS with respect to α). Also it can be easily verified that at $n = \frac{\alpha}{\alpha-1}$, the inequality holds for every α . Therefore, for any given α value, there exists some $n_{\min} \geq 0$ such that this inequality is satisfied for all $n \in [n_{\min}, \frac{\alpha}{\alpha-1}]$. As per Theorem 6.1, the condition for existence of an equilibrium with all players having nonzero frequency in game G' is $n \leq \frac{\alpha}{\alpha-1}$.

So all the conditions for an equilibrium of game *G* are satisfied if $\max(2, n_{\min}) \le n \le \min(N-1, \frac{\alpha}{\alpha-1})$. Therefore, any equilibrium of game *G'* where all players have nonzero frequency is also an equilibrium of game *G* where all the remaining players have zero frequency and vice versa. The players in game *G'* can be chosen in $\binom{N}{n}$ ways. Therefore, we have proved that in a symmetric N-player game, for any integer *n* such that $\max(2, n_{\min}) \le n \le \min(N-1, \frac{\alpha}{\alpha-1})$, there exist exactly $\binom{N}{n}$ asymmetric equilibria such that exactly *n* players have nonzero frequency. To show that there exists at least one such integer *n*, consider two cases. If $\alpha > 1.5$, then it is easy to verify that the inequality (93) is always satisfied for n = 2. If $\alpha \le 1.5$, then we see that (93) is satisfied by $n = \frac{1}{\alpha-1} = \frac{\alpha}{\alpha-1} - 1$. In either case, $\frac{\alpha}{\alpha-1} > 2$ is always satisfied. So there always exists some such *n*. The frequency of each player with nonzero frequency equals $\frac{\alpha Mp}{n} \frac{n-1}{n^2}$.

Theorem A.8 Among all equilibria with exactly n players $(n \le N)$ having nonzero frequency, the total frequency is maximum for the symmetric equilibrium (Same as Theorem 6.4 in the main text).

Proof As per Theorem 6.2, any possible asymmetric equilibria with exactly n players having nonzero frequency must involve at least one player without excess seating capacity. Let player i be such a player with nonzero frequency and without excess seating capacity at equilibrium. So the effective competitor frequency y must be less than y_{cr} and

$$x_i > x_{\rm cr} = \frac{M}{S} \left(1 - \frac{C}{\alpha p S} \right) > \frac{M}{2S}$$
(94)

Therefore, each such player must carry at least $\frac{M}{2}$ passengers. Therefore, at equilibrium there can be at most one such player. So each of the remaining n-1 players has excess capacity. Using the same argument as the one used in proving Theorem 6.2, we can prove that each player with nonzero frequency and excess capacity will have equal frequency at equilibrium. Let us denote the equilibrium frequency of the sole player without excess capacity by x_1 and that of each of the remaining players as x_2 . We will denote the equilibrium market share of the player without excess capacity as *l*. Therefore, the total frequency under the asymmetric equilibrium equals,

$$(n-1)x_{2} + x_{1} = \frac{\alpha Mp}{C} \frac{(n-1)x_{2}^{\alpha}}{(n-1)x_{2}^{\alpha} + x_{1}^{\alpha}} \left(1 - \frac{x_{2}^{\alpha}}{(n-1)x_{2}^{\alpha} + x_{1}^{\alpha}}\right) + \frac{M}{S} \frac{x_{1}^{\alpha}}{(n-1)x_{2}^{\alpha} + x_{1}^{\alpha}}$$
$$= \frac{\alpha Mp}{C} (1-l) \left(1 - \frac{1-l}{n-1}\right) + \frac{M}{S}l$$
(95)

Let us assume that there exists an asymmetric equilibrium where the total frequency is greater than that under the corresponding n-symmetric equilibrium, which equals $\frac{\alpha Mp n-1}{C n}$. This condition translates into

$$\frac{\alpha Mp}{C}(1-l)\left(1-\frac{1-l}{n-1}\right) + \frac{M}{S}l > \frac{\alpha Mp}{C}\frac{n-1}{n}$$
(96)

which further simplifies to

$$nl(5 - n - 2l) > 2$$
 (97)

But we know that $n \in \mathbb{I}^+$, $n \ge 2$ and $l \ge \frac{1}{2}$. So 5 - n - 2l > 0 only if $n < 5 - 2l \le 4$. So n = 2 or n = 3. For n = 2, the conditions for existence of type BC equilibrium in the two-player symmetric case require $\frac{\alpha PS}{C} \le 2$, which contradicts our assumption. For n = 3, we need some *l* such that

$$3l^2 - 3l + 1 < 0 \tag{98}$$

which is true if and only if

$$3(l-0.5)^2 + 0.25 < 0 \tag{99}$$

which is also impossible. Thus our assumption leads to a contradiction. So we have proved that among all equilibria with exactly *n* players ($n \le N$) having nonzero frequency, the total frequency is maximum for the symmetric equilibrium.

References

- 1. Adler N (2001) Competition in a deregulated air transportation market. Eur J Oper Res 129 (2):337–345
- 2. Aguirregabiria V, Ho CY (2012) A dynamic oligopoly game of the US airline industry: estimation and policy experiments. J Econometrics 168(1):156–173
- 3. Ball M, Barnhart C, Dresner M, Neels K, Hansen M, Odoni A, Peterson E, Sherry L, Trani A, Zou B (2010) Total delay impact study: a comprehensive assessment of the costs and impacts of flight delay in the United States, NEXTOR
- Barnett A (2008) Is it really safe to fly. In: Chen ZL, Raghavan S (eds) Tutorials in operations research: state-of-the-art decision making tools in the information-intensive age, INFORMS, Hanover, pp 17–30
- 5. Barnett A, Paull G, Iaedeluca J (2000) Fatal US runway collisions over the next two decades. Air Traffic Control Q 8(4):253–276
- 6. Baseler R (2002) Airline fleet revenue management: design and implementation. In: Jenkins D (ed) Handbook of airline economics, 2nd edn. Aviation Week, Washington, DC
- 7. Belobaba P (2009) The airline planning process. In: Belobaba P, Odoni A, Barnhart C (eds) The global airline industry. Wiley, New York, pp 153–181
- Belobaba P (2009) Overview of airline economics, markets and demand. In: Belobaba P, Odoni A, Barnhart C (eds) The global airline industry. Wiley, New York, pp 47–71
- 9. Binggeli U, Pompeo L (2006) Does the s-curve still exist? Techical report, McKinsey & Company
- 10. Bonnefoy p (2008) Scalability of the air transportation system and development of multiairport systems: a worldwide perspective. Thesis, Massachusetts Institute of Technology, Cambridge
- Brander JA, Zhang A (1993) Dynamic oligopoly behaviour in the airline industry. Int J Ind Organ 11(3):407–435
- 12. Brueckner JK (2010) Schedule competition revisited. J Transp Econ Policy 44(3):261-285
- 13. Brueckner JK, Flores-Fillol R (2007) Airline schedule competition. Rev Ind Organ 30 (3):161–177
- Brueckner JK, Van Dender K (2008) Atomistic congestion tolls at concentrated airports? the internalization debate. J Urban Econ 64(2):288–295
- 15. Bureau of Transportation Statistics (BTS) (2010) TranStats. www.transtats.bts.gov. Accessed 15 Feb 2010
- Button K, Drexler J (2005) Recovering costs by increasing market share: an empirical critique of the S-curve. J Trans Econ Policy 39(3):391–404
- 17. Daniel J, Harback K (2008) (When) Do hub airlines internalize their self-imposed congestion delays? J Urban Econ 63:583–612
- Dobson G, Lederer PJ (1993) Airline scheduling and routing in a hub-and-spoke system. Transp Sci 27(3):281–297
- Goodman D, Mandayam N (2000) Power control for wireless data. IEEE Pers Commun 7 (2):48–54
- Hansen M (1990) Airline competition in a hub-dominated environment: an application of noncooperative game theory. Transp Res Part B Methodol 24(1):27–43
- Hong S, Harker PT (1992) Air traffic network equilibrium: toward frequency, price and slot priority analysis. Transp Res Part B: Methodol 26(4):307–323
- 22. Kahn AE (1993) Change, challenge, and competition: A review of the airline commission report. Regulation 3:1–10
- 23. Koutsoupias E, Papadimitriou C (2009) Worst-case equilibria. Comput Sci Rev 3(2):65-69
- 24. Norman VD, Strandenes SP (1994) Deregulation of Scandinavian airlines: a case study of the Oslo-Stockholm route. In: Krugman P, Smith A (eds) Empirical studies of strategic trade policy. NBER Books, National Bureau of Economic Research, Cambridge

- 25. O'Connor WE (2001) An introduction to airline economics. Greenwood Publishing Group, Westport
- Pels E, Nijkamp P, Rietveld P (2000) Airport and airline competition for passengers departing from a large metropolitan area. J Urban Econ 48(1):29–45
- Rosen JB (1965) Existence and uniqueness of equilibrium points for concave N-Person games. Econometrica 33(3):520–534
- 28. Roughgarden T (2005) Selfish routing and the price of anarchy. The MIT Press, Cambridge
- Saraydar C, Mandayam N, Goodman D (2001) Pricing and power control in multicell wireless data network. IEEE J Sel Areas Commun 19(10):1883–1892
- Saraydar C, Mandayam N, Goodman D (2002) Efficient power control via pricing in wireless data network. IEEE Trans Commun 50(2):291–303
- Simpson RW (1970) A market share model for US domestic airline competitive markets. Massachusetts Institute of Technology, Flight Transportation Laboratory, Cambridge
- 32. Taneja NK (1968) Airline competition analysis. MIT Flight Transportation Laboratory, Cambridge
- 33. Taneja NK (1976) The commercial airline industry. DC Heath, Lexington
- Transport Canada (2000) National civil aviation safety committee. Sub-committee on runway incursions. Final report 14 Sept 2000, Montreal, Canada
- Vaze V, Barnhart C (2012) Modeling airline frequency competition for airport congestion mitigation. Transp Sci 46(1):1–24
- 36. Wei W, Hansen M (2005) Impact of aircraft size and seat availability on airlines' demand and market share in duopoly markets. Transp Res Part E Logist Transp Rev 41(4):315–327
- Wei W, Hansen M (2007) Airlines' competition in aircraft size and service frequency in duopoly markets. Transp Res Part E Logist Transp Rev 43(4):409–424
- Wilke S, Majumdar A (2012) Critical factors underlying airport surface accidents and incidents: a holistic taxonomy. Airport Manage 6(2):170–190

A Simulation Game Application for Improving the United States' Next Generation Air Transportation System NextGen

Ersin Ancel and Adrian Gheorghe

Abstract Societies around the world depend on the proper functioning of various infrastructures. However, changes in technology, societal needs/expectations, political shifts, and environmental concerns cause infrastructure systems to underperform (i.e., congestion, energy shortage, air transportation delays, etc.). In order to accurately plan the next generation infrastructure systems, understanding the interactions between technical, political, and economic factors as well as stakeholders are of paramount importance. The current research pursued the development and deployment of a simulation game which aimed to serve as a venue to generate and evaluate data for next generation infrastructure development efforts. The problem domain was selected as the Next Generation Air Transportation System (NextGen) transition environment. The complex and stakeholder-rich environment of NextGen provided an accurate test-bed for the sociotechnical system transformation, highlighting the interaction of variables like system capacity, safety, public demand, and stakeholder behavior with diverging agendas under various world scenarios.

Keywords Serious gaming • Game theoretic • Infrastructure planning • Expert elicitation • Air transportation

1 Introduction

1.1 Future of Critical Infrastructures

Infrastructures are defined as facilities and systems that have strong links to economic development and public sector involvement. They provide the underpinnings of the nation's defense, a strong economy, and health and safety. Following the

A. Gheorghe e-mail: agheorgh@odu.edu

© Springer International Publishing Switzerland 2015 K. Hausken and J. Zhuang (eds.), *Game Theoretic Analysis of Congestion, Safety and Security*, Springer Series in Reliability Engineering, DOI 10.1007/978-3-319-13009-5_8

E. Ancel $(\boxtimes) \cdot A$. Gheorghe

Old Dominion University, Norfolk, VA, USA e-mail: eance001@odu.edu

9/11 terrorist attacks, the term "critical infrastructure" became of interest with seemingly similar definition to that of infrastructures. Critical infrastructures are "national infrastructures [...] so vital that their incapacity or destruction would have a debilitating impact on the defense or economic security of the United States" [43]. Infrastructures like energy, banking/finance, transportation, water, emergency services, government, health services, food/agriculture, telecommunications information networks, national monument icons among others are considered critical for various criteria such as national defense, economic security, public health and safety as well as national morale. Today's critical infrastructures are large-scale systems comprised of several elements and forces, involving various stakeholders, technologies, policies, and social factors [14, 28]. Since infrastructures involve both social (people, rules, regulations, etc.) and technical components (roads, bridges, power stations, etc.), they are also referred as sociotechnical systems [3].

In recent years, various infrastructures started to undergo a series of structural changes in order to respond to increased performability, sustainability, and environmental efficiency demands [6]. However, modernization of these infrastructures is being held back for reasons besides economics [15]. Because infrastructures like power, transportation, and communication contain multidimensional complexity and are essentially stable, transforming the existing systems to more efficient alternatives is challenging [41]. The planning and implementation phases of such infrastructure transitions require close monitoring of performance parameters like safety, efficiency, and sustainability. Ensuring that infrastructure transition reveals a safer and more sustainable system has become a major challenge for society [10, 26]. The design requirements in infrastructure systems are dynamically affected by the presence of various competitive stakeholders, shifts in public perceptions, and changes in the political environment. The issues associated with the infrastructure transition prohibit the use of traditional professional knowledge which often provides a narrow, technical rationality. The problem solving strategies provided by the mathematics and physics alone are too limited in scope and do disregard the presence of competing frameworks (actors, stakeholders, etc.) [25].

1.2 Simulation Gaming to Assist Infrastructure Transitions

Due to the challenges highlighted in the previous section, the need to comprehend infrastructures at the societal level and understand technical, political, and economic factors' interaction becomes more and more prominent [19]. Managing the complexity and uncertainty associated with infrastructure systems requires understanding and teaching the internal logic of the system [25]. However, teaching the dynamic infrastructure behavior (whether the basics to students or the implications of new policy measures to policy makers and strategists) is proven to be difficult [46].

Also, solely employing past strategies and historical data regarding previous infrastructure systems are no longer adequate for next generation infrastructure systems design because (1) previous systems evolved via incremental changes

which lead them to be unsustainable (i.e., congestion, energy shortage, air transportation delays, etc.), and (2) previous infrastructures were made to last and were robust but also are resistant to change that can causes challenges [10, 14, 19]. Additionally, the increased presence of societal aspects in the sociotechnical system structure causes complications in understanding and foreseeing solutions. The evolutionary nature of infrastructure systems and the ever-changing societal dynamics make every problem essentially unique, rendering historical data somehow ineffective [5, 21, 40].

The research presented in this chapter concentrates on the development and the execution of a simulation gaming application to assist the infrastructure transition process. Besides employing historical data, the insight obtained from a well planned and executed simulation gaming could potentially assist the infrastructure planning and transition process. The complex stakeholder interactions as well as counter-intuitive behavior of social systems (e.g., policy resistance) could be identified via simulation games. Consequently, the simulation gaming approach could benefit the policy makers and strategists to be more involved with the dynamic infrastructure transition behavior, yielding to more informed decision-making processes [25, 46]. The feasibility of such an application is investigated by developing a simulation game representing the transition process of the current National Airspace System (NAS) into the Next Generation Air Transportation System (NextGen) environment. The goal of the game is to identify stakeholder interactions and to simulate the fluctuations in aviation safety during the NextGen transition process.

1.3 Simulation Gaming and Policy Gaming Overview

Over the last few decades, practitioners and management scholars increasingly criticized the conventional strategy making methods, arguing that rapidly changing environments require emerging and creative approaches. Serious gaming (simulation game or gaming, used interchangeably within the text) discipline is found to be increasingly useful within the mainstream strategy literature involved with former strategy making approaches [16]. A definition of simulation gaming is given as a representation of a set of key relationships and structure elements of a particular issue or a problem environment, where the behavior of actors and the effects of their decision are a direct result of the rules guiding the interaction between these actors [47].

Serious gaming is an activity where two or more independent decision makers seek to achieve their objectives within a limited context: "The participants (or the players) of the game perform a set of activities in an attempt to achieve goals in a limiting context consisting of constraints and of definitions of contingencies [18]." The common point on each simulation game is that reality is simulated through the interaction of role players using nonformal symbols as well as formal, computerized submodels when necessary. This approach allows the group of participants to create and analyze future worlds they are willing to explore [48].

Simulation games have many different forms and aim to provide insights for various goals. This chapter is focused on the policy gaming exercises which can carry various objectives like understanding system complexity, improving communication, promoting individual and collective learning, creating consensus among players, and motivating participants to enhance their creativity or collaboration. The policy gaming exercises allow the exploration of complex interaction of social aspects among participants [4, 8, 16].

2 Research Design and Game Setup

As previously discussed, this chapter presents the development of a simulation gaming exercise for an infrastructure transformation problem. The history of infrastructure development shows that the majority of the challenges associated with transitions are related to social aspects, rather than technology related issues. The pressure from various stakeholders with different agendas renders the infrastructure transition rather challenging. Meijer [29] argues that simulation gaming helps researchers study trust and cheating within complex supply networks. This feature of simulation gaming was also a key element in identifying hidden agendas and negotiations among various NAS stakeholders. The simulation game is designed to capture the components of an infrastructure transition process including drivers for change (new technologies, congestion, decay, efficiency, and reliability, changing needs, etc.), constraints (existing structure, cost, environmental, social, and political impacts and externalities) as well as context (government intervention, stakeholder actions, social factors, economic and political opportunities including developing new standards and protocols) [19].

The proposed game consists of three phases: *pre-gaming*, *gaming*, and *post-gaming*. The pre-gaming phase consists of collection of the gaming variables depending on the modeled infrastructure or system. Such variables include scenarios, stakeholders and their interactions, and historical data regarding the system and information on the parameter(s) upon which the success of the transition process will be measured. The identification and selection of an appropriate simulation mechanism is also identified in this phase. Depending on the application, a computer-based simulation can be used to evaluate risk or reliability of an infrastructure system or keep track of generation capacity or throughput of a certain utility. The game development phase is an iterative process where versions are often tested by playing with several groups followed by fine-tuning.

The *gaming* phase includes the execution of the gaming exercise with the participation of experts. The game usually starts with the presentation of the scenario to the participants. Participants are asked to perform according to their predetermined roles. Considering the new information they have been presented, participants are asked to make collective decisions about the investigated parameters. The decisions are taken as the input variables for the computer assisted simulation mechanism where initial conditions for the next step are calculated in an iterative manner. The presence of participants (preferably experts or actual stakeholders of the system) and their social values, norms, and beliefs provide the realistic input for the social interaction and the decision-making process.

The *post-gaming* phase involves data collection and analysis which surfaced during the gaming cycle. At this level, elicited data are arranged and presented back to the participants for further analysis and feedback.

3 Problem Domain: Next Generation Air Transportation System Safety

The infrastructure transformation problem presented in the previous section is demonstrated by modeling the transition from the current airspace system into the NextGen environment. The goal of the modeling effort is to demonstrate various challenges of transformation of sociotechnical systems and to simulate the air transportation system safety throughout the transition process. This section provides an overview of the modeled NextGen characteristics followed by the game details, of data requirements, game rules, scenarios, stakeholders, and their interactions. Aside from the test runs, the game was played once by the SMEs representing the stakeholders. The data collection mechanism and the results of the gaming exercise were also presented in this section.

3.1 NextGen Overview

The United States' (NAS) is a vast, multilayered array of operations covering virtually everything involving air transportation. With well over 800 million passengers, NAS requires input from more than 15,000 air traffic controllers to assist 590,000 pilots on board 239,000 aircraft that take off and land at 20,000 U.S. airports. Within recent years, delays have heavily impacted passenger travel, and they are forecasted to be even higher in the future as the demand for air transportation is expected to increase [12, 44].

The NextGen consists of a system-wide upgrade of the current NAS. The Joint Planning and Development Office (JPDO), which is a cooperative partnership between public and private stakeholders, is charged with developing concepts, architectures, roadmaps, and implementation plans for transforming the current NAS into the NextGen [23]. The NextGen goals include flying an increased number of passengers and cargo more safely, precisely, and efficiently, while using less fuel and decreasing environmental impact [13]. During the next two decades, it is expected that the system will provide two to three times the current air vehicle operations and will be agile enough to accommodate a changing fleet that includes unmanned aircraft systems (UASs), and space vehicles while maintaining a safe airspace [22].

However, the complex nature of the NAS, combined with numerous operational and management challenges, threatens the NextGen efforts. The reports from the Office of the Inspector General (OIG) reveal that the Federal Aviation Administration (FAA) is facing difficulties in developing a strategy to engage stakeholders, not to mention managing and integrating multiple NextGen efforts [45]. Also, challenges like multidimensional research and development along with complex software development, workforce changes, mixed equipage, and policy issues need addressing.

3.2 Gaming Overview

The game is aimed to determine the NAS safety values during the transition process by examining the chief safety related NextGen enablers and technologies against a set of predetermined scenarios. The high-level overview of the NextGen transition game is given in Fig. 1. During the pre-game phase, information regarding NextGen environment components such as scenarios, realistic timelines, stakeholders, and technologies are gathered via the literature review and expert opinion. The gaming phase includes the iterative decision-making cycles and the dynamic calculation of concerned variables (i.e., NAS safety value) at each time step (defined as one year) using the risk simulation mechanism called Rapid Risk Assessment Model (RRAM). The post-gaming phase includes data collection and analysis. The data obtained from the gaming exercise include the behavior of the stakeholders, the 2010–2025 dynamic aviation risk values under varying scenarios and the ranking of the NextGen enabler and technology alternatives.



Fig. 1 Serious gaming design-high-level gaming architecture

3.3 Pre-gaming Phase

3.3.1 Data Requirements

In order to setup the boundary conditions and support the decision-making process, various data sources have been used throughout the game development. Due to the nature of the problem at hand, a combination of numerical and elicited data has been used in different parts of the game. Historical data concerning the current aviation accident rates and fatalities for FAA Part 121 are taken from NTSB general aviation statistics [37] and the Aviation Accident and Incident Data System. In both datasets, the average of the past 10 years was considered.

Besides the incident/accident-related data, current airline and airport financial data are extracted to create a realistic baseline for the gaming activity. For that purpose, the eight largest airports¹ and ten airline companies² are identified and their financial data are selected as initial conditions for the game. For the airline activities, average operating revenues (passenger enplanements, transportation revenues, baggage fees, and reservation cancellation fees) as well as operating expenses (flying operations, maintenance, and specific fuel cost) for years 2005–2009 were used to obtain operating profit/loss values. Airline data were obtained from Research and Innovative Technology Administration (RITA) Bureau of Transportation Statistics (BTS) website.³

The airports are selected according to the passenger enplanement in 2009. The hub airports and the financial data are obtained from the FAA Compliance Activity Tracking System (CATS)—Summary Report 127.⁴ The total operating income is calculated using the aeronautical and non-aeronautical revenues and expenses. The financial data for the airline companies and hub airports are handled collectively, and their 5 year averages (2005–2009) are adopted as the initial conditions for the gaming exercise.

Finally, information regarding the future NextGen enablers and technologies and their attributes are collected. By their nature, technologies planned for future infrastructure systems and their parameters carry uncertainties (e.g., advantages of a certain technology in 15 years, cost-benefit values, etc.). Consequently, the financial and technical (schedule and safety impact) data was obtained via FAA/JPDO as well as expert opinions.

¹ Airports considered in the calculation include Hartsfield-Jackson Atlanta (ATL), Chicago O'Hare (ORF), Los Angeles (LAX), Dallas Fort Worth (DFW), Denver (DEN), John F. Kennedy (JFK), George Bush Intercontinental-Houston (IAH), and Las Vegas-McCarran International (LAS).

 $^{^2}$ At the time of writing, the selected ten airline companies (Delta, Southwest, United, U.S., Northwest, JetBlue, Continental, American, Alaska, and Airtran) represented around 70 % of all NAS enplanements in the domestic market.

³ http://www.transtats.bts.gov/, Retrieved July 3rd, 2014.

⁴ http://cats.airports.faa.gov/Reports/reports.cfm, Retrieved July 3rd, 2014.

3.3.2 Game Rules

Each serious game is a dedicated simulation exercise, specifically tailored and developed for the problem at hand. The actual run of the serious game is a collective and interactive process, designed by the very owners of the problem since participants can dynamically change or create rules, as opposed to formal, blackbox simulation exercises [16]. The identification of the game rules plays a very important role. Rigid and rule-based gaming works well for well-structured environments like military gaming where specific rule-sets, formalized by mathematical and/or computational methods exist. The rigid-type rule-sets are successful when the problem at hand is well defined and understood. On the other hand, in social arenas with public and intense stakeholder interactions where firm rules do not exist, free-form gaming is more suitable. In this type of approach, positions, objects, and rules can be challenged, modified, and improved by players during game play as they see fit [27]. Also, Dormans [7] argues that in order to effectively communicate the modeled system, game's perspective should be focused on maximizing the effectiveness of game mechanics by emphasizing the expressive power of relatively simple game mechanics.

Since the primary goal of this game is to provide insights into future NAS safety and data gathering regarding future systems, a combination of rigid and free-form gaming rules was found most suitable. The cases where the participants are actual stakeholders in the aerospace industry, the common ground rules are usually well known. Nonetheless, a number of basic stakeholder rules are determined to help guide the stakeholder interactions throughout the gaming process (Fig. 2).



Fig. 2 Basic stakeholder rules

3.3.3 Stakeholders

One of the most productive outcomes of the policy gaming exercises is the participants' interaction with the problem at hand. As Duke [9] argues, real-world complex problems often include a sociopolitical context, created by the idiosyncratic or irrational players present during the decision-making process.

The NextGen transition process involves multiple stakeholders with diverse agendas that can directly or indirectly affect the outcome. In order to model such a dynamic environment, a simplified list of involved stakeholders and respective objectives were identified. Collectively, each stakeholder desires a safer and more profitable NAS environment; however, individual objectives and agendas often inhibit the overarching goal. The interested parties along with their primary (and often conflicting) goals are given in Table 1. However, all stakeholders are bounded to the scenario variables.

Government, FAA, and Military Stakeholders. The government, FAA, and military stakeholders are grouped together for simplification purposes. This group represents the "big brother" role over the airports and airlines. The government is responsible for determining tax values for various areas such as income, environmental and security, along with aviation fuel tax. The FAA's role as the enforcer of aviation safety is also controlled by this stakeholder based on the dynamic yearly risk levels. They also have the ability to spare funds for assisting airlines and airports in purchasing large-ticket items such as ADS-B and Data Link enablers. The final task of this stakeholder is to reflect the military agenda based on the scenario presented (i.e., the adjustments required for UAS integration to NAS).

Corporate Airlines. Individual corporate airlines are represented as a single entity, assuming they have similar goals. This group determines yearly average ticket prices, reservation and cancellation fees, and passenger baggage fees. Their finances are directly affected by the strategies followed by other players, e.g., airports collect landing fees and the government collects various taxes. Airlines are also affected by the predetermined scenario mandating aircraft jet fuel before taxes or global terrorist threats. Airlines are encouraged to engage in coalitions with other stakeholders and invest funds in NextGen enablers and technologies because they are the primary beneficiaries of the increase in NAS capacity. Corporate airlines are expected to reflect their expenses in passenger ticket prices and fees; however, the general public stakeholder can react to increased ticket prices by choosing other modes of transportation.

Stakeholder name	Objectives
Government, FAA, and military	Safety, protect, and nurture aviation industry, throughput
Corporate airlines	Market share, profit, throughput, safety
Airport operators	Throughput, revenue neutral, safety
Public	Affordable and safe transportation

Table 1 Stakeholders listand their primary objectives

Airport Operators. For simplification purposes, the airport stakeholder represents the main hub airports in the continental United States. Like the corporate airlines stakeholder, airport operators interact with other players in determining their strategies and pricing such as aeronautical and non-aeronautical fees. Aeronautical fees include airport landing fees that a passenger facility charges that are billed to airline companies. Non-aeronautical fees include parking fees, concession fees, airport shop rental fees, etc. Since NextGen enablers allow airports to increase their landing capacities, players representing the airports are inherently motivated to invest in these technologies as well.

General Public. The general public stakeholder indirectly works with the game master in order to determine the "actual" air transportation capacity. This stakeholder reviews the information regarding various forms of transportation and determines if he/she agrees with the projected air traffic capacity. The public stakeholder adjusts the air transportation capacity by comparing alternative methods of transportation (automobile and high speed train) in predetermined routes. At the end of each time step (i.e., simulated year), the public stakeholder decides whether to agree or adjust the projected air transportation capacity from -10 to 10% with 5 % intervals. The general public stakeholder can reflect upon the increased air transportation costs that were decided by airport and airline stakeholders.

3.3.4 Scenarios

The gaming exercise requires a dynamic environment to enable participants (or players) to interact with each other. The dynamic scenario enables game facilitators or interested parties to evaluate various scenarios and extract the collective response from all the stakeholders. The scenarios presented in the gaming application are adapted from a workshop conducted by the National Research Council (NRC) [35]. A modified version of the NRC study scenarios is given in Table 2. Four scenarios considered for the gaming activity include "Pushing the Envelope," "Grounded," "Regional Tensions," and "Environmental Challenged." Each year (or time step) has its distinctive scenario with specific attributes and players are expected to adjust their strategies accordingly. It is possible to select and experiment with other scenarios, sources, or combinations.

3.4 Gaming Phase

The gaming phase consists of the actual human-in-the-loop data gathering effort and is adopted from a policy gaming platform developed by Geurts et al. [16]. A modified version of the play sequence is supported by the risk simulation mechanism and COTS software in order to accommodate the NextGen safety framework (Fig. 3). The process is initiated by the presentation of the game to the stakeholders including the game rules, overall NextGen goals, and available

Table 2	Scenario environment	by year					
Year	Scenario name	U.S. eco- nomic state	Demand for aeronau- tics services	Global secu-	Government	Base fuel	Anticipated passenger
2010	Pushing the	Strong	High growth	Low	Low	\$2.136	100
2011	envelope Pushing the	Strong	High growth	Low	Low	\$2.136	120
2012	Pushing the envelope	Strong	High growth	Low	Low	\$2.136	135
2013	Pushing the envelope	Strong	High Growth	Low	Low	\$2.136	150
2014	Pushing the envelope	Strong	High growth	Low	Low	\$2.136	165
2015	Grounded	Strong	High growth	High	High	\$2.136	50
2016	Grounded	Strong	Low growth	High	High	\$2.138	60
2017	Regional tensions	Strong	High growth	High	High	\$2.140	165
2018	Regional tensions	Weak	High growth	High	High	\$2.140	175
2019	Regional tensions	Weak	High Growth	High	High	\$2.140	185
2020	Regional tensions	Weak	High growth	High	High	\$2.140	195
2021	Regional tensions	Weak	High growth	High	High	\$2.155	205
2022	Regional tensions	Weak	Low growth	High	High	\$2.155	215
2023	Environmentally challenged	Weak	Low growth	High	High	\$2.170	220
2024	Environmentally challenged	Weak	Low growth	High	High	\$2.180	225
2025	Environmentally challenged	Weak	Low growth	High	High	\$2.180	230



Fig. 3 NextGen safety risk assessment gaming sequence overview

resources. Participants are then presented with their respective roles, resources, and goals. Usually, a dry-run is favorable in order to ensure participants' understanding of the game rules and expectancies.

The gaming sequence starts with participants briefing regarding environmental variables such as economic state, global security threat, and values like NAS capacity and fuel prices. Then, participants are asked to review their resources and consider acquiring potential enablers that can help reach their safety/capacity goals. Following discussions within each participant group, their decisions are gathered and inputted into the risk simulation mechanism where next year's updated NAS risk values are calculated. Along with updated scenario variables, the NAS risk values constitute the initial conditions for the next iteration step. The game is iterated until the desired year is reached. The gaming simulation is concluded with debriefing, discussions, and feedback. Detailed gaming steps are given below.

- 1. The game master/facilitator announces the variables of the specific calendar year including the anticipated air transportation capacity, political, economic, social environments, and the untaxed fuel price as given in Table 2.
- 2. According to predictions for the upcoming year, participants experiment with their variables and simulate their budgets using the provided personalized Excel spreadsheets, allowing them to determine the funds that can be used for NextGen enablers.
- 3. The participants are given 5 min to discuss the enabler acquisition strategy and possible coalitions. To help the process, airline and airport stakeholders are provided with LDW models where they can experiment with the most desirable enabler selection.

- 4. The elicitation round starts with the government's announcement of taxes for the coming year, in accordance with the political and economic environment (government participation level and U.S. economic state).
- 5. The participants representing airport authorities announce their landing fees and concession fees, along with the enablers they are willing to purchase this year.
- 6. The airline stakeholders announce their variables: passenger ticket fees, reservation cancellation fees, baggage fees, and the planned NextGen enabler acquisitions.
- 7. Another 5 min are allowed for stakeholders to discuss their fees among them before they are announced to the game facilitator.
- 8. Using the input provided by the participants, the game facilitator runs the simulation mechanism to reveal the risk values for the specific year.
- 9. Once the "new air transportation environment" is revealed, the general public stakeholder examines the cost for various modes of transportation along with the safety of air travel and determines the final air transportation capacity by adjusting the previously announced anticipated capacity. Adjustments can be done from -10 to 10 % change with 5 % increments.
- 10. With the actual passenger capacity determined, stakeholder budgets are adjusted, and the following year's variables are stated by the game facilitator, and next round is initiated beginning with Step 1.

3.4.1 Risk Simulation Mechanism

In order to represent dynamic risk values with the NAS, a comprehensive, yet intuitive risk calculation method named (RRAM⁵) was developed and integrated into the game. The RRAM tracks the desired parameter (i.e., air transportation safety) and employs the traditional risk construct defined as the product of probability of an accident and its respective consequences [2]. The RRAM calculates accident probabilities and consequences separately and the resultant risk is demonstrated on a risk matrix (Fig. 4). The following sections highlight the components of RRAM.

Consequences. The consequences (the *x*-axis of the risk matrix, seen in Fig. 4) of aviation accidents are based on *fatalities*, considering that the ultimate goal of NextGen related safety efforts within JPDO is concerned with saving human lives. The consequences are estimated as a product of various components comprised of (a) baseline fatality rate of Federal Aviation Regulations (FAR) Part 121 aviation [37], (b) air traffic density rate (function of *t*), and (c) presence of the correcting or mitigating factors regarding the survivability rate in accident scenarios.

$$C_i = F_{\text{baseline}} \times \delta_i \times n_{\text{survival rate},i} \tag{1}$$

⁵ The RRAM was developed based on hazardous material handling, storage, processing, and transportation study conducted by a United Nations joint consortium [20].

Severity	<0.30	0.30	0.38	0.45	0.53	0.60	0.68	0.75	0.83	>0.9
Likelihood	Minimal	Mi	nor	Major		Hazardous		Catastrophic		N Scale
Frequent										N ≤ 3
Probable										3 < N ≤ 5
Bemote		л.								5 < N ≤ 6
										6 < N ≤ 7
Extremely										7 < N ≤ 8
Remote										8 < N ≤ 9
Extremely Improbable										N > 9

Fig. 4 Risk matrix, adapted from FAA [11]

where:

C_i	Consequences at time, t
F _{baseline}	Average fatality rate for FAR Part 121
δ_i	NAS air traffic density at time, t
n _{survival rate,i}	Crash Survivability correcting factors (Table 3)

The crash survivability correcting factors (i.e., fire/smoke mitigation, survivability of aircraft structures, and accident response procedures) are adopted from the National Science and Technology Council and are provided in Table 3. The formula is developed to estimate accident fatalities per 100,000 flight hours.

Presence of cor- recting factors	1: Enhanced post- impact fire/smoke mitigation (-15 %)	1 + 2: Improved crash survivability of aircraft structures (-15 %)	1 + 2 + 3: Improved evacuation and accident response procedures (-15 %)
Initial value: 1.0	0.85	0.7	0.55

 Table 3 Crash survivability correcting factors [36]

Probabilities. The accident probabilities are estimated via Probability Number Method (PNM) where the probability of a certain accident happening is calculated via a dimensionless "probability number," N, which is then transformed to actual probabilities. The relationship between the probability and N is given via $N = |\log_{10} P|$. The y-axis of the risk matrix, seen in Fig. 4, consists of accident likelihood or probability given in N numbers.

The probability number is updated according to the presence of various correcting factors. In the context of NAS safety, the correcting factors are enablers and mitigation technologies, obtained from NextGen JPDO's Avionics Roadmap [24] and subject matter experts (SMEs). The tools, methods, and programs considered below do not constitute an exhaustive list; however, efforts from both NASA and FAA-guided programs are included. The accident probabilities are calculated via the equation below.

$$N_i = N_i^* + n_{\rm rs} + n_{\rm asr} + n_{\rm icing} + n_{\rm ac} + n_{\rm Wxa} + n_{\rm turb}$$
(2)

where:

 N_i Calculated probability number for the system at time = t N^* The average probability number for the current NAS setup Correction parameter for runway safety and collision avoidance $n_{\rm rs}$ Correction parameter for aircraft systems reliability technologies nasr Correction parameter for icing mitigation technologies nicing Correction parameter for airborne collision avoidance n_{ac} Correction parameter for weather avoidance precautions n_{Wxa} Correction parameter for turbulence (wake) avoidance solutions n_{turb}

The calculated probability number (N_i) is updated at each time frame and is used as the initial average probability value (N_i^*) for the next time step.

Probability, Consequence Definitions, and Risk Matrix Thresholds. The NAS safety risk is obtained by the combination of accident probabilities and consequences and displayed on the risk matrix adapted from the FAA's Safety Management System Manual [11]. This graphical means of determining risk levels is chosen since the game aims to calculate the likelihood (probability) and the severity (consequences) for each risk independently where the risk is the product of these two (Fig. 4).

According to FAA's Risk matrix, the *x*-axis provides the severity of a threat, where the *y*-axis determines the probability of such threat to occur. The intersection of the two axes provides the risk, shown via a colored traffic light matrix. The "traffic light" approach is taken where the red areas demonstrate the unacceptable risk areas, caused by an event carrying catastrophic consequences, major consequences with a high likelihood value. The yellow and green areas signify the medium and low risk levels, respectively. The definitions of the severity and likelihood axes are given in the following tables (Tables 4 and 5) within the context of NextGen safety assessment game.

Consequences	Minimal	Minor	Major	Hazardous	Catastrophic
\rightarrow	5	4	3	2	1
Fatalities/100.000FH	0	0.291	0.436	0.582	0.727
Normalized	0	0.25	0.5	0.75	1

 Table 4
 Consequences definitions

The consequences or severity levels are defined based on the FAA risk definitions, using five-point Likert scale, ranging from *minimal* to *catastrophic*. The historical average of 0.291 fatalities per 100,000 flight hours (obtained from the NTSB website for 2000–2009 timeframe) is assumed to be a minor risk that the aviation industry inherently carries. Since there have been years without any fatalities within the FAR Part 121, the lower end of the axis is assigned as "0." The upper end of the scale designates the worst-case scenario where there are no crash survivability efforts in 2025 with NAS air traffic density is 2.5 times the current density. Within this assumption, the threshold value for catastrophic consequences can be shown as: $C_{2025} = 0.291 \times 2.5 \times 1 = 0.727$. The remainder of the table is normalized to reflect minor, major, and hazardous consequences (Table 4).

Similarly, the probability axis values range from *frequent* to *extremely improbable* along with their respective probabilities of occurrence and probability numbers are given in Table 5. Since the game is focused on estimating the overall NAS accidents causing fatalities, the FAA's quantitative probability definition for NAS systems and

Probability ↓	NAS systems and ATC opera- tional (quantitative)	Probability of occurrence	Probability number N
Frequent A	Probability of occurrence per operation/operational hour is equal to or greater than 1×10^{-3}	$P \ge 1 \times 10-3$	$N \leq 3$
Probable B	Probability of occurrence per operation/operational hour is less than 1×10^{-3} , but equal to or greater than 1×10^{-5}	$1 \times 10^{-3} > P \ge 1 \times 10^{-5}$	$3 < N \leq 5$
Remote C	Probability of occurrence per operation/operational hour is less than or equal to 1×10^{-5} , but equal to or greater than 1×10^{-7}	$1 \times 10^{-5} > P \ge 1 \times 10^{-7}$	5 <i>< N</i> ≤ 7
Extremely remote D	Probability of occurrence per operation/operational hour is less than or equal to 1×10^{-7} , but equal to or greater than 1×10^{-9}	$1 \times 10^{-7} > P \ge 1 \times 10^{-9}$	7 < N ≤ 9
Extremely improbable E	Probability of occurrence per operation/operational hour is less than 1×10^{-9}	$P < 1 \times 10^{-9}$	<u>N</u> > 9

 Table 5
 Probability definitions

ATC Operational definitions are adopted. The baseline accident rates (0.208/100,000 FH for 2000–2009) indicate an initial average of N = 5.681 using (designated as Remote) before any NextGen-related technology implementation is present. The current aviation statistics indicate that the initial risk value is identified in the "Low Risk" area, located at the intersection of severity 4 and likelihood C.

3.4.2 Accident and Enabler Categories

By definition, number of accidents and their consequences become greater as the NAS capacity increases. In order to ensure the risk level remains at acceptable limits, participants are required to invest in various accident mitigation technologies and/or enablers. The enabler categories were selected based on the National Transportation Safety Board (NTSB) aviation accident statistics and potential future accident areas with the introduction of increased traffic within the FAA Part 121–Commercial Air Carrier Category. The categories include:

- runway safety and collision avoidance, including runway capacity and visibility
- *aircraft systems reliability*, including propulsion health, airframe health, and software health management systems
- *icing mitigation*, including airframe aerodynamic modeling and atmospheric modeling
- airborne collision, including Near mid-air collision (NMAC) and loss of separation
- weather (thunderstorm) avoidance, including thunderstorm and visibility
- turbulence avoidance as in-flight turbulence and wake turbulence.

For each accident category, several enablers or mitigation technologies developed under various programs within NASA's Aviation Safety and FAA [33, 44]. These mitigation technologies and their respective information on cost, operational timeline, and content were obtained from the (JPDO) and other sources [13, 23, 24, 30–32]. The enabler categories and the respective technologies/methods are limited to safety related areas; technologies associated with increased capacity or reduced environmental impact goals are not within the scope of this research. A comprehensive list of included technologies and respective cost, timeline, and operational benefits can be found at Ancel [1].

Within the gaming cycle, the selection of the enablers is done by the participants of the relevant stakeholder groups. Participants decide on the timeline and collaborations regarding the adoption of the predetermined enablers under several categories. Participants are asked to evaluate enabler benefits, costs, mixed equipage risk and implementation timeline, then review their budget and plan for the near future in order to make the decision about when to "acquire" the enablers and how to construct collaborations whenever it is possible. During this process, the participants were encouraged to use Logical Decisions for Windows® (LDW) software in order to support enabler selection process by dynamically adjusting utilities to determine the ranking. Figure 5 provides a snapshot of enabler ranking for Airlines stakeholder.



Fig. 5 Airlines enablers LDW snapshot

3.5 Post-gaming: Data Collection, Analysis, and Observations

The proposed game setup presented was executed several times for testing and calibration purposes among different groups. However, the data presented below was generated in a session played on February 14, 2011 with contributions from aviation professionals working at NASA Langley Research Center.

One of the most tangible outcomes of the gaming exercise is the 2025 NAS safety values with respect to the FAA's Risk Matrix (Fig. 6) acceptability measures along with the intermediate risk values during the transitional technology implementation phase. The cumulative effect of various safety-related technological implementations within the NextGen operating environment help decision makers to define technologies or areas that require further analysis and understanding. Also, throughout the gaming effort, discussions and possible negotiations within the opposing parties are important findings that can lead to different constructive problem-solving approaches.

Severity	<0.30	0.30	0.38	0.45	0.53	0.60	0.68	0.75	0.83	>0.9
Likelihood	Minimal	Mi	Minor		Major		Hazardous		strophic	N Scale
Frequent										N ≤ 3
Probable				2014	4					3 < N ≤ 5
Remote	2010			20	25					5 < N ≤ 6 6 < N ≤ 7
Extremely	2015	2017-	2020							7 < N ≤ 8
Remote										8 < N ≤ 9
Extremely Improbable										N > 9

Fig. 6 Evolution of NAS safety values with time

3.5.1 Data Collection Mechanism

In order to aggregate and process the data, the *serious gaming platform* presented within the previous chapter is coupled with the *data aggregation platform*, a designated, comprehensive Excel® file assigned to calculate and communicate the dynamic NAS Risk values and other statistics among players and facilitators. The data aggregation platform contains all the financial relationships, accident statistics, and risk assessment model calculations necessary to generate interim safety values and other statistics (Fig. 7).

3.5.2 Generic Scenarios and Sample Data

The gaming session is initiated with "Pushing the Envelope" scenario starting from 2010 until 2015, followed by "Grounded" for 2 years. The game then assumes 5 years of "Regional Tensions" starting in 2017 and is finalized by 4 years of "Environmentally Challenged" scenario, from 2022 until 2025 (see Table 2 for scenario parameters).



Fig. 7 Data elicitation mechanism

Pushing the Envelope. This scenario phase represents a rapid growth in the aviation industry, starting with the game year 2010, in-line with the FAA's FY2010 expectations. Within the scenario, the increase in air transportation capacity pushes the accident likelihood toward "probable" where accident severity remains with "minor" consequences. The scenario also depicts a continuously growing strong

economy and a liberal trade policy environment, allowing stakeholders to regulate the market. During this timeframe, stakeholders are required to invest in transportation infrastructure components like ADS-B initiation, Data Link setup, and many other enablers to accommodate the anticipated increase in air travel.

Grounded. In years 2015 and 2016, the air transportation capacity is largely hampered by a scenario-driven series of terrorist attacks. This scenario was generated by the NRC study from 1997, somehow portraying the September 2011 events. Within the NRC study, terrorist attacks are caused by large gap between the income levels and living standard of developed nations compared to second or third world countries. The scenario for these 2 years is called "Grounded" where air travel is no longer safe and a sharp decline in the air transportation capacity takes place.

Regional Tensions. This scenario represents a changing global scenario where harmonious globalization is no longer available. Although demand for aeronautics products and services is back up, increased oil cost deeply affects airline companies. Due to the initial NextGen enabler investments, the NAS safety values are better compared to baseline 2010 levels with less likelihood of accident. For the years 2017–2020, the increase in air transportation capacity does not deteriorate NAS safety. Even with a considerable terrorist attack risk, air transportation stays rather stable and safe.

Environmentally Challenged. Initiated in 2022, this scenario simulates a very CO_2 conscious world where carbon-based fuel usage is very limited and resources are costly. High fuel prices are reflected to all available transportation modes. The NAS safety values start to migrate toward the unacceptable areas due to increased capacity levels, but the unfavorable economic environment prevents further capacity growth, and final air transportation safety values stay within the acceptable limits. During the gaming exercise, at the end of year 2025, the likelihood of an accident stayed within the "remote" area; however, the accident consequences migrated toward "major" category due to increased aircraft capacity (Fig. 6).

3.5.3 Sample Stakeholder Specific Variables

Government Stakeholder Variables. The government variables (income tax, environmental tax, security tax, and fuel tax) are given in Fig. 8. As expected, during the "Pushing the Envelope" era (2010–2014), the U.S. economy is strong and tax rates are relatively low, since there are no terrorist or environmental concerns, there is no taxation on these areas. With the introduction of the "Grounded" scenario, air transportation industry faced a steep increase in security and income taxes in order to compensate for elevated global terror risk and declining economic status. During the "Regional Tensions" era, the security threat remains stable, with constant increase in income taxes and a slight increase in environmental taxes. Due to the decline in U.S. economic competitiveness and the disruption global structure, starting from year 2017, the government started to collect taxes from air transportation stakeholders.



Fig. 8 Government variables-taxes and total fuel price

Airport Stakeholder. The main operating revenue for the airport operators are aeronautical revenues (passenger airline landing fees, terminal arrival fees, rents and utilities), and non-aeronautical revenues (terminal food and beverage, retail stores, and duty free). Figure 9 outlines the airport variables and the capacity change. As anticipated, during the competitive air transportation environment (2010–2014), airport charges are rather constant, and they are generating low income. During the "Grounded" era, the fees climb in order to compensate for increased governmental taxes and increased NextGen-related expenses. From 2017 until 2022, airports raise fees steadily mostly since the air transportation remains the main choice of transportation in the United States. Due to the competition between the participants, the airport stakeholder increased the landing fees toward the end of the game when the air transportation capacity reached around 185 % of the 2010 values.

Corporate Airlines Stakeholder. In order to compensate for the large acquisitions mandated by the FAA, the corporate airlines raised all of their fees throughout the game (Fig. 10). One prominent observation that surfaced during the gaming exercise was the increase of ticket prices when the airline expenses (i.e., taxes, fuel, etc.) are elevated. However, even when taxes are back to their normal values, the airlines did not reflect the relief in their fees, which is in accord with a real-world environment. Like the airports stakeholder, corporate airlines had to lower their



Fig. 9 Airport variables-landing fees and concession and parking fees versus capacity



Fig. 10 Corporate airlines variables-ticket fees, cancellation fees, and baggage fees

ticket fees when the general public stakeholder reacted and adjusted the passenger capacity. During that time, the raw ticket fee was decreased from \$205 to \$179; however, it climbed back up to \$209 once passenger capacity recovered. Throughout the game, passengers experienced a more than \$50 increase in ticket prices (\$377 compared to \$325 in 2010, after the government taxes are reflected). Baggage fees were increased from \$25 to \$31 while reservation cancellation fees went up from \$150 to \$198 over the course of 16 years.

Airline companies are highly susceptible to jet fuel price, especially for the future capacity values reaching almost three times the current state. Although passenger capacity reached 240 % of 2010 values and ticket prices were increased more than 15 %, baggage fees more than 25 %, and reservation fees more than 30 %, corporate airlines still stayed below the profit margin level experienced at years 2017 and 2019.

General Public Stakeholder. The general public indirectly decides air travel passenger capacity at the end of each time-step by comparing the transit time and cost for the two predetermined routes. These one-way routes are the 228 mile Washington, DC (Union Station) to New York, NY (Penn Station) and 437-mile Washington, DC (Union Station) to Boston, MA (South Station) routes. As of March 2011, these two routes are the only two high-speed rail routes existing in the United States (*Acela Express* by Amtrak⁶). The three modes of transportation considered are rail, automobile, and air transportation. Travel times and costs include to and from the train stations and airports. Also, an automobile with 25 mpg, \$3.113/gallon national gas average was also assumed.

The Fig. 11 shows cost and transit times for the three modes of transportation with respect to the simulation year for the first configuration, from Washington, DC to New York. For this particular trip setup, driving is the lowest cost option where the rail and air options are converging toward the end of the scenario timeframe. With the introduction of future high-speed rail systems, it is assumed that rail prices will rise in order to compensate for increased infrastructure investments while transit times will be faster. Rail option is consistently faster than the other two methods, where automobile and flight transit times vary based on the scenario parameters.

The second configuration variables are given in Fig. 12. Due to increased travel distance, the train mode is not considerably cheaper than the air mode, but it is still much slower.

The air transportation mode provides the fastest service with the highest cost until around the year 2020, when High Speed Rail infrastructure starts to offer faster service times. By the end of the simulation, transit times for air and train modes of transportation are comparable, and costs for both of the modes are on the rise.

The general public stakeholder participant adjusted air transportation capacity on six occasions throughout the game (Table 6). The 2015 terrorist attacks hampered air transportation capacity 10 % more than anticipated; however, even with the same terror risk, in the following year, the perceived terror risk was lower than

⁶ http://www.amtrak.com/home, Retrieved July 3rd, 2014.



Fig. 11 General public announcement variables (configuration 1)



Fig. 12 General public announcement variables (configuration 2)

projected. Even with the higher transportation costs and slower travel speeds, the air transportation mode was adjusted by the general public stakeholder and reached 240 % of the 2010 passenger capacity. This resulted in over 1.5 billion passengers in the NAS.
Year/scenario	Capacity adjustment (%)	Reason provided
2014/Pushing the envelope	+5	Strong U.S. Economy + relatively inexpensive transportation fees
2015/Grounded	-10	Perceived terror risk higher than anticipated
2016/Grounded	+10	Ongoing perceived terror risk, lower than anticipated
2022/Environmen- tally challenged	+5	No increase on air transportation fees
2023/Environmen- tally challenged	-5	High speed rail presence and the increase air travel cost
2024/Environmen- tally challenged	+10	Capacity increase in response to steep decrease in air transportation fees previous year
2025/Environmen- tally challenged	+10	Continuing satisfaction from air transportation services

 Table 6
 Public intervention values and provided reasons

3.5.4 Enabler Acquisition Timeline and Surfaced Strategies

The selected 18 enablers from seven different categories are all implemented within the first 3 years of the game timeline. Since gaming participants were experts in the aviation field, they were aware of the necessity of key enablers like ADS-B and Data Link, along with other safety-related enablers (Table 7). The acquisition cost for the ADS-B and the Data Link were collected by increased ticket fees and other fees charged by the airline companies. The corporate airline stakeholder compensated for the majority of the widespread application of ADS-B technology which is the main enabler for many NextGen technologies. The Data Link acquisition was realized by the government contribution, around 25 % of the Data Link acquisition cost over the 11 years.

4 Bridging the Gap Between Game Theory and Gaming

This section briefly makes the first step in bridging the gap between game theory and gaming. That is, we cursorily stipulate the NextGen problem in game theoretic propositions. For our four kinds of players, we propose 4, 2, 2, and 1, respectively, strategies where the player can ordinally choose a high versus low value, i.e.,

Government, FAA, and Military Stakeholder: High taxes versus low taxes, enforce versus not enforce aviation safety, spare funds versus not spare funds, and reflect versus not reflect the military agenda.

Corporate Airlines: Determine high versus low prices for tickets, cancellation, baggage; offer customers high versus low quality.

Table 7Enabler acquisition timeline and s	trategies					
Enabler package definition/denomination	Enabler acquisition planned	Enabler acquisition completed	Primary cost bearer	Other cost bearers	Total cost	Coalition surfaced
ADS-B	2011	2021	Airline: \$2,486 M	Airport: \$33 M	\$2,519M	Airport and airline
Data link	2012	2021	Airline: \$3,410 M	Gov't: \$1,100 M	\$4,510M	Airline and Gov't
Capacity/safety related runway enablers/ R1	2012	2019	Airport	N/A	\$78M	N/A
Runway visibility/R2	2010	2020	Airport	N/A	\$452M	N/A
Collision—NMAC/C1	2011	2017	Airport	N/A	\$98M	N/A
Collision—loss of separation/C2	2013	2018	Airport	N/A	\$570M	N/A
A/C powerplant/AC1	2010	2023	Airline	N/A	\$112M	N/A
A/C structures/AC2	2010	2023	Airline	N/A	\$112M	N/A
A/C Systems/AC3	2010	2023	Airline	N/A	\$154M	N/A
Icing—structures/I1	2010	2023	Airline	N/A	\$234M	N/A
Icing—engine/I2	2010	2023	Airline	N/A	\$392M	N/A
Weather	2012	2017	Airline	N/A	\$60M	N/A
Weather-visibility/W2	2010	2016	Airline	N/A	\$70M	N/A
Turbulence-in flight/Tl	2011	2018	Airline	N/A	\$336M	N/A
Turbulence—ground wake/T2	2010	2020	Airport	N/A	\$132M	N/A
Enhanced post-impact fire/smoke miti- gation/S1	2013	2020	Airline	N/A	\$400M	N/A
Improved crash survivability of aircraft structures/S2	2013	2020	Airline	N/A	\$400M	N/A
Improved evacuation and accident response procedures/S3	2013	2020	Airport	N/A	\$240M	N/A

A Simulation Game Application for Improving the United States ...

Table 8 2×2 Game between General Public and Corporate Airlines when Government, FAA, andMilitary Stakeholder choose high taxes and Airport Operators choose high aeronautical and non-
aeronautical fees

		Corporate airlines	
		High prices Low prices	
General public	General public Travel frequently		Intermediate, low
	Travel infrequently	Low, low	Intermediate, low

Airport Operators: Charge high versus low aeronautical fees (including airport landing fees), charge high versus low non-aeronautical fees (parking fees, concession fees, airport shop rental fees, etc.)

General Public: Choose to travel frequently or infrequently (i.e., substituting or not with other transportation methods).

It is beyond the scope of this chapter to provide a full mathematical formulation of this problem, but we here provide five 2×2 payoff matrices [39].

Table 8 shows the game between the General Public and Corporate Airlines when the Government, FAA and Military Stakeholder choose high taxes and Airport Operators choose high aeronautical and non-aeronautical fees. The payoff before the comma is to the row player (General Public), and the payoff after the comma is to the column player (Corporate Airlines). First (lower left cell), when the General Public travels infrequently and the Corporate Airlines choose high prices, they both earn low payoff. Second (upper left cell), when the General Public travels frequently and the Corporate Airlines keep the high prices, the General Public still earns low payoff (because of the high prices), while the Corporate Airlines earn intermediate payoff (because of the larger travel volume). Third (upper right cell), when the General Public travels frequently and the Corporate Airlines choose low prices, the payoffs are reversed. That is, the General Public benefits from the low prices and earns intermediate payoff, while the Corporate Airlines suffer from the low prices and earn low payoff. Finally (lower right cell), when the General Public travels infrequently and the Corporate Airlines choose low prices, the General Public still earns intermediate payoff (because of the low prices), while the Corporate Airlines earn low payoff (because of the low prices). With these payoffs, which are merely ordinal at two levels, the General Public is indifferent between frequent and infrequent travel, while the Corporate Airlines prefer high prices when the General Public travels frequently, and is otherwise indifferent. This gives a Nash [34] equilibrium with high prices, shown in bold in Table 8. A Nash [34] equilibrium is a state of affairs from which no player prefers to deviate unilaterally.

Table 9 shows the game between the General Public and Corporate Airlines when the Government, FAA, and Military Stakeholder choose high taxes and Airport Operators choose high aeronautical and non-aeronautical fees. This causes seven payoffs in Table 8 to increase while the General Public's payoff from infrequent travel with high prices remains low. In Table 9 the Corporate Airlines are indifferent between high and low prices for any given strategy by the General

 Table 9
 2×2 Game between General Public and Corporate Airlines when Government, FAA, and

 Military Stakeholder choose low taxes and Airport Operators choose low aeronautical and non-aeronautical fees

		Corporate airlines	
		High prices Low prices	
General public	Travel frequently	uently Intermediate, high high, high	
	Travel infrequently	Low, intermediate	High, intermediate

Table 102×2 Game between Airport Operators and Government, FAA, and Military StakeholderCorporate Airlines when Corporate Airlines choose low prices and General Public choose to travelfrequently

		Airport operators	
		High prices	Low prices
Government, FAA, and military stakeholder	High taxes	High, intermediate	High, low
	Low taxes	Intermediate, high	Intermediate, intermediate

Public, while the General Public prefers frequent travel when prices are high. This gives a Nash [34] equilibrium with frequent travel, shown in bold in Table 9.

Another 2×2 game can be observed between the Airport Operators on the one hand and the Government, FAA and Military stakeholders (Government for short) on the other hand. Table 10 shows the game between these two stakeholders when the Corporate Airlines choose to charge low prices while the general public travels infrequently. First (lower left cell), the Government earns intermediate and the Airport Operators earn high payoff when the General Public travels frequently. Second (upper left cell) quadrant shows intermediate payoff for the Airport Operators since despite the large travel volume, the high Government taxes still hampers profit whereas the Government enjoys high payoff given they chose high taxes. Third (upper right cell), choosing low prices and high Government taxes, the Airport Operators' payoff is low, but the Government enjoys high tax income revenue and thus high payoff. Finally (lower right cell), although the Airport Operators and Government choose low prices, high travel volume yields to an intermediate payoff to both players due to the frequent General Public travel. In Table 10, the Government prefers high taxes regardless which strategy the Airport Operators choose, and the Airport Operators prefer high prices regardless which strategy the Government choose. This gives one unique Nash [34] equilibrium with high taxes and high prices in the upper left cell, shown in bold in Table 10.

Table 11 shows the game between the Government and Airport Operators when the Corporate Airlines choose high prices while the General Public travels infrequently. Consequently, the cells given in Table 11 show lower payoffs than the values in Table 10. The presence of infrequent air travel affects both the Airport Operators and the Government equally, but additionally, high prices set by the

 Table 11
 2×2 Game between Airport Operators and Government, FAA, and Military Stakeholder

 Corporate Airlines when Corporate Airlines choose high prices and General Public choose to
 travel infrequently

		Airport operators	
		High prices	Low prices
Government, FAA, and military stakeholder	High taxes	Intermediate, low	Intermediate, low
	Low taxes	Low, intermediate	Low, low

Corporate Airlines further affect Airport Operators payoffs. For that reason, in all cells except for the lower left, the Airport Operators have low payoff. In the lower left cell, the Airport Operators earn intermediate payoff due to relaxed Government taxation and set high prices.

The Government payoffs in all cells were degraded by one level (i.e., high was degraded to intermediate, and intermediate to low) due to low air traffic volume. Also, it is possible to observe that the Government payoffs shown in Tables 10 and 11 were only affected by the air traffic volume and not by the Airport Operator selections. However, the Government as well as air traffic volume and Corporate Airlines selections affected the Airport Operators payoffs. In Table 11, as in Table 10, the Government prefers high taxes regardless which strategy the Airport Operators choose. However, the Airport Operators are indifferent between high and low prices when the Government choose high taxes, and prefers high prices when the Government chooses low taxes. As in Table 10, this gives one unique Nash [34] equilibrium with high taxes and high prices in the upper left cell, shown in bold in Table 11. However, the Nash equilibrium is less stable and easily perturbed since the Airport Operators are indifferent between high and low prices in the Nash equilibrium.

Six 2×2 games can be specified between any two players out of a total of four players, i.e., 12, 13, 14, 23, 24, and 34. If four different conditions are assumed for the other players, for each of the four games, we get a total of 24 2×2 games. Furthermore, games can be set up between subplayers of the four kinds of players, e.g., government against military stakeholders, individual corporate airlines against each other, of different segments of the general public against each other. Future research should develop payoff functions for the various players to scrutinize the equilibria more thoroughly.

5 Conclusions and Discussions

In order to assist the planning of infrastructure transitions, the current research pursued the development and deployment of a simulation game to serve as a venue to generate and evaluate data that may complement historical data. The simulation game was applied to the NextGen framework which involves the transformation of current US airspace system. The data extracted from a gaming run played by SMEs identified that under varying scenario parameters, the aviation safety values remained within the acceptable region during the NextGen transition. Also, it was observed that discussions surfaced throughout the game allowed SMEs to experience the issues associated with NextGen due including technical, financial, and social challenges. The SMEs believed that a more comprehensive game design with detailed stakeholder interactions might reveal strategies that can assist planning and implementation of infrastructure transitions. The gaming exercise can potentially be applied to any large-scale infrastructure system to generate preliminary data regarding the systems' characteristics such as capacity, power generation, throughput, risk evaluation/acceptance, resources prioritization, identification of potential future issues (strategic behavior, public perception, etc.), while considering both social and technical aspects of the system. The linkage between game theory and gaming was briefly considered by setting up some simple 2×2 games between a subset of the four kinds of players. Future research should extend that linkage.

5.1 Limitations of the Study

Unlike computer models or other hard-science alternatives, the development, execution, and validation phases of the research design require extensive SME contribution. For that reason, the gaming environment delineates the physical presence of all the prominent stakeholders of that particular infrastructure system on several occasions, which creates logistical challenges such as travel, scheduling, cost, etc.

Besides logistics limitations, the interdisciplinary nature of simulation and gaming, in most cases, limits the use of this approach for educational/training purposes in businesses. Researchers believe that further work must be performed to theorize and establish serious gaming as a field of study; whether simulation and gaming is a beneficial tool or an academic field is still an uncertainty [42].

The accurate representation of a large and complex system, fusing multiple perspectives and multiple disciplines, has proven challenging in practice since the selection of scenario elements and the actual composition of the scenario are based primarily on the game developer's perspective. Similar to the scenario construct, the abstraction of the elements of the reference system and translating them to the model poses a challenge. Consequently, the NextGen game demonstrated in this chapter contains major assumptions and simplifications applied to the game rules, players, and their interactions, which limits the scope of the project. For instance, the consequence calculations within RRAM were limited to fatalities whereas other consequences such as accidents, serious injuries, hull losses, and extended accident damage to the surrounding environments and associated costs were excluded. Also, since the timeframe was over 15 years, the technologies, their implications, benefits, and costs were not well defined and demanded simplifications.

public stakeholder was represented by an artificial participant whereas it might be possible to use alternative methods to elicit the behavior of this stakeholder (realtime crowdsourcing, polls or other methods). However, given the dynamic structure of the game where the feedback from the public stakeholder is considered by other stakeholders in adjusting their strategies, the use of alternative aimed to elicit public behavior might be limited.

Finally, due to lack of resources and time, the gaming exercise was only executed once and resulting data was presented to demonstrate the game outputs. In order to generate useful data, the game needs to involve more players and to be executed several times where a statistically meaningful database can be obtained. Similarly, since the validation of the gaming activity also relies on SME input, limited number of game runs prevented a thorough validation.

5.2 Validation

The validity of the model presented in this chapter was investigated using definitions given by Greenblat [17] and Peters et al. [38]. The validation parameters given by Greenblat include face validity, empirical validity, and theoretical validity whereas Peters et al. [38] discuss the validity of games from various aspects such as psychological reality, structural validity, process validity, and predictive validity. The face validity (similar to psychological reality defined by Peters et al.) refers to the realistic gaming environment experienced by the participants. For a game to be valid, the environment must portray similar characteristics to the reference system. The empirical validity suggested by Greenblat designates the closeness of the game structure to the reference system. However, Peters et al. separate this concept into two sections: structural validity (including the game structure, theory and assumptions) and process validity (concerning the information/resource flows, actor interactions, negotiations, etc.). For the simulation to be valid, all the elements of the game (actors, information, data, laws, norms, etc.) should be isomorphic, meaning the elements and relations do not necessarily have to be identical but should be able to demonstrate congruency between them. Finally, the last element covered by both definitions is related to the *theoretical* (or *predictive*) validity: the models' ability to reproduce historical outcomes or predict the future, and conform to existing logical principles.

The current research relies heavily on subjective assessments obtained from experts at all levels (pre-gaming, gaming, and post-gaming phases) including the validation of the simulation game mostly due to lack of predictable data regarding the future states of the current NAS. Consequently, the *theoretical* (or *predictive*) element of the validation is challenging to obtain simply because the outcomes of the game (i.e., 2025 NAS safety values) cannot be put to test. The *empirical* (*structural/process*) as well as *face* (*psychological reality*) validity assessments were obtained by the SMEs both during and after the gaming exercise via the validation questionnaire. The questionnaire is aimed to acquire SME's ranking of

the various validation categories including elements from face and structural validity. However, due to the limitations highlighted in the previous section, statistical analysis of questionnaire results were not performed due to lack of sufficient data points.

Acknowledgments The authors would like to thank the Subject Matter Experts at System Analysis and Concepts Directorate (SACD) at NASA Langley Research Center for their guidance with the development and execution of the gaming exercise. The authors thank the editor for aid with Sect. 4.

References

- 1. Ancel E (2011) A systemic approach to next generation infrastructure data elicitation and planning using serious gaming methods. Doctoral Dissertation, Old Dominion University, Norfolk, VA
- 2. Bedford T, Cooke R (2001) Probabilistic risk analysis foundations and methods. Cambridge University Press, United Kingdom
- 3. Bekebrede G (2010) Experiencing complexity: a gaming approach for understanding infrastructure systems. Doctoral Disseration, Delft University of Technology, the Netherlands
- 4. Bekebrede G, Mayer I, van Houten SP, Chin R, Verbraeck A (2005) How serious are serious games? Some lessons from infra-games. In: Proceedings of digital games research association (DiGRA) conference: changing views—worlds in play, Vancouver, Canada, 16–20 June 2005
- 5. Brewer G (2007) Inventing the future: scenarios, imagination, mastery and control. Sustain Sci 2:159–177
- Chappin EJL, Dijkema GPJ (2008) On the design of system transitions: Is transition management in the energy domain feasible? Paper presented in IEEE international engineering management conference, Estoril, Portugal, 28–30 June 2008
- 7. Dormans J (2011) Beyond iconic simulation. Simul Gaming 42:610-631
- 8. Duke RD (1974) Gaming: the future's language. Wiley, New York
- 9. Duke RD (1980) A paradigm for game design. Simul Gaming 11:364-377
- Elzen B, Wieczorek A (2005) Transitions towards sustainability through system innovation. Technol Forecast Soc 72:651–661
- 11. FAA (2008) Air traffic organization safety management system manual, Washington, DC
- 12. FAA (2009) FAA's NextGen implementation plan 2009, Washington, DC
- 13. FAA (2010) FAA's NextGen implementation plan 2010, Washington, DC
- 14. Frantzeskaki N, Loorbach D (2008) Infrastructures in transition: role and response of infrastructures in societal transitions. Paper presented in international conference on infrastructure systems, NGInfra Foundation, Rotterdam, Netherlands, 10–12 Nov 2008
- Geels FW (2005) Processes and patterns in transitions and system innovations: refining the coevolutionary multi-level perspective. Technol Forecast Soc 72:681–696
- Geurts JLA, Duke RD, Vermeulen PAM (2007) Policy gaming for strategy and change. Long Range Plan 40:535–558
- 17. Greenblat CS (1975) From theory to model to gaming-simulation: a case study and validity test. In: Greenblat CS, Duke R (eds) Gaming-simulation: rationale, design, and applications, Wiley, New York
- Greenblat CS, Duke R (1975) Gaming-simulation: rationale, design and applications. Wiley, New York
- Hansman RJ, Magee C, de Neufville R, Robins R, Roos D (2006) Research agenda for an integrated approach to infrastructure planning, design and management. Int J Critical Infrastruct 2:146–159

- 20. International Atomic Energy Agency (1996) Manual for the classification and prioritization of risks due to major accidents in process and related industries, Vienna
- Janssen M, Chun SA, Gil-Garcia JR (2009) Building the next generation of digital government infrastructures. Gov Inf Q 26:233–237
- 22. JPDO (2007) Concept of operations for the next generation air transportation system (NextGen), Washington, DC
- 23. JPDO (2008) Next generation air transportation system integrated work plan: a functional outline v1.0, Washington, DC
- 24. JPDO (2008) NextGen avionics roadmap v1.0, Washington, DC
- 25. Klabbers JHG (1996) Problem framing through gaming: learning to manage complexity, uncertainty, and value adjustment. Simul Gaming 27:74–92
- Luna-Reyes LF, Zhang J, Gil-García JR, Cresswell AM (2005) Information systems development as emergent socio-technical change: a practice approach. Eur J Inf Syst 14:93–105
- 27. Mayer I (2009) The gaming of policy and the politics of gaming: a review. Simul Gaming 40:825–862
- Mayer I, Bockstael-Blok W, Valentin EC (2004) A building block approach to simulation: an evaluation using CONTAINERS ADRIFT. Simul Gaming 35:29–52
- 29. Meijer S (2009) The organisation of transactions: studying supply networks using gaming simulations, Wageningen Academic Publishers, the Netherlands
- 30. NASA (2009) Integrated intelligent flight deck technologies: technical plan summary, Washington, DC
- 31. NASA (2009) Integrated resilient aircraft control: stability, manueverability, and safe landing in the presence of adverse conditions, Washington, DC
- 32. NASA (Integrated vehicle health management: automated detection, diagnosis, prognosis to enable mitigation of adverse events during flight, Washington, DC
- 33. NASA (2010) FY 2010 aeronautics budget estimate, Washington, DC
- 34. Nash J (1951) Non-cooperative games. Ann Math 54(2):286-295
- 35. National Research Council (1997) Maintaining US leadership in aeronautics: scenario-based strategic planning for NASA's aeronautics enterprise. National Academy Press, Washington, DC
- 36. National Science and Technology Council (2010) National aeronautics research and development plan, Washington, DC
- NTSB (2010) Aviation accident statistics. http://www.ntsb.gov/data/aviation_stats.html. Accessed 3 July 2014
- 38. Peters V, Vissers G, Heijne G (1998) The validity of games. Simul Gaming 29:20-30
- 39. Rapoport A, Guyer M (1966) A taxonomy of 2 × 2 games. Gen Syst 11:203-214
- 40. Rittel HWJ, Webber MM (1973) Dilemmas in a general theory of planning. Policy Sci 4:155–169
- 41. Roos D, de Neufville R, Moavenzadeh F, Connors S (2004) The design and development of next generation infrastructure systems. In: Proceedings of the IEEE international conference on systems, man and cybernetics. The Hague, Netherlands, pp 4662–4666, 10–13 Oct 2004
- 42. Shiratori R (2005) Toward a new science of simulation and gaming: ISAGA and the identity problem of simulation and gaming as an academic discipline. In: Shiratori R, Arai K, Kato F (eds) Gaming, simulations and society: research scope and perspective. Springer, Tokyo, pp 3–8
- 43. US Congressional Research Service (2003) Critical infrastructures: what makes an infrastructure critical? The Library of Congress, Washington, DC, Jan 29 2013
- 44. US Department of Transportation (2010) Budget estimates FY2011. Federal Aviation Administration, Washington, DC
- 45. US Department of Transportation (2010) Memorandum: timely actions needed to advance the next generation air transportation system, Washington, DC

- 46. Vries LJD, Subramahnian E, Chappin EJL (2009) A power game: simulating the long-term development of an electricity market in a competitive game. Paper presented at the NGInfra developing 21st century infrastructure networks, Chennai, India 9–11 Dec 2009
- 47. Wenzler I (2005) Simulations and social responsibility: why should we bother? In: Shiratori R, Arai K, Kato F (eds) Gaming, simulations and society: research scope and perspective. Springer, Tokyo, pp 139–148
- Wenzler I, van Muijen S (2009) Simulation games for managing change. Paper presented at the 40th annual conference of international simulation and gaming association, Singapore, June 29–3 July 2009

A Congestion Game Framework for Emergency Department Overcrowding

Elizabeth Verheggen

Abstract Hospitals often manage capacity and resource constraints by different strategies implemented at their system access points. Emergency departments are key portals where timely access to care is a crucial quality of service and safety metric. Individuals vying for both urgent and nonurgent care seek these services analogous to the Tragedy of the Commons archetype. In a commons, a resource is used as if it belonged to everyone. Competition for a finite, decentralized, and shared resource risks its depletion as individuals optimize their own objectives while impacting the choices of others. As a result, overall system performance degrades. Ambulance diversion, extensive wait times and patient elopements, referred to as left without being seen, epitomize overutilization and inefficient load balancing. Traditionally, many hospitals were able to build their way out of congestion. Adding capacity, however, is at odds with concerted efforts to reign in the costs of health care. In an effort to break with this tradition, we exploited insights from game theory to inform the development of policies for more effective capacity management related to emergency department use, and to highlight related challenges. We examined emergency department overcrowding within the framework of a congestion game, the El Farol Bar Game and its variants, which illustrate the Tragedy of the Commons. In a series of agent-based simulations of the games, we found no statistically significant difference between the predictions of two games and our empirical observations during our most congested time periods of nonurgent patient attendance. Given the new competitive social context of real-time publicly advertised door-to-doctor wait times, and the implications that burgeoning information technologies have for the strategies invoked by providers and patients, it seems a bar might be the best metaphor to understand emergency department congestion.

Safety and Security, Springer Series in Reliability Engineering, DOI 10.1007/978-3-319-13009-5_9

E. Verheggen (🖂)

Lehigh Valley Health Network, Allentown, PA, USA e-mail: Elizabeth.Verheggen@lvhn.org

E. Verheggen University of South Florida College of Medicine, Tampa, FL, USA

[©] Springer International Publishing Switzerland 2015 K. Hausken and J. Zhuang (eds.), *Game Theoretic Analysis of Congestion*,

Keywords Emergency department overcrowding \cdot Tragedy of the commons \cdot Game theory \cdot Congestion games \cdot El Farol Bar Game \cdot Minority game

1 Introduction

How do individuals with low-severity medical needs utilize an emergency department (ED) as wait time information becomes publicly available? We explore the possibilities for their arsenal of strategic behavior in this work. The El Farol Bar Game, also known as the Santa Fe Bar problem, introduced by Arthur [1], provides the context for examining how individuals optimize their own behavior in crowding situations that conflict with the collective social optimization of a decentralized resource. It is intuitive that congestion affects the efficient use of a health system's ED "commons." The majority of studies of ED crowding have been done in the US, however, around the world, the disequilibrium between supply and demand in EDs with different care delivery and payment systems has become a public health problem. Fifteen countries outside of the US have been studied in the recent literature to understand the causal mechanisms of worsening rates of congestion. These include: Australia, Canada, Denmark, Finland, France, Germany, Hong Kong, India, Iran, Italy, The Netherlands, Saudi Arabia, Catalonia (Spain), Sweden, and the United Kingdom [2]. International overcrowding studies address a variety of endpoints, including operational efficiency, safety, and quality, as well as health outcomes. Various methods have been used, including: econometric analyses (time series regression); operations research and management (simulation, system dynamics and queuing theory); computational models (agent-based modeling); and mixed models (qualitative survey and quantitative analysis). A European perspective discussing the multifactorial nature of this worldwide phenomenon suggests a direct link to public funding of hospital beds, staffing, and community care facilities as main contributors to crowding [3]. A 40 % increase in annual ED visits over 12 years in Taiwan, motivated research to predict visitor volume to address the mismatch between ED capacity and input, throughput, and output factors of ED operations management [4]. Patient elopements, commonly referred to as left without being seen (LWBS) in the ED literature, have been studied in the UK [5, 6]Canada [7], Switzerland [8] Pakistan [9], and Australia [10, 11]. Use of the ED for nonurgent complaints has been addressed considering general practitioner characteristics [12, 13] in The Netherlands and Italy. Operations research methods based on a system dynamics simulation were implemented to improve patient flow in a Malaysian ED [14]. Canadian studies have used agent-based modeling (ABM) to address over taxed EDs functioning in a compromised situation of many competing priorities, [15] as well as load balancing optimization via "crowdinforming" simulations in [16]. Three Swedish EDs were modeled to improve waiting management using a grounded theory mixed models approach, [17] while six Swedish EDs were studied in [18] in an effort to build a registry for quality of care measures which improved patient throughput and presentation times. In the UK, the 4 h NHS standard for emergency care was investigated [19]. In a related inquiry to our work considering published wait times, a small sample of ED patients was surveyed to confirm a limited awareness of published waits for two Canadian EDs [20]. Other research from a Canadian perspective studied ED crowding along the lines of appropriate measurements, threats to patient dignity and its power as an agent of change, and improved turnaround time for emergency medical technicians in their patient transfer activity [21-23]. ED presentations studied in Saudi Arabia concluded that age and presentation times were significantly associated with ED waits [24]. Wait times, length of stay (LOS) concerns, and the nearly 10,000 patients treated daily in EDs in Denmark, Sweden, and Norway prompted changes for the emergency medicine specialty [25]. Finally, the general issue of waiting time, which has become a major public health policy concern in many of the countries in the 34 member Organization for Economic Cooperation and Development (OECD), was investigated from the viewpoint of competition and waiting times in hospital markets [26]. Among their numerous findings for many countries' hospital density scenarios, the authors find that for countries where waiting times are excessively high and prices are too low, such as UK, Norway, and Finland, hospital competition will reduce waiting times and increase activity and welfare for the case with a sufficiently small demand segment. While this work focuses on a typical US ED, it has implications for EDs wherever they operate, because characteristics of the Tragedy of the Commons pervade most emergency systems which operate at or beyond capacity [27]. Because the US healthcare system is poised to undergo dramatic changes, its emergency services make an interesting case study of a complex adaptive system (CAS) at the edge of chaos [28]. Metaphorically, its "geologic" landscape resembles plate tectonics, with several fault lines: a new business model, accurate costs, population centricity, and implementation of The Affordable Care Act (ACA). On the surface, a value-based (pay-for-performance) model will replace the conventional fee-for-service model. The mythical relationship between hospital charges (reimbursement) and accurate economic cost will eventually be dispelled in order to effectively compete in the new business model. The pay-for-performance paradigm will test healthcare systems' understanding of quantitative performance and quality measures, costs, and risk. New Centers for Medicare and Medicaid Services (CMS) incentives, tying reimbursement to outcomes, are driving commercial payers to design new payment models incentivized by quality and efficiency goals [29]. Additionally, when accurate costs are tied to efficiency, the healthcare business model will need to incorporate nonhealthcare business strategies in the new competitive environment. One of these, advertising wait times to ensure market share, will put congestion metrics on the customer's radar. Macroeconomic forces will goad the transformation from hospital and inpatient centric to preventative ambulatory and population-centric care. Finally, the expansion of insurance coverage under the ACA will likely swell the publicly insured Medicaid population. The seismic pressure of healthcare's plate tectonics, shifting simultaneously over the next several years, may challenge the efficient and safe delivery of high quality care, hence new paradigms to model and predict the system's behavior are in order. Nonlinear systems, complexity, game theory, and econophysics approaches successfully applied in other sectors hold promise.

We present our work in the remainder of this chapter. In Sect. 2, we discuss the issues of ED overcrowding, utilization, and negative primary care sentiment. Section 3 outlines preliminary analytic studies that defined the "left without being seen" population as the focus for our game theory models. The congestion game framework, outlining the El Farol Bar, El Farol Network Congestion and Minority Game approaches is described in Sect. 4. Results of agent-based simulations of the games are presented in Sect. 5, with statistical analysis of empirical observations and simulations in Sect. 6. Future work covering game extensions, the Braess Paradox, learning schemes, and uncertainty measures are explored in Sect. 7. We conclude our work in Sect. 8.

2 Background and Motivation

The literature on ED utilization problems is rife with the topic of overcrowding and its consequences. Wait times in the US have increased between 1997 and 2004 by 36 % [30]. EDs are currently dealing with the uncertainty of whether it will get worse [31]. A well-documented subject, emergency service congestion has been discussed across a spectrum of ED endpoints, including operations, staffing, process reengineering, crowding measures, interventions, and the opportunity loss of boarding admitted patients [32-37]. Several studies have documented myriad adverse patient quality and safety outcomes. These encompass research which attributed 5 % greater odds of dying after being admitted in a congested ED to adverse cardiovascular outcomes for chest pain patients [38–46]. EDs are a primary gateway for hospital systems, accounting for as much as 50 % of hospitals' nonobstetric acute-care admissions and 60 % of Medicare admissions [47, 48]. With 24 h availability, sophisticated medical equipment, specialized manpower and procedures, and high fixed costs, EDs are also a microcosm of the larger hospital system—a hospital within a hospital. Legal requirements such as the Emergency Medical Treatment and Labor Act (EMTALA) require that EDs treat on the basis of need rather than on the ability to pay, thus the uninsured have been more likely to turn to the ED for care rather than a physician's office [49]. Moreover, the US Congressional Budget Office projects that while the Patient Protection and ACA will provide an additional 32 million people with health insurance by 2019, 23 million people will still remain uncovered [50]. Their ranks will include those who choose to remain uninsured as well as undocumented immigrants. Primary care shortfalls have not been accounted for, as was the case with Massachusetts healthcare reform in 2006. Without a commensurate increase in primary care physicians, already constrained EDs might continue to be the preference for the burgeoning Medicaid population and the uninsured [51]. The low socioeconomic status (SES) population, even in countries other than the US with near universal insurance coverage, displays a pattern of seeking low value care-underuse of primary care and overuse of hospital-based services [52]. These learned behaviors might also contribute to ED overcrowding in the near term and aggravate the struggle to maintain high quality and safety service levels. The ED utilization lessons from the statewide Massachusetts healthcare reform in 2006 may serve to guide strategies for already constrained EDs [53–55]. Theoretical aspects of the ACA, targeted to reduce low value care by encouraging primary care outside of the ED setting, may not map seamlessly to practical aspects in the national setting. Therefore, chronic ED crowding due to the traditional components of ED congestion, namely the preferences of low SES, nonurgent medical needs, and "frequent flyer" patients, may be exacerbated by the shifting tectonic plates of the ACA, population health management, and the safety and quality driven payfor-performance paradigm on the horizon.

3 Related ED Congestion Work

Several preliminary investigations regarding ED congestion at our facility are briefly surveyed here, providing the contextual backdrop for recasting our problem in a game theoretic framework. Our study was designed at Lehigh Valley Health Network, an eastern Pennsylvania USA tertiary medical center with 1,011 licensed acute-care beds and 5 campuses and EDs. As Pennsylvania's first Level 1 Trauma Center, it has a catchment area of approximately 4,500 square miles and a population base of over two million people. It remains the only Trauma Center with additional qualifications for pediatric trauma and the region's only children's ED. It is a teaching hospital, academically affiliated with the University of South Florida. It is also part of an interstate consortium, AllSpire Health Partners, comprised of seven hospital systems totaling 31 hospitals, the largest healthcare consortium at the time of its formation in the US. Our institutional review board approved the study reported here, which pertains to the adult ED at one of our campuses which had over 110,000 annual ED presentations.

With an increasing rate of ambulance diversion, an understanding of the system congestion, in both the ED and the inpatient setting was analyzed using discrete event simulation for process modeling and a popular machine learning algorithm, a decision tree classification model. These models were used to guide policy decisions aimed at reducing ED congestion tactically, focusing on ambulance diversion as an operational definition and proxy measure of crowding. We briefly outline key points which informed our investigation of game theory approaches; approaches we felt had the potential for incorporating elements of the patient's perspective for understanding the causes of congestion and potential solutions, guided by strategic rather than tactical considerations.

Data was SQL queried and joined from three real-time electronic tracking systems: the ED's electronic database, T-System, which tracks all patients registered in the ED, the patient logistics database or hospital bed board Teletracker, and the perioperative tracking database, Amelior ORTracker, an RFID tracking system

for surgical patients. Demographic retrospective data was queried from a data warehouse where GE Centricity electronic medical record data is uploaded. Ambulance diversions in response to congestion are logged as hours of diversion status, therefore all data reflected hourly units. Three years of data (2010–2013) comprised of over 300,000 attendances were compiled for time-stamped hourly variables, yielding 26,280 hourly instances. In an effort to organize our research agenda, we followed the standard conceptual model of ED crowding based on input, throughput, and output variables [56]. This operations management framework facilitates understanding the ED as a CAS because it highlights system coupling by outlining the numerous sub-systems and transitions. More importantly, in terms of a CAS, it is a lens of the numerous interactions between the coupled dynamical subsystems. Input variables included ED self and ambulance arrivals, census, emergency severity index (ESI) or priority of patients in the ED, referred to as a triage number on a 1–5 scale of decreasing severity, and transfer patients. Transfer patients reflect the utilization rate of an ED as a referral site for other providers such as skilled nursing facilities, urgent care centers, and others. Throughput variables were comprised of standard ED process steps, including doorto-triage, door-to-room, door-to-doctor, door-to-test, and door-to-departure for nonadmitted patients, as well as the time from disposition decision until departure for admitted patients, and the rate of inpatient admissions. Output variables were acute hospital LOS, observational LOS, elective surgery same day admissions, direct transfers from other facilities, intensive care unit occupancy, inpatient unit occupancy, total admissions, discharges, and their hourly standard deviations. We also developed a novel output variable, a congestion volatility index (VIX), which characterized intra-house transfers as an indicator of congestion accommodation. While the movement or transfer of an inpatient reflects standard hospital processes, such as transferring out of a critical care unit to a medical or surgical unit in a stepdown process, it also reflects hospital administrators' preparations to create space or buffer anticipated arrivals before a bottleneck occurs. Lastly, we included patients who arrive at the ED and register, but leave before completing treatment. This elopement is commonly called left without being seen, hereafter LWBS. A small but related category of patients who leave after starting treatment, but before treatment is complete (LBTC) was not included in this study. Similar to ambulance diversion, LWBS is related to all categories of the conceptual model and is a consequence of congestion. Descriptive statistics and data cleansing were performed using SQL commands for the joined databases. Reconciliation of the electronic data in the form of aberration detection methods and multivariate outlier analysis was done in conjunction with empirical observations of ED operations during all shifts over a 3-month period by the author. The final database of 26,280 records defined by the 25 operational variables was analyzed using the computer applications ARENA® for simulations [57] and CART for decision trees, a form of supervised pattern recognition [58]. More specialized audiences seeking additional information about these methods can refer to more extensive descriptions of discrete event simulations and classification and regression tree design in several engineering and machine learning references [59–61].

The results of our simulations confirmed congestion bottlenecks due to a lack of consideration of coupled systems. In particular, the pattern recognition models' high prediction accuracy for ambulance diversion also yielded scores of variable importance, with our VIX, LWBS, and LOS variabilities as the top performing variables. Both models quantitatively confirmed the temporal nature of our congestion problem, previously understood only qualitatively. Inpatient congestion variables on Mondays, which formed the output operational variables in the conceptual ED model, were the primary culprit in turning away ED bound ambulances for nonsevere medical conditions.

An unsurprising result, Mondays enjoy a unique position in many business sectors which reflect seasonality, calendar, and day-of-the-week effects, the stock market being a common example. The Monday syndrome in emergency care is also well documented. In the US, LWBS patients were found to occur 1.4 times more frequently from Monday through Wednesday in a study where extensive wait times and nonurgent LWBS were investigated [62]. The phenomenon is not unique to the US. In the study of six EDs in Sweden, Monday was the busiest day and Saturday the least busy [18]. Similar results have been reported in the UK, where both a Monday effect and evening hours with longer wait times was found to increase the probability of LWBS [6]. In another large UK study of over 15 million ED attendances, of which 11.7 % were inappropriate or nonurgent, a Monday effect was also found [6, 63]. We also found a Monday effect for recurrent ED attendances. During the study period, 100,779 patients made 205,953 ED presentations to the smaller of our suburban EDs. Over 70 % were repeat attendances. At our larger suburban ED, over 65 % were repeat attendances out of 305,890 total attendances by 164,172 patients. At our urban ED 50,070 patients made 121,431 attendances, with over 78 % recurrent attendances. Our recurrent attenders comprised 39, 35, and 46 % of all ED patients and accounted for 70, 65, and 78 % of ED visits, respectively. Additionally, 47-58 % of our recurrent ED attendances occurred on Mondays for our urban and suburban sites. Attendance frequency details are listed in Table 1 and the Beta probability distribution (0.5 + 81 * Beta)(0.392, 20.2) of the corresponding attendance frequency versus count (x, y) is shown is Fig. 1.

Based on our findings for Mondays, we further developed our classification tree by reducing the input space and measuring the prediction results using the receiver

Attendances	Suburban ED 1 (%)	Urban ED (%)	Suburban ED 2 (%)
>50	0.4	0.4	0.4
26–50	1.7	1.4	1.1
>25	2.1	1.8	1.5
7–25	18.2	25.5	14.8
2-6	50.1	50.7	48.8
1	29.6	22.0	34.9

Table 1 Frequency of EDattendances

Fig. 1 ED attendance probability distribution—count versus attendance frequency (x,y)

operating characteristic (ROC) summary statistic [61]. An ROC curve is graphed as the false alarm rate (1-specificity) on the abscissa, against sensitivity on the ordinate axis for all possible thresholds of a binary classification task. The area under the ROC curve, referred to as AUC or A_7 , represents the discrimination ability of test. Perfect discrimination is represented by 1.0, while 0.5 denotes no discrimination or randomness. A₇ was calculated for daily models of ambulance diversion prediction by the hour. Our best classifiers, defined by $A_7 > 0.85$, achieved performances of 0.87, 0.88, and 0.97 for predicting ambulance diversion with the smallest subset of variables-LOS, VIX, and LWBS. Figure 2 graphically displays a subset of the average congestion variables for Mondays which led to ambulance diversion, depicting a typical example of all variables relating to congestion. Included were: hour of the day, inpatient hospital occupancy (referred to as house congestion), ICU occupancy, ED census and admission rate, the elopement of patients who register (check in) but do not wait for emergency service or "left without being seen" (LWBS), the percentage of patients present with higher triage status (1-3), and LWBS wait time. Our ED capacity is delineated by adult, children, trauma, and surge (capacity for sudden excess demand), accomodating103 patients. From Fig. 2, the notated red or largest sphere, indicating the largest LWBS count, occurred when the ED census reached 62 patients. We will return to this characteristic of ED crowding in our game theory models.

Our goal to preserve ED service quality and safety by targeting the most explanatory variables for ambulance diversion, in the nascent context of population health, motivated examining the LWBS population as an additional congestion indicator. Although our VIX and LOS were equally important, our LWBS rate had increased over the study period. While it is not surprising that both ambulance diversion and elopements are interrelated signs of constrained resources, ED elopements are also concerns in their own right. We analyzed our LWBS population using a bootstrapped cumulative summation (CUSUM) or Change Point Detection, a standard process control engineering technique to statistically characterize the timing of elopements as the new marketing environment unfolded [64]. While Change Point Detection is typically used as a surveillance method for disease outbreaks because it can detect subtle shifts from the process mean and provide estimates of when the change occurred and its magnitude, it is also useful



Fig. 2 Decision tree results as a graphic of congestion resulting in ambulance diversion

to characterize any rare event in time series data. When augmented with statistical resampling or bootstrap methods, where random iterations of the original dataset are generated, confidence levels are obtained for every change in the mean or variation detected. In the interval studied, hospital EDs were starting to scale-up competition for market share with other EDs, urgent care centers, and retail health clinics, similar to market entry strategies in the commercial or industrial and business sectors. In Fig. 3, statistically significant increases in elopements occurred throughout 2012, while a steady rate was maintained in 2013.

Elopements were further characterized by payer status. As previously mentioned regarding the newly insured through the provisions of the ACA healthcare reform, an analysis by payer categories (self-pay, uninsured, commercial, Medicare, and Medicaid) was of interest from various perspectives. At this juncture, we wanted to understand whether any potential shifts in payer status over time were associated with the LWBS congestion indicator. Although no significant changes in total payer categories were noted, a LWBS increase (38.6 %) in the 2010–2013 timeframe occurred for the Medicaid population. A large corpus of literature is testament to the safety and quality issues associated with the LWBS population, with LWBS used as a proxy indicator of ED performance and overcrowding, as a safety net concern, and for its association with low-income and poorly insured patients [5, 9, 10, 65, 66]. To further characterize these elopements, we analyzed the population by ESI



Fig. 3 Statistical change point detection of LWBS

to understand where their presentations occurred in the severity spectrum, with ESI 1 and 2 representing urgent conditions and 3–5 as nonlife threatening conditions. The majority of our patients who did not wait for service had nonurgent severity indices at ED registration.

Given the congestion which resulted in ambulance diversions, we assumed a component of the LWBS phenomenon was dissatisfaction with the wait times. In light of the new strategy by US hospitals to advertise ED wait times coincident with increased LWBS, we decided on a figure-ground shift approach. Investigating the social coordination background could elucidate a variety of congestion phenomena, including how people decide to use an often congested ED for nonurgent conditions, as well as decision making under uncertainty at the server-diverting ambulances. Recent marketing of ED wait times might have contributed to the changes we were experiencing [67–70]. Advertisements using various real- or neartime venues such as text messaging, smartphone applications, Internet sites, and billboards were introduced in our ED catchment area during 2011. In general, while it seems appealing from the customer viewpoint to acquire wait information, it is seldom clear whether it is measurably beneficial. Consider how people decide while holding in a call center scenario with an estimate of the wait time-"to balk or to renege: that is the question." Furthermore, a Canadian study found low awareness of published online ED wait time data [20], complicating prediction. Similarly, little is known about the benefits to the server regarding publicly posted queue information in the healthcare sector. However, even the US government has launched an app, ER Wait Watcher, based on data gathered by the Centers for Medicare and Medicaid Services (CMS) in conjunction with Google. The CMS collects care measures from US hospitals, who are financially incentivized to voluntarily provide average wait times [71]. It pairs the data with Google real-time travel estimates based on traffic conditions, in theory eliminating the need for guesswork regarding

265

waiting. Still, there are conflicting opinions in the grab for market share. These span positive outcomes for relieving ED congestion by more effectively distributing patients across several servers, to negative consequences, such as discrepancies in the definition of waiting time that individual hospitals report [72–74]. Additionally, the Federal Communications Commission, which oversees the Lifeline program to ensure that low-income Americans are provided with telephone service, expanded the program to include smartphone upgrades during 2011-2012. In the context of our Medicaid population's increase in elopements, the likely expansion of this population under the ACA, the 23 million who will remain uninsured, and the old habits that die hard (people who prefer the ED to primary care), we find ourselves on new turf. It seemed plausible that the confluence of enhanced flow of real-time wait information, regardless of the payer status of the user of the data, as well as previous wait experiences would provide an interesting backdrop to understand how people decide to go to the ED for nonurgent care sans crowding. Fortunately, the same kind of problem has been examined by economists, computer scientists, and econophysicists, although not in an ED, but in a bar.

4 Game Theory Models

Game theory models cooperation and conflict between intelligent and rational decision makers. A game is simply the description of strategic interactions, which includes constraints on the players' actions and interests. Nash equilibrium is a central tenet of decentralized decision making. A Nash equilibrium is a solution in which no player can improve her selfish objective by unilaterally changing her strategy. In a Nash equilibrium, each player has chosen a strategy which is a best response to other players' choices. While equilibrium solutions are stable, they can be less efficient than centralized optimal solutions which minimize the cost to the entire system. For the problem we consider here-the ED self-referral system as a population of agents who have coupled dynamics-the El Farol Bar Game, also known as the Santa Fe Bar problem, provides an interesting and relevant game theoretic framework. Additionally, extensions of the game allowed us to make comparisons for our congestion problem. In the remainder of this section, we develop the central tenets of the El Farol Bar Game, discuss its interpretation as a congestion game, a Network Congestion game, and as an abstraction of the El Farol Bar Game, the Minority Game. All modeling was done using the agent-based modeling (ABM) platform, NetLogo [75–78]. In contrast with our previous discrete event simulations, ABM is decentralized-there is no definition of global system behavior. Rather, behavior is defined at the individual level (bottom-up modeling) and the global system behavior emerges as a result of individuals following behavior rules, interacting and communicating with and in their environment. The NetLogo implementation facilitated modeling speed.

4.1 The El Farol Bar Game

Introduced by Arthur [1] to investigate how to model bounded rationality in economics, the game is used to model a system of agents who inductively adapt their beliefs in the aggregate environment they create when they cannot explicitly coordinate with each other directly or through some external mechanism. Inspired by the El Farol bar in Santa Fe, New Mexico, which used to host Irish music every Thursday night, the original problem was set-up as follows:

N people decide independently each week whether to go to a bar that offers entertainment on a certain night. For correctness, let us set N at 100. Space is limited, and the evening is enjoyable if things are not too crowded—specifically if fewer than 60 percent of the possible 100 are present. There is no sure way to tell the numbers coming in advance; therefore a person or an agent goes (deems it worth going) if he expects fewer than 60 to show up or stays home if he expects more than 60 to go [1].

Arthur's simulations showed that while the mean attendance was near the capacity level of the resource at 60, the artificial agents did not effectively coordinate because the variation around the capacity fluctuated widely. It is straightforward to make the connection to the ED and other common pool resources. The bar's capacity is a resource which is subject to congestion, therefore the El Farol Bar Game is a stylized version of many problems in economics which examine efficient exploitation of common public resources. As we noted earlier, the ED can be cast in the framework of the Tragedy of the Commons—a decentralized system where individual users compete to maximize their individual, local payoffs which impact the system-wide objective, often causing performance degradation of the global system. Depending on the problem domain, centralized systems can be impractical or costly to implement. Therefore, solutions to the problems of decentralized decision making—individuals minimizing their own costs of receiving service in the context of the centralized objective of minimizing the total cost of providing service to everyone—are often sought.

We describe central tenets of the problem and refer the interested reader to the extensive literature on the subject spanning a variety of perspectives [79–87]. In our explanation, we closely follow the derivations for El Farol and its relative the Minority Game, which are outlined with clarity in [83–86]. The words player and agent are used interchangeably in our discussions.

4.1.1 The El Farol Stage Game

We introduce the one-shot or El Farol stage game, with the definition of the bar's threshold, *T*, a nonzero integer, and a population of *N* people who are bar goers. They decide independently to go to the bar. Going to the bar is worthwhile if it is not too crowded, otherwise they prefer to stay home. Specifically, the bar is considered crowded when n (where $n \le N$) people attend and the threshold is exceeded. Likewise, it is enjoyable when the attendance is at or below the threshold. Given

a set of actions, $x = \{0, 1\}$, where 1 denotes attending and 0 denotes staying home, the general form of the payoff function is given by:

$$U(x,n) = \begin{cases} u_1 \text{ if } x = 1 \text{ and } n < T \\ u_2 \text{ if } x = 1 \text{ and } n \ge T \\ u_3 \text{ if } x = 0 \end{cases}$$
(1)

There are different payoffs depending on the results of each game round. The conditional payoffs which occur according to the state of the bar are as follows: a player receives the highest score, u_1 , when going to the bar and finding that it is not crowded, while the lowest score, u_2 , or a penalty is awarded when attending the bar and discovering it is crowded. Staying at home nets an unconditional payoff, u_3 , a medium number of points. This is expressed in the payoff matrix in Table 2. For ease of interpretation, we substitute high, medium, and low (H, M, L) for u_1, u_2, u_3 . The payoffs are strictly ordered such that H > M > L.

There are three types of Nash equilibria of the El Farol stage game: pure strategy, symmetric mixed strategy, and asymmetric mixed strategy. The number of pure strategy Nash equilibria, where all players play a pure strategy, with $N \in \mathbb{N}$ players and a threshold capacity of $T \in \mathbb{N}$ is:

$$\binom{N}{T} = \frac{N!}{T!(N-T)!}$$
(2)

There is also a unique symmetric mixed strategy equilibrium, where all players play the same mixed strategy defined by (p, [1-p]), where p is the probability of going to the bar and [1-p] denotes the probability of staying home:

$$\frac{M-L}{H-L} = \sum_{m=0}^{T-1} {N-1 \choose m} p^m [1-p]^{N-1-m}$$
(3)

Furthermore, the number of asymmetric mixed strategy equilibria, where some players play a pure strategy and the remaining play a mixed strategy, is countable.

We conclude our discussion noting that the Nash criterion for an equilibrium in many cases, as in the El Farol bar stage game, is a weak criterion—many strategy combinations are admissible as equilibria, far beyond the number of "equilibria" that actually realize in the data we see. Naturally, this is disadvantageous for prediction purposes. It would be difficult to know which Nash equilibrium will occur in the data when there are a countably infinite number possible. In order

Table 2 El Farol Bar Stage Game payoff matrix		Attendance	
Sume pujon maan	Action	Crowded	Uncrowded
	Attend	L	Н
	Stay home	М	М

to correct for the large number of Nash equilibria, which rises as N increases, a different equilibrium criterion, usually a refinement of the Nash criterion but more strict so as to reduce the number of possible equilibria is a remedy. Furthermore, when a game of both coordination and differentiation is repeatedly played by many players, clear predictions from game theory are not forthcoming. For this reason, simulation models of the underlying game afford the next best approach to actual experiments, so that different simulated agents can use different strategies as the game is repeated. Allowing the players to use learning rules accomplishes this end. Reinforcement learning, essentially trial and error learning based on interactions with an environment, and Arthur's inductive learning are two possibilities.

4.1.2 The El Farol Game

In the El Farol Bar problem, Arthur notes that psychologists agree on the tendency for people to develop internal models, mental schemes, and behavior rules from which to search for patterns in order to make decisions in complicated situations. They are essentially creating forecasts. If their forecast predicts low attendance, they will go. The converse holds true. The problem with these mental schemes and hypotheses is that no forecast can be employed by everyone and simultaneously be accurate [80, 86]. To see this, consider the case where a forecast based on past attendances predicts high attendance for the coming week. Then all potential attendees using the forecast will stay home. This response invalidates the forecast, and exemplifies the problem that no deductive solution can be found. To model bounded rationality, Arthur assumed individuals had a set of hypotheses or forecasting models to base their decisions on. Therefore in his simulations, agents predict using deterministic inductive reasoning: if they predict attendance will be less than 60, they go to the bar; if they predict it will be greater than 60, they will stay home. The predictive heuristics used could be anything from moving or simple averages, a mirror of the last week, a period two cycle detector, and so forth. They were ranked according to their accuracy at the end of the week to use as a predictor for the coming week.

Consider each agent, k, to have a set of predictors, s^k which they score and rank in a weekly period t, for their upcoming decision given the attendance history, h, over the last d weeks where $d(h_{t-1})D$ and D is the set of all possible attendance profiles for the last d weeks. Armed with the publicly posted actual attendance number of the most recent week, a score is assigned to each model, $U_t(s^k)$. It is computed from the weighted average of the score of the same model in the previous week and the absolute difference between the forecasting model's last prediction, $s^k(d(h_{t-1}))$ and the current realized attendance y_t in period t. Formally, the derivation can be expressed based on [86, 87], where λ is a number strictly between zero and one: A Congestion Game Framework for Emergency Department Overcrowding

$$U_t(s^k) = \lambda \ U_{t-1}(s^k) + (1-\lambda) |s^k(d(h_{t-1})) - y_t|$$
(4)

For each week, the forecast which yields the highest score becomes the agent's predictor and she goes or stays accordingly. The actual attendance is publicly available. The agents realize their payoffs, update the score for their forecast models, and choose the predictor for the coming week's decision to go to the bar or not.

4.2 El Farol as a Congestion Game

As a repeated simultaneous move game, the El Farol Bar problem describes N players with identical preferences attempting to coordinate their actions so as to maximize local payoffs subject to the crowding externality of a public common resource. Without the ability to communicate, they must independently decide that when going, they will receive a payoff greater than what they would have received had they remained at home. Similarly, when choosing to stay, they would receive a payoff greater than what they would have realized had they not stayed away. Analogously, market entry games, interpreted as truncated two-stage games, depict players in the first stage simultaneously choosing whether to enter or not [86]. The entrants' payoffs are determined from these actions in the second stage. We find the market entry interpretation of El Farol intuitively appealing because of the marketlike dynamics of competing healthcare services in our current climate, where consumers' service provision and providers' resource allocation behaviors suffer conflict. Competing EDs with various express care configurations, retail and convenient care clinics, urgent care or "emergi-centers," as well as consumers strategizing how to individually utilize these resources for their nonsevere conditions, resemble market operations. There is some dissimilarity with the market entry interpretation to the El Farol bar problem, however, because the entrants' payoffs are related to the bar attendance discontinuously, unlike markets where the payoffs are negatively related to the number of entrants continuously. Notwithstanding, the healthcare market's uniqueness encompasses numerous other economic distortions worthy of novel modeling approaches.

4.3 Congestion Games

Other research casts the El Farol bar problem as a congestion game. Congestion games and selfish routing problems are advantageous frameworks for decentralized systems and are often studied from a multidisciplinary setting of game theory, computer science, operations research, transportation engineering, and communication network routing. Congestion games, as the name implies, model resource competition where players simultaneously attempt to utilize a resource, with the

payoff a function of the congestion level of the players allocating the common resource. In this framework, when a player unilaterally switches strategy, the change in her payoff or cost, such as a delay, is the same as the change in a global objective. Congestion games were first introduced by Rosenthal [88] and have been extensively studied [87–97]. They were later generalized to a class of games called potential games, which incorporate information about Nash equilibria in a single real-valued potential function over the strategy space [95]. Rosenthal's original congestion game applications, road networks and factory production, have since been extended to many analogous problems in areas spanning engineering, biology, and network routing design.

We overview congestion games as our perspective for recasting the El Farol bar problem as the El Farol Network Congestion Game for ED crowding. For further details regarding congestion game derivations, the reader is referred to [96] and the references therein. Their work includes the mathematical relationship between congestion games and potential games in an analogous decentralized problem of spectrum sharing in a wireless communication system.

Congestion games are a class of games given by the tuple $(\mathcal{N}, \mathcal{R}, (\sum_i)_{i \in N_i}, (g_r)_{r \in \mathbb{R}})$, where $\mathcal{N} = \{1, 2, ..., N\}$ denotes a set of users, $\mathcal{R} = \{1, 2, ..., R\}$ denotes a set of resources, and $\sum_i \subset 2^{\mathcal{R}}$ is the strategy space of player *i*. The payoff function associated with resource *r* is $g_r : \mathbb{N} \to \mathbb{Z}$. The payoff g_r is a function of the total number of users using resource *r* and is nonincreasing in the congestion. A player tries to maximize (minimize) its total payoff (cost) which equals the sum total of payoff over all resources involved by its strategy. The strategy profile is: $\sigma = (\sigma_1, \sigma_2, ..., \sigma_N)$, where $\sigma_i \in \sum_i$. Hence user *i*'s total payoff is given by:

$$g^{i}(\sigma) = \sum_{r \in \sigma_{i}} g_{r}(n_{r}(\sigma))$$
(5)

where $n_r(\sigma)$ is the total number of users for resource *r* under the strategy profile σ . Rosenthal's potential function $\varphi : \sum_1 \times \sum_2 \times \cdots \times \sum_n \to Z$ is given by:

$$\Phi(\sigma) = \sum_{r \in R} \sum_{i=1}^{n_r(\sigma)} g_r(i) \tag{6}$$

$$=\sum_{i=1}^{N}\sum_{r\in\sigma_{i}}g_{r}\left(n_{r}^{i}(\sigma)\right)$$
(7)

The second equality is the result of exchanging the two sums. The number of players using resource *r* whose indices do not exceed *i* in the set $\{1, 2, ..., i\}$ is given by $n_r^i(\sigma)$. When player *i* unilaterally moves from strategy σ_i to the profile σ'_i the potential changes by

A Congestion Game Framework for Emergency Department Overcrowding

$$\Delta \Phi \left(\sigma_i \to \sigma'_i \right) = \sum_{r \in \sigma'_i, r \notin \sigma_i} g_r(n_r(\sigma) + 1) - \sum_{r \in \sigma_i, r \notin \sigma'_i} g_r(n_r(\sigma))$$
$$= \sum_{r \in \sigma'_i} g_r(n_r(\sigma')) - \sum_{r \in \sigma_i} g_r(n_r(\sigma))$$
$$\Delta \Phi \left(\sigma_i \to \sigma'_i \right) = g^i \left(\sigma^{-i}, \sigma'_i \right) - g^i (\sigma^{-i}, \sigma_i),$$
(8)

where, for resources that are used by both strategies, σ_i and σ'_i , there is no change in their total number of users, yielding the second equality. Here, σ^{-i} denotes a vector of strategies for all players other than player *i*, which is held fixed since we are only considering a deviation in player i's strategy. This result can be obtained by exchanging the two sums in the potential definition (7) and by assuming the N-th player is considered. An intuitive interpretation is presented in [97], conveying that the gain (loss) as a result of any player's unilateral move is exactly the same as the gain (loss) in the potential. The potential can be viewed as a global objective function. Because the potential of any strategy profile is finite, then every sequence of improvement steps is also finite. Known as the finite improvement property (FIP), the improvement steps converge to a pure strategy Nash Equilibrium, which is a local maximum (minimum) point of the potential function Φ , defined as a strategy profile where changing one coordinate cannot result in greater value of Φ . Therefore unilateral improvement steps, regardless of sequence, converge to a pure strategy Nash equilibrium—a local optimum of the global objective function given by $\Phi()$, the exact potential function:

$$g^{i}\left(\sigma^{-i},\sigma_{i}^{'}\right) - g^{i}\left(\sigma^{-i},\sigma_{i}\right) = \Phi\left(\sigma^{-i},\sigma_{i}^{'}\right) - \Phi(\sigma^{-i},\sigma_{i})$$

$$\tag{9}$$

A congestion game is thus an exact potential game, and every potential game may be converted into an equivalent congestion game [95].

In the introduction we conjectured how game theory could be used to highlight how prospective ED customers might be optimizing their own behavior in crowding situations, which in turn conflicts with the collective, social optimization of our decentralized resource. By the FIP for congestion games, advanced through its relationship with the potential function, an understanding of the stable state of a system is within reach. Hence, the appeal of the congestion game framework is that it provides a context to examine how individual and self-interested, local optimizing behaviors collectively may result in a globally optimized solution for a complex system in periods of common resource scarcity.

4.4 El Farol Network Congestion Game

The interpretation of the El Farol bar problem as a congestion game, where each player's payoff depends on the number of other players vying to utilize the same resource, has been investigated in the literature [77, 79, 80]. In one congestion game framework, Arthur's deterministic inductive reasoning in the original problem was ported to a stochastic setting, and an adaptive learning algorithm was implemented. The agents need not make explicit predictions about the bar based on the decisions of others; rather, they use their own experiences [80]. This distinction, where the authors used the Internet as example of a public resource which can become congested due to overutilization, is interesting because it accounts for coordination failure. The lack of coordination is due to the agents' uncertainty regarding the actions of other agents. Algorithmically, the adaptive strategy resembles habit formation. This feature may be relevant for our problem domain, where ED utilization patterns for nonurgent medical care have been ingrained in some populations, as mentioned in the introduction. Additionally, by using the stochastic version of an adaptive learning rule, the authors observed that the characteristics of the equilibria observed depend on the nature of information available to agents. Limiting the information leads to successful coordination and a Pareto efficient equilibrium, while providing more leads to an inefficient outcome. Recall our earlier discussion regarding the current transition of hospitals to market entrants with strategies for garnering share based on their quality, safety, and performance metrics, specifically in the form of wait times. The adaptive El Farol congestion game may shed some light on how much information is too much information. Having developed the congestion game foundation for this distinctive interpretation of the El Farol bar problem, we refer readers interested in more complete derivations to the original work [79, 80], highlighting the salient features of the game in what follows.

Consider that over time, agents adapt their probability of attending the bar based on a history of their own experiences—they cannot know ahead of time whether they will be satisfied with the service level (the bar will not be too crowded) because they cannot know the actions of everyone else. If an agent attends p percent of the time, then the agent will go more often (increase p slightly) if the bar is not too crowded. Alternatively, if the bar is crowded, the agent will decrease p. Furthermore, consider M agents competing for the N spaces at the bar. Then the probability the *i*th agent attends is p_i . Let k be the iteration, N(k) the number of agents attending at time k, and μ the characteristic parameter which defines how much each agent changes p_i in response to new information, and $p_i(k)$ the instantaneous value of p_i at time k. Let

$$N(k) = \sum_{i=1}^{M} x_i(k),$$
(10)

where the $x_i(k)$ are independent Bernoulli random variables which are one with probability $p_i(k)$ and zero otherwise. The probabilities then follow in (11)

$$p_{i}(k+1) = \begin{cases} 0, & \text{if } p_{i}(k) - \mu(N(k) - \mathcal{N})x_{i}(k) < 0\\ 1, & \text{if } p_{i}(k) - \mu(N(k) - \mathcal{N})x_{i}(k) > 1\\ p_{i}(k) - \mu(N(k) - \mathcal{N})x_{i}(k), & \text{otherwise} \end{cases}$$
(11)

At each time k, the agent flips a biased coin, attending with probability $p_i(k)$ in [0, 1]. When the agent attends, $p_i(k)$ is increased proportionally to N(k) - N if the bar is uncrowded and conversely, decreasing proportionally to N(k) - N if the bar is crowded. For not attending, we have $x_i(k) = 0$ and $p_i(k+1) = p_i(k)$.

4.5 Minority Game

The Minority Game is a symmetric version of the El Farol Bar Game. Where Arthur assumed patrons would enjoy the bar if 60 % or less attended, the Minority Game simplifies the cutoff to 50 %. Introduced in the physics community by Challet and Zhang [98–100], it is well documented in the econophysics literature to study market entry and financial market interactions [101–114]. Similar to the El Farol problem, it can be used to model agents with bounded rationality competing for scarce resources and is a simple version of a congestion game. As a repeated game, N agents must decide between two actions, such as in the El Farol bar problem, to go to the bar or stay home, with N odd so as to identify the minority action. Therefore, agents trying to distinguish themselves from the crowd must choose between two alternatives (0 and 1) independently, and those who take the minority action, receive one positive payoff unit. Denoting the *i*th player's action at time *t* by $a_i(t) \in \{0, 1\}$, the total action of all players at time *t* is then computed as:

$$A(t) = \sum_{i=1}^{N} \{2a_i(t) - 1\}$$
(12)

The winning group $W(t) \in \{0, 1\}$ at time t is then defined as:

$$W(t) = H[-A(t)] \tag{13}$$

where *H* is the Heaviside function. For each player, their decision about the action in the next round is based on the last *m* outcomes of the game, illustrated by: $\{W(t-m), \ldots, W(t-1)\}$. Based on the common knowledge of the past record, all players can exploit the global information $\mu(t) \in \{0, \ldots, 2^m - 1\}$, a decimal representation of the binary vector of the last *m* rounds. Players who take the minority decision gain one point at the end of each round, and vice versa. Players decide on their actions by adopting strategies. A strategy is a map from all possible *m*-step histories onto the binary set, mapping $\mu(t)$ to an action a(t). Each player is equipped with a pool of *S* strategies which do not necessarily differ. Strategies assigned to them at the beginning of the game are limited, corresponding to their bounded rationality, and all players are initiated with a random global history their pool of strategies. The strategy *s* of player *i* is given by s_i , where $s \in \{1, ..., S\}, i \in \{1, ..., N\}$. Assuming the players can remember the correct minority decisions of the last *m* rounds, they can resolve 2^m possible histories $\mu(t)$ and therefore the strategy space comprises different strategies. Players decide which of the *S* strategies to play in the next round by keeping track of the *virtual* points that each of its own strategies would have gained if it had been used from the beginning of the game forward. The points are only virtual because they record the merit of a strategy regardless of whether or not it was used in the round. The virtual score U_{is} of the strategy *s* of player *i* is updated by

$$U_{i,s}(t+1) = U_{i,s}(t) - \text{sign}[(2a_{i,s}(t) - 1)A(t)]$$
(14)

At every round, the player then selects the strategy with the highest virtual score. When there is a tie among possible strategies, the player randomly selects one of them. A diagram of the minority game concept with N = 103, depicting the difference between an efficient outcome, which has a smaller difference between the majority and minority (51:52), versus an inefficient outcome or underutilized resource (1:102) as well as an example of strategy table with m = 3 are shown in Fig. 4 and Table 3.



Fig. 4 Minority game diagram—a smaller difference between the majority and minority groups is a better result

for $m = 3$	History	Prediction
101 m - 5	{0,0,0}	1
	{0,0,1}	0
	{0,1,0}	0
	{0,1,1}	1
	{1,0,0}	1
	{1,0,1}	0
	{1,1,0}	1
	{1,1,1)	0

5 Results

Three series of agent-based model were implemented in NetLogo. The NetLogo agent-based modeling environment is an open source software tool with an easy to use graphical user interface. All games used in this study are available in the library of models, including the source code [75]. Standardized descriptions of ABM that promote modeling protocols for social science simulations are provided in [115]. We followed the systems engineering approach for discrete event simulations in [59], which similarly follows the perspective of overview, design concepts, and details in [115]. Additional information regarding design protocol which we found helpful can be found in the NetLogo implementation of EL Farol detailed in [116]. The ED studied has a capacity of 81 beds. Surge and trauma comprise an additional 22 beds for a total capacity of 103. A pseudocolor heat map representation was graphed which depicts census or occupancy by color, with orange indicating the most occupied and dark blue the least occupied. The Monday heat map of attendances is shown in Fig. 5. Compared with all other days of the week, the rate of ambulance diversion on Mondays was the greatest at 38 %. Additionally, the LWBS rate for nonurgent ESI care, triage 4 and 5 levels, coincided with peak ambulance diversions, commencing at an ED census of 55 patients during the interval of noon to midnight. The time interval of greatest congestion is represented by the orange color and its variations.

We parameterized the game theory models based on these empirical insights from our ED, given that the results of a decision tree learning algorithm can be decomposed into a set of rules and parameters. For the El Farol Bar and the El Farol Network Congestion games we set N = 81, our capacity without trauma and surge, because the trauma patients have dedicated bays based on their severity and surge in general reflects mass casualty scenarios.

For the completely deterministic El Farol Bar Game, where the agents choose whether or not to attend based on their best current (active) predictor, Arthur's model is modified in the NetLogo implementation following the work defined in [81]. In this interpretation, a prediction strategy is a list of weights determining how each agent believes that each time period of past data will affect the attendance prediction for the current week. In our simulations, each agent was assigned five



Fig. 5 ED census heat map

prediction strategies and decided which to use by assessing which one would have performed the best had it been used in the preceding weeks. We chose five heuristically, based on our observations of repeated ED attendances. The potential strategies are distributed randomly to the agents upon initialization, and at each time step the agent utilizes only one strategy, based on its previous predictive power for the attendance. The length of the attendance history that the agents can use to evaluate a strategy (memory size) was varied from 1 to 5. The overcrowding threshold was set at 55.

In the El Farol Network Congestion Game, the frequency update was varied from 1 to 7. It represents the value to update the attendance frequency based on a crowded or uncrowded condition such that agents can utilize a different time span until their next visit. This is called the phase in the algorithm derivation, meaning that agents will have a different amount of time (1-7 days) until their next visit [80]. Each agent's phase value decreases as the model runs until it drops below 1 and the agent goes to the bar. Once there, it changes the attendance frequency based on the equilibrium value, set at 55. If the bar is crowded, the attendance frequency is increased by the frequency update value, meaning that it will wait longer for a future visit. The adaptive strategy also defines a dead zone below the point at which the bar becomes crowded. In this transitional attendance zone between uncrowded and crowded, the agents do not care. Thus it is an interval where agents neither increase or decrease their attendance, avoiding the knife-edge scenario of requiring exactly 60 agents in the zone where the bar is considered neither overcrowded or undercrowded, as in the original problem. Our dead zone was set at seven agents, based on the observed variance of our LWBS indicator of congestion. Therefore if the attendance falls below the difference between the equilibrium value and the dead zone, in our case 55 - 7 = 48, the bar is not too crowded and the agent will decrease the attendance frequency by the frequency update value. When the attendance falls within the interval of the dead zone, the agent does not change attendance frequency.

For the generalized El Farol problem, the Minority Game, each agent uses a finite set of strategies based on the past record, however in this case the record is not an actual attendance count, but rather a record of which group, 0 or 1, was in the minority. Each of the agents began with a score of 0, and both the standard ED capacity and trauma/surge were modeled, using N = 81 and N = 103. The length of history (memory) used by the agents for predictions as well as strategies per agent were varied from 2 to 5.

Simulations were repeated for 10 runs using 10,000 iterations. Representative examples of typical runs are shown as time series in Figs. 6, 7, and 8, where graphs, consistent across all simulations, depict iterations 9,000–10,000 without the initial transient behavior. We compare the variance of attendances as agents strategize in the three games using process control charts, a standard framework in a hospital performance setting, to visualize three standard deviations above and below the mean (6 σ).

As previous studies have shown, the El Farol Bar Game fluctuates around the equilibrium, never settling down [80, 81, 82]. In the original simulations by Arthur, the mean attendance oscillates around 60, but varied above 70 and below 50. In our simulations, we duplicated these results with the parameter settings for our ED instead of the bar. The attendance time series is shown in Fig. 6, with a mean of 58 attendees (57.503), and a standard deviation of 7.980. The 3 σ upper and lower control limits show that 99 % of the population can be expected to fall within the range 33–81. The majority of our triage 4 and 5 elopements commence at an ED census (attendance) of 55. As in any goods and service sector, variance is also a critical performance measure in the hospital setting. Therefore, given the variance of the standard El Farol Bar model, predictions regarding the behavior of prospective ED nonurgent patients would be of limited utility for resolving our congestion problem. While it may not be possible to satisfy all people all of the time, we might end up satisfying only some people, some of the time.

In the El Farol Network Congestion Game model, the attendance time series yielded a mean of 50 (50.694), and a standard deviation of 1.312, with 99 % of the attendees between the upper control limit of 54.629 and lower control limit 46.759





Fig. 7 El Farol Network Congestion Game model time series 9,000–1,000 iterations

as shown in Fig. 7. We note that our assignment for the dead zone value, 48, falls within this interval and it correlates with the yellow portions in the heat map, Fig. 5. In this interval, our ED nonurgent care operations are functioning at 60 % capacity (48/81), not so crowded so as to cause diversion or elopements, and not that uncrowded so as to waste resources. The prospects for anticipating that the predictive collective behavior of patrons will enable our system to self-organize into a globally efficient state look more promising with this model.

The agents' inductive learning dynamics in the Minority Game model yield predictions of attendance in between the previous models. The time series attendance in Fig. 8 shows a mean fluctuating around 47 (46.549) by a standard deviation of 1.999. While we could expect that 99 % of our attendees will fall between 53 and 40, and the ED congestion problems at 55 would not be evident, we would experience more idle capacity when the attendance falls to 40, operating at 49 % capacity, compared to the network congestion model which is lower bounded at 58 % capacity (47/81). The Minority Game, however, has other interesting properties which we also examined.



Given our predilection for controlling variation in health care, and our VIX intrahouse transfer variable's importance in our ambulance diversion model, the variance of attendance in the Minority game under the influence of various system parameters was of particular interest. With the minority rule of the game, the number of winners is always smaller than the number of losers. Using "0" and "1" for the actions, the attendance $\mathcal{A}(t)$ averaged over time is N/2, hence the variance or volatility, is given by: $\sigma^2 = \langle \mathcal{A}^2 \rangle - \langle \mathcal{A} \rangle^2$ or σ^2 / N . Thus the volatility is a measure of global efficiency in terms of cooperation and global gain, with higher variance leading to larger aggregate loss. It was previously observed that the ratio given by $\alpha = 2^m/N$ is the most important control parameter [112], and σ^2/N depends only on this ratio. Figure 9 shows normalized volatility σ^2/N as a function of α on a log-log scale for different numbers of players with S = 2. As α is decreased, the minimum of σ^2/N at intermediate α is found, and at small α the volatility of the game is large: $\sigma^2/N > 1$ as $\alpha \to 0$. This exemplifies a collective behavior which is worse than random or a crowded regime. It can be seen that for large α , $\sigma^2/N \approx 1$. This reflects the situation where the agents are making cointoss decisions, the random choice limit, and has been attributed primarily as due to the collective motion of whole crowds of agents making the same decision and of the corresponding anti-crowds [102, 106, 113]. When α increases, the volatility starts to decrease up to a certain critical point, α_c . This part of the game is a better than random regime, and it is here that coordination among the players is improved. The volatility achieves the minimum value for $\alpha = \alpha_c$, which is smaller than the coin-toss limit. At α_c a transition between a symmetric and asymmetric phase of the game occurs. In earlier work regarding this phenomenon, the authors concluded that for $\alpha < \alpha_c$, it is not possible to discern from the last m bits of the history whether one of the two possible actions was more likely to win [112]. Because both actions are equally probable, this is an unpredictable phase. Thus the symmetric phase has high volatility and low predictability. They also found that for $\alpha > \alpha_c$ the probability of winning is not equal when analyzing the last m history bits. This asymmetric phase is predictable with low volatility. This noteworthy feature is not due to players exploiting previous outcomes of the game, as the players are selfishly maximizing their own local objective and consequently approach global efficiency, as shown in Fig. 9.







6 Empirical Statistics and Simulation Results

In order to compare each model and analyze the statistical difference between the models and our empirical observations, we fitted a series of ten probability distributions to all the simulation results and computed their squared error. In addition to the error for each distribution we also computed a series of nonparametric statistical tests. All tests were done at the 0.05 level. The Kolmogorov–Smirnov or KS test was used to compare the fitted distribution with a known distribution to determine similarity, shape, and position. The Chi-square statistic tested the null hypothesis that there was no significant difference between the expected and observed result. The Kruskal–Wallis test of the null hypothesis that the true location parameters (expected values and medians) were the same for each distribution was also computed. The simulations were then compared to the hourly empirical distributions for all hours in the congestion interval of noon to midnight on Mondays using the same statistical tests.

The EFBP results follow a triangular distribution, with 0.118 least squared error. The error range for all fitted distributions was on [0.118, 0.156]. Both KS and Chi-square were not significant at the 0.05 level (p = 0.781 and 0.055), indicating that the EFBP distribution and the triangular distribution do not differ significantly. Our hourly empirical data is Beta distributed, with 0.0134 least squared error for the fitted distribution in the range of [0.0134, 0.0579]. Both KS and Chi-square were not significant, showing no difference between the empirical and Beta distribution (p = 0.0912, and 0.151). A comparison of the EFBP distribution to the empirical distribution for each of the 12 hours on Mondays was further verified using the KS and Kruskal-Wallis tests and was significant, confirming the distributions of the two samples are different and a nonzero difference in the location parameters exists (p = 0.003 and 0.001). The MG results follow a beta distribution, with 0.0378 least squared error. The error range for all fitted distributions was on [0.0378, 0.233]. The KS and Chi-square were not significant (p = 0.462 and 0.063), therefore the MG distribution and the beta distribution do not differ significantly. The comparison of the MG distribution to the empirical distribution using the KS and Kruskal-Wallis tests was not significant, verifying that both follow the same distribution and the difference in the location parameters was not statistically significant (p = 0.059 and 0.626). The EFNC results are also beta distributed, with 0.0184 least squared error. The error range for all fitted distributions was on [0.0184, 0.13]. Both KS and Chi-square were not significant (p = 0.767 and 0.134), indicating that the EFNC distribution and the beta distribution do not statistically differ at the 0.05 level. The KS and Kruskal-Wallis tests were not significant, confirming that the distributions of the EFNC simulations and our observations were not different and the difference in the location parameters was insignificant (p = 0.287 and 0.414).
7 Discussion and Future Directions

Based upon the lack of a significant difference between our observations and the agent-based simulations of the MG and EFNC games, both games provide insights for the ED congestion we originally sought to address. The novelty of the El Farol Bar game as a complex systems approach facilitated by simulating a repeated game primarily served as a developmental path to the EFNC and MG. Additional work exploring machine learning techniques to model human induction, similar to work of Fogel [81] with genetic algorithms, is an area we would like to explore for this model. Similarly, the Kolkata Paise Restaurant Problem (KPR) is a variant of the EFBP that promises insights for ongoing work in the catchment area of the hospital ED modeled here [117, 118]. The KPR resource utilization problem swaps Santa Fe for Kolkata and has a similarly interesting storyline. In Kolkata where paise is the smallest Indian coin, these inexpensive and fixed rate restaurants, some ranked better than others, were frequented by laborers. Walking to a restaurant and finding it crowded meant missing lunch. Walking to the next restaurant meant reporting back late from lunch. In both the KPR and EFBP the number of players is macroscopically large, however, the number of choices is also macroscopically large in the KPR problem compared with only two in the EFBP. The history of choices made by different players for different restaurants is available to all players, as in the EFBP. For the situation where choices for a single restaurant on any evening are made by more than one player, one is randomly selected from them and served food (payoff = 1), while the others lose (payoff = 0). Hence, while each player gains a point (payoff) if her choice of the restaurant any evening is unique (not made by other players on the same evening), the resource utilization is maximized when each restaurant is chosen by at least one player. An extension of this context to include ranks has been suggested by the authors [118]. In the ranked scenario, hospitals are provided in every locality but their local patients may choose hospitals of better rank in other localities. This sets up a competition with the local patients of the perceived higher ranked hospitals. Consequences, such as the lack of timely treatment or the underutilization of some hospitals similarly address the coordination and congestion attributes of the EFBP. In our setting, this model has the potential to offer insights for the nonurgent care population who can choose between our EDs, our competitors' EDs, and the expanding network of walk-in urgent care or 'emergi-centers.' The rank analogy for our ED patients could be described by their perceptions of rank derived from the annually determined US News and World Report score, which is widely advertised, as well as the publicly advertised wait times for service currently being marketed by competing EDs.

The network congestion model, with aspects of adaptive learning, dovetails with our observations of patient characteristics in various population subgroups seeking nonurgent emergency service. However, the value of wait time information in complete or partial states will likely be addressed more fully using a mixed models approach in a prospective study design, given that we have no direct evidence regarding how prospective patients may or may not be using the publicly posted wait times to adaptively learn when the ED will be least crowded. While ambulance diversions were our original motivation to characterize crowding, the nonurgent medical conditions of the growing "left without being seen" population became a clarion call for deciphering how people make decisions. For the former case, the server (ED) makes the decision regarding congestion and reduces the queue via diversions. For the latter, the customers choose. While the decision hierarchy of the service provider would be interesting through the lens of game theory, we have focused our initial inquiries on the customer's perspective, captured in the "left without being" seen variable. The provider as a market entrant is a secondary goal. In the network model, the authors' adaptive scheme was compared with habit formation and reinforcement learning. Based on their Internet example, they realistically suggest that agents might observe the congestion level when they log on. Even though people might prefer to access the Internet at a particular time, if they develop a habit of logging on at the same time, they are signaling others to avoid these times. This type of implicit coordination can act to smooth demand [80]. An analogous situation might also occur in an ED setting. Smoothing patient flow has in fact become the mantra of health care operations leaders and consultants, although toolkits proffered are generally devoid of game theory approaches. Similarly, health care's new marketing strategy conundrum—how to navigate the pros and cons quagmire of ED wait time advertisements-might be resolved from a game theory vantage point. As an illustration of the effects of providing similar information, the network congestion framework had the capability of using complete or partial information regarding the resource. In the complete information scenario, Arthur's original information structure is maintained, such that agents can update their attendance probabilities at every iteration, whether they attended or not. The price was increased variance. Our network congestion simulations were run choosing the partial information setting. In this scenario with less information, agents coordinate their behavior so that the system achieves a Pareto efficient outcome. Along these lines, simulations in transportation routing yield similar system behavior. When agents pursue a strict best-response strategy which results in a change of route, the simulations showed that system-wide performance degradation occurred if more than 25 % of the drivers used traffic congestion information [119]. Will this situation isomorphically map to our setting as ER Wait Watcher is launched? It would be instructive to direct further inquiries along these lines of heterogeneity versus homogeneity of information for our congestion problem. Moreover, the strategies invoked by the provider and those invoked by the consumer are really two sides of the same coin-congestion is likely best mitigated when both adaptively learn to modify their behavior along some complementary trajectory. We reserve these contemplations regarding the value of providing advanced queuing information, the type of information, and the task of distinguishing balking and reneging in the "left without being seen" population for ongoing work combining operations research and machine learning methods. Exploring perspectives from a variety of analogous problem domains regarding the value and impact of sharing lead time information will continue to motivate our future work [120–128].

The network theme within our large physical framework and network congestion game insights lends itself to other extensions. There are other areas outside of transportation and computer networks where the counterintuitive phenomenon of adding more capacity for self-interested agents choosing the quickest route may reduce efficiency. Wardrop's User Equilibrium Principle states that all users traveling between a given origin and destination incur the same cost, and it is not greater than that which would be incurred by traversing any of the unused routes [129]. The Braess Paradox shows that a Wardrop equilibrium is not minimizing the congestion in a network, and depicts the consequences of self-centered users. inefficiency, and paradoxical behavior that a network designer or service provider must consider when designing or upgrading a network [129–133]. Our hospital structure, in particular, is networked and bears similarities to other large-scale networks like the Internet, transportation, electric grid, or social networks. Common to all of them are highly heterogeneous users interacting in a decentralized environment. Thus the complexity of agents interacting in a networked environment raises congestion concerns for understanding optimal flows and system efficiency whether agents are people, data packets, vehicles, or electrons. Analogous situations have been reported in the healthcare settings. An epidemic game of infection was studied in [134], where each agent could choose either preemptive vaccination, laissez-faire (the opposite strategy), or self-protection. Vaccination is expensive but would provide 100 % immunity against infection. Laissez-faire does not have a cost, but carries with it the greatest potential exposure to infection. Self-protection, such as face masks, restricted travel, and increased hygiene carries a moderate cost. The authors found that when the chance of self-protection succeeding was low, most individuals either chose vaccination or laissez-fair, given the negligible benefit of choosing self-protection. When the chance of self-protection succeeding was high, almost everyone chose either self-protection or laissez-fair, given the minimal benefit to pay for the expensive vaccine. In this scenario, the proportion of the population that becomes infected declines. With moderately low levels of the chance of self-protection succeeding, however, individuals begin to notice the benefits of self-protection. Given its lower cost, it is more widely chosen, but its inferior effectiveness produces the effect of reduced vaccinations. Therefore, the number of infected people actually begins to increase with the increased effectiveness of infection protection. Hence, increased disease prevention paradoxically increases the population of infected people just as increasing roads or bridges can increase congestion and delays. Another healthcare example of the Braess paradox at work explored the conflicting objectives of maximum social and individual utility related to the spread of disease in agent-driven contagion dynamics through transportation networks [135]. While game theory shows that the best options for individual users yield a Nash equilibrium, the lack of a social optimum due to selfish behavior in decentralized networks has an efficiency measure referred to as the price of anarchy-the ratio of the total cost of the Nash equilibrium to the total cost of the social optimum. The price of anarchy was addressed in another healthcare setting for the problem of allowing individuals to choose healthcare providers, where choice in a system which copes with demand was found to have a negative effect [136]. We plan to extend our work to quantify efficiency along these lines in future work.

Further testing of the predictive capabilities of the underlying games in an agentbased simulation approach will continue in the next phase of our research. Considering extensions of the Minority Game and its prediction capabilities, it has been shown that because the information common to all agents rather than the feedback of their actions determines the game dynamics, a decoupling of strategies leads to pockets of predictability [137]. Market predictability has been previously demonstrated using agent-based modeling and the Minority game on both real financial time series and simulated market data [138, 139]. We briefly introduced phase transitions or volatility insights in our work to similarly set the stage for related future directions. Markets can be described as operating in two phases. One relates to an equilibrium market and the other to mismatched supply and demand or out-ofequilibrium market. The market is unpredictable when there are no market opportunities in the equilibrium phase. Market opportunities exist, however, when there is excess demand, and in this nonequilibrium phase predictability becomes possible. These ideas underpinned research in a fast moving consumer goods market, where the Minority game was applied to predict the timing of promotional actions for four distinct markets, and provided notable improvements over random guessing in both single agent and multiagent realizations [140]. The results are particularly impressive from a time interval perspective, having achieved a longterm forecast of 6-8 weeks. Time series methods to predict volume or load, as in patient visits to an ED, are less effective for long-term forecasts, a requirement often desired in the healthcare setting for staffing schedules which are based on a 2-month time horizon. We have developed computational models along these lines, and look forward to comparing our previously developed wavelet methods with a similar Minority game approach. Predictive success was also achieved by reverse engineering real-world financial time series [141]. Using Minority game variants and agent-based modeling with a genetic algorithm, the authors demonstrated through a validation test set the applicability of the method for prediction of complex systems with a multi agent structure.

In terms of learning schemes and wait time information characteristics, we foresee a need to explore uncertainty from several perspectives in the decision making process more rigorously. In the real-world setting of publicly available congestion and wait time information, the role of human judgment in determining risk, dealing with uncertainty, and understanding imprecision and ambiguity are important considerations. As we consider model extensions which take these concepts into account, we note the distinction between risk and uncertainty that was formalized decades ago by the economist Frank Knight [142]. Interpretations of his book *Risk, Uncertainty, and Profit* can be found in [143]. The authors clarify the interesting evolution of his distinction between risk and uncertainty that was interpreted over time from the view of subjective probabilities versus objective probabilities to the more recent challenge of asymmetric information theory. Following the traditional interpretation of Knight's risk and uncertainty distinction, risk applies to situations where the outcome is unknown, but the odds can be

accurately measured. Uncertainty governs situations where all of the necessary information cannot be known in order to accurately determine odds. Therefore, for the person weighing the risk of congestion, risk is a single probability distribution governing a single uncertain outcome. Given this interpretation, people would be using a single probabilistic model of past ED attendances to forecast future attendances. For the situation where Knight's uncertainty regarding partial knowledge is considered, the traditional interpretation that multiple probability distributions can govern a single uncertain outcome can be adopted. Additionally, in the multiple probability distributions interpretation, uncertainty can be distinguished in terms of ambiguity and model uncertainty. Ambiguity entails multiple priors but no central, focal model. Uncertainty about probability, subjective expected utility (SEU) as a theory of choice, and its challenge in the Ellsberg thought experiment outlining unambiguous versus ambiguous probability is outlined in [144]. Explaining that ambiguity as uncertainty about probability stems from missing information that is relevant and could be known, the authors distinguish both ambiguity over probability and ambiguity over outcomes, defining ambiguity averse or risk averse, respectively. Research regarding the role of ambiguity in the decision-maker's choices has been investigated from a variety of perspectives that explore the limitations of the SEU theory of decision making under uncertainty and extensions which capture aspects of ambiguity. Ambiguity attitudes in relation to variational preferences were studied in [145]. Biseparable preferences, which include the SEU extensions Choquet expected utility (CEU) maximization (nonadditive probabilities) and multiple priors preferences, also known as "maxmin" expected utility (MEU) theory [146], were studied in [147]. The separation between ambiguity, as a characteristic of a decision-maker's subjective beliefs, and ambiguity attitude, which characterizes the decision-maker's tastes is modeled in [148]. While decision theoretic extensions were beyond the scope of our initial investigations, ongoing work that addresses reliable ambiguity measures would provide insights regarding wait time information and decisions about congestion tolerance. A mixed models approach using qualitative survey methods to underpin a quantitative model, such as the KPR extension, is a possibility.

Alternatively, uncertainty can be addressed in terms of model uncertainty, which considers multiple distributions that are centered on a particular model. It is derived from robust control theory. For an in-depth discussion of robustness and control theory the reader is referred to [149], where the authors clarify that when a decision maker is trying to make optimal decisions when her model is correct, standard control theory is used. Alternatively, robust control theory informs how to make good enough decisions when her model only approximates the correct model. A review of classical control theory with an emphasis on the role of feedback to reduce the effects of uncertainty in various systems is also addressed in the framework of model uncertainty, could entail simulations where agents do not know which model correctly describes attendances and they would need to have the capacity to develop a rudimentary approximation of the underlying data generating process of congestion to consider perturbations of their model when making their

decisions. Learning mechanisms which facilitate understanding model uncertainty could be explored from this viewpoint.

While Bayesian statistical methods applied to problems where uncertainty and information can be expressed by a probability distribution are the ideal way to address randomness, uncertainty reflects more than randomness. Uncertainty can take many forms: fuzziness or the lack of definite or sharp distinctions; discord or disagreement in choosing among several alternatives; ambiguity or one-to-many relationships; nonspecificity or two or more alternatives left unspecified [151]. Modeling subjective judgments made by decision makers coping with uncertainty often has less to do with the stochastic nature of uncertainty than it does with vagueness in human thought. We quote the work of Luce and Raiffa in Games and Decisions, who in 1957 commented: "So the second change we propose in utility theory is to admit we shall be dealing with fuzzy subjective probabilities, not sharp objective ones" [152]. The mathematical formalism of "fuzzy," however, did not take on its distinction as a multi-valued logic complementary with two-valued Boolean logic until 1965 when Lotfi A. Zadeh published his seminal work "Fuzzy Sets" [153]. Zadeh's fuzzy logic has a precise math definition; the fuzziness describes real-world applications permeated with the many forms of uncertainty. While fuzzy logic and probability are different, they are complementary ways of expressing uncertainty and both can be used to represent subjective belief. As discussed earlier regarding SEU, probability theory uses the concept of subjective probability which answers the question: how probable do I think that a variable is in a set? Fuzzy set theory uses the concept of set membership functions, and asks a different question: how much or to what degree is a variable in a set? A possibility measure is derived for this question, which is different from a probability measure. Membership functions (or the values false and true) operate over the range of real numbers [0.0, 1.0], and can handle partial truth by degree of membership in a fuzzy set. Classical Boolean logic relies on one of only two values: 0 (nonmembership) or 1 (membership). The emergence of knowledge-based systems in Artificial Intelligence created the impetus for handling a wider scope of uncertainty and fueled the need to address the limitations of probability theory as a satisfactory model of subjective uncertainty given empirical deviations of the SEU model. These limitations are formulated in [154] and can be summarized as five main areas. The need for a reference set of exhaustive and mutually exclusive elementary events which mentally may change such that events are imprecisely perceived. The additive rule is insufficiently flexible for the way humans handle grades of uncertainty. Probability theory cannot model weak states if knowledge regarding the uncertainty about some event is only loosely related to the uncertainty about the contrary event. Total ignorance (i.e., when the probabilities are unknown) cannot be expressed by a probability measure. The uncertainty numbers people use are not reliable, but rather are fuzzy probabilities. An in-depth discussion of these limitations can be found in [155]. For an overview of modern methods for uncertainty modeling the interested reader can find detailed treatments in [156-164]. Insights regarding the equivalence of fuzzy sets to random sets or loci of two-point conditional probabilities can be found in [165]. An engineering reliability motivation for investigating imprecise probabilities from the viewpoint of random sets and fuzzy sets as consonant random sets is discussed in [166], as well as the coexistence of fuzzy and random sets for modeling perception-based information gathering processes or coarsening schemes in [167]. We plan to explore ambiguity and vagueness in several facets of the wait time information, learning, and decision processing using possibility theory and fuzzy logic in extensions of our work. Q-learning, a close contender to human learning, has also been previously shown to converge to the Nash equilibrium in Minority games of increasing difficulty [168]. It has previously been used in a Minority game analog, a market selection game [169]. Fuzzy Q-learning has also been explored in related research [170]. We also envision exploring social network analysis (SNA) and its impact on learning to quantify negative ambulatory care sentiment, specifically as it relates to patients of low SES, who prefer the ED over various primary care options for routine care.

8 Conclusion

The Institute of Medicine's Crossing the Quality Chasm set a foundation for the goals of health care, indicating that it should be "safe, effective, patient-centered, timely, efficient, and equitable" [171]. Our game theoretic approach is a different tenor for health care, shedding new light on three of these tenets—safety, timeliness, and efficiency. The literature on forward looking research on the science of ED crowding calls for comparative effectiveness of ED interventions based on modeling several waiting time, quality, and satisfaction measures [172]. Variance is a critical measure in health care, because it is fundamental to understanding and maintaining safety and quality service measures, particularly for the ED. Therefore, we have examined the games emphasizing their variance and volatility measures in relation to our real-world observations. Both the Minority and El Farol Network Congestion games predict attendances which were shown to be a statistically significant fit to our empirical attendances. Extensions of this developmental phase of game theory models are planned, in line with a comparative effectiveness approach which includes operations research, machine learning, and systems science methods.

Ultimately, the ED is a CAS; therefore achieving predictability so as to avoid congestion will likely receive continued priority in the evolving healthcare landscape. In order to "nudge" the system along a path toward efficient states, quantifying characteristics of congestion with decentralized decision makers using game theory and agent-based models has helped us pin down the possibilities for improvements a little more tightly. We have tried to convey some well-vetted analogs of systems, networks, and applications with common themes: repeated scenarios of many agents (people), whose interactions (games) produce complex adaptive feedbacks, while using memory and learning, to compete for a common resource (Tragedy of the Commons). The direction for our future work is to expand the armamentarium of approaches along these themes, and quantify the interplay between information and congestion in an ED. It is our hope that with the addition of game theory models, predictions from a variety of approaches for a complex system can be compared and applied where they are most appropriate. In conclusion, we note that progress in science is often punctuated by thought experiments simple stories and toy models that capture a salient point which opens the mind to new ways of thinking. To Newton's apple, Schrödinger's cat, and Archimedes' bath, among others, we would add the El Farol Bar. New ways of thinking about healthcare's problems need to supplant often used conventional wisdom. Overcrowding problems would be a good place to start. Congestions games and related concepts, such as the counterintuitive notion that building our way of congestion can make things worse, have fruitfully informed other analogous problem domains. Adding ED beds would thus be of little utility if Yogi Berra's nonsequitur begins to take on new meaning in the medical commons: "Nobody goes there anymore. It's too crowded." Therefore, we have found the bar to be a most instructive paradigm to investigate crowding problems in health care and plan for the future.

References

- 1. Arthur WB (1994) Inductive reasoning and bounded rationality (the El Farol problem). Am Econ Rev 84:406–411
- Pines M, Hilton JA, Weber EJ, Alkemade AJ, Shabanah HA, Anderson PD, Bernhard M, Bertini A, Gries A, Ferrandiz S, Kumar VA, Harjola V-P, Hogan B, Madsen B, Mason S, Öhlén G, Rainer T, Rathlev N, Revue E, Richardson D, Sattarian M, Schull MJ et al (2011) International perspectives on emergency department crowding. Acad Emerg Med 18 (12):1358–1370
- Jayaprakash N, O'Sullivan R, Bey T, Ahmed SS, Lotfipour S (2009) Crowding and delivery of healthcare in emergency departments: the European perspective. West J Emerg Med 10 (4):233–239
- Chen CF, Ho WH, Chou HY, Yang SM, Chen IT, Shi HY (2011) Long-term prediction of emergency department revenue and visitor volume using autoregressive integrated moving average model. Comput Math Meth Med 2011:395690. doi:10.1155/2011/395690
- Clarey AJ, Cooke MW (2012) Patients who leave emergency departments without being seen: literature review and English data analysis. Emerg Med J 29:617–621
- 6. Goodacre S, Webster A (2005) Who waits longest in the emergency department and who leaves without being seen? Emerg Med J 22:93–96
- Goldman RD, Macpherson A, Schuh S, Mulligan C, Pirie J (2005) Patients who leave the pediatric emergency department without being seen: a case-control study. Can Med Assoc J 172(1):39–43
- Grosgurin O, Cramer B, Schaller M, Sarasin FP, Rutschmann OT (2013) Patients leaving the emergency department without being seen by a physician: a retrospective database analysis. Swiss Med Wkly 143(w13889):1–5
- 9. Fayyaz J, Khursheed M, Umer M, Mehmood A (2013) Missing the boat: odds for the patients who leave ED without being seen. BMC Emerg Med 13(1):1–9
- Kennedy M, MacBean C, Brand C, Sundararajan V, Taylor DM (2008) Review article: leaving the emergency department without being seen. Emerg Med Australa 20:306–313
- 11. Martin M, Champion C, Kinsman L, Masman K (2010) Mapping patient flow in a regional Australian emergency department: a model driven approach. Int Emerg Nurs 19:75–85

- Lega F, Mengoni A (2008) Why non-urgent patients choose emergency over primary care services? empirical evidence and managerial implications. Health Policy 88:326–338
- 13. van Moll Charante EP, ter Riet G, Bindels P (2008) Self-referrals to the A&E department during out-of-hours: patients motives and characteristics. Patient Educ Couns 70:256–265
- Ahmad N, Ghani NA, Kamil AA, Tahar RM (2012) Emergency department problems: a call for hybrid simulation. In: Proceedings of the world congress on engineering vol III WCE 2012
- 15. Laskowski M, McLeod RD, Friesen MR, Podaima BW, Alfa AS (2009) Models of emergency departments for reducing patient wait times. PLoS ONE 4(7):e6127
- 16. Neighbour R, Oppenheimer L, Mukhi S, Friesen MR, McLeod RD (2010) Agent based modeling of "crowdinforming" as a means of load balancing at emergency departments. Public Health Inform 2(3):e4
- 17. Burstrom L, Starrin B, Engstrom ML, Thulesius H (2013) Waiting management at the emergency department—a grounded theory study. BioMed Cent Health Serv Res 13:95
- Ekelund U, Kurland L, Eklund F, Torkki P, Letterstal A, Lindmarker P, Castren M (2011) Patient throughput times and inflow patterns in Swedish emergency departments. A basis for ANSWER, A National Swedish Emergency Registry. Trauma Resusc Emerg Med 19 (37):1–10
- 19. Higginson I (2012) Emergency department crowding. Emergency Medicine Journal 29:437-443
- Yip A, McLeod S, MacRae A, Xie B Influence of publicly available online wait time data on emergency department choice in patients with noncritical complaints. CJEM 14(4):233–242
- Forster AJ (2005) An agenda for reducing emergency department crowding. Ann Emerg Med 45(5):479–481
- 22. Mah R (2009) Emergency department overcrowding as a threat to patient dignity. Can J Emerg Med Care 11(4):365–369
- Segal E, Verter V, Colacone A, Afilalo M (2006) The in-hospital interval: a description of EMT time spent in the emergency department. Prehospital Emerg Care 10:378–382
- Elkum N, Fahim M, Shoukri M, Al-Madouj A (2009) Which patients wait longer to be seen and when? A waiting time study in the emergency department. East Mediterr Health J 15 (2):416–424
- 25. Hallas P, Ekelund U, Bjørrnsen L, Brabrand M (2013) Hoping for a domino effect: a new specialty in Sweden is a breath of fresh air for the development of Scandinavian emergency medicine. Scand J Trauma Resusc Emerg Med 21(26):1–3
- Brekke KR, Siciliani L, Straume OR (2008) Competition and waiting times in hospital markets. J Public Econ 92:1607–1628
- 27. Lewis RJ (2004) Academic emergency medicine and the 'tragedy of the commons'. Acad Emerg Med 11:423-427
- Smith M, Feied C (1999) The emergency department as a complex system. New England Complex Systems Institute, accessed Jul 1 2013. http://www.necsi.org/docs/miscellaneous/ Misc/complexity%20necsi%20paper-02f.pdf. Accessed 19 July 2014
- Kaplan RS, Porter ME (2011) How to solve the cost crisis in health care. Harv Bus Rev 89:46–52
- Wilper AP, Woolhandler S, Lasser KE, McCormick D, Cutrona SL, Bor DH, Himmelstein DU (2004) Waits to see an emergency department physician: U.S. trends and predictors, 1997–2004. Health Aff 27:w84–w95
- Armour S (2013) Under Obamacare, emergency rooms may get even more crowded. Bloomberg BusinessWeek. http://www.businessweek.com/articles/2013-05-30/underobamacare-emergency-rooms-may-get-even-more-crowded. Accessed 19 July 2014
- 32. Falvo T, Grove L, Stachura R, Vega D, Stike R, Schlenker M, Zirkin W (2007) The opportunity loss of boarding admitted patients in the emergency department. Acad Emerg Med 14:332–337

- Pines JM, Pilgrim RL, Schneider SM (2011) Practical implications of implementing emergency department crowding interventions: summary of a moderated panel. Acad Emerg Med 18:1278–1282
- Foley M, Kifaieh N, Mallon WK (2011) Financial impact of emergency department crowding. West J Emerg Med 12:192–197
- 35. Eitel DR, Rudkin SE, Malvehy MA, Killeen JP, Pines JM (2008) Improving service quality by understanding emergency department flow: a white paper and position statement prepared for the American Academy of Emergency Medicine. J Emerg Med 38:70–79
- Schweigler LM, Desmond JS, McCarthy ML, Bukowski KJ, Ionides EL, Younger JG (2009) Forecasting models of emergency department crowding. Acad Emerg Med 16:301–308
- McCarthy ML, Ding R, Pines JM, Zeger SL (2011) Comparison of methods for measuring crowding and its effects on length of stay in the emergency department. Acad Emerg Med 18:1269–1277
- Sun BC, Hsia RY, Weiss RE, Zingmond D, Liang L-J, Han W, McCreath H, Asch SM (2013) Effect of emergency department crowding on outcomes of admitted patients. Ann Emerg Med 61(6):605–611
- 39. Pines JM, Pollack CV, Diercks DB, Chang AM, Shofer FS, Hollander JE (2009) The association between emergency department crowding and adverse cardiovascular outcomes in patients with chest pain. Acad Emerg Med 16:617–625
- 40. Sprivulis PC, Da Silva JA, Jacobs IG, Frazer ARL, Jelinek GA (2006) The association between hospital overcrowding and mortality among patients admitted via Western Australian emergency departments. Med J Aust 184:208–221
- 41. Chalfin DB, Trzeciak S, Likourezos A, Baumann B, Dellinger RP (2007) Impact of delayed transfer of critically ill patients from the emergency department to the intensive care unit. Crit Care Med 35:1477–1483
- 42. Pines JM, Hollander JE, Localio AR, Metlay JP (2006) The association between emergency department overcrowding and hospital performance on antibiotic timing for pneumonia and percutaneous intervention for myocardial infarction. Acad Emerg Med 13:873–878
- Schull MJ, Vermeulen M, Slaughter G, Morrison L, Daly P (2004) Emergency department crowding and thrombolysis delays in acute myocardial infarction. Ann Emerg Med 44:577–585
- 44. Pines JM, Hollander JE (2008) Emergency department crowding is associated with poor care for patients with severe pain. Ann Emerg Med 51:1–5
- 45. Bernstein SL, Aronsky D, Duseja R, Epstein S, Handel D, Hwang U, McCarthy M, McConnell KJ, Pines JM, Rathlev N, Schafermeyer R, Zwemer F, Schull M, Asplin BR et al (2009) The effect of emergency department crowding on clinically oriented outcomes. Acad Emerg Med 16:1–10
- 46. US Government Accountability Office (2009) Hospital emergency departments: crowding continues to occur and some patients wait longer than recommended time frames. GAO Report 1–52
- 47. Pennsylvania Patient Safety Advisory (2010) Managing patient access and flow in the emergency department to improve patient safety. PA Patient Safety Auth 7:123–134
- Health Grades (2012) Emergency Medicine in American Hospitals. Health Grades Inc Denver CO pp 1–10
- 49. Simonet D (2009) Cost reduction strategies for emergency services: insurance role, practice changes and patients accountability. Health Care Anal 17:1–19
- Marco CA, Moskop JC, Schears RM, L'Hommedieu Stankus J, Bookman KJ, Padela AI, Baine J, Bryant E (2012) the ethics of health care reform: impact on emergency medicine. Acad Emerg Med 19:461–468
- 51. Smulowitz PB, Lipton R, Wharam JF, Adelman L, Baugh CW, Schuur JD, Liu SW, McGrath ME, Liu B, Sayah A, Burke MC, Pope JH, Landon BE (2011) Emergency department utilization after the implementation of Massachusetts health reform. Ann Emerg Med 58:225–234

- 52. Kangovi S, Barg FK, Carter T, Long JA, Shannon R, Grande D (2013) Understanding why patients of low socioeconomic status prefer hospitals over ambulatory care. Health Aff 32:1196–1203
- Weissman JS, Bigby JA (2009) Massachusetts health care reform-near-universal coverage at what cost? NEJM 361:2012–2015
- American College of Emergency Physicians (2014) Massachusetts health care: ED utilization increasing. http://www.acep.org. Accessed 19 July 2014
- 55. Long SK, Stockley K (2009) Emergency department visits in Massachusetts: who uses emergency care and why? Urban Institute, Washington
- Asplin BR, Magid DJ, Rhodes KV, Solberg LI, Lurie N, Camargo CA Jr (2003) A conceptual model of emergency department crowding. Ann Emerg Med 42:173–180
- 57. ARENA® Simulation Software version 14.5 © 2013 Rockwell Automation, Inc
- 58. CART Salford Predictive Modeler® version SPM7 ©1984-2012 Salford Systems
- 59. Kelton WD, Sadowski RP, Sadowski DA (2006) Simulation with Arena, 5th edn. McGraw-Hill Professional, New York
- 60. Breiman L, Friedman J, Stone C, Olshen R (1984) Classification and regression trees. In: Wadsworth Statistics/Probability Series. Chapman and Hall/CRC
- 61. Webb A (2002) Statistical pattern recognition. Wiley, New York
- Vieth T, Rhodes K (2006) The effect of crowding on access and quality in an academic ED. Am J Emerg Med 24(7):787–794
- McHale P, Wood S, Hughes K, Bellis MA, Demnitz U, Wyke S (2013) Who uses emergency departments inappropriately and when—a national cross-sectional study using a monitoring data system. BMC Med 11:258
- 64. Moustakides George V, Polunchenko Aleksey S, Tartakovsky Alexander G (2011) A numerical approach to performance analysis of quickest change-point detection procedures. Stat Sin 21(2):571–596
- 65. Hsia RY, Asch SM, Weiss RE, Zingmond D, Liang L-J, Han W, McCreath H, Sun BC (2011) Hospital determinants of emergency department left without being seen rates. Ann Emerg Med 58:24–32
- 66. Parekh KP, Russ S, Amsalem DA, Rambaran N, Wright SW (2013) Who leaves the emergency department without being seen? a public hospital experience in Georgetown, Guyana. BMC Emerg Med 13:1–6
- Neergaard L (2010) Some ERs post wait times by text, billboard. AP, MSNBC. http://www. msnbc.msn.com/id/38820121/ns/health-health_care/t/some-ers-post-wait-times-textbillboard. Accessed 19 July 2014
- NJ hospital to show ER wait times on billboards. AP, MSNBC. http://www.msnbc.msn.com/ id/42557663/ns/health-health_care/t/nj-hospital-show-er-wait-times-billboards. Accessed 12 April 2011. (Retrieved 19 July 2014)
- 69. O'Reilly KB. Posting emergency wait times: good marketing or good medicine? American medical news. Amednews.com. http://www.ama-assn.org/amednews/2010/10/11/prl21011. htm. Accessed 19 July 2014
- Niska R, Bhuiya F, Xu J (2010) National hospital ambulatory medical care survey: 2007 emergency department summary. Natl Health Stat Report 26:1–31
- Groeger L (2013) How long will you wait at the emergency room? ProPublica. http://www. propublica.org/article/how-long-will-you-wait-at-the-emergency-room. Accessed 19 July 2014
- American College of Emergency Physicians (2012) Publishing wait times for emergency department care: an information paper, pp 1–5
- 73. Hemaya SA, Locker TE (2012) How accurate are predicted waiting times, determined upon a patient's arrival in the emergency department? Emerg Med J 29:316–318
- 74. Jouriles N, Simon EL, Griffin P, Williams CJ, Haller NA (2013) Posted emergency department wait times are not always accurate. Acad Emerg Med 20:421–423

- 75. Wilensky U (2004) NetLogo minority game model. Center for connected learning and computer-based modeling, Northwestern Institute on Complex Systems, Northwestern University, Evanston, IL. http://ccl.northwestern.edu/netlogo/models/MinorityGame
- Wilensky U (1999) NetLogo. Center for connected learning and computer-based modeling, Northwestern Institute on Complex Systems, Northwestern University, Evanston IL. http:// ccl.northwestern.edu/netlogo
- 77. Wilensky U (2003) NetLogo El Farol network congestion model. Center for Connected learning and computer-based modeling, Northwestern Institute on Complex Systems, Northwestern University, Evanston, IL. http://ccl.northwestern.edu/netlogo/models/ ElFarolNetworkCongestion
- Rand W, Wilensky U (2007) Netlogo El Farol model. Center for connected learning and computer-based modeling, Northwestern Institute on Complex Systems, Northwestern University, Evanston, IL. http://ccl.northwestern.edu/netlogo/models/ElFarol
- Bell AM, Sethares WA (2001) Avoiding global congestion using decentralized adaptive agents. IEEE Trans Sig Proc 49:2873–2879
- Bell AM, Sethares WA, Bucklew JA (2003) Coordination failure as a source of congestion in information networks. IEEE Trans Sig Proc 51:875–885
- Fogel DB, Chellapilla K, Angeline PJ (1999) Inductive reasoning and bounded rationality reconsidered. IEEE Trans Evol Comp 3:142–146
- Casti JL (1996) Seeing the light at El Farol: a look at the most important problem in complex systems theory. Complexity 1:7–10
- 83. Challet D, Marsili M, Ottino G (2004) Shedding light on El Farol. Phys A 332:479-482
- 84. Franke R (2003) Reinforcement learning in the El Farol model. J Econ Behav Org 51:367–388
- 85. Marsili M, Challet D, Zecchina R (2000) Exact solution of a modified El Farol's bar problem: efficiency and the role of market impact. Phys A: Stat Mech Appl 280:522–553
- Whitehead D (2008) The El Farol bar problem revisited: reinforcement learning in a potential game. ESE Discuss Pap 186:1–30
- Zambrano E (2004) The interplay between analytics and computation in the study of congestion externalities: the case of the El Farol problem. J Pub Econ Theory 6:375–395
- Rosenthal RW (1973) A class of games possessing pure-strategy Nash equilibria. Int J Game Theory 2:65–67
- Meyers CA, Schulz AS (2008) The complexity of congestion games. Massachusetts Institute of Technology, Cambridge, pp 1–16
- Facchini G, van Megen F, Borm P, Tijs S (1997) Congestion models and weighted Bayesian potential games. Theory Dec 42:193–206
- 91. Sandholm W (2001) Potential games with continuous player sets. J Econ Theory 97:81-108
- 92. Ui T (2000) A Shapley value representation of potential games. Games Econ Behav 31:121–135
- Voorneveld M (1997) Equilibria and approximate equilibria in infinite potential games. Econ Lett 56:163–169
- 94. Voorneveld M, Borm P, van Megen F, Tijs S, Facchini G (1999) Congestion games and potentials revisited. Int Game Theory Rev 1:283–299
- 95. Monderer D, Shapley LS (1996) Potential games. Games Econ Behav 14(1):124-143
- Liu M, Wu Y (2008) Spectrum sharing as congestion games. Communication control and computing, pp 1146–1153
- 97. Vocking B, Aachen R (2006) Congestion games: optimization in competition. In: Proceedings of the 2nd algorithms and complexity in durham workshop, pp 1–12
- Challet D, Zhang YC (1997) Emergence of cooperation and organization in an evolutionary game. Phys A 246:407–418
- 99. Marsili M, Challet D, Zecchina R (2000) Exact solution of a modified El Farol's bar problem: efficiency and the role of market impact. Phys A: Stat Mech App 280:522–553
- 100. Challet D, Marsili M, Zecchina R (1999) Statistical mechanics of systems with heterogeneous agents: minority games. Phys Rev Lett 84:1824

- 101. Johnson NF, Jefferies P, Hui PM (2003) Financial market complexity. Oxford University Press, Oxford
- 102. Johnson NF, Hart M, Hui PM, Zheng DL (2000) Trader dynamics in a model market. Int J Thoer App Fin 3:443–450
- 103. Jefferies P, Hart ML, Hui PM, Johnson NF (2001) From market games to real-world markets. Euro Phys J B-Cond Matter Comp Sys 20:493–501
- Kalinowski T, Schulz HJ, Briese M (2000) Cooperation in the minority game with local information. Phys A 277:502–508
- 105. Chmura T, Pitz T (2006) Successful strategies in repeated minority games. Phys A 363:477–480
- Hart M, Jefferies P, Johnson NF, Hui PM (2000) Crowd-anticrowd model of the minority game. Phys A 298:537–544
- 107. Bottazzi G, Devetag G (2004) Coordination and self-organization in minority games: experimental evidence. The Complex Dynamics of Economic Interaction. Springer, Heidelberg, pp 283–300
- 108. Kurihara S, Fukuda K, Hirotsu T, Akashi O, Sato S, Sugawara T (2004) How does collective intelligence emerge in complex environment? The Complex Networks of Economic Interactions Lecture Notes in Economics and Mathematical Systems Essays in Agent-Based Economics and Econophysics, vol 567, pp 279–289
- Johnson NF, Hart M, Hui PM (1999) Crowd effects and volatility in markets with competing agents. Phys A 269:1–8
- 110. Bottazzi G, Devetag G, Dosi G (2001) Adaptive learning and emergent coordination in minority games. Simul Model Pract Theory 10:321–347
- 111. Beier R, Czumaj A, Krysta P, Vocking B (2004) Computing equilibria for congestion games with (im)perfect information. In: Proceedings of the 15th annual ACM-SIAM symposium on discrete algorithms. Society for Industrial and Applied Mathematics, pp 746–755
- 112. Savit R, Manuca R, Riolo R (1999) Adaptive competition, market efficiency, and phase transitions. Phys Rev let 82(10):2203–2206
- 113. Johnson NF, Jefferies P, Hui PM (2003) Financial market complexity. Oxford University Press, Oxford
- 114. Galla T, Mosetti G, Zhang Y (2006) Anomalous fluctuations in minority games and related multi-agent models of financial markets, pp 1–26. arXiv:http://arxiv.org/abs/arXiv.physics/ 0608091 [physics.soc-ph]
- 115. Grimm V, Berger U et al (2010) The ODD protocol: a review and first update. Ecol Model 221:2760–2768
- 116. Garofalo M (2006) Modeling the 'El Farol Bar Problem' In NetLogo. Preliminary Draft, Dexia Bank Belgium
- 117. Chakraborti A, Challet D, Chatterjee A, Marsili M, Zhang YC, Chakrabarti BK (2013) Statistical mechanics of competitive resource allocation, pp 1–24. arXiv:http://arXiv.org/abs/ 1305.2121 [physics.soc-ph]
- 118. Chakrabarti AS, Chakrabarti BK, Chatterjee M, Mitra M (2009) The Kolkata paise restaurant problem and resource utilization. Phys A 388:2420–2426
- 119. Mahmassani HS, Jayakrishnan R (1991) System performance and user response under realtime information in a congested traffic corridor. Transp Res A 25A(5):293307
- 120. Kitamura R, Nakayama S (2007) Can travel time information really influence network flow? implications of the minority game. Transp Res Rec 2010:14–20
- 121. Dobson G, Pinker E (2006) The value of sharing lead time information. IIE Trans $38{:}171{-}183$
- 122. Chorus C, Arentze T, Molin E, Timmermans HJP (2006) The value of travel information: decision strategy-specific conceptualizations and numerical examples. Transp Res Part B 40:504–519
- 123. Levinson D (2003) The value of advanced traveler information systems for route choice. Transp Res Part C 11:75–87

- 124. Allon G, Bassamboo A (2011) The impact of delaying the delay announcements. Oper Res 59(5):1198–1210
- 125. Debo L, Veeraraghavan S (2010) Prices and congestion as signals of quality. Chicago booth working paper
- 126. Armony M, Shimkin N, Whitt W (2009) The impact of delay announcements in many-server queues with abandonment. Oper Res 57(1):66–81
- 127. Cui S, Veeraraghaven S (2014) Blind queues: the impact of consumer beliefs on revenues and congestion. JEL classification: D60, D80, L10 working papers series. http://dx.doi.org/ 10.2139/ssrn.2196817. Accessed 19 July 2014
- 128. Garcia D, Archer T, Moradi S, Ghiabi B (2012) Waiting in vain: managing time and customer satisfaction at call centers. Psychology 3(2):213–216
- 129. El Azouzi R, Altman E, Pourtallier O (2005) Braess paradox and properties of wardrop equilibrium in some multiservice networks. In: Haurie A, Zaccour G (eds) Dynamic games, theory and applications. Springer, New York, pp 57–77
- 130. Mondragon RJ (2006) Optimal networks, congestion and Braess' paradox. In: Proceedings from the 2006 workshop on interdisciplinary systems approach in performance evaluation and design of computer and communications systems, ACM, pp 1–9
- 131. Chudak F, Eleuterio VDS (2006) The traffic equilibrium problem. IFOR Miteilungen, pp 1-4
- 132. Nagurney A, Boyce D (2005) Preface to "on a paradox of traffic planning". Trans Sci 39:4435–4445
- 133. Steinberg R, Zangwill WI (1983) The prevalence of Braess' paradox. Trans Sci 17:301-318
- 134. Zhang H-F, Yang Z, Wu Z-X, Wang BH, Zhou T (2013) Braess's paradox in epidemic game: better condition results in less payoff. Sci Rep 3(3292):1–8
- 135. Nicolaides C, Cueto-Felgueroso L, Juanes R (2013) The price of anarchy in mobility-driven contagion dynamics. J R Soc Interface 10(87):1742–5662
- 136. Knight V, Harper P (2013) Selfish routing in public services. Eur J Oper Res 230(1):122-132
- 137. Andersen J, Sornette D (2005) A mechanism for pockets of predictability in complex adaptive systems. Europhys Lett 70(5):697–703
- 138. Johnson N, Lamper D, Jeffries P, Hart M, Howison S (2001) Application of multi-agent games to the prediction of financial time series. Phys A 299:222–227
- 139. Lamper D, Howison S, Johnson N (2002) Predictability of large future changes in a competitive evolving population. Phys Rev Lett 88(017902):U190–U192
- Groot R, Musters P (2004) Minority game of price promotion in fast moving consumer goods markets. Phys A 350:553–557
- 141. Wiesinger J, Sornette D, Satinover J (2010) Reverse engineering financial markets with majority and minority games using genetic algorithms. Swiss Finance Institute Research Paper Series
- 142. Knight FH (1921) Risk, uncertainty and profit. Houghton, Mifflin
- 143. Langlois R, Cosgel M (1993) Frank Knight on risk uncertainty, and the firm: a new interpretation. Economic inquiry, vol XXXI, pp 456–465
- 144. Camerer C, Weber M (1992) Recent developments in modeling preferences: uncertainty and ambiguity. J Risk Uncertainty 5:325–370
- 145. Maccheroni F, Marinacci M, Rustichini A (2006) Ambiguity aversion, robustness, and the variational representation of preferences. Econometrica 74(6):1447–1498
- 146. Gilboa I, Schmeidler D (1989) Maxmin expected utility with a non-unique prior. J Math Econ 18:141–153
- 147. Ghirardato P, Marinacci M (2002) Ambiguity made precise: a comparative foundation. J Econ Theory 102:251–289
- Klibanoff P, Marinacci M, Mukerji S (2005) A smooth model of decision making under ambiguity. Econometrica 73(6):1849–1892
- 149. Hansen LP, Sargent TJ (2001) Robust control and model uncertainty. Econ Rev 91(2):60-66
- 150. Å ström KJ (2000) Model uncertainty and robust control. COSY workshop—ESF course Valencia Spain

- 151. Klir G, Yuan B (1995) Fuzzy sets and fuzzy logic: theory and applications. Prentice-Hall, New York
- 152. Luce RD, Raiffa H (1957) Games and decisions: introduction and critical survey. Dover Publications, New York
- 153. Zadeh LA (1965) Fuzzy sets. Inf Control 8:338-353
- 154. Dubois D, Prade H (1988) Decision evaluation methods under uncertainty and imprecision. In: Kacprzyk J, Fedrizzi M (eds) Combining fuzzy imprecision with probabilistic uncertainty in decision making, lecture notes in economics and mathematical systems, vol 310. Springer, Berlin, pp 48–65
- 155. Dubois D, Prade H (1986) Recent models of uncertainty and imprecision as a basis for decision theory. Towards less normative frameworks. In: Hollnagel E, Mancini G, Woods DD (eds) Intelligent decision support in process environments, NATO-ASI Series F.21,1986. Springer, Berlin, pp 3–24
- 156. Dymowa L (2011) Soft computing in economics and finance intelligent systems reference library, vol 6. Springer, Berlin, pp 41–105
- 157. Dubois D, Prade H (1980) Fuzzy sets and systems: theory and applications. Academic Press, London
- 158. Zadeh LA (1978) Fuzzy sets as a basis for a theory of possibility. Fuzzy Sets Syst 1:3-28
- 159. Dubois D, Prade H (1995) Possibility theory as a basis for qualitative decision theory. In: Proceedings of the 14th IJCAI95: 1924–1930. Morgan Kaufmann, San Francisco
- 160. Klir GJ (2005) Uncertainty and information. Foundations of generalized information theory. Wiley, Hoboken
- 161. Dubois D, Prade H, Smets P (1996) Representing partial ignorance. IEEE Trans Syst Man Cybern Part A Syst Hum 26(3):361–377
- 162. Dubois D, Prade H (1988) Possibility theory. Plenum Press, New York
- 163. Yager RR (1979) Possibilistic decision making. IEEE Trans Syst Man Cybern 9:338-392
- 164. Kosko B (1990) Fuzziness vs. Probability. Int J General Syst 17:211-240
- 165. Kosko B (1994) Fuzzy systems as universal approximators. IEEE Trans Comput 43 (11):1329–1333
- 166. Bernardini A, Tonon F (2009) Extreme probability distributions of random sets, fuzzy sets, and p-boxes. Int J Reliab Saf 3(1/2/3):57–78
- 167. Nguyen HT (2005) Fuzzy and random sets. Fuzzy Sets Syst 156 (2005):349-356
- 168. Catteeuw D, Manderick B (2009) Learning in the time-dependent minority game. In: GECCO '09 Proceedings of the 11th annual conference companion on genetic and evolutionary computation conference: late breaking papers: 2011–2016
- Ishibuchi H, Sakamoto R, Nakashima T (2003) Learning fuzzy rules from iterative execution of games. Fuzzy Sets Syst 135:213–240
- 170. Ishibuchi H, Nakashima T, Miyamoto H, Oh C (1997) Fuzzy Q-learning for a multi-player -30 repeated game. In: Proceeding of 6th IEEE international conference on fuzzy systems, pp 1573–1579
- 171. Institute of Medicine (2001) Crossing the quality chasm: a new health system for the 21st century. Washington DC National Academies Press, Washington
- 172. Pines J, Yealy D (2009) Advancing the science of emergency department crowding: measurement and solutions. Ann Emerg Med 54(4):511–513

Author Index

A

Aachen, R., 270, 271 Abbosh, A., 146, 149, 150 Abdel-Rahim, A., 116 Abido, M.A., 94 Adelman, L., 258 Adjerohb, D.A., 5 Adler, N., 177, 178 Afilalo, M., 257 Aguirregabiria, V., 127, 178 Ahmad, N., 256 Ahmadi, H., 85 Ahmed, I., 5 Ahmed, S.S., 256 Akashi, O., 273 Akella, A., 28 Akyildiz, F., 5 Alagoz, F., 5, 14 Albrecht, P., 101 Alfa, A.S., 256 Aliprantis, C.D., 116, 132, 155 Alkemade, A.J., 256 Allan, R., 101 Allon, G., 282 Al-Madouj, A., 257 Altman, E., 283 Amelin, M., 85 Ammar, R.A, 147 Amsalem, D.A., 263 Ancel, E., 219, 235 Andersen, J., 284 Anderson, P.D., 256 Andersson, G, 86-88, 109 Angeline, P.J., 266, 275, 277, 281 Ansari, J., 9 Antoniou, J., 27, 29, 33, 33-38, 40-45, 47, 159 Aram, A., 33 Archer, T., 282 Arentze, T., 282

Armony, M., 282 Armour, S., 258 Aronsky, D., 258 Arthur, W.B., 256, 266 Asch, S.M., 258, 263 Asplin, B.R., 258, 260 Å ström, K.J., 285 Axelrod, R.M., 33

B

Bacci, V., 9 Baillo, A., 87 Baine, J., 258 Bajcsy, R., 4 Balakrishnan, A., 116 Balakrishnan, H., 5, 14 Baldick, R., 86, 88 Ball, M., 174 Ball, M.O., 121 Barg, F.K., 259 Barnett, A., 175 Barnhart, C., 173, 174, 177, 178 Basar, T., 32 Baseler, R., 176 Baskar, S., 91, 94 Bassamboo, A., 282 Baugh, C.W., 258 Baumann, B., 258 Bean, N., 154, 159, 161, 165 Bedford, T., 231 Beier, R., 273 Bekebrede, G., 220, 222 Bell, A.M., 266, 268, 272, 276, 277, 282 Bell, M.G.H., 114, 122, 159, 168 Bellana, S.K., 14 Bellis, M.A., 261 Belobaba, P., 174, 176 Bennett, A.C., 57 Berger, U., 275, 287

© Springer International Publishing Switzerland 2015 K. Hausken and J. Zhuang (eds.), *Game Theoretic Analysis of Congestion, Safety and Security*, Springer Series in Reliability Engineering, DOI 10.1007/978-3-319-13009-5 Bernhard, M., 256 Bernstein, S.L., 258 Bertini, A., 256 Bey, T., 256 Bhavaraju, M., 101 Bhuiya, F., 264 Bier, V.M., 57, 85, 117, 142, 143, 145-147, 158 Bigby, J.A., 259 Billinton, R., 101 Bindels, P., 256 Binggeli, U., 177 Bjørrnsen, L., 257 Bockstael-Blok, W., 220 Bompard, E., 85 Bonnefoy, P., 174 Bookman, K.J., 258 Boot, W.R., 56 Bor, D.H., 258 Borenstein, S., 84 Borm, P., 32, 270 Bottazzi, G., 273 Boyce, D., 283 Brabrand, M., 257 Braess, D., 154 Brand, C., 256, 263 Brander, J.A., 178 Breiman, L., 260 Brekke, K.R., 257 Brewer, G., 221 Briese, M., 273 Brueckner, J.K., 178, 198 Bruso, K., 146, 147, 155 Bryant, E., 258 Bucklew, J.A., 266, 268, 272, 276, 277, 282 Bukowski, K.J., 258 Burke, M.C., 258 Burstrom, L., 256 Bushnell, J., 84 Butenko, S., 132 Button, K., 176, 177, 203

С

Callaway, E., 10 Camargo, C.A. Jr, 260 Camerer, C., 285 Campbell, A.T., 4 Cannon-Bowers, J.A., 157 Cappanera, P., 153, 159, 160 Caromi, R., 146, 149, 150 Carter, T., 259 Casti, J.L., 266, 277 Castillo, E., 97 Castren, M., 256, 261 Castrillo, D.P., 155 Catteeuw, D., 287 Cayirci, E., 5 Chae, C.B., 51 Chakrabarti, A.S., 281 Chakrabarti, B.K., 281 Chakrabarti, S.K., 116, 132, 155 Chakraborti, A., 281 Chalfin, D.B., 258 Challet, D., 266, 273, 281 Champion, C., 256 Chan, Y., 113, 114, 116, 118-120, 123-125, 127, 131, 133, 134, 141-144, 146, 147, 149, 150, 153-155, 159, 165 Chandrakasan, A.P., 5, 14 Chang, A.M., 258 Chang-Albitres, C.M., 132 Chappin, E.J.L., 220, 221 Chari, M.K., 120 Chatterjee, A., 281 Chatterjee, M., 17, 281 Chaturvedi, K.T., 94 Chellapilla, K., 266, 275, 277, 281 Chen, A., 114, 127 Chen. C.F., 256 Chen, C.T., 98 Chen, I.T., 256 Chen, P., 10 Chen, S., 2, 3 Chen, S.S., 85 Chin, R., 222 Chin, Y.Z., 57 Chinery, R., 85 Chiravath, S.S., 57 Chiti, F., 9 Chmura, T., 273 Cholvi, V., 32 Chong, S.Y., 33 Chorus, C., 282 Chou, H.Y., 256 Chudak, F., 283 Chun, S.A., 221 Church, R.L., 153 Clarey, A.J., 256, 263 Colacone, A., 257 Colbourn, C.J., 120, 121 Cominetti, R., 167 Conejo, A.J., 86, 87, 96, 97, 101 Connors, S., 220

Contreras, J., 86 Cooke, M.W., 256, 263 Cooke, R., 231 Correa, J.R., 167 Correia, P., 85 Cosgel, M., 284 Cox, L.A. Jr, 122 Cramer, B., 256 Cresswell, A.M., 220 Croxton, K.L., 153, 159, 161 Cruz, A., 89 Cueto-Felgueroso, L., 283 Cui, S., 282 Culler, D.E., 4 Cutrona, S.L., 258 Czumaj, A., 273

D

Da, Silva, J.A., 258 Dai, J.G., 150, 159, 160 Daly, P., 258 Dananjayan, P., 10, 14 Daneshmand, M., 3 Daniel, J., 198 Daniels, M., 114 DaSilva, L.A., 32 Dayoub, I., 10 de la Torre, S., 86 de Neufville, R., 220-222 Debbah, M., 32 Debo, L., 282 Dehkordi, H, 33-35, 40-45, 47 Del Vecchio, J.R., 146, 147, 155 Dellinger, R.P., 258 Demirkol, I., 5, 14 Demnitz, U., 261 Desmond, J.S., 258 DesRoches, S., 115, 132 Dessante, P., 87 Devetag, G., 273 Devore, J.L., 56 Diercks, D.B., 258 Diestel, R., 30 Dietz, D.C., 150, 153 Dijkema, G.P.J., 220 Ding, R., 258 Ding, Y., 57 Ding, Z., 91, 94 Dobson, G., 177, 178, 282 Dong, S., 85 Dormans, J., 226 Dorton, S.L., 57

Dosi, G., 273 Downward, A., 167 Dresner, M., 174 Drexler, J., 176, 177, 203 Dubois, D., 286 Duke, R., 221 Duke, R.D., 221, 222, 226–228 Duseja, R., 258 Dymowa, L., 286

Е

Eckel, N., 56 Ee, C.T., 4 Eisenman, S.B., 4 Eitel, D.R., 258 Ekelund, U., 256, 257, 261 Eklund, F., 256, 261 El Azouzi, R., 283 Eleuterio, V.D.S., 283 El-Hoiydi, A., 5 Elkum, N., 257 Elzen, B., 220, 221 Engstrom, M.L., 256 Epstein, S., 258 Erosy, C., 5, 14 Estrin, D., 4 Evans, J.R., 149

F

Facchini, G., 270 Fahim, M., 257 Falvo, T., 258 Fayyaz, J., 256, 263 Fei, X., 152 Feied, C., 257 Feldman, R.M., 132 Fernadez, F.R., 117 Fernandez, A., 32 Ferrandiz, S., 256 Fleming, A., 56 Fletcher, B.A., 85 Flores-Fillol, R., 178 Fogel, D.B., 266, 275, 277, 281 Foley, M., 258 Fonzone, A., 159 Ford, L.R. Jr, 128, 148 Forster, A.J., 257 Franceschetti, M., 32 Frank, H., 114 Franke, R., 266 Frantzeskaki, N., 220, 221 Frazer, A.R.L., 258

Friedman, J., 260 Friesen, M.R., 256 Frisch, I.T., 114 Fu, F., 85 Fukuda, K., 273 Fulkerson, D.R., 128, 148

G

Galla, T., 273 Gamberini, R., 57 Gan, D., 101 Garcia, A., 156, 159, 162, 165, 167 Garcia, D., 282 García-Bertrand, R., 86, 87, 96, 101 Garofalo, M., 275 Gasson, S., 56 Gaukler, G.M., 57 Geels, F.W., 220 Gendron, B., 153, 159, 161 Georgilakis, P., 101 Gesbert, D., 51 Geurts, J.L.A., 221, 222, 226, 228 Ghani, N.A., 256 Gheorghe, A., 219 Ghiabi, B., 282 Ghirardato, P., 285 Gilboa, I., 285 Gil-García, J.R., 220, 221 Gilliam, R.R., 57 Glachant, J.M., 87 Goldman, R.D., 256 Gong, D.W., 91, 94 Goodacre, S., 256, 261 Goodman, D., 178 Govindan, R., 4 Grande, D., 259 Greenblat, C.S., 221, 250 Gries, A., 256 Griffin, P., 265 Grigg, C., 101 Grimm, V., 275, 287 Groeger, L., 264 Groot, R., 284 Groot Bruinderink, R., 155 Grosgurin, O., 256 Gross, G., 85 Grossman, W.M., 33 Grove, L., 258 Gubbins, K.E., 57 Guha, R.K., 32 Guikema, S.D., 132, 167 Guyer, M., 246

H

Haapola, J., 9, 10, 14 Haimes, Y.Y., 122, 124 Hallas, P., 257 Haller, N.A., 265 Han, W., 258, 263 Han, Z., 32 Han, Zhu., 32 Handel, D., 258 Hansen, L.P., 285 Hansen, M., 174, 177, 178 Hansman, R.J., 220-222 Haphuriwat, N., 57 Harback, K., 198 Harjola, V.-P., 256 Harker, P.T., 178 Harms, D.D., 120 Harper, P., 284 Hart, M., 273, 279, 284 Hart, M.L., 273 Hassan, J., 31 Hatziargyriou, N.D., 101 Hausken, K., 85, 117, 142, 145-147, 158 Head, E., 116, 123, 142, 143 Health, Grades., 258 Heath, R.W. Jr, 51 Heidemann, J., 4 Heijne, G., 250 Heinzelman, W.B., 5, 14 Heinzow, T., 87 Helman, U., 86, 91 Hemaya, S.A., 265 Hicks, I.V., 132 Higginson, I., 257 Hilton, J.A., 256 Himmelstein, D.U., 258 Hirotsu, T., 273 Hjorungnes, A., 32 Ho, C.Y., 178 Ho, W.H., 256 Hobbs, B.F., 86, 88, 91 Hogan, B., 256 Hollander, J.E., 258 Holmgren, J.Å., 167 Holyland, A., 114 Hong, S., 117, 178 Howison, S., 284 Hsia, R.Y., 258, 263 Hsieh, C.-C., 150 Hu, J., 159 Hu, Y., 3 Huang, G.L., 149

Hughes, K., 261 Hui, K.-P., 154, 159, 161, 165 Hui, P.M., 273, 279 Hung-po, C., 88, 91 Hwang, C.L., 99 Hwang, S.A., 85 Hwang, U., 258

I

Iaedeluca, J., 175 Ibrahim, S., 147 Iida, Y., 114 Intanagonwiwat, C., 4 Ionides, E.L., 258 Ishibuchi, H., 287

J

Jacobs, I.G., 258 Janakiram, M., 135 Janssen, M., 221 Jayakrishnan, R., 282 Jayaprakash, N., 256 Jefferies, P., 273, 279, 284 Jelinek, G.A., 258 Jenelius, E., 167 Jiao, L.W., 85 Johnson, B., 116 Johnson, N., 284 Johnson, N.F., 273, 279 Johnson, W., 56 Jordan, K.E., 85 Jorswieck, E., 29 Jouriles, N., 265 Juanes, R., 283 Judd, G., 28

K

Kahn, A.E., 176, 177, 203 Kalinowski, T., 273 Kalpana, G., 10 Kalyoncu, H., 116 Kamil, A.A., 256 Kang, D.H., 132 Kangovi, S., 259 Kanturska, U., 159 Kaplan, R.S., 257 Kaplan, S.M., 85 Karamchandani, N., 32 Karray, F., 154 Kasu, S.R., 14 Katz, R.H., 14 Kelton, W.D., 57, 260, 275 Kemfert, C., 87 Kendall, G., 33 Kennedy, M., 256, 263 Keseric, N., 87 Keskinocak, P., 135 Khursheed, M., 256, 263 Kifaieh, N., 258 Killeen, J.P., 258 Kinsman, L., 256 Kirschen, D., 84, 91 Kitamura, R., 282 Klabbers, J.H.G., 220, 221 Kleinrock, L., 7 Klemperer, P.D., 86 Klibanoff, P., 285 Klir, G., 286 Klir, G.J., 286 Knight, F.H., 284 Knight, V., 284 Kongsomsaksakul, S., 127 Kosko, B., 286 Kountouris, M., 51 Koutsoupias, E., 198, 201 Kovenock, D., 157, 159, 163, 166 Kraetzl, M., 154, 159, 161, 165 Krafft, J.A., 86-88, 109 Kramer, A.F., 56 Krause, T., 86-88, 109 Krishnamurthy, L., 4 Kroese, D.P., 154, 159, 161, 165 Kronik, L., 57 Krugler, P.E., 132 Krysta, P., 273 Kuik, O., 87 Kumar, A., 33, 85 Kumar, C., 14 Kumar, V.A., 256 Kunz, F., 85 Kurihara, S., 273 Kurland, L., 256, 261 Kwait, K.A., 17 Kwak, K.S., 14, 85

L

L'Hommedieu Stankus, J., 258 Lai, L.L, 84 Lai, Y.-S., 127, 133, 134 Lam, W.H.K., 114, 131 Lamper, D., 284 Landfeldt, B., 31 Landon, B.E., 258 Langlois, R., 284 Larsson, E., 29 Laskowski, M., 256 Lasser, K.E., 258 Lazarou, G.Y., 7 Lederer, P.J., 177, 178 Lee, A.H., 87, 99 Lee, E.S., 98 Lee, K.H., 86-88, 109 Lega, F., 256 Lesani, H., 85 Leshem, A., 32 Lesta, V.P., 29, 33-38, 40-45, 47 Letterstal, A., 256, 261 Levinson, D., 131, 282 Lewis, R.J., 257 Li, B., 3, 4 Li, C., 57 Li, D., 32 Li, J., 7 Li, Z., 123, 132 Liang, L.-J., 258, 263 Libman, L., 27, 29, 33-38, 40-45, 47 Likourezos, A., 258 Lim, C., 142, 144, 146, 147 Lin, C.Y., 87, 99 Lin, M.-H., 150 Lin, Q., 3, 4, 15 Lin, S., 85 Lin, W., 150, 159, 160 Linderhof, V., 87 Lindmarker, P., 256, 261 Lipton, R., 258 Lise, W., 87 Liu, B., 87, 258 Liu, J., 32 Liu, M., 270 Liu, S.W., 258 Liu, Y., 86-88, 109 Liu, Z., 97 Lleras, G., 114 Localio, A.R., 258 Locker, T.E., 265 Lolli, F., 57 Long, J., 85 Long, J.A., 259 Long, S.K., 259 Loorbach, D., 220, 221 Lopez, L., 32 Losada, C., 150, 159, 160, 164 Lotfipour, S., 256 Lou, Y., 121, 159 Lownes, N.E., 147

Luan, F., 87 Luce, R.D., 286 Luna-Reyes, L.F., 220 Luo, M., 85 Luo, X., 85 Lurie, N., 260 Lyle, D., 116, 123, 142, 143

М

Ma, D., 85 MacBean, C., 256, 263 Maccheroni, F., 285 Machado, R., 15 MacKenzie, A.B., 32 Macpherson, A., 256 MacRae, A., 257, 264 Maddox, M.W., 57 Madsen, B., 256 Magee, C., 220-222 Magid, D.J., 260 Magnanti, T.L., 153, 159, 161 Mah, R., 257 Mahmassani, H.S., 282 Mahonen, P., 9, 10, 14 Majumdar, A., 175 Maku, T., 135 Malaviya, N., 3 Mallon, W.K., 258 Malvehy, M.A., 258 Mandayam, N., 178 Mandelbaum, S.A., 57 Manderick, B., 287 Manuca, R., 273, 279 Manzini, R., 57 Marco, C.A., 258 Marín, J., 89 Marinacci, M., 285 Marsh, A., 114, 118, 119, 123, 125 Marsili, M., 256, 266, 273, 281 Masman, K., 256 Mason, S., 256 Mayer, I., 220, 222, 226 McCarley, J.S., 56 McCarthy, J., 146, 147, 165 McCarthy, M., 258 McCarthy, M.L., 258 McConnell, K.J., 258 McCormick, D., 258 McCreath, H., 258, 263 McGrath, M.E., 258 McHale, P., 261 McLaughlin, B.J., 115, 132

McLay, L.A., 56 McLeod, R.D., 256 McLeod, S., 257, 264 Mehmood, A., 256, 263 Mehta, S., 14 Meijer, S., 222 Mengoni, A., 256 Merrick, J.R.W., 56 Metlay, J.P., 258 Meyer, M.A., 86 Meyers, C.A., 270 Milano, F., 86, 87, 96, 97, 101 Minero, P., 32 Mínguez, R., 97 Mira, P., 127 Mirchandani, P., 116 Mitra, M., 281 Miyamoto, H., 287 Moavenzadeh, F., 220 Molin, E., 282 Monderer, D., 270, 271 Mondragon, R.J., 283 Moradi, S., 9, 282 Morrison, L., 258 Mosetti, G., 273 Moskop, J.C., 258 Moustakides George, V., 262 Mukerji, S., 285 Mukhi, S., 256 Mulligan, C., 256 Muñoz, A., 89 Murillo-Sánchez, C.E., 101 Murray-Tuite, P.M., 152 Murrell, S.D., 115, 132 Musa, J.D., 57 Musters, P., 284 Myerson, R.B., 57

Ν

Nagurney, A., 283 Nakashima, T., 287 Nakayama, S., 282 Nash, J., 34, 246–248 Nash, J.F., 34 Natarajan, H.P., 116 Neels, K., 174 Neergaard, L., 264 Neighbour, R., 256 Nesa Sudha, M., 3 Nguyen, H.T., 287 Ni, Y.X., 85 Nicolaides, C., 283 Nijkamp, P., 177, 178 Niska, R., 264 Nojima, N., 114 Norman, V.D., 178

0

O'Brien, L., 154 O'Connor, W.E., 176, 177, 203 O'Dea, B., 10 O'Hanley, J.R., 150, 159, 160, 164 O'Mahony, M., 154 O'Neill, R.P., 88, 91 O'Reilly, K.B., 264 O'Sullivan, R., 256 Odoni, A., 174 Oh, C., 287 Öhlén, G., 256 Oloomi-Buygi, M., 98, 101 Olshen, R., 260 Oman, P., 116 Oppenheimer, L., 256 Orfanos, G.A., 101 Östling, R., 87 Ottino, G., 266

P

Padela, A.I., 258 Page, L.B., 120 Pandit, M., 94 Pang, J.S., 86, 91 Papadimitriou, C., 198, 201 Papadopoulou-Lesta, V., 27 Papaefthymiou, G., 97 Parekh, K.P., 263 Park, S., 14 Patek, S.D., 156, 159, 162, 165, 167 Patino, D., 114 Paull, G., 175 Peleg, B., 32, 33 Pels, E., 177, 178 Pendergraft, D.R., 57 Peng, M., 5 Peng, Z., 85 Pennsylvania, 258 Periyasamy, R., 1 Perry, J.E., 120 Perumal, D., 1 Peters, V., 250 Peterson, E., 174 Pettigrew, R., 135 Philpott, A.B., 167 Pickett, K.W., 132 Pilgrim, R.L., 258 Pines, J., 287

Pines, J.M., 258 Pines, M., 256 Pinker, E., 282 Pirie, J., 256 Pitsillides, A., 27, 29, 33-35, 40-45, 47, 159 Pitsillides, S., 29, 35-38 Pitz, T., 273 Podaima, B.W., 256 Pollack, C.V., 258 Polunchenko Aleksey, S., 262 Pomalaza-Raez, C., 10, 14 Pompeo, L., 177 Pope, J.H., 258 Porter, K., 85 Porter, M.E., 257 Pourtallier, O., 283 Prade, H., 286 Provan, J.S., 120, 121 Puerto, J., 117

R

Radhakrishnan, S., 4 Raghunathan, V., 14 Rahimi-Kian, A., 81, 86-88, 98, 109, 101 Raiffa, H., 286 Rainer, T., 256 Raines, R., 120, 123, 124, 127, 131, 133, 134, 142, 144, 146, 147, 155 Raja, P., 10, 14 Rajasekaran, S., 147 Rakshit, S., 32 Rambaran, N., 263 Ramos, A., 87 Rand, W., 265 Rapoport, A., 246 Rathlev, N., 256, 258 Ratliff, J., 15 Rausand, M., 114 Regattieri, A., 57 Revue, E., 256 Rhodes, K., 261 Rhodes, K.V., 260 Richardson, D., 256 Riet, G., 256 Rietveld, P., 177, 178 Riihijarvi, J., 9 Riolo, R., 273, 279 Rittel, H.W.J., 221 Rivier, M., 87 Roberson, B., 157, 159, 163, 166 Robertson, C.V., 57 Robins, R., 220-222 Robinson, A.R., 150, 153

Roos, D., 220–222 Rosen, J.B., 195, 197 Rosenthal, R.W., 270 Rothkopf, M.H., 88, 91 Roughgarden, T., 198 Rouvaen, J.M., 10 Roxy, P., 56 Rudkin, S.E., 258 Russ, S., 263 Rustichini, A., 285 Rutschmann, O.T., 256

S

Saad, W., 32 Saaty, T.L., 99, 107 Saber, A.Y., 91, 101 Sadowski, D.A., 260, 275 Sadowski, R.P., 57, 260, 275 Saguan, M., 87 Sahraei-Ardakani, M., 86-88, 109 Saizer, T., 51 Sakamoto, R., 287 Salas, E., 157 Salehizadeh, M.R., 81, 98, 101 Samuelson, L., 157 Sanchez-Silva, M., 114 Sánchez-Úbeda, E.F., 89 Sancho, N.G.F., 149 Sandholm, W., 270 Sankarasubramaniam, Y., 5 Sankur, B., 116 Saranga, V., 4 Sarasin, F.P., 256 Saraydar, C., 178 Sargent, T.J., 285 Sarkar, S., 33 Satinover, J., 284 Sato, S., 273 Sattarian, M., 256 Savit, R., 273, 279 Sayah, A., 258 Scaparra, M.P., 150, 153, 159, 160, 164 Schafermeyer, R., 258 Schaller, M., 256 Schavland, J., 120, 123, 124, 127, 131, 133, 134, 142, 144, 146, 147, 155 Schears, R.M., 258 Schlenker, M., 258 Schmeidler, D., 285 Schmocker, J.-D., 159 Schneider, S.M., 258 Schuh, S., 256 Schull, M., 258

Schull, M.J., 256, 258 Schulz, A.S., 270 Schulz, H.J., 273 Schurgers, C., 14 Schuur, J.D., 258 Schweigler, L.M., 258 Segal, E., 257 Sengupta, S., 17 Seshan, S., 28 Sethares, W.A., 266, 268, 272, 276, 277, 282 Shabanah, H.A., 256 Shakir, M., 5 Shannon, R., 259 Shapira, Y., 57 Shapley, L.S., 270, 271 Sharma, D., 147 Shebli, F., 10 Shelby, Z., 9, 10, 14 Shelfer, K.M., 56 Sheremeta, R.M., 157, 159, 163, 166 Sherry, L., 174 Shi, H.Y., 256 Shier, D., 121 Shih, H.S., 98 Shimkin, N., 282 Shiratori, R., 249 Shofer, F.S., 258 Shoukri, M., 257 Shrader, S., 57 Shustorovich, E., 57 Shyur, H.J., 98 Siciliani, L., 257 Sikdar, B., 6 Simon, E.L., 265 Simonet, D., 258 Simpson, R.W., 176 Singh, C., 33, 101 Singh, S.N., 85 Sinha, K., 156, 159, 162, 165, 167 Sinha, K.C., 123, 132 Sirisena, H., 31 Slaughter, G., 258 Smets, P., 286 Smith, J.C., 142, 144, 146, 147 Smith, M., 257 Smith, R.E., 132 Smulowitz, P.B., 258 Sohraby, K., 3, 4 Solberg, L.I., 260 Son, Y.S., 86, 88 Song, C., 55 Song, F., 149 Sornette, D., 284

Sprivulis, P.C., 258 Srivastava, L., 94 Srivastava, M.B., 14 Srivastava, S.C., 85 Stachura, R., 258 Starrin, B., 256 Steenkiste, P., 28 Steinberg, R., 283 Stemm, M., 14 Steuer, R., 134, 156 Steuer, R.E., 124 Stier-Moses, N.E., 167 Stike, R., 258 Stirling, W.C., 135 Stockley, K., 259 Stone, C., 260 Strandenes, S.P., 178 Straume, O.R., 257 Strbac, G., 84, 91 Su, W., 5 Subramahnian, E., 221 Sudholter, P., 32, 33 Sugawara, T., 273 Sun, B.C., 258, 263 Sun, H., 85 Sun, L., 15 Sun, L.-J., 3, 4 Sundarapandian, V., 57 Sundararajan, V., 256, 263 Supekar, N.S., 85 Suris, J.E., 32 Swets, N.B., 57 Szeto, W.Y., 154, 155, 159, 162, 164

Т

Tadelis, S., 57 Tahar, R.M., 256 Taneja, N.K., 176 Tang, W.H., 114 Tartakovsky Alexander, G., 262 Taylor, D.M., 256, 263 Tekinay, S., 15 Telang, N.G., 85 Ter, 256 Tessema, B., 97 Thulesius, H., 256 Tijs, S., 155, 270 Timmermans, H.J.P., 282 Tobagi, F.A., 7 Torkki, P., 256, 261 Trani, A., 174 Trzeciak, S., 258 Tung, L.-W., 116

U

Ui, T., 270 Umer, M., 256, 263

v

Valarmathi, M.L., 3 Valentin, E.C., 220 van de Nouweland, A., 32 van den Nouweland, A., 155 Van Dender, K., 198 van der Sluis, L., 97 van Golstein Brouwers, W., 32, 155 van Houten, S.P., 222 van Hove, J.C., 146, 149, 150 van Megen, F., 270 van Moll Charante, E.P., 256 van Muijen, S., 221 Vargas, L.G., 99, 107 Varma, S., 3 Vaze, V., 173, 177, 178 Veeraraghavan, S., 282 Vega, D., 258 Veit, D.J., 86-88, 109 Venayagamoorthy, G.K., 91, 101 Ventosa, M., 87 Verbraeck, A., 222 Verheggen, E., 255 Vermeulen, M., 258 Vermeulen, P.A.M., 221, 222, 226, 228 Verter, V., 257 Vidhya, J., 10 Vidoni, E.D., 56 Vieth, T., 261 Visalakshi, S., 91, 94 Vissers, G., 250 Vocking, B., 270, 271, 273 Voorneveld, M., 270 Vries, L.J.D., 221

W

Wakabayashi, H., 114 Wan, C.Y., 4 Wang, B.H., 283 Wang, C., 3, 4 Wang, Q., 147 Wang, R., 15 Wang, R.-C., 3, 4 Wang, W., 5 Wang, X., 32, 55, 76, 85 Wang, X.D., 149 Washburn, A., 117 Webb, A., 260, 262 Webber, M.M., 221 Weber, E.J., 256 Weber, M., 285 Webster, A., 256, 261 Wei, W., 177, 178 Weidlich, A., 86-88, 109 Weiss, R.E., 258, 263 Weissman, J.S., 259 Wen, F.S., 85 Wenzler, I., 221 Westin, J., 167 Wettstein, D., 155 Wharam, J.F., 258 Whitehead, D., 266, 268, 269 Whitt, W., 282 Wickens, C.D., 56 Wieczorek, A., 220, 221 Wiesinger, J., 284 Wilensky, U., 265, 272, 275 Wilke, S., 175 Williams, C.J., 265 Wilper, A.P., 258 Wolfram, S., 56 Won, J., 154 Wong, P., 101 Woo, A., 4 Wood, K., 117 Wood, S., 261 Woolhandler, S., 258 Wright, S.W., 263 Wu, F., 85 Wu, F.F., 86-88, 109 Wu, J., 85 Wu, Y., 270 Wu, Z.-X., 283 Wyke, S., 261

Х

Xia, S., 135 Xie, B., 257, 264 Xu, G., 114, 131 Xu, J., 264 Xu, Y., 32

Y

Yadav, R., 3 Yager, R.R., 286 Yaghmaee, M.H., 5 Yan, Z., 85 Yang, C.N., 87, 99 Yang, H., 6, 114 Yang, N., 2, 3 Yang, Q., 85 Yang, S.M., 256 Yang, Z., 283 Yao, X., 33 Ye, W., 4 Yealy, D., 287 Yim, E., 114, 118, 119, 123, 125 Yip, A., 257, 264 Yoon, D., 57 Yoon, K.P., 99 Younger, J.G., 258 Yuan, B., 286

Z

Zadeh, L.A., 286 Zakeri, G., 167 Zambrano, E., 266, 268, 270 Zangwill, W.I., 283 Zecchina, R., 266, 273 Zeger, S.L., 258 Zehavi, E., 32 Zeltyn Marmor, Y.N., 57 Zeng, M., 87 Zhang, A., 178 Zhang, H.-F., 283 Zhang, J., 87, 220 Zhang, L., 121, 131, 159 Zhang, Y., 91, 94, 273 Zhang, Y.C., 273, 281 Zhang, Y.P., 85 Zhang, Z., 87 Zheng, D.L., 273, 279 Zheng, T., 4 Zhou, T., 283 Zhou, Z., 127 Zhuang, J., 55, 57, 76, 85, 117, 142, 143, 145-147, 158 Zimmerman, R.D., 101 Zingmond, D., 258, 263 Zirkin, W., 258 Zou, B., 174 Zou, X., 85 Zwemer, F., 258

Subject Index

A

Air transportation, 223, 228, 230, 231, 238–240, 242 Airline competition, 178, 179, 203, 204 Airline probability, 183 Airport congestion, 174, 175, 179, 198, 202–204

B

Best response, 175, 178, 181–185, 192, 195, 203, 207, 208, 210, 211

С

Congestion games, 269–271 Cooperation, 28–33, 37–40, 47, 51 Cooperative game, 32, 33, 35

D

Degree of inefficiency, 198 Dense Wi-Fi access points, 51

Е

El Farol Bar Game, 265, 266, 273, 275, 277 Electricity market, 85–87, 89, 91 Emergency department overcrowding, 255 Energy, 87, 91 Expected maximum-flow, 118 Expert elicitation, 224

G

Game theory, 28, 32–34, 55, 57, 79, 244–248, 258, 265, 268, 271, 282, 283 Generators, 81–98, 100, 101, 103–105, 107–109 Graph theory, 34 Graphical game, 34, 38, 51

I

Imperfect screening, 58, 69 Independent system operator, 81, 84 Information value theory, 166 Infrastructure planning, 221

L

Leader-follower game, 85, 87, 101, 108, 109

М

Minority game, 265, 266, 273, 274, 277–279, 284, 287 Multicriteria decision-making, 135

N

Nash equilibrium, 116, 126, 131 Network reliability, 114, 116, 117, 120, 121, 127 Network security, 116, 122, 136, 142, 143 Non-cooperative game, 117, 126, 128

P

Posturing, 141 Power systems, 85, 87, 97, 108 Price of anarchy, 175, 179, 198, 201–203

S

Security screening policy, 55 Serious gaming, 221, 249 Shapley value, 155 Stochastic network, 114, 118, 120, 122, 125, 148–150, 153

Т

Tragedy of the commons, 255, 257, 266

© Springer International Publishing Switzerland 2015 K. Hausken and J. Zhuang (eds.), *Game Theoretic Analysis of Congestion, Safety and Security*, Springer Series in Reliability Engineering, DOI 10.1007/978-3-319-13009-5 Transmission congestion management, 87, 88, 95, 97, 108, 109 Two-stage queueing network, 55

U

Unmanaged wireless deployment, 28

W

Waiting time, 55-57, 60, 61, 66, 67, 79