

Essays in Operations and Inventory Management

A THESIS

SUBMITTED TO THE FACULTY OF THE GRADUATE SCHOOL OF THE
UNIVERSITY OF MINNESOTA

BY

Yimin Yu

IN PARTIAL FULFILLMENT OF THE REQUIREMENTS

FOR THE DEGREE OF

Doctor Of Philosophy

Adviser: Saif Benjaafar

August, 2009

© Yimin Yu 2009

ALL RIGHTS RESERVED

Acknowledgements

I am indebted to many people who guided, supported, amused, entertained, encouraged, and motivated me during my Ph.D study at the University of Minnesota. Without their assistance and support this thesis would not have been possible.

First of all I would like to thank my supervisor, Saif Benjaafar, without whose expert guidance, financial support, feedback, patience and all-around great supervising skills this thesis would not be attempted.

Sincere thanks must go to Bill Cooper and Bharath Rangarajan, who have taken time to comment on my work, send me papers, or otherwise assist. Among my fellow students, I wish to thank Dan Zhang, Lei Wang, Tingliang Huang, Le Li, Yu Wang, Fei Li, Dong Xiao, Xi Chen and Ang Liu for the good memories we shared.

Nearer to my heart, I wish to thank my beloved parents and my brother for their endless support. Finally, I owe my present wonderful life to my wife and my daughter, without whom I could never have done any of this.

Abstract

This thesis consists of two essays. The first essay is on capacity pooling and cost sharing among independent firms in the presence of congestion. We analyze the benefit of production/service capacity sharing for a set of independent firms. Firms have the choice of either operating their own production/service facilities or investing in a facility that is shared. Facilities are modeled as queueing systems with finite service rates. Firms decide on capacity levels (the service rate) to minimize delay costs and capacity investment costs possibly subject to service level constraints. If firms decide to operate a shared facility they must also decide on a scheme for sharing the costs. We formulate the problem as a cooperative game and identify a cost allocation that is in the core. The allocation rule charges every firm the cost of capacity for which it is directly responsible, its own delay cost, and a fraction of buffer capacity cost that is consistent with its contribution to this cost. In settings where unit delay costs are private information, the cooperative capacity sharing game becomes embedded with a non-cooperative information reporting game. We show how a cost allocation rule can be designed to induce all firms to report truthfully this information. Moreover, we show that, under this allocation rule, truth telling is a dominant strategy, with each firm reporting truthfully its private information regardless of the reporting decisions of other firms.

The second essay is on a customer-item decomposition approach to inventory problems. We consider inventory systems with periodic review, correlated, non-stationary stochastic demand and correlated, non-stationary stochastic and sequential leadtimes. We treat systems with both single and multiple stages. We use the customer-item decomposition approach to decompose the associated inventory control problem into sub-problems, each involving a single customer-item pair. We then formulate each subproblem as an optimal stopping problem. We use properties that arise from this formulation to show that the

optimal policy is a state-dependent base-stock policy and to show, for the case of positive demand, that the optimal policy can be obtained via an algorithm whose complexity is polynomial in the length of the planning horizon. We also use the formulation to construct myopic heuristics which lead to explicit solutions for the optimal policy in the form of a critical fractile. We characterize conditions under which the myopic heuristics are optimal. We show how the results can be extended to systems with advance demand information and batch ordering.

Contents

Acknowledgements	i
Abstract	ii
List of Tables	vi
List of Figures	vii
1 Introduction	1
2 Capacity Pooling and Cost Sharing among Independent Firms in the Presence of Congestion	4
2.1 Introduction	4
2.2 Related Literature	9
2.3 Systems without Capacity Sharing	12
2.4 Capacity Sharing with Full Information	17
2.4.1 Capacity Optimization	19
2.4.2 Cost Sharing	22
2.5 Capacity Sharing with Incomplete Information	25
2.6 Extensions to Systems with General Demand and Processing Times	33
2.7 Summary and Concluding Comments	38

3	A Customer-Item Decomposition Approach to Stochastic Inventory Systems with Correlation	42
3.1	Introduction	42
3.2	Related Literature	46
3.3	Systems with a Single Stage	49
3.4	Systems with Multiple Stages in Series	63
3.5	Myopic Policies	70
3.5.1	Systems with a Single Stage	71
3.5.2	Systems with Multiple Stages	75
3.6	Extensions	78
3.6.1	Inventory Systems with Advanced Demand Information	78
3.6.2	Inventory systems with batch ordering	80
3.7	Conclusion	85
4	Future Research and Other Completed and Ongoing Research Projects	97
4.1	Possible Extensions to Capacity Pooling among Independent Firms	97
4.2	Possible Extensions of the Customer-Item Decomposition Approach to Other Stochastic Inventory Systems	99
4.3	Other Completed and Ongoing Research Projects	100
4.3.1	Inventory Systems with Concave Ordering Costs	100
4.3.2	Optimal Incentives in Supply Chain with Asymmetric Information	101
4.3.3	A Marginal Revenue Approach to Airline Seats Allocation Problems	102
4.3.4	Optimal Control of Make-to-Stock Queues	102
	Bibliography	103

List of Tables

3.1	Solution times (in seconds) for varying planning horizons	95
3.2	Solution times (in seconds) for varying demand distributions	96

List of Figures

3.1	An illustration of the demand and supply processes	53
-----	--	----

Chapter 1

Introduction

This thesis consists of three parts. The first two parts, which correspond to Chapters 2 and 3, describe two of completed research. The third part, which is included in Chapter 4, discusses other completed research projects and possible future research directions. Chapters 2 and 3 are independent, self-contained, and deal with separate topics. We expect future research to build on both of these independent research streams.

The following is a brief synopsis of Chapters 2 and 3. Chapter 2 deals with a problem involving capacity pooling and cost sharing among independent firms in the presence of congestion. We analyze the benefit of production/service capacity pooling for a set of independent firms. Firms have the choice of either operating their own production/service facilities or investing in a facility that is shared. Facilities are modeled as queueing systems with finite service rates. Firms decide on capacity levels (the service rate) to minimize delay costs and capacity investment costs subject to service level constraints. If firms decide to operate a shared facility they must also decide on a scheme for sharing the costs. We formulate the problem as a cooperative game and identify a cost allocation that is in the core (i.e., a single facility shared by all the firms is the best arrangement). The allocation rule charges every firm the cost of capacity for which it is directly responsible, its

own delay cost, and a fraction of buffer capacity (which is equal to the difference between total capacity and total demand rate) cost that is consistent with its contribution to this cost. In settings where unit delay costs are private information, the cooperative capacity sharing game becomes embedded with a non-cooperative information reporting game. We show how a cost allocation rule can be designed to induce all firms to report truthfully this information. Moreover, we show that, under this allocation rule, truth telling is a dominant strategy, with each firm reporting truthfully its private information regardless of the reporting decisions of other firms. We also analyze the benefit of pooling under a queueing model with general demand and processing times. We identify a condition under which pooling is beneficial. This chapter is based on a paper co-authored with Saif Benjaafar and Yigal Gerchak and currently in its second review with the journal *Manufacturing & Service Operations Management*.

Chapter 3 deals with a stochastic inventory control problem. We consider inventory systems with periodic review, correlated, non-stationary stochastic demand and correlated, non-stationary stochastic and sequential leadtimes. We treat systems with both single and multiple stages. We use the customer-item decomposition approach to decompose the associated inventory control problem into sub-problems, each involving a single customer-item pair. We then formulate each subproblem as an optimal stopping problem. We use properties that arise from this formulation to show that the optimal policy is a state-dependent base-stock policy and to show, for the case of positive demand, that the optimal policy can be obtained via an algorithm whose complexity is polynomial in the length of the planning horizon. We also use the formulation to construct myopic heuristics which lead to explicit solutions for the optimal policy in the form of a critical fractile. We characterize conditions under which the myopic heuristics are optimal. We show how the results can be extended to systems with advance demand information and batch ordering. This chapter

is based on a paper jointly authored with Saif Benjaafar.

Chapter 4 discusses potential directions for future research and other completed research projects. With regard to capacity pooling, we plan to extend the analysis to (1) more complex queueing systems such as systems with parallel servers or servers in a network configuration and (2) to systems with both cooperation and competition. With regard to the inventory control problem, we also expect to extend the analysis to more complex systems, such as assembly systems and distribution systems. We also plan to investigate the usefulness of the decomposition approach to other applications, such as revenue management. In this chapter we also provide a brief description other finished projects:(1) inventory systems with concave ordering costs;(2) optimal incentives in supply chains with asymmetric information;(3) a marginal approach to airline seat allocation problems.

Chapter 2

Capacity Pooling and Cost Sharing among Independent Firms in the Presence of Congestion

2.1 Introduction

Capacity sharing refers to the fulfillment of demand that arises from multiple sources from a single facility instead of facilities dedicated to each demand source. In a system without capacity sharing, each dedicated facility fulfills its own demand relying solely on its capacity. It has long been known that capacity sharing can be beneficial when demand is random. This benefit can be in the form of improved quality of service with the same amount of capacity or in the form of less capacity needed to provide the same quality of service. Capacity sharing can also be beneficial when there are economies of scale associated with acquiring capacity or fulfilling demand. These benefits have been shown to be true for various forms of capacity, including manufacturing, service, and inventory.

Capacity sharing has been studied mostly in situations where a single firm owns all the capacity in the system, has full information, and has responsibility for serving all the demand. This firm makes the decision about whether or not to share capacity and how much capacity to acquire. In this paper, we consider a system with n *independent* firms, each facing its own demand and each having the option of either operating its own independent facility or joining some or all the other firms in a shared facility. The firms may vary in their demand levels and in their tolerance for capacity shortage. They may also possess private information regarding their costs or service level requirements which they may not report truthfully. If some or all of the firms decide to share capacity, they must also decide on how to allocate the cost of the shared facility. They must do so in a *fair* manner that prevents any of the firms from defecting and perhaps sharing a facility with a subset of the firms or staying on their own. Hence, firms that contribute more to the cost of the shared facility (because of their higher usage of capacity or lower tolerance for capacity shortage) are expected to pay a greater share of total cost. In the presence of private information, the cost allocation scheme should also mitigate the possibility of firms not truthfully disclosing their private information.

Capacity sharing among independent firms is increasingly common in the manufacturing, service, and public sectors. In manufacturing, there are numerous instances of independent firms sharing the same production facilities. For example, several car manufacturers, such as Toyota and General Motors, share final assembly plants. Similar sharing arrangements can be found in the electronics industry where firms share printed circuit board assembly facilities or semiconductor manufacturing plants. In the service sector, the sharing of service facilities is also common. For example, airlines share check-in counters, maintenance facilities, and reservation systems; hospitals share expensive diagnostic equipment, laboratory facilities, and in some cases medical specialists and surgeons. In the public sector,

the sharing of resources between independent government entities is also widespread. For instance, local governments in rural communities share fire and police departments, 911 call centers, and other social services. In various sectors, independent organizations are increasingly sharing infrastructure resources such as telecommunication networks, computer services, basic manufacturing resources, and facilities for handling back-office operations. In this paper, we are motivated by such settings. That is, we are motivated by settings where firms have the option of sharing generic infrastructure resources whose capacity can be easily scaled and which can be accessed with equal efficiency by all firms.

Capacity sharing among independent firms¹ raises several important questions. For example, is capacity sharing always beneficial to all firms? Does it always lead to a reduction in total capacity in the system? How should capacity costs be allocated among the different firms? Is capacity sharing among all the firms the best arrangement or would sharing among smaller subsets of the firms be more beneficial to particular firms? Can capacity sharing be beneficial when firms do not report truthfully private information, especially when this information is used in determining capacity levels and cost allocation? Is it possible to induce firms, via cost allocation, to disclose truthfully their private information? If so, would such cost allocation ensure that all firms continue to benefit from capacity sharing?

In this paper, we address these and other related questions for a specific setting. We consider applications where facilities can be modeled as queueing systems. Demand for each firm consists of an independent stream of customers (or orders) that arrive continuously over time with random inter-arrival times. Customers are processed at each facility one at a time with stochastic service time. The capacity at each facility is determined by the rate at which

¹We use the term *independent firms* broadly to include economic entities who are independently owned and also entities, such as sub-divisions within a single firm, who may have a common owner. What is important to our analysis is that these entities are empowered to make independent decisions that minimize their individual costs.

customers can be processed. Because customers are processed one at a time and because customer arrivals and processing times are random, congestion arises and customers can experience delay prior to processing (if a customer arrives and finds the service facility busy, the customer must wait for service). Each firm can install and operate its own facility where its customers are processed. Firms make decisions about how much service capacity to acquire in order to minimize two types of costs, delay cost due to customers spending time at the facility prior to completing service and capacity investment cost, subject to a constraint on the amount of delay customers experience. Alternatively, firms may choose to collectively operate a shared facility. In that case, in addition to determining the optimal amount of capacity (taking into account the reported delay costs and service levels of all the firms), we must also determine how the corresponding cost must be allocated.

The main contributions of our paper are summarized below.

- We provide a framework for modeling capacity sharing in queueing systems with independent firms. We consider systems with and without full information. In systems with full information, parameters of all the firms are common knowledge; in systems with incomplete information, unit delay costs are private information to each firm. To our knowledge, our paper is among the first to model the issue of cooperation and capacity sharing in a queueing context and to do so for systems with and without complete information.
- We formulate capacity sharing as a cooperative game and show that for systems where facilities are modeled as $M/M/1$ queues (queues with Poisson arrivals and exponential service times) the *core* of the game is non-empty. That is, there always exists a cost allocation rule for which all the firms are better off than under any other alternative sharing arrangement, including being on their own.

- In systems with full information, we identify a simple and easy to implement allocation rule with desirable properties that is in the core. The allocation rule charges every firm the cost of capacity for which it is directly responsible, its own delay cost, and a fraction of buffer capacity cost that is consistent with its contribution to this cost.
- In systems with incomplete information, firms may act strategically and misreport their unit unit delay costs. This leads to a non-cooperative “information reporting” game embedded in the cooperative “capacity sharing” game. We show that allocation rules that are in the core under full information could lead to significant misreporting of private information in the presence of incomplete information.
- Although firms can act strategically when they possess private information, we show how a cost allocation rule can be designed to induce all firms to truthfully report their private information and to do so regardless of the reporting decisions of other firms. That is, our proposed allocation rule is *incentive-compatible* with truth-telling being a *dominant strategy*. Moreover, we show that our proposed allocation rule is in the core.
- We extend our treatment beyond the M/M/1 queue framework and discuss the extent to which our results continue to hold in more general settings.

Our choice of modeling framework is in part motivated by the fact that types of service and manufacturing facilities can be viewed as queueing systems. There is a rich literature that takes this modeling view (see Sections 2.2 and 2.3 for further discussion). Surprisingly very little of this literature addresses the issue of cooperation and capacity sharing when there are independent firms. Therefore, we view our paper as a step toward a more comprehensive examination of the issue of cooperation in queueing systems, whether it arises in services, manufacturing, or elsewhere. We also view it as a contribution, in the form

of a potentially rich application domain, to the literature on cooperative games with and without full information.

The rest of the paper is organized as follows. In Section 2.2, we provide a brief review of related literature. In Section 2.3, we treat the case with no capacity sharing. In Section 2.4, we analyze capacity sharing when there is full information. In Section 2.5, we consider the case of capacity sharing with incomplete information. We extend the analysis to systems with more general arrival and service processes in Section 2.6. In Section 2.7, we discuss additional extensions and offer concluding comments.

2.2 Related Literature

There is a rich literature on capacity *pooling* in queueing systems, with applications ranging from manufacturing and service operations to telecommunications systems to computer networks. This literature can be classified broadly as relating to either the *pooling of service rates* or the *pooling of servers*. Server rate pooling refers to the consolidation of multiple servers into a single one with a faster rate (e.g., N servers, each with service rate μ and demand rate λ , are replaced by a single server with service rate $N\mu$ and demand rate $N\lambda$). Server pooling on the other hand refers to placing multiple servers in a single facility from which all demand streams are served (e.g., N single server queues are replaced by a single multi-server queue with N servers and a demand rate $N\lambda$).

Kleinrock (1976) discusses various examples of both types of pooling. Stidham (1970) considers a design problem where the decision variables are the number of parallel servers and the service rate of each server. Smith and Whitt (1981) and Benjaafar (1995) show that server pooling, when the number of servers is exogenously determined, is beneficial as long as all customers have identical service time distributions. Buzacott (1996) considers the pooling of N servers in series, with each server dedicated to one task, into N parallel servers,

with each server carrying out all the tasks. Mandelbaum and Reiman (1998) consider the pooling of general Jackson networks into single server queues with phase-type service time distributions.

Tekin et al. (2004) use approximations to evaluate the benefit of partitioning servers in multiple pools instead of a single large one. Sheikhzadeh et al. (1998), Gurumurthi and Benjaafar (2004) and Jordan et al. (2005) study the *chaining* of servers, where each server can process customers from two customer streams and each customer can be routed to two servers. They show that in systems with homogeneous demand rates and service time requirements, chaining can achieve most of the benefits of total server pooling; see also Hopp et al. (2004), Iravani et al. (2004), Bassamboo et al. (2008), Aksin and Karaesmen (2008), Wallace and Whitt (2005) and the references therein. These papers belong to the growing literature on queueing systems with server flexibility (or cross-training); see Jouini et al. (2008), Aksin et al. (2005) and Koole and Pot (2005) for recent reviews.

The treatment in this paper is different from the above literature in four important aspects. First, we do not assume that there is a single decision maker that determines whether or not to pool. Instead, we consider multiple firms that decide independently on either operating their own facilities or sharing one with other firms (pooling here does not imply a merger however). Second, we do not assume that service capacity is exogenously given. We allow for this to be an outcome of an optimization carried out by the firms either individually or jointly. Third, we are concerned with identifying cost allocation schemes under which all firms prefer a single shared facility to any other capacity sharing arrangement, including remaining on their own, Fourth, we allow for the possibility of private information regarding delay costs or service levels and for the possibility of firms not reporting this information truthfully.

The literature dealing with capacity sharing in the context of independent firms is

limited. Gonzalez and Herrero (2004), and also Garcia-Sanz et al. (2007), consider a special case of the M/M/1 model we consider. However in both cases, they do not optimize capacity (before or after pooling), do not consider delay costs, and assume truthful reporting of all information. In our case, the presence of delay costs significantly complicates the process of cost allocation since we seek allocations that could allow for each firm to absorb its own cost of delay. We also consider systems where firms might have private information. Anily and Haviv (2008) treat a related M/M/1 model where the issue is how to allocate delay cost to ensure that the allocation is in the core. However, in their case, capacity is exogenously determined so that capacity cost is not included. They also assume full information regarding all parameters.

Dewan and Mendelson (1990) consider the problem faced by the manager of a service facility with multiple users. The manager decides on the capacity of the service facility, which is modeled as a single server queue, and on the price to charge each user. The prices affect the demand rates of the users with higher prices resulting in lower demand rates. This is different from our setting where the demand rates are exogenous. Also, in their case, customers do not have the option of operating independent facilities or forming coalitions.

Our work is of course related to the vast literature on cooperative game theory and, more broadly, the economics of coalition formation and joint ventures; see Moulin (1995) for a general introduction to the topic. Some of this literature has focused on cooperation involving sequencing and scheduling; see for example Moulin and Stong (2002), Maniquet (2003), and Katta and Sethuraman (2006). This literature sometimes refers to these problems as queueing problems. However, they typically involve a finite population of customers who simultaneously arrive to the system, and therefore are not concerned with steady state behavior and congestion in the way that we are in this paper. In Operations Management,

there is growing literature that applies cooperative game theory to joint ordering problems, particularly in the context of economic order quantity models (see Anily and Haviv (2007), Dror and Hartman (2007) and the many references therein), economic lot sizing models (see for example van den Heuvel (2007) and Chen and Zhang (2006), among others), and news-vendor models (see Muller et al. (2002), Nagarajan and Sošić (2007), Kemahlioglu-Ziya (2004), Chen and Zhang (2007), and Hanany and Gerchak (2008) and the references therein). In these papers, a major focus is proving the existence of a core allocation under the assumption that the players involved in the cooperative game have no bargaining power and the objective of each coalition is to optimize the total expected value for the coalition. Consistent with this literature, we also assume that firms have no bargaining power and the objective of each coalition is to select a capacity level that minimizes the total expected cost for that coalition.

Finally, we should note that one could view the decision to invest in a shared facility (instead of dedicated facilities) as a decision by the corresponding firms to outsource. There is a rich literature on outsourcing and procurement, including for settings where the outsourcing supplier is modeled as a queueing system; see for example, Cachon and Harker (2002), Allon and Federgruen (2006), Gans and Zhou (2007), and Benjaafar et al. (2007). In general, this literature does not deal with cost allocation or coalition formation.

2.3 Systems without Capacity Sharing

Consider a system consisting of a set $\mathcal{N} = \{1, \dots, n\}$ of n firms. Firm i , $i \in \mathcal{N}$, faces an independent demand stream with customers arriving according to a Poisson process with rate λ_i (we treat more general arrival processes in Section 2.6). When firms operate independently, each firm invests in a separate service facility and chooses a certain level of capacity in the form of a service rate. We refer to this scenario as the scenario without

capacity sharing. Once the facilities are built, each firm serves its customers from its own facility one at a time on a first-come, first-served (FCFS) basis. We assume service times are independent and identically distributed random variables denoted by X_i where X_i is of the form Y/μ_i and Y is a random variable that is exponentially distributed with a mean equal to 1. Hence, service time is also exponentially distributed with mean $E[X_i] = 1/\mu_i$. The parameter μ_i , ($\mu_i > 0$) is a scaling parameter that corresponds to the service rate or capacity.

The random variable Y can be viewed as the work content associated with each customer. We assume that work content is homogeneous across firms. Given the exponential nature of both customer inter-arrival times and service times, each firm behaves like an M/M/1 queue. There is significant literature on the economics of queues in *competitive* settings that primarily focuses on the M/M/1 queue (and where the service rate is the decision variable); see Hassin and Haviv (2003) for a review of that literature and see Cachon and Harker (2002), Cachon and Zhang (2007), Benjaafar et al. (2007), and Allon and Federgruen (2007), among many others, for example applications. Our treatment of the M/M/1 queue is consistent with assumptions made in that literature and can be viewed as complementing it for cooperative settings.

We assume that service rate can be varied continuously and that firms incur a capacity cost c per unit of service rate per unit time. This is justified in settings where capacity can be continuously scaled over a sufficiently large interval (e.g., the speed of computing facilities, the bandwidth of communication networks, or the throughput of production lines). It is also consistent with treatments elsewhere in the literature (see for example Kalai et al. (1992), Mendelson and Whang (1990), Ha (2001), Allon and Federgruen (2007, 2008), Cachon and Zhang (2007), and the vast literature reviewed therein). It is also consistent with the significant literature on capacity planning, as noted recently by Bassamboo et al.

(2008). The assumption of linear capacity cost implies that there are neither economies nor diseconomies of scale. This is an important case that has been widely studied in the literature (see Allon and Federgruen (2007, 2008), Dewan and Mendelson (1990), Stidham (1992), Cachon and Harker (2002), and Bassamboo et al. (2008) among others), leads to tractable analysis, and provides a useful benchmark for other cost structures.

We assume that the demand rate for each firm is known. This of course does not mean that demand is deterministic. Inter-arrival times between consecutive customers are stochastic. Therefore, the number of customers that arrive over a given period of time is random. The assumption of known demand rate is consistent with most of the existing literature on capacity planning in queueing systems (and indeed in most of the queueing literature); see for example Kleinrock (1976), Cachon and Harker (2002), Bassamboo et al. (2008), and Allon and Federgruen (2007, 2008), among many others.²

The objective of each firm is to minimize its capacity investment while limiting the amount of delay its customers experience. Limiting customer delay can be achieved by enforcing a service level constraint or by associating a cost with the amount of delay customers experience. A service level constraint may take several forms, including a constraint on the probability of customer delay not exceeding a specified threshold, or a constraint on expected delay not exceeding a certain maximum amount. Service level constraints are managerial decisions that typically reflect either a position in the marketplace that a firm would like to take or contractual obligations that a firm has negotiated with its customers.

Delay costs can reflect either direct or indirect costs. Direct costs are penalties incurred by the firm due to delays experienced by its customers (for example, payments to customers

²It is possible to consider systems where the demand rates themselves are random (e.g., when the demand is modulated by another process). However, depending on the assumptions made regarding this modulating process, the analysis could become significantly less tractable and we leave this as a potential area for future research.

to compensate for the total time they spend in the system) or indirect costs due to loss of customer goodwill. Hence, delay costs are not unlike backorder costs, common in inventory settings (Zipkin 2000). Delay costs may also reflect the cost of work-in-process accumulation when there is a physical product released to the queue with the arrival of each customer, as in many manufacturing applications. The use of delay costs and service levels are both common in the literature; see for example Dewan and Mendelson (1990), Mendelson and Whang (1990), Ha (1998, 2001), Allon and Federgruen (2007, 2008) and the references therein.

In this paper, we limit our analysis to the case where service level is expressed in terms of a probability that total delay in the system (time in the queue + time in service) for each customer in steady state does not exceed a specified threshold. We also limit ourselves to the case where a unit delay cost h_i is incurred for each unit of time a customer spends in the system (time either in the queue or in service in steady state) and the objective is to minimize the long run expected delay cost.

Let $z_i(\mu_i)$ denote the expected total cost incurred by firm i given a service rate μ_i (for stability, we assume that $\lambda_i/\mu_i < 1$). Let W_i , a random variable, denote the total time a customer of firm i spends in the system (customer delay) and $P(W_i \leq w_0)$ the probability that customer delay does not exceed w_0 where $w_0 \geq 0$. The problem faced by firm i can then be stated as follows

$$\text{Minimize } z_i(\mu_i) = c\mu_i + \frac{h_i\lambda_i}{\mu_i - \lambda_i} \quad (2.1)$$

subject to

$$P(W_i \leq w_0) = 1 - e^{-(\mu_i - \lambda_i)w_0} \geq \alpha_i, \quad (2.2)$$

and

$$\lambda_i/\mu_i \leq 1. \quad (2.3)$$

The objective function in the above optimization problem consists of two terms: a capacity cost term and a delay cost term, where the decision variable is the capacity level of firm i as determined by the service rate μ_i . The formulation captures two important special cases: (1) the case where $\alpha_i = 0$ for all $i \in \mathcal{N}$ and (2) the case where $h_i = 0$ for all $i \in \mathcal{N}$. The first corresponds to a pure cost-based formulation with no constraints on service levels, while the second corresponds to a service level-based formulation with no delay costs. In the absence of service level constraints, the optimal capacity level μ_i^* can be obtained from the first order condition of optimality, since z_i is convex in μ_i , as

$$\mu_i^* = \lambda_i + \sqrt{\frac{h_i \lambda_i}{c}}. \quad (2.4)$$

In systems with service level constraints but no delay costs, the optimal capacity level is given by the smallest μ_i that satisfies inequality (2.2). This leads to the following optimal capacity level

$$\mu_i^* = \lambda_i + \frac{\ln(\frac{1}{1-\alpha_i})}{w_0}. \quad (2.5)$$

In both cases, the optimal capacity is the sum of two components. The first corresponds to the demand rate, λ_i (since all demand must be satisfied) while the second corresponds to *buffer* capacity that increases in either the ratio $\frac{h_i \lambda_i}{c}$ or the service level α_i . The expressions in equations (2.4) and (2.5) are not new. Similar expressions have been derived elsewhere; see for example Kleinrock (1976), Allon and Federeguren (2008) and Hassin and Haviv (2003).

In the general case, with both delay costs and service level constraints, the optimal capacity level is given by

$$\mu_i^* = \lambda_i + \eta_i, \quad (2.6)$$

where

$$\eta_i = \max\left\{\frac{\ln(\frac{1}{1-\alpha_i})}{w_0}, \sqrt{\frac{h_i \lambda_i}{c}}\right\}. \quad (2.7)$$

Substituting μ_i^* in (2.1), we obtain the optimal expected cost for firm i as

$$z_i^* = c(\lambda_i + \eta_i) + \frac{h_i \lambda_i}{\eta_i}. \quad (2.8)$$

This leads to a total system cost of $z_{1,\dots,n}^* = \sum_{i \in \mathcal{N}} z_i^*$. In systems where $\sqrt{\frac{h_i \lambda_i}{c}} \geq \frac{\ln(\frac{1}{1-\alpha_i})}{w_0}$ for all $i \in \mathcal{N}$, the optimal cost simplifies to

$$z_i^* = c\lambda_i + 2\sqrt{h_i \lambda_i c}. \quad (2.9)$$

This leads to a total system cost, $z_{1,\dots,n}^*$, given by

$$z_{1,\dots,n}^* = c \sum_{i \in \mathcal{N}} \lambda_i + 2 \sum_{i \in \mathcal{N}} \sqrt{h_i \lambda_i c}. \quad (2.10)$$

In the case of identical firms, with $\lambda_i = \lambda$ and $h_i = h$ for all $i \in \mathcal{N}$, the optimal total cost in (2.10) reduces to

$$z_{1,\dots,n}^* = cn\lambda + 2n\sqrt{h\lambda c}, \quad (2.11)$$

and the total capacity in the system to

$$\sum_{i \in \mathcal{N}} \mu_i^* = n \left(\lambda + \sqrt{\frac{h\lambda}{c}} \right). \quad (2.12)$$

As we can see, both the optimal cost and the optimal buffer capacity in the system increase linearly in the number of firms n . Similar observations can be made for systems in which $\sqrt{\frac{h_i \lambda_i}{c}} \leq \frac{\ln(\frac{1}{1-\alpha_i})}{w_0}$. That is, in this case too, both the optimal cost and the optimal buffer capacity in the system increase linearly in n when the firms have identical cost, service level, and demand parameters.

2.4 Capacity Sharing with Full Information

In this section, we consider the scenario where the firms decide to form a coalition and invest in a single shared facility (a joint venture) from which the demand of all the firms

will then be satisfied. We assume that the rules governing the joint venture (as negotiated by members of the coalition) require that the choice of capacity, in the form of a service rate, for the shared facility takes into account the demand levels of each member of the coalition, their delay costs, and their service level requirements. In particular, we assume that the service rate is chosen by the managers of the joint venture so that it minimizes the total cost for the coalition (the sum of expected delay costs experienced by customers of all the firms and the cost of capacity) and satisfies all service level constraints. We assume that all members of the coalition are truthful in their reporting of their demand rates, delay costs, and service levels. In Section 2.5, we consider the case where firms act strategically and may misreport some of this information. There may of course be firms who are unwilling to share any information, truthfully or not. These firms will not be allowed to participate in the coalition. We assume throughout that, although independent, the firms are not competitors so that their demands are exogenously determined and are not affected by decisions made by any of the firms.

The assumption of full information applies to settings where the information is public and can be independently verified by all the firms. For example, delay penalties and service level guarantees could be publicly advertised by the firms themselves as part of their marketing strategy. In some cases, delay penalties and service levels may also adhere to well-known industry standards. In settings where delay costs are directly incurred by the shared facility (e.g., the shared facility is responsible for handling delay penalty payments to the customers), firms would also need to provide the pooled facility with the correct delay costs. Similarly, service levels must be known to the shared facility if contractual agreements with the customers regarding service levels are handled directly by the shared facility. Demand rates are in most cases verifiable since demand would eventually be satisfied from the shared facility. Firms can be induced to disclose their true demand rates by imposing

high penalties if the originally reported rates are higher than the realized rates (measured over a sufficiently long period of time) once the facility is in operation. The assumption of full information is of course applicable to the case where the firms are all subsidiaries of a single large firm.

We refer to the service rate in the shared facility from which the demand of all firms is satisfied as $\mu_{\mathcal{N}}$ (from heretofore, we shall index parameters associated with a set of firms with the name of that set while parameters associated with individual firms with the name of the firm). Because the superposition of independent Poisson processes is also a Poisson process, the demand process at the shared facility is Poisson with rate $\sum_{i \in \mathcal{N}} \lambda_i$. Similarly, because the work content for each customer regardless of its firm is exponentially distributed, the processing time at the shared facility is a random variable $X_{\mathcal{N}} = Y/\mu_{\mathcal{N}}$ with the exponential distribution and mean $1/\mu_{\mathcal{N}}$. We assume that customers regardless of their firm affiliation are served in a FCFS fashion. Hence, the system with the shared facility behaves again as an M/M/1 queue.

2.4.1 Capacity Optimization

We assume that the terms of the joint venture between the participating firms in the coalition require that the shared facility invests in capacity so as to minimize the total cost to the coalition while satisfying the service level constraint of each firm. The total cost to the coalition consists of the sum of capacity cost and expected delay cost (experienced by customers of all the firms over the long run). Satisfying the service level constraints of all the firm requires satisfying the highest of these service level constraints. If we let $z_{\mathcal{N}}(\mu_{\mathcal{N}})$ denote total system cost and let $W_{\mathcal{N}}$, a random variable, refer to customer delay, then the capacity optimization problem can be stated as follows:

$$\text{Minimize } z_{\mathcal{N}}(\mu_{\mathcal{N}}) = c\mu_{\mathcal{N}} + \frac{\sum_{i \in \mathcal{N}} h_i \lambda_i}{\mu_{\mathcal{N}} - \sum_{i \in \mathcal{N}} \lambda_i} \quad (2.13)$$

subject to

$$P(W_{\mathcal{N}} \leq w_0) = 1 - e^{(\mu_{\mathcal{N}} - \sum_{i \in \mathcal{N}} \lambda_i)w_0} \geq \alpha_{\mathcal{N}}, \quad (2.14)$$

and

$$\sum_{i \in \mathcal{N}} \lambda_i / \mu_{\mathcal{N}} \leq 1, \quad (2.15)$$

where $\alpha_{\mathcal{N}} = \max(\alpha_1, \dots, \alpha_n)$. Then, the optimal capacity is given by

$$\mu_{\mathcal{N}}^* = \sum_{i \in \mathcal{N}} \lambda_i + \eta_{\mathcal{N}}, \quad (2.16)$$

where

$$\eta_{\mathcal{N}} = \max\left\{\frac{\ln\left(\frac{1}{1-\alpha_{\mathcal{N}}}\right)}{w_0}, \sqrt{\frac{\sum_{i \in \mathcal{N}} h_i \lambda_i}{c}}\right\}. \quad (2.17)$$

Similar to the distributed system, the optimal capacity consists of two components. The first corresponds to the total demand rate, while the second to buffer capacity which, in this case, increases in either the sum of the ratios $\frac{h_i \lambda_i}{c}$ or the maximum service level $\alpha_{\mathcal{N}}$.

The following theorem shows that by investing in a shared facility, the firms are able to reduce total cost in the system while investing in less capacity.

Theorem 2.4.1 $z_{\mathcal{N}}^* \leq z_{1, \dots, n}^*$ and $\mu_{\mathcal{N}}^* \leq \sum_{i=1}^n \mu_i^*$, where $z_{\mathcal{N}}^*$ is the optimal cost in the shared facility.

Proof. To prove that $\sum_{i \in \mathcal{N}} \mu_i^* \geq \mu_{\mathcal{N}}^*$, note that

$$\begin{aligned} \sum_{i \in \mathcal{N}} \mu_i^* &= \sum_{i \in \mathcal{N}} \max\left\{\frac{\ln\left(\frac{1}{1-\alpha_i}\right)}{w_0}, \sqrt{\frac{h_i \lambda_i}{c}}\right\} \geq \max\left\{\frac{\ln\left(\frac{1}{1-\alpha_{i_{max}}}\right)}{w_0}, \sqrt{\frac{h_{i_{max}} \lambda_{i_{max}}}{c}}\right\} + \sum_{i \in \mathcal{N}, i \neq i_{max}} \sqrt{\frac{h_i \lambda_i}{c}} \\ &\geq \max\left\{\frac{\ln\left(\frac{1}{1-\alpha_{\mathcal{N}}}\right)}{w_0}, \sqrt{\frac{\sum_{i \in \mathcal{N}} h_i \lambda_i}{c}}\right\} = \mu_{\mathcal{N}}^*, \end{aligned}$$

where $i_{max} \in \{i : \alpha_i = \max(\alpha_1, \dots, \alpha_n)\}$. In order to prove that $z_{\mathcal{N}}^* \leq z_{1, \dots, n}^*$, we distinguish two cases.

(1) $\frac{\ln(\frac{1}{1-\alpha_{\mathcal{N}}})}{w_0} \leq \sqrt{\frac{\sum_{i \in \mathcal{N}} h_i \lambda_i}{c}}$: In this case, we have

$$\begin{aligned} z_{\mathcal{N}}^* &= \frac{\sum_{i \in \mathcal{N}} h_i \lambda_i}{\sqrt{\frac{\sum_{i \in \mathcal{N}} h_i \lambda_i}{c}}} + c \sum_{i \in \mathcal{N}} \lambda_i + c \sqrt{\frac{\sum_{i \in \mathcal{N}} h_i \lambda_i}{c}} \leq \sum_{i \in \mathcal{N}} \left(\frac{h_i \lambda_i}{\sqrt{\frac{h_i \lambda_i}{c}}} + c \lambda_i + c \sqrt{\frac{h_i \lambda_i}{c}} \right) \\ &\leq \sum_{i \in \mathcal{N}} \left(c(\lambda_i + \eta_i) + \frac{h_i \lambda_i}{\eta_i} \right) = z_{1, \dots, n}^*. \end{aligned}$$

The last inequality is due to the fact that $\sqrt{\frac{h_i \lambda_i}{c}}$ is the optimal solution for the optimization problem $\max_{x>0} \{ \frac{h_i \lambda_i}{x} + cx \}$.

(2) $\frac{\ln(\frac{1}{1-\alpha_{\mathcal{N}}})}{w_0} > \sqrt{\frac{\sum_{i \in \mathcal{N}} h_i \lambda_i}{c}}$: In this case, we have

$$\begin{aligned} z_{\mathcal{N}}^* &= \frac{\sum_{i \in \mathcal{N}} h_i \lambda_i}{\frac{\ln(\frac{1}{1-\alpha_{\mathcal{N}}})}{w_0}} + c \sum_{i \in \mathcal{N}} \lambda_i + c \frac{\ln(\frac{1}{1-\alpha_{\mathcal{N}}})}{w_0} \\ &\leq \frac{h_{i_{max}} \lambda_{i_{max}}}{\frac{\ln(\frac{1}{1-\alpha_{i_{max}}})}{w_0}} + c \lambda_{i_{max}} + c \frac{\ln(\frac{1}{1-\alpha_{i_{max}}})}{w_0} + \sum_{i \in \mathcal{N}, i \neq i_{max}} \left(\frac{h_i \lambda_i}{\sqrt{\frac{h_i \lambda_i}{c}}} + c \lambda_i + c \sqrt{\frac{h_i \lambda_i}{c}} \right) \\ &\leq \sum_{i \in \mathcal{N}} \left(c(\lambda_i + \eta_i) + \frac{h_i \lambda_i}{\eta_i} \right) = z_{1, \dots, n}^*. \end{aligned}$$

The last inequality is due to the fact that for firm i_{max} , we have $\frac{h_{i_{max}} \lambda_{i_{max}}}{\eta_{i_{max}}} + c(\lambda_{i_{max}} + \eta_{i_{max}}) = \frac{h_{i_{max}} \lambda_{i_{max}}}{\frac{\ln(\frac{1}{1-\alpha_{i_{max}}})}{w_0}} + c(\lambda_{i_{max}} + \frac{\ln(\frac{1}{1-\alpha_{i_{max}}})}{w_0})$ and $\sqrt{\frac{h_i \lambda_i}{c}}$ is the optimal solution for the optimization problem $\max_{x>0} \{ \frac{h_i \lambda_i}{x} + cx \}$. ■

The potential magnitude of the savings from capacity sharing can be more easily seen in a system with identical firms where $\alpha_i = \alpha$, $h_i = h$, and $\lambda_i = \lambda$ for all $i \in \mathcal{N}$. Consider the case where $\sqrt{\frac{h\lambda}{c}} \geq \frac{\ln(\frac{1}{1-\alpha})}{w_0}$. This leads to $\mu_{\mathcal{N}}^* = n\lambda + \sqrt{\frac{nh\lambda}{c}}$, $z_{\mathcal{N}}^* = cn\lambda + 2\sqrt{cnh\lambda}$, and $E(W_{\mathcal{N}}^*) = \sqrt{\frac{c}{nh\lambda}}$ from which we can observe that both buffer capacity and expected delay, and consequently delay cost, are reduced by a factor of a square root of n (relative to those observed in the case of no capacity sharing). In the case where $\sqrt{\frac{nh\lambda}{c}} \leq \frac{\ln(\frac{1}{1-\alpha})}{w_0}$, we have $\mu_{\mathcal{N}}^* = n\lambda + \frac{\ln(\frac{1}{1-\alpha})}{w_0}$, $z_{\mathcal{N}}^* = c(n\lambda + \frac{\ln(\frac{1}{1-\alpha})}{w_0}) + \frac{nh\lambda w_0}{\ln(\frac{1}{1-\alpha})}$, and $E(W_{\mathcal{N}}^*) = \frac{w_0}{\ln(\frac{1}{1-\alpha})}$. Here, the

magnitude of savings on capacity is even larger with buffer capacity reduced by a factor of n , but expected delay remains unchanged from the case without capacity sharing.

2.4.2 Cost Sharing

We have so far showed that capacity sharing is system-optimal. However, whether or not it is also optimal for individual firms depends on how the cost of the shared facility is allocated among the firms. We assume that each firm incurs its own delay cost and pays a fraction of capacity cost. A firm would prefer the shared facility if the sum of its share of capacity cost and its long run expected delay cost is lower than the cost it would incur without capacity sharing. Moreover, in many settings, the choice is not just between a single facility shared among all firms or facilities operated individually by each firm. There may instead be a range of facility sharing options. For example, a firm may find it more advantageous to share capacity with only a subset of the firms. This could lead firms to form groupings around multiple smaller shared facilities. A single shared facility would be preferred by all firms only if there exists a cost allocation under which the firms are better off than under any other capacity sharing arrangement, including operating individual facilities. Hence, it is desirable that the cost allocation for the shared would be designed so that it deters firms from breaking away and engaging in other facility sharing arrangements.

The problem of determining whether or not there exists a cost allocation scheme under which firms prefer to share a single facility to any other facility sharing configuration can be formulated as a *cooperative game* among the independent firms in the set \mathcal{N} . Consistent with standard terminology from cooperative game theory, let us refer to the subset of firms $\mathcal{J} \subseteq \mathcal{N}$ as *coalition* \mathcal{J} and to the set \mathcal{N} , the largest coalition, as the *grand coalition*. A cooperative game is then defined by a characteristic function which specifies the value associated with each coalition \mathcal{J} . In our context, this corresponds to the total expected

cost associated with a subset of firms \mathcal{J} sharing a single facility. We refer to this cost as $z_{\mathcal{J}}^*$, where $z_{\mathcal{J}}^* \equiv z_{\mathcal{J}}(\mu_{\mathcal{J}}^*)$. A vector $\phi = (\phi_1, \dots, \phi_n)$ is called an allocation rule if ϕ_i corresponds to the portion of total expected cost in the grand coalition that is incurred by firm i . If $\sum_{i=1}^n \phi_i = z_{\mathcal{N}}^*$, then the allocation rule is said to be efficient. An allocation rule is said to be individually rational if $\phi_i \leq z_i^*$ and to be stable for a coalition J if $\sum_{i \in \mathcal{J}} \phi_i \leq z_{\mathcal{J}}^*$. An allocation is said to be a member of the core if it satisfies the following inequalities:

$$\sum_{i \in \mathcal{J}} \phi_i \leq z_{\mathcal{J}}^*, \quad \forall \mathcal{J} \subseteq \mathcal{N}, \quad (2.18)$$

and

$$\sum_{i \in \mathcal{N}} \phi_i = z_{\mathcal{N}}^*. \quad (2.19)$$

When an allocation rule is in the core, no subset of players would want to secede from the grand coalition and form smaller coalitions, including being on their own. Hence the existence of an allocation rule that is in the core (the core is non-empty) is sufficient in our context to show that it is optimal for all the firms to share a single facility. This single facility is a superior arrangement to any other arrangement that may involve a set of partially pooled facilities shared among multiple subsets of the firms.

In addition to the requirement of being in the core, it is desirable for an allocation rule to be perceived as *fair*. In general, a fair allocation is one that assigns a higher portion of total cost to firms whose membership in the coalition contribute more to total cost. In particular, everything else being equal, firms with higher demand rates, higher delay costs, or higher service levels should pay a greater portion of total cost. In what follows, we show that a relatively simple allocation rule has both the properties of being in the core and satisfying the above intuitive notions about fairness (for a more extensive discussion of fairness in cost allocation rules see Moulin 1995).

Consider the following cost allocation rule:

$$\phi_i = \frac{h_i \lambda_i}{\eta_{\mathcal{N}}} + c \lambda_i + \gamma_i, \quad (2.20)$$

where

$$\gamma_i = \frac{h_i \lambda_i}{\sum_{i \in \mathcal{N}} h_i \lambda_i} c \eta_{\mathcal{N}} \quad \text{if} \quad \frac{\ln(\frac{1}{1-\alpha_{\mathcal{N}}})}{w_0} \leq \sqrt{\frac{\sum_{i \in \mathcal{N}} h_i \lambda_i}{c}} \quad (2.21)$$

and, otherwise (if $\frac{\ln(\frac{1}{1-\alpha_{\mathcal{N}}})}{w_0} > \sqrt{\frac{\sum_{i \in \mathcal{N}} h_i \lambda_i}{c}}$),

$$\gamma_i = \begin{cases} c \eta_{\mathcal{N}} - c \sqrt{\frac{\sum_{i \in \mathcal{N}, i \neq i_{max}} h_i \lambda_i}{c}} & \text{if } i = i_{max}, \text{ and} \\ c \frac{h_i \lambda_i}{\sum_{i \in \mathcal{N}, i \neq i_{max}} h_i \lambda_i} \sqrt{\frac{\sum_{i \in \mathcal{N}, i \neq i_{max}} h_i \lambda_i}{c}} & \text{if } i \neq i_{max}, \end{cases} \quad (2.22)$$

with again $i_{max} \in \{i : \alpha_i = \max(\alpha_1, \dots, \alpha_n)\}$. Under the above allocation rule, each firm (1) incurs its own delay cost, $\frac{h_i \lambda_i}{\eta_{\mathcal{N}}}$ and (2) a portion of total capacity cost, $c \lambda_i + \gamma_i$. The portion of total capacity cost has itself two parts: (a) an amount proportional to the firm's demand rate that can be directly attributed to each firm (this amount corresponds to the minimum cost needed to satisfy demand from this firm) and (b) a portion of the cost of buffer capacity. This portion is non-decreasing in the demand rate, delay cost, and service level of each firm. If $\frac{\ln(\frac{1}{1-\alpha_{\mathcal{N}}})}{w_0} \leq \sqrt{\frac{\sum_{i \in \mathcal{N}} h_i \lambda_i}{c}}$, this fraction is proportional to the firms' demand-weighted delay costs. If $\frac{\ln(\frac{1}{1-\alpha_{\mathcal{N}}})}{w_0} > \sqrt{\frac{\sum_{i \in \mathcal{N}} h_i \lambda_i}{c}}$ (the case where the service level constraint is more restrictive), firm i_{max} determines the service level requirement for the entire system. Therefore, it is treated differently to ensure that it is allocated a portion of the cost that is sufficiently high so that other firms do not break away from the coalition. This allocation appears to be consistent with those observed in practice, where combinations of volume based and capacity/service level based fees are common; see for example Gans and Zhou (2003, 2007) and Aksin et al. (2008).

Theorem 2.4.2 *The cost allocation rule $\phi = (\phi_1, \dots, \phi_n)$ as specified in (2.20) – (2.22) is in the core. That is, under this cost allocation, no subset of the firms in \mathcal{N} has an incentive to secede from the grand coalition.*

Proof. We distinguish two cases here.

(1) $\frac{\ln(\frac{1}{1-\alpha_N})}{w_0} \leq \sqrt{\frac{\sum_{i \in \mathcal{N}} h_i \lambda_i}{c}}$: First note that $z_{\mathcal{J}}^* \geq c \sum_{i \in \mathcal{J}} \lambda_i + 2\sqrt{c \sum_{i \in \mathcal{J}} h_i \lambda_i}$. Since $\sum_{i \in \mathcal{J}} \phi_i - \left[c \sum_{i \in \mathcal{J}} \lambda_i + 2\sqrt{c \sum_{i \in \mathcal{J}} h_i \lambda_i} \right] = 2 \left[\sum_{i \in \mathcal{J}} h_i \lambda_i \sqrt{\frac{c}{\sum_{i \in \mathcal{N}} h_i \lambda_i}} - \sqrt{c \sum_{i \in \mathcal{J}} h_i \lambda_i} \right] \leq 0$, we have $z_{\mathcal{J}}^* \geq \sum_{i \in \mathcal{J}} \phi_i, \forall \mathcal{J} \subseteq \mathcal{N}$. It follows that the allocation rule is in the core.

(2) $\frac{\ln(\frac{1}{1-\alpha_N})}{w_0} > \sqrt{\frac{\sum_{i \in \mathcal{N}} h_i \lambda_i}{c}}$: For coalition $\mathcal{J} \subseteq \mathcal{N} \setminus \{i_{max}\}$, we have $z_{\mathcal{J}}^* \geq c \sum_{i \in \mathcal{J}} \lambda_i + 2\sqrt{c \sum_{i \in \mathcal{J}} h_i \lambda_i}$. Since

$$\begin{aligned} c \sum_{i \in \mathcal{J}} \lambda_i + 2\sqrt{c \sum_{i \in \mathcal{J}} h_i \lambda_i} &= \sum_{i \in \mathcal{J}} \left[\frac{h_i \lambda_i}{\sqrt{\sum_{i \in \mathcal{J}} h_i \lambda_i}} + c \lambda_i + c \frac{h_i \lambda_i}{\sum_{i \in \mathcal{J}} h_i \lambda_i} \sqrt{\frac{\sum_{i \in \mathcal{J}} h_i \lambda_i}{c}} \right] \\ &\geq \sum_{i \in \mathcal{J}} \left[\frac{h_i \lambda_i}{\sqrt{\sum_{i \in \mathcal{N}, i \neq i_{max}} h_i \lambda_i}} + c \lambda_i + c \frac{h_i \lambda_i}{\sum_{i \in \mathcal{N}, i \neq i_{max}} h_i \lambda_i} \sqrt{\frac{\sum_{i \in \mathcal{N}, i \neq i_{max}} h_i \lambda_i}{c}} \right] \\ &\geq \sum_{i \in \mathcal{J}} \phi_i, \end{aligned}$$

we have $z_{\mathcal{J}}^* \geq \sum_{i \in \mathcal{J}} \phi_i^*, \forall \mathcal{J} \subseteq \mathcal{N} \setminus \{i_{max}\}$. If $i_{max} \in \mathcal{J}$, then we have $\max_{i \in \mathcal{J}} \alpha_i = \alpha_N$, and hence $z_{\mathcal{J}}^* = \frac{\sum_{i \in \mathcal{J}} h_i \lambda_i}{\frac{\ln(\frac{1}{1-\alpha_N})}{w_0}} + c \sum_{i \in \mathcal{J}} \lambda_i + c \frac{\ln(\frac{1}{1-\alpha_N})}{w_0} \geq \sum_{i \in \mathcal{J}} \phi_i, \forall \mathcal{J} : i_{max} \in \mathcal{J}$. Consequently, the allocation is in the core. \blacksquare

2.5 Capacity Sharing with Incomplete Information

In this section, we consider the case where unit delay costs are private information to each firm. Hence, firms could act strategically and misreport this information if doing so is individually beneficial. In other words, each firm i makes a decision about what unit delay cost \hat{h}_i to report, where \hat{h}_i can be different from the true value h_i . A firm makes this decision knowing that the reported information will be used to determine the corresponding optimal capacity level. As in Section 2.4.2, each firm incurs in the long run two costs: (1) a private expected delay cost and (2) a fraction of the total capacity cost, where the latter

is determined by the cost allocation rule. For tractability, we restrict our treatment to the case of pure delay costs (i.e., no service level constraints). Treating systems where both service levels and delay costs appears difficult in the presence of incomplete information. Moreover, it is arguably more important to focus on the case where unit delay costs, and not service level requirements, are private information. The misreporting (at least under-reporting) of service level requirements is less plausible since the only guarantee a firm has that its service level would be fulfilled is to truthfully report it.

We assume the following sequence of events. First, the cost allocation rule, which may depend on the reported information, is announced. Second, the firms report unit delay costs taking into account the announced allocation rule. Third, based on the reported information, the optimal capacity is selected. Finally, firms incur cost based on the realized delay and their share of the capacity cost. Hence, for a given cost allocation rule, the problem faced by the firms can be viewed as a noncooperative information reporting game where the strategy set for each firm consists of the reported values of unit delay costs and service levels. The objective of each firm is to report values that minimize its total expected cost given the reported values by the other firms.³

The presence of private information raises several important questions. Would the cost allocation rule described in Section 2.4.2 lead to misreporting of private information? If so, is it possible to design an alternative cost allocation rule that does lead to truth telling? Would such a cost allocation be in the core and would it preserve desirable fairness properties? In this section, we provide answers to these and other related questions.

We assume that, given reported unit delay costs $(\hat{h}_1, \dots, \hat{h}_n)$, total capacity is determined

³To formulate a game with incomplete information, as for example in a Nash-Bayes game, one would typically require additional assumptions regarding what each firm might know about the private information of other firms. However, as we shall see in Theorem 2.5.1, we do not need to specify such assumptions, as the described cost allocation rule leads to truth-telling being a dominant strategy.

so as to minimize the sum of expected delay costs (based on the reported information) and investment capacity cost. The corresponding capacity selection problem can be stated as follows:

$$\text{Minimize } \hat{z}_{\mathcal{N}}(\mu_{\mathcal{N}}) = c\mu_{\mathcal{N}} + \frac{\sum_{i \in \mathcal{N}} \hat{h}_i \lambda_i}{\mu_{\mathcal{N}} - \sum_{i \in \mathcal{N}} \lambda_i}. \quad (2.23)$$

The optimal capacity, which we denote by $\mu_{\mathcal{N}}^*(\hat{h}_1, \dots, \hat{h}_n)$, is then given by

$$\mu_{\mathcal{N}}^*(\hat{h}_1, \dots, \hat{h}_n) = \sum_{k=1}^n \lambda_k + \sqrt{\frac{\sum_{k=1}^n \hat{h}_k \lambda_k}{c}}, \quad (2.24)$$

and the resulting expected delay cost experienced by firm i is given by $\frac{h_i \lambda_i}{\mu_{\mathcal{N}}^*(\hat{h}_1, \dots, \hat{h}_n) - \sum_{k=1}^n \lambda_k}$.

Let us first consider the cost allocation rule discussed in Section 2.4.2 where firm i , for all $i \in \mathcal{N}$, incurs privately its delay cost, the cost of capacity for which it is directly responsible, and a fraction of buffer capacity cost that is proportional to its demand-weighted unit delay cost. Consequently, given reported unit delay costs $\hat{h}_1, \dots, \hat{h}_n$, the long run expected cost incurred by firm i is

$$\phi_i(\hat{h}_i, \hat{h}_{-i}) = \frac{h_i \lambda_i}{\sqrt{\frac{\sum_{k=1}^n \hat{h}_k \lambda_k}{c}}} + c\lambda_i + \frac{\hat{h}_i \lambda_i}{\sum_{k=1}^n \hat{h}_k \lambda_k} \sqrt{\frac{\sum_{k=1}^n \hat{h}_k \lambda_k}{c}}, \quad (2.25)$$

where $\hat{h}_{-i} = (\hat{h}_1, \dots, \hat{h}_{i-1}, \hat{h}_{i+1}, \dots, \hat{h}_n)$ denotes the set of unit delay costs of firms other than firm i .

Examining the above cost function, we can see that the expected cost of firm i is affected by both its true unit delay cost as well as the one it reports. It is also affected by the unit delay cost reported by other firms. For example, by under-reporting (reporting a lower unit delay cost than its true one), a firm could benefit by incurring a smaller fraction of buffer capacity cost. However, it could also incur a higher delay cost because less capacity could be installed. The extent to which a firm benefits from misreporting depends on the reporting decisions of other firms.

Given that firms $j \neq i$ report unit delay costs \hat{h}_{-i} , firm i would choose to report unit delay cost $\hat{h}_i^*(\hat{h}_{-i})$ that minimizes its total expected cost. Noting that the expected cost

function $\phi_i(\hat{h}_i, \hat{h}_{-i})$ is convex in \hat{h}_i , the optimal reported costs $\hat{h}_i^*(\hat{h}_{-i})$ can be obtained from the first order condition of optimality as

$$\hat{h}_i^*(\hat{h}_{-i}) = \max\{0, h_i - \frac{2 \sum_{j \neq i} \hat{h}_j \lambda_j}{\lambda_i}\}. \quad (2.26)$$

As we can see, firm i would always under-report its true delay cost regardless of the reporting decision of other firms. This is the case, which is perhaps surprising, even if other firms are truthful in their reporting. The under-reporting appears due to the proportionality in how buffer capacity cost is allocated among the firms, making it more advantageous for firms to always under-report and reduce their share of buffer capacity than over-report (or report truthfully) and reduce their delay cost. The under-reporting can be significant, leading firm i in some cases to even report a unit delay cost of zero. Thus, the cost allocation rule of Section 2.4.2 is not incentive compatible. In the remainder of this section, we turn our attention to constructing an allocation rule that is incentive compatible.

Consider the cost allocation rule under which the long run expected cost for firm i is given by the following

$$\begin{aligned} \phi_i(\hat{h}_1, \dots, \hat{h}_n) &= \frac{h_i \lambda_i}{\mu_{\mathcal{N}}^*(\hat{h}_1, \dots, \hat{h}_n) - \sum_{k \in \mathcal{N}} \lambda_k} + \sum_{j \neq i} \frac{\hat{h}_j \lambda_j}{\mu_{\mathcal{N}}^*(\hat{h}_1, \dots, \hat{h}_n) - \sum_{k \in \mathcal{N}} \lambda_k} \\ &\quad + c \mu_{\mathcal{N}}^*(\hat{h}_1, \dots, \hat{h}_n) - p_i(\hat{h}_{-i}), \end{aligned} \quad (2.27)$$

where $p_i(\hat{h}_{-i})$ is a positive function that depends only on the reported unit delay costs of firms $j \neq i$. The first term on the right-hand side of (2.27) corresponds to the private expected delay cost of firm i while the remaining terms constitute its share of capacity cost. Then, we can show that the following result holds.

Theorem 2.5.1 *Under the allocation rule defined in (2.27),*

$$\phi_i(h_i, \hat{h}_{-i}) \leq \phi_i(\hat{h}_i, \hat{h}_{-i}), \forall \hat{h}_{-i}.$$

That is, the allocation rule is incentive compatible with truth telling being a dominant strategy for each firm.

Proof. First note that the choice $p_i(\hat{h}_{-i})$ does not affect firm i 's choice of \hat{h}_i since it does not depend on \hat{h}_i . Given other firms' reported information, \hat{h}_{-i} , firm i reports a unit delay cost that minimizes its allocated cost $\phi_i(\hat{h}_i, \hat{h}_{-i})$, which can be rewritten as

$$\phi_i(\hat{h}_i, \hat{h}_{-i}) = \frac{h_i \lambda_i + \sum_{j \neq i} \hat{h}_j \lambda_j}{\mu_{\mathcal{N}}^*(\hat{h}_1, \dots, \hat{h}_n) - \sum_{k=1}^n \lambda_k} + c \mu_{\mathcal{N}}^*(\hat{h}_1, \dots, \hat{h}_n) - p_i(\hat{h}_{-i}) \quad (2.28)$$

Hence, firm i would like capacity to be set as $\sqrt{\frac{h_i \lambda_i + \sum_{j \neq i} \hat{h}_j \lambda_j}{c}} + \sum_{i=1}^n \lambda_i$, as the above capacity is the unique minimizer of its expected total cost. However, this is possible only if firm i reports its true delay cost h_i , regardless of whether other firms report their information truthfully or not. In other words, truth reporting is a dominant-strategy. It is also the unique strategy that minimizes the expected total cost of each firm, by virtue of the fact that the capacity optimization problem in (2.23) admits the unique minimizer given in (2.24). \blacksquare

Although the above cost allocation rule is incentive-compatible, it must also be efficient so that the sum of the allocated costs equals the total actual cost incurred by the system. This means that we must have $\sum_{i \in \mathcal{N}} \phi_i(\hat{h}_i, \hat{h}_{-i}) = z_{\mathcal{N}}^*(\mu_{\mathcal{N}}(\hat{h}_1, \dots, \hat{h}_n))$, where

$$z_{\mathcal{N}}^*(\mu_{\mathcal{N}}(\hat{h}_1, \dots, \hat{h}_n)) = \frac{\sum_{k=1}^n h_k \lambda_k}{\mu_{\mathcal{N}}^*(\hat{h}_1, \dots, \hat{h}_n) - \sum_{k=1}^n \lambda_k} + c \mu_{\mathcal{N}}^*(\hat{h}_1, \dots, \hat{h}_n), \quad (2.29)$$

which corresponds to the actual total cost in the system. Noting that

$$\sum_{i \in \mathcal{N}} \phi_i(\hat{h}_i, \hat{h}_{-i}) = z_{\mathcal{N}}^*(\mu_{\mathcal{N}}(\hat{h}_1, \dots, \hat{h}_n)) + (n-1) \hat{z}(\mu_{\mathcal{N}}(\hat{h}_1, \dots, \hat{h}_n)) - \sum_{i \in \mathcal{N}} p_i(\hat{h}_{-i}),$$

where

$$\hat{z}(\mu_{\mathcal{N}}(\hat{h}_1, \dots, \hat{h}_n)) = \frac{\sum_{k=1}^n \hat{h}_k \lambda_k}{\mu_{\mathcal{N}}^*(\hat{h}_1, \dots, \hat{h}_{-i}) - \sum_{k=1}^n \lambda_k} + c \mu_{\mathcal{N}}^*(\hat{h}_1, \dots, \hat{h}_n),$$

we can see that the functions $p_i(\hat{h}_{-i})$ must satisfy the equality

$$\sum_{i \in \mathcal{N}} p_i(\hat{h}_{-i}) = (n-1) \hat{z}_{\mathcal{N}}^*(\mu_{\mathcal{N}}^*(\hat{h}_1, \dots, \hat{h}_n)). \quad (2.30)$$

In what follows, we describe how the functions p_i can be constructed to satisfy the above condition. First, note that $\hat{z}_{\mathcal{N}}^*(\mu_{\mathcal{N}}^*(\hat{h}_1, \dots, \hat{h}_n)) = 2\sqrt{c \sum_{k=1}^n \hat{h}_k \lambda_k} + c \sum_{k=n}^n \lambda_k$. Next, suppose that $h_i \geq a$ and $\lambda_i \geq b$ for some i and some positive a and b (we rule out the trivial cases where $h_i = 0$ and $\lambda_i = 0$ for all $i \in \mathcal{N}$). This would then ensure the existence of a Taylor expansion for $\sqrt{c \sum_{k=1}^n \hat{h}_k \lambda_k}$. In particular, we have

$$\sqrt{c \sum_{k=1}^n \hat{h}_k \lambda_k} = \sqrt{cab} + \sum_{m=1}^{\infty} \frac{f^{(m)}(cab)}{m!} (c \sum_{k=1}^n \hat{h}_k \lambda_k - cab)^m,$$

where $f(x) = \sqrt{x}$ and $f^{(m)}(x)$ is the m -th derivative of f evaluated at x . Define now the function $p_i(\hat{h}_i)$ as

$$\begin{aligned} p_i(\hat{h}_{-i}) &= 2 \sum_{m=1}^{\infty} \frac{f^{(m)}(cab)}{m!} (c^m \sum_{j, k \neq i} \sum_{l=0}^m \binom{m}{l} (\hat{h}_k \lambda_k)^l (\hat{h}_j \lambda_j)^{m-l} \\ &+ c^m \sum_{k \neq i} \sum_{l=1}^{m-1} \binom{m}{l} (\hat{h}_k \lambda_k)^l (-ab)^{m-l} + (-cab)^m) + 2\sqrt{cab} + \sum_{j \neq i} c \lambda_j. \end{aligned} \quad (2.31)$$

Define $b_i = (c^m \sum_{j, k \neq i} \sum_{l=0}^{m-1} \binom{m}{l} (\hat{h}_k \lambda_k)^l (\hat{h}_j \lambda_j)^{m-l} + c^m \sum_{k \neq i} \sum_{l=1}^{m-1} \binom{m}{l} (\hat{h}_k \lambda_k)^l (-ab)^{m-l} + (-cab)^m)$, which is the expansion of $(c \sum_{i=1}^n \hat{h}_i \lambda_i - cab)^m$ by excluding the terms containing \hat{h}_i . It is easy to verify that $\sum_{i=1}^n b_i = (n-1)(c \sum_{i=1}^n \hat{h}_i \lambda_i - cab)^m + (-cab)^m$. Therefore

$$\begin{aligned} \sum_{i=1}^n p_i(\hat{h}_{-i}) &= 2 \sum_{i=1}^n \sum_{m=1}^{\infty} \frac{f^{(m)}(cab)}{m!} b_i + 2n\sqrt{cab} + \sum_{i=1}^n \sum_{j \neq i} c \lambda_j \\ &= 2(n-1) \sqrt{\sum_{i=1}^n \hat{h}_i \lambda_i} + (n-1) \sum_{i=1}^n \lambda_i + 2 \sum_{m=1}^{\infty} \frac{f^{(m)}(cab)}{m!} (-cab)^m + 2\sqrt{cab} \end{aligned}$$

But $2 \sum_{m=1}^{\infty} \frac{f^{(m)}(cab)}{m!} (-cab)^m + 2\sqrt{cab} = 0$. It follows that $\sum_{i=1}^n p_i(\hat{h}_{-i}) = (n-1) \hat{z}_{\mathcal{N}}^*(\hat{h}_1, \dots, \hat{h}_n)$.

Hence, the cost allocation function specified by (2.27) and (2.31) is efficient. Note that under this cost allocation, each firm continues to incur its own delay cost, the cost of capacity

for which it is directly responsible, and a fraction of buffer capacity cost that is decreasing in its own unit delay cost and demand rate.

What now remains to show is that if the above cost allocation rule is always used to allocate cost in any coalition, then firms would prefer the grand coalition to any other coalition. First note that if the cost allocation rule, as specified in (2.27) and (2.31), is applied to every coalition \mathcal{J} , then a firm $i \in \mathcal{J}$ would report truthfully its unit delay cost and incurs expected cost

$$\phi_i(h_i, h_{-i}|\mathcal{J}) = \frac{\sum_{j \in \mathcal{J}} h_j \lambda_j}{\mu_{\mathcal{J}}^* - \sum_{k=1}^n \lambda_k} + c\mu_{\mathcal{J}}^* - p_i(h_{-i}|\mathcal{J}), \quad (2.32)$$

where $p_i(h_{-i}|\mathcal{J})$ is defined similarly as $p_i(h_{-i})$ for the set of firms in \mathcal{J} .

Theorem 2.5.2 *The cost allocation rule specified in (2.32) is in the core. That is, $\phi_i(h_i, h_{-i}|\mathcal{N}) < \phi_i(h_i, h_{-i}|\mathcal{J})$ for any subset \mathcal{J} of \mathcal{N} .*

Proof. Noting that

$$\begin{aligned} \phi_i(h_i, h_{-i}|\mathcal{N}) &= 2 \sum_{m=1}^{\infty} \frac{f^{(m)}(cab)}{m!} (c^m \sum_{j=1, j \neq i}^n \sum_{l=1}^{m-1} \binom{m}{l} (h_j \lambda_j)^{m-l} (h_i \lambda_i)^l + c^m \sum_{l=1}^{m-1} \binom{m}{l} (-ab)^{m-l} (h_i \lambda_i)^l \\ &\quad + c^m (h_i \lambda_i)^m) + c\lambda_i, \end{aligned}$$

and

$$\begin{aligned} \phi_i(h_i, h_{-i}|\mathcal{J}) &= 2 \sum_{m=1}^{\infty} \frac{f^{(m)}(cab)}{m!} (c^m \sum_{j \in \mathcal{J}, j \neq i} \sum_{l=1}^{m-1} \binom{m}{l} (h_j \lambda_j)^{m-l} (h_i \lambda_i)^l + c^m \sum_{l=1}^{m-1} \binom{m}{l} (-ab)^{m-l} (h_i \lambda_i)^l \\ &\quad + c^m (h_i \lambda_i)^m) + c\lambda_i, \end{aligned}$$

we have

$$\begin{aligned} \phi_i(h_i, h_{-i}|\mathcal{N}) - \phi_i(h_i, h_{-i}|\mathcal{J}) &= 2 \sum_{m=1}^{\infty} \frac{f^{(m)}(cab)}{m!} \left(c^m \sum_{j \notin \mathcal{J}} \sum_{l=1}^{m-1} \binom{m}{l} (h_j \lambda_j)^{m-l} (h_i \lambda_i)^l \right) \\ &= 2 \sum_{j \notin \mathcal{J}} \left(\sqrt{cab + c(h_i \lambda_i + h_j \lambda_j)} - \sqrt{cab + ch_i \lambda_i} - \sqrt{cab + ch_j \lambda_j} \right). \end{aligned}$$

As $cab \rightarrow 0$, we have

$$\begin{aligned} & 2 \sum_{j \notin \mathcal{J}} \left(\sqrt{cab + c(h_i \lambda_i + h_j \lambda_j)} - \sqrt{cab + ch_i \lambda_i} - \sqrt{cab + ch_j \lambda_j} \right) \\ \rightarrow & 2 \sum_{j \notin \mathcal{J}} \left(\sqrt{c(h_i \lambda_i + h_j \lambda_j)} - \sqrt{ch_i \lambda_i} - \sqrt{ch_j \lambda_j} \right) < 0. \end{aligned}$$

Since we can make cab as small as we want, consequently, the cost allocation scheme defined by ϕ_i is in the core. ■

The above results are rather remarkable. Not only is it possible to design a cost allocation rule to induce all firms to reveal their true unit delay costs, and for this to be a dominant strategy, but this allocation scheme also ensures that all firms prefer the grand coalition. Furthermore, the cost allocation ensures that each firm incurs its own delay cost and the cost of capacity for which it is directly responsible, with the allocated cost being increasing in each firm's unit delay cost and demand rate. That is, the desirable fairness properties observed in the allocation of Section 2.4.2 continue to hold.

Intuitively, the above cost allocation scheme is truth revealing because each firm, in addition to incurring its own delay cost, incurs a fraction of the delay costs of all other firms. This, coupled with the fact that capacity investment decisions are optimized based on reported information, leads firms to prefer truth telling. In this sense, the cost allocation rule can be viewed as an example of a Groves mechanism from the theory of mechanism design; see for example, Groves (1973, 1976) and Groves and Loeb (1979). In fact, just like a Groves mechanism, the cost allocation rule has broad applicability for any queueing system where the expected delay is well defined and for which the p_i functions can be specified; see the Appendix for details.

2.6 Extensions to Systems with General Demand and Processing Times

In this section, we briefly discuss systems where the customer inter-arrival times and processing times are not necessarily exponentially distributed. Our objective here is not to provide a comprehensive analysis, which is outside the scope of this paper, but rather to offer preliminary insights into the impact of relaxing assumptions made so far and the extent to which results we obtained under these assumptions would continue to hold. Exact analysis for general systems is difficult. Therefore, to obtain these preliminary insights, we rely throughout on approximations that have been extensively used in the literature. For simplicity, we also restrict our treatment to the case of pure delay costs, although the analysis can be extended to systems with service level constraints.

We consider systems where customer inter-arrival times for each firm $i \in \mathcal{N}$ are independent and identically distributed (i.e., arrivals form a renewal process) with mean $1/\lambda_i$ and coefficient of variation c_{a_i} . Customer processing times are independent, identically distributed, and described by a random variable of the form Y/μ , where Y has a mean equal to one and coefficient of variation c_s . The parameter μ is again a scaling factor that corresponds to the service rate. In systems without capacity sharing, each independent facility can thus be modeled as a $GI/G/1$ queue. To obtain an explicit expression for the expected delay cost, we rely on an approximation that is asymptotically correct when the demand rates are high (i.e., when $\lambda_i \rightarrow \infty$). In particular, we approximate the expected number of customers at firm i , given capacity level μ_i , as follows

$$E[Q_i(\mu_i)] \approx \sigma_i^2 \frac{\lambda_i}{\mu - \lambda_i},$$

where $\sigma_i = \sqrt{\frac{c_{a_i}^2 + c_s^2}{2}}$. Motivation and supporting arguments for this approximation can be found in Harrison (1985) and more recently in Bassombo et al. (2008) and the references

therein. The problem faced by each firm can then be restated as

$$\text{Minimize } z_i(\mu_i) \approx \frac{h_i \lambda_i}{\mu_i - \lambda_i} \frac{c_{a_i}^2 + c_s^2}{2} + c\mu_i.$$

This leads to an optimal capacity given by $\mu_i^* = \lambda_i + \sigma_i \sqrt{\frac{h_i \lambda_i}{c}}$, and corresponding optimal cost

$$z_i^* = c\lambda_i + 2\sigma_i \sqrt{ch_i \lambda_i}. \quad (2.33)$$

Bassombo et al. (2008) show that this capacity is asymptotically optimal when the demand rate is high (i.e., $\lambda_i \rightarrow \infty$). It also reduces to the optimal capacity for the $M/M/1$ case (in that case, $\sigma_i = 1$). Note that the above expressions capture now explicitly the effect of both demand and processing time variability.

For systems with capacity sharing, the analysis is more complicated since the superposition of renewal processes is not necessarily a renewal process. To handle this difficulty, we approximate superposed renewal processes by a renewal process whose coefficient of variation is obtained via a two-moment approximation, see Albin (1984) and Whitt (1982). In particular, we approximate the arrival process to a facility shared by the N firms by a renewal process with rate $\sum_{i \in \mathcal{N}} \lambda_i$ and coefficient of variation $c_{a_N}^2 = \sum_{i \in \mathcal{N}} \frac{\lambda_i c_{a_i}^2}{\sum_{i=1}^n \lambda_i}$. For the case of full information, the capacity optimization problem can be stated as follows:

$$\text{Minimize } z_{\mathcal{N}}(\mu_{\mathcal{N}}) \approx c\mu_{\mathcal{N}} + \sigma_{\mathcal{N}}^2 \frac{\sum_{i \in \mathcal{N}} h_i \lambda_i}{\mu_{\mathcal{N}} - \sum_{i \in \mathcal{N}} \lambda_i},$$

where $\sigma_{\mathcal{N}} = \sqrt{\frac{c_{a_N}^2 + c_s^2}{2}}$. Hence, the optimal capacity is given by $\mu_{\mathcal{N}}^* = \sum_{i \in \mathcal{N}} \lambda_i + \sigma_{\mathcal{N}} \sqrt{\frac{\sum_{i \in \mathcal{N}} h_i \lambda_i}{c}}$.

and the optimal cost by

$$z_{\mathcal{N}}^* = c \sum_{i \in \mathcal{N}} \lambda_i + 2\sigma_{\mathcal{N}} \sqrt{c \sum_{i \in \mathcal{N}} h_i \lambda_i}. \quad (2.34)$$

Observation 2.6.1 *Capacity sharing can lead to higher total cost in the system. That is, it is possible to have $z_{\mathcal{N}}^* > \sum_{i \in \mathcal{N}} z_i^*$.*

Proof. Comparing the optimal costs in (2.33) and (2.34), we can see that $z_{\mathcal{N}}^* \leq \sum_{i \in \mathcal{N}} z_i^*$ does not always holds. To see that, let $\lambda = \lambda_i$. Then, in order to have $z_{\mathcal{N}}^* \leq \sum_{i \in \mathcal{N}} z_i^*$ we must have

$$\frac{1}{n} \sum_{i=1}^n \sigma_i^2 \sum_{i=1}^n h_i \leq \left(\sum_{i=1}^n \sqrt{h_i \sigma_i^2} \right)^2.$$

But if $n = 2$, $\lambda_1 = \lambda_2$, $\sigma_1 = 1$, $\sigma_2 = 0$, $h_1 = 0$, $h_2 = 1$, we have

$$\frac{1}{n} \sum_{i=1}^n \sigma_i^2 \sum_{i=1}^n h_i = \frac{1}{2} > 0 = \left(\sum_{i=1}^n \sqrt{h_i \sigma_i^2} \right)^2,$$

which is a counterexample. Note that in this counterexample, firms, when they are on their own, do not need any buffer capacity as either the unit delay cost is zero or variability is zero. When the firms share the same facility, the overall variability is positive and, therefore, there is congestion which leads to delay costs being incurred by customers of firm 2. ■

Although capacity sharing is not always beneficial, it is still possible to identify plausible ranges of parameter values for which capacity sharing is beneficial. The following result describes such a setting.

Theorem 2.6.2 *Capacity sharing is beneficial if, for each pair of firms i and j , $h_i \geq h_j$ and only if $\sigma_i \geq \sigma_j$. In other words, capacity sharing is beneficial if firms with higher delay costs also have higher demand variability.*

Proof. To show that $z_{\mathcal{N}}^* \leq \sum_{i=1}^n z_i^*$, it suffices to show that

$$\sum_{i=1}^n \sqrt{c h_i \lambda_i \sigma_i^2} \geq \sqrt{\frac{\sum_{i=1}^n \lambda_i \sigma_i^2}{\sum_{i=1}^n \lambda_i}} \sqrt{c \sum_{i=1}^n h_i \lambda_i}.$$

Taking the square of both sides of the inequality, it is enough to show that

$$\sum_{i=1}^n h_i \lambda_i \sigma_i^2 + 2 \sum_{j \neq i} \sqrt{h_i \lambda_i \sigma_i^2 h_j \lambda_j \sigma_j^2} \geq \frac{\sum_{i=1}^n \lambda_i \sigma_i^2 \sum_{i=1}^n h_i \lambda_i}{\sum_{i=1}^n \lambda_i}.$$

But to prove $\sum_{i=1}^n \lambda_i \left(\sum_{i=1}^n h_i \lambda_i \sigma_i^2 + 2 \sum_{j \neq i} \sqrt{h_i \lambda_i \sigma_i^2 h_j \lambda_j \sigma_j^2} \right) \geq \sum_{i=1}^n \lambda_i \sigma_i^2 \sum_{i=1}^n h_i \lambda_i$, it is sufficient to show that $\sum_{i=1}^n \lambda_i \sum_{i=1}^n h_i \lambda_i \sigma_i^2 \geq \sum_{i=1}^n \lambda_i \sigma_i^2 \sum_{i=1}^n h_i \lambda_i$. Note that

$$\sum_{i=1}^n \lambda_i \sum_{i=1}^n h_i \lambda_i \sigma_i^2 = \sum_{i=1}^n h_i \lambda_i^2 \sigma_i^2 + 2 \sum_{j \neq i} h_i \lambda_i \lambda_j \sigma_i^2$$

and

$$\sum_{i=1}^n \lambda_i \sigma_i^2 \sum_{i=1}^n h_i \lambda_i = \sum_{i=1}^n h_i \lambda_i^2 \sigma_i^2 + 2 \sum_{j \neq i} h_i \lambda_i \lambda_j \sigma_j^2.$$

Now let $a(i, j) = h_i \lambda_i \lambda_j \sigma_i^2$, $a(j, i) = h_j \lambda_j \lambda_i \sigma_j^2$, $b(i, j) = h_i \lambda_i \lambda_j \sigma_j^2$ and $b(j, i) = h_j \lambda_j \lambda_i \sigma_i^2$.

Then,

$$a(i, j) + a(j, i) - b(i, j) - b(j, i) = (h_i - h_j) \lambda_i \lambda_j (\sigma_i^2 - \sigma_j^2).$$

Hence if we have $(h_i - h_j)(\sigma_i^2 - \sigma_j^2) \geq 0$, we must also have

$$\sum_{j \neq i} h_i \lambda_i \lambda_j \sigma_i^2 \geq \sum_{j \neq i} h_i \lambda_i \lambda_j \sigma_j^2 \quad \text{and} \quad \sum_{i=1}^n \lambda_i \sum_{i=1}^n h_i \lambda_i \sigma_i^2 \geq \sum_{i=1}^n \lambda_i \sigma_i^2 \sum_{i=1}^n h_i \lambda_i,$$

which completes the proof. ■

Interestingly, the condition in Theorem 2.6.2 is independent of the firms' demand rates. In particular, if $\sigma_i^2 = \sigma_j^2$ for all i and j , then capacity sharing is always beneficial. The case of exponential inter-arrival times and processing of course satisfies this condition, but it is also satisfied by a broader class of problems. The condition in Theorem 2.6.2 is also satisfied if firms have the same unit delay cost with $h_i = h_j$ for all i and j or if the ratio h_i/σ_i is the same for all $i \in \mathcal{N}$. These results make intuitive sense, firms that have high demand variability but low delay costs (or high delay cost but low demand variability) could get away with little buffer capacity. However, this ceases to be the case if firms with high demand variability but low delay costs share the same facility with firms with high delay cost but low demand variability. The condition in Theorem 2.6.2 can be viewed as a requirement that firms be relatively "alike" in the sense that the magnitude of their unit delay costs are consistent with the magnitude of their variability parameters.

In systems where unit delay costs are private information, it is always possible to design a cost allocation rule under which all firms that decide to share a single facility would truthfully report their private information. In particular, we can show that an allocation rule of the same form as the one we considered in Section 2.5.2 is incentive compatible with truth telling being a dominant strategy for each firm.

Theorem 2.6.3 *Consider the cost allocation rule under which the expected cost of firm i is given by*

$$\begin{aligned} \phi_i(\hat{h}_1, \dots, \hat{h}_n) = & \frac{\sigma_{\mathcal{N}}^2 h_i \lambda_i}{\mu_{\mathcal{N}}^*(\hat{h}_1, \dots, \hat{h}_n) - \sum_{k \in \mathcal{N}} \lambda_k} + \sum_{j \neq i} \frac{\sigma_{\mathcal{N}}^2 \hat{h}_j \lambda_j}{\mu_{\mathcal{N}}^*(\hat{h}_1, \dots, \hat{h}_n) - \sum_{k \in \mathcal{N}} \lambda_k} + c \mu_{\mathcal{N}}^*(\hat{h}_1, \dots, \hat{h}_n) \\ & - p_i(\hat{h}_{-i}). \end{aligned} \quad (2.35)$$

where $p_i(\hat{h}_{-i})$ is defined by

$$\begin{aligned} p_i(\hat{h}_{-i}) = & 2\sigma_{\mathcal{N}} \sum_{m=1}^{\infty} \frac{f^{(m)}(cab)}{m!} (c^m \sum_{j, k \neq i} \sum_{l=0}^m \binom{m}{l} (\hat{h}_k \lambda_k)^l (\hat{h}_j \lambda_j)^{m-l} \\ & + c^m \sum_{k \neq i} \sum_{l=1}^{m-1} \binom{m}{l} (\hat{h}_k \lambda_k)^l (-ab)^{m-l} + (-cab)^m) + 2\sigma_{\mathcal{N}} \sqrt{cab} + \sum_{j \neq i} c \lambda_j. \end{aligned} \quad (2.36)$$

Then, under this cost allocation, all firms report truthfully their unit delay costs. Moreover, truthful reporting is a dominant strategy.

A proof for the above result is similar to the proof of Theorem 2.5.1 and is therefore omitted for brevity. The functions p_i can be constructed in a similar fashion to the one described in Section 2.5 to obtain an efficient allocation. Furthermore, we can show that under some conditions (e.g., when the ratio $h_i/\sigma_i^2 = h_j/\sigma_j^2$ for all $i, j \in \mathcal{N}$) the resulting cost allocation rule is in the core with all firms preferring the grand coalition to any other capacity sharing arrangements.

We conclude this section by noting that the results of this section highlight the fact that, although capacity sharing among all firms may not be always beneficial, it may be so

among subsets of the firms that satisfy certain conditions. For these firms, the results we obtained in this paper can be used as a basis for determining capacity and for allocating cost.

2.7 Summary and Concluding Comments

In this paper, we presented models to study the benefit of capacity sharing among independent firms. We formulated the capacity sharing problem as a cooperative game. We showed that capacity sharing can lead to significant savings in total system cost. However, we found that whether or not firms choose to share capacity, and with whom, depends on how the associated cost are allocated. It also depends on whether or not firms are truthful about their private information. We showed that it is possible to design a cost allocation rule that induces all firms to report truthfully their private information and for this allocation to be in the core. We showed that there exists settings for which capacity sharing among all the firms may not be beneficial. In such settings, capacity sharing among subsets of the firms with similar characteristics may still be beneficial.

We view the results of this paper as a step toward a better understanding of the issue of cooperation among independent firms via capacity sharing in the presence of congestion. The results identify some of the important mechanisms that might be needed to make capacity sharing desirable and to mitigate the effect of incomplete information. Much more work obviously remains to be done. Please refer to the Chapter 4 of this thesis for possible extensions.

Acknowledgments: We thank seminar participants at Northwestern University, University of Wisconsin, and Lehigh University for useful feedback and comments. We particularly thank Albert Ha, Eran Hanany, Donald Hausch, Moshe Haviv, the Associate Editor, and

two anonymous referees for insightful comments and suggestions on earlier versions of the paper.

Appendix

In this appendix, we show how the incentive compatible allocation rule described in Section 2.4.2 can be extended to settings much more general than the M/M/1 setting of Section 2.5. To illustrate, we consider the case where facilities are modeled as GI/G/1 queues as in Section 2.6 (the applicability of the allocation rule is however significantly broader). Let $E[W_{\mathcal{N}}(\mu_{\mathcal{N}})]$ refer to expected delay when firms in coalition \mathcal{N} share a single facility with capacity level $\mu_{\mathcal{N}}$. Given reported unit delay costs $(\hat{h}_1, \dots, \hat{h}_n)$, let $\mu_{\mathcal{N}}^*(\hat{h}_1, \dots, \hat{h}_n)$ denote the capacity level that minimizes the expected cost

$$\hat{z}_{\mathcal{N}}(\mu_{\mathcal{N}}) = \sum_{i \in \mathcal{N}} \hat{h}_i \lambda_i E[W_{\mathcal{N}}(\mu_{\mathcal{N}})] + c\mu_{\mathcal{N}}.$$

Consider now the cost allocation rule under which the long run expected cost of firm i is given by

$$\begin{aligned} \phi_i(\hat{h}_1, \dots, \hat{h}_n) = & h_i \lambda_i E[W_{\mathcal{N}}(\mu_{\mathcal{N}}^*(\hat{h}_1, \dots, \hat{h}_n))] + \\ & \sum_{j \neq i} \hat{h}_j \lambda_j E[W_{\mathcal{N}}(\mu_{\mathcal{N}}^*(\hat{h}_1, \dots, \hat{h}_n))] + c\mu_{\mathcal{N}}^*(\hat{h}_1, \dots, \hat{h}_n) - p_i(\hat{h}_{-i}), \end{aligned} \quad (2.37)$$

where $p_i(\hat{h}_{-i})$ is a positive function that depends only on the reported unit delay costs of firms $j \neq i$. It is easy to verify that the cost of firm i would be minimized if the capacity level for the shared facility is set equal to $\mu_{\mathcal{N}}^*(\hat{h}_1, \dots, \hat{h}_{i-1}, h_i, \hat{h}_{i+1}, \dots, \hat{h}_n)$. However, this is possible only if firm i reports truthfully its unit delay cost h_i . Hence, the allocation rule specified in (2.37) is incentive-compatible with truth telling being a dominant strategy. The above arguments are similar to those used to construct a Groves mechanism; see for example Groves (1973, 1976) and Groves and Loeb (1979). We briefly summarize the key steps below. In a typical application of a Groves mechanism, the capacity of a public good is determined based on the reported valuation of this public good by the users. Each user's valuation depends on her type, which is private information. Let θ_i denote the true type

of user i and $\hat{\theta}_i$ be the reported type which may be different from the true one (type plays the role of unit delay cost in our application). User i derives private utility $v(x, \theta_i)$ from the public good if the capacity level of the public good is x (private utility plays the same role as expected delay cost in our application). Under a Groves mechanism, the amount of capacity of the public good is chosen based on the reported types so as to maximize $\sum_{i \in \mathcal{N}} v(x, \hat{\theta}_i) - g(x)$, where $v(x, \hat{\theta}_i)$ is the estimated utility user $i \in \mathcal{N}$ would derive from capacity level x based on her reported type and $g(x)$ is the cost of investing in capacity level x . Let $x^*(\hat{\theta}_1, \dots, \hat{\theta}_n)$ denote this capacity level. Consider a mechanism under which each user is charged a fee such that user i 's overall utility is given by

$$v(x^*(\hat{\theta}_1, \dots, \hat{\theta}_n), \theta_i) + \sum_{j \neq i} v(x^*(\hat{\theta}_1, \dots, \hat{\theta}_n), \hat{\theta}_j) - g(x^*(\hat{\theta}_1, \dots, \hat{\theta}_n)) - p_i(\hat{\theta}_{-i}),$$

where $p_i(\hat{\theta}_{-i})$ is a function that depends on the reported information of firms other than firm i . Then, using arguments similar to the ones used for our application, it is easy to verify that such a mechanism would induce users to report truthfully their type and to do so regardless of the reporting decisions of other firms.

Chapter 3

A Customer-Item Decomposition Approach to Stochastic Inventory Systems with Correlation

3.1 Introduction

A common assumption in the inventory literature is that demands and leadtimes are independent across review periods and independent of external events. In practice, this is rarely the case. Demand is often correlated over time. Factors that influence an increase or decrease in demand in one period often persist over the next several periods. In fact, the presence of time-correlation is what makes demand forecasting a meaningful activity. Demand is also often correlated with external events. For example, inflation levels, levels of employment, and interest rates tend to significantly affect consumption and, consequently, demand for most products. Similarly, correlated leadtimes are common in practice. For example, a long leadtime from a supplier in one period often leads to long leadtimes from

that same supplier in the next few periods (if that supplier is backlogged or experiencing production difficulties). Also, external events, such as disruptions to the supply network, often affect the length of leadtimes over multiple periods.

The relatively limited treatment of correlation in the literature appears to be due to the mathematical and computational intractability of the problem. Traditional approaches to the problem, which rely on dynamic programming and use aggregate information about demand and inventory, tend to lead to problems with too many dimensions, making them difficult to analyze and solve. In this paper, we explore an alternative approach based on the *customer-item decomposition method*. Under this approach, demand is viewed as consisting of a stream of individual customers arriving over time, with a batch of such customers arriving in each period. Inventory (both in-stock, on-order, and yet to be ordered) is also viewed as a stream of individual items that are delivered over time, with a batch of such items delivered in each period. Each item is then matched to a customer. This reduces the inventory control problems to making decisions about when to place an order for the item destined to a particular customer. This decision can be made, in some cases, by taking into account only the marginal cost associated with each customer-item pair. In turn, this decomposes the problem into independent sub-problems, one for each single customer-item pair. As we show in this paper, such decomposition can greatly simplify the inventory control problem, lead to computationally-efficient solution algorithms, and provide a more direct approach to characterizing the structure of the optimal policy and the effectiveness of certain heuristic policies.

The origin of the approach appears to date back to Axsäter (1990) who uses a customer-item decomposition in the analysis of a two-stage distribution system that operates under a heuristic policy of one-for-one replenishment with fixed base-stock levels. Axsäter (1993) extends this treatment to systems with batch-ordering. Katircioglu and Atkins (1998)

used a decomposition approach to study a continuous review system with arbitrary order inter-arrival time distributions with increasing failure rates and argue that a time-delayed base-stock policy is optimal. Janakiraman and Muckstadt (2005) use the decomposition method to analyze a serial system. In their case, the system consists of only two stages with independent demands and fixed leadtimes but with capacity constraints. de Albeniz and Lago (2007) study the optimality of myopic policies in the context of a single stage system and provide under some conditions closed-form expressions for the corresponding ordering decisions. Muharremoglu and Tsitsiklis (2008) use customer-item decomposition to study a serial system under periodic-review with Markov-modulated demand.

In this paper, we use the customer-item decomposition method to study stochastic inventory systems under general assumptions regarding demand and leadtimes. In particular, we consider systems where both demands and leadtimes are stochastic and correlated. Correlation can be across periods or with an external process consisting of a vector of external indicators. The correlation models we consider incorporate most models treated in the literature, including Markov-modulated processes, as special cases. In the systems we study, we allow for the cost parameters to be non-stationary and to vary arbitrarily from period to period. We also treat systems with both single and multiple stages, as well as generalizations of these systems, including systems with advance demand information and batch ordering. In each case, we provide a formulation that leads to efficient solution methods, characterize the structure of the optimal policy, provide simple characterizations of myopic policies and identify conditions under which these policies are optimal.

The contributions of this paper are threefold. First, we study problems that are largely intractable, both in terms of computing the optimal solution and in characterizing the structure of the optimal policy, using a traditional dynamic programming approach. Second, we contribute to the development of the customer-item decomposition method. We do so by

introducing the notions of customer arrival and delivery time distributions which simplify how the marginal costs associated with each customer-item pair are computed. We then formulate the inventory control problem for each customer-item pair as an *optimal stopping problem*. This is significant because *submodularity* properties of the cost function that arise from this formulation simplify the method for characterizing the structure of the optimal policy. We show how the optimal stopping formulation also leads in the important case of strictly positive demand to polynomial time algorithms for computing the optimal policy. This is not possible using other decomposition formulations, such as the one in Muharremoglu and Tsitsiklis (2008)(see Section 3.2 for further detail), and of course not possible using the traditional dynamic programming method.

Third, we show how the marginal cost formulation can be used to construct myopic heuristics. These myopic heuristics lead to explicit solutions for item ordering times in the form of critical fractiles (a so-called *news-vendor* solution). We do so for both single and multiple stage systems. For systems with multiple stages, we provide easy to verify conditions under which a myopic newsvendor solution is optimal at each stage. To our knowledge, we are the first to provide such conditions for problems with multiple stages. Shang and Song (2003) do propose similar myopic news-vendor solutions, which they show to be close to the optimal solution for the special case of systems with i.i.d demand, infinite horizon, and under the average cost criterion. However they do not characterize the optimality conditions for such solutions.

Another important contribution of our paper is in showing how the analysis can be extended to even more complex systems, such as those with advance demand information (ADI) and batch ordering. For systems with batch ordering, our results generalize the results in Chen (2000) who considers serial systems with i.i.d demand and fixed leadtimes. In our case, we consider systems with general demand correlation and sequential stochastic

and correlated leadtimes. More significantly, we characterize again, in closed-form, myopic policies for these systems and specify conditions under which the myopic policies are optimal. To our knowledge, these results are also the first of their kind in the literature.

The rest of this paper is organized as follows. In Section 3.2, we provide a brief review of related literature. In Section 3.3, we describe our approach in the context of systems with a single stage. In section 3.4, we extend our analysis to systems with multiple stages. In section 3.5, we discuss one-period look-ahead policies (the myopic heuristics). In Section 3.6, we discuss extensions to systems with advance demand information and systems with batch ordering. In Section 3.7, we provide concluding comments.

3.2 Related Literature

Our paper is related to two streams in the inventory literature, one dealing with demand correlation and the other with leadtime correlation. Existing literature that deals with demand correlation focuses on either characterizing the structure of optimal policies in specific settings or evaluating the performance of heuristics (sub-optimal policies). A sub-stream within this literature considers systems with Markov-modulated demand, where demand is stochastic and affected (modulated) by an exogenous Markov process; see for example, Song and Zipkin (1993, 1996a, 1996b), Chen and Song (2001), Muharremoglu and Tsitsiklis (2008) and the references therein.

Another sub-stream within this literature considers demand correlation across periods (time-correlated demand) for the purpose of modeling inventory systems with forecasting. This includes treatments with simple time-series models such as the order-one autoregressive process AR(1) in Scarf (1959, 1960), Johnson and Thompson (1975), Erkip *et al.* (1990), and Lee *et al.* (1997, 1999), among others, or more general models such as the Martingale model of forecast evolution (MMFE) introduced by Graves *et al.* (1986, 1998) and Heath

and Jackson (1994), and the random walk model of Graves (1999) and Lee and Whang (1998). Dong and Lee (2003) analyze the optimal policy for multi-echelon systems under a Martingale model of forecast evolution and derive a lower bound for the optimal echelon base stock level. Levi *et al.* (2005) consider a general model for demand correlation and then use it to propose a heuristic with certain performance guarantees.

Papers that study heuristics include the early papers by Veinott (1963, 1965a and 1965b) and Ignall and Veinott (1969), and more recent ones by Iida and Zipkin (2006), Dong and Lee (2003), Lu *et al.* (2006), among several others. An important focus of this literature is the evaluation of so-called myopic heuristics where in each period the objective is to minimize the expected cost for that period, ignoring the effect on the cost of future periods. Myopic policies have been shown to be optimal in some cases (See Karlin 1960, Veinott 1965a and 1965b, Johnson and Thompson 1975, Iida and Zipkin 2006, Lu *et al.* 2006, de Albeniz and Largo 2007). Variations on these policies have also been shown to provide useful bounds (Levi *et al.* 2005). Papers that study heuristics in the context of systems with multiple stages include Chen and Zheng (1994), Gallego and Zipkin (1999), Shang and Song (2003), Dong and Lee (2003), and Levi *et al.* (2005). Shang and Song (2003) derive critical-fractile lower bounds and upper bounds on the optimal *echelon* base-stock policies for systems with independent and identically distributed demand. In this paper, we also derive a critical-fractile upper bound but do so for more general systems with time-varying cost parameters and generally correlated demand. Furthermore, we characterize conditions under which the critical-fractile upper-bound is optimal.

The literature on correlated stochastic leadtimes is relatively limited. Song and Zipkin (1996b) studied a single stage system with Markov modulated leadtimes and Muharremoglu and Tsitsiklis (2008) studied a serial system with Markov modulated leadtimes. The optimal policy is shown to be a base-stock policy with the base-stock levels dependent on the state

of the modulating Markov process. Markov modulated leadtimes is a special case of the leadtime correlation we consider in this paper.

Among existing the papers that use the customer-item decomposition approach, ours is most closely related to Muharremoglu and Tsitsiklis (2008). An important difference is that our paper treats systems with more general (than Markov modulated) models of demand and leadtime correlation. This generalization is not trivial and it is not obvious how the approach in Muharremoglu and Tsitsiklis (2008) could be extended to non-Markovian settings¹. Our formulation of the inventory control problem as an optimal stopping problem is also different and leads to a different and more direct method for proving the structure of the optimal policy. Our associated solution algorithm is different and leads to a polynomial time algorithm for systems with strictly positive demand. Muharremoglu and Tsitsiklis do not study heuristics, as we do, and do not consider systems with advance demand information and batch ordering. Finally, our treatment allows for the cost parameters to be time-dependent, even in the case of infinite horizon, while Muharremoglu and Tsitsiklis assume stationary parameters.

¹The optimality proof in Muharremoglu and Tsitsiklis utilizes the notions of *decomposable systems* and *decomposable policies*. In particular they define the state of the system in period t by the vector $x_t = (x_t^1, x_t^2, \dots)$ where x_t^i is associated with subsystem i and $x_t^i = k$ denotes that customer i is the k th future customer to arrive, among those that have not arrived yet prior to period t , given the available information at time t . As stated by Muharremoglu and Tsitsiklis (2008), they assume that the “the various state components x_t^1, x_t^2, \dots evolve independently of each other, the only coupling arising through exogenous processes s_t and w_t .” However, when demand is correlated across time, the realized demand in one period affects the demand distribution in all future periods. Therefore, the evolution of component x_t^i depends on the entire history of all components $x_{t-1}^1, \dots, x_{t-1}^{i-1}, x_{t-2}^1, \dots, x_{t-2}^{i-1}, \dots$. Therefore, a single variable x_t^i is not sufficient to describe the state of subsystem i . Similarly, in the presence of time-coorelation, the notion of decomposable policies (as defined in Muharremoglu and Tsitsiklis) is no longer valid since the optimal policy for one subsystem in a given period depends on the the state of other subsystems in previous periods.

Finally, our paper is also related to de Albeniz and Largo (2007) which focuses on the study of myopic policies using the customer-item decomposition approach. They obtain results similar to ours for systems with single stage. These results appear to have been independently developed from ours, although an earlier version of our paper (Yu and Benjaafar 2006) seems to predate theirs. The focus of our paper is of course significantly broader, as we also study the structure of optimal policies, treating both single and multiple stage systems as well as systems with batch ordering and advance demand information. In addition, we also characterize the computational complexity of the corresponding optimal solution approach.

3.3 Systems with a Single Stage

Consider a single stage inventory problem with multiple periods, stochastic demands, and stochastic sequential lead times. Let $\mathbf{D} = (D_0, \dots, D_T)$ be a sequence of random variables corresponding to demand in each period and $\mathbf{E} = (\mathbf{E}_0, \dots, \mathbf{E}_T)$ be a sequence of external discrete random vectors where \mathbf{E}_t corresponds to the vector of observable indicators at the beginning of period t . Let $H_{t-1} = \{\mathbf{e}_0, \dots, \mathbf{e}_t, d_0, \dots, d_{t-1}\} \in \mathcal{H}_{t-1}$ be the realized values of the indicator variables from period 0 to period t and of demand from period 0 to period $t-1$ where \mathcal{H}_{t-1} is the set of all possible such realizations. Demand in each period t may depend on past demand realizations and the realizations of the indicator variables. Therefore, at the beginning of each period, the distribution of demand of all future periods is updated based on the information provided by H_{t-1} , where H_{t-1} consists of the entire history of the demand process up to time t . Note that \mathbf{E}_t is observable at the beginning of period t and can be used in making decisions in period t . Therefore, we also define \mathcal{H}_{-1} which consists of all possible realizations of the random vector \mathbf{E}_0 .

Inventory is replenished from an outside supplier with a random leadtime L_t , so that an

item ordered in period t is received in period $t + L_t$. We assume that the outside supplier has ample stock. We assume that leadtimes are sequential so that $t + L_t(\omega) < s + L_s(\omega)$ for all $s > t$ and for all sample paths ω for the random variables L_0, \dots, L_T ; in other words, there is no order crossing. We also assume that leadtimes are finite and bounded so that there exists $r_{max} > 0$ such that $L_t(\omega) < r_{max}$ for all ω and all periods t . In each period t , in addition to observing the information set H_{t-1} for demand, we also observe an information set $\hat{H}_{t-1} = \{\hat{\mathbf{e}}_0, \dots, \hat{\mathbf{e}}_t, l_0, \dots, l_{t-1}\} \in \hat{\mathcal{H}}_{t-1}$ for lead time, where $\hat{\mathbf{e}}_t$ is the realized value of the vector of random variables $\hat{\mathbf{E}}_t$ for the set of leadtime indicators in period t and l_t is the realized value of leadtime in period t and $\hat{\mathcal{H}}_{t-1}$ is the set of all possible such realizations. For notational convenience, we also define the sets $\mathbf{L} = (L_0, \dots, L_T)$ and $\hat{\mathbf{E}} = (\hat{\mathbf{E}}_0, \dots, \hat{\mathbf{E}}_T)$. We allow the intersection of $\hat{\mathbf{E}}_t$ and \mathbf{E}_t to be non-empty (the demand and leadtime may depend on the same set of indicators). We let $F_t = H_t \cup \hat{H}_t$, so that \mathcal{F}_t corresponds to the set of all possible realizations of H_t and \hat{H}_t . We make the following additional assumptions:

- The evolution of \mathbf{E}_t and $\hat{\mathbf{E}}_t$ is independent of \mathbf{D} and \mathbf{L} . That is, $P(\mathbf{E}_t = \mathbf{e}_t | H_{t-2}) = P(\mathbf{E}_t = \mathbf{e}_t | \mathbf{e}_0, \dots, \mathbf{e}_{t-1})$ and $P(\hat{\mathbf{E}}_t = \hat{\mathbf{e}}_t | \hat{H}_{t-2}) = P(\hat{\mathbf{E}}_t = \hat{\mathbf{e}}_t | \hat{\mathbf{e}}_0, \dots, \hat{\mathbf{e}}_{t-1})$ for all $H_{t-2}, \hat{H}_{t-2} \in \mathcal{F}_{t-2}$.
- $E[D_{t'} | H_t] < \infty$ for each period t and $t' > t$ under each possible H_t .
- The demand process and leadtimes are independent of inventory policies and inventory status, including the current inventory level or the number of orders outstanding.

Our setting and our model of demand and leadtime correlation generalizes most known models of correlation in the literature.

Demand is satisfied from on-hand inventory, if any is available; otherwise it is backlogged. In each period, the decision maker must decide on how many units to order to minimize the expected discounted cost over the entire planning horizon. There are three types of

cost: (1) an ordering cost c_t per unit ordered in period t , (2) a holding cost h_t per unit of inventory held in period t and (3) a backorder cost b_t for each order that stays backlogged in period t . Note that we allow for the cost parameters, h_t, b_t and c_t , to be time-varying. All costs for period t are incurred at the end of period t . Also, without loss of generality, we assume that $h_T = b_T = 0$.

We assume that there is no speculative motivation for holding inventory or backloging orders. Therefore, we assume that the following conditions on the cost parameters are satisfied.

Assumption 1 : $\beta^t c_t + \beta^{t+r} h_{t+r} > \beta^{t+1} c_{t+1}$, and $\beta^t c_t < \beta^{t+1} c_{t+1} + \beta^{t+r} b_{t+r}$ for $t = 0, 1, \dots, T$ and $r = 0, 1, \dots, r_{max}$, where $\beta \in (0, 1]$ is the discount factor.

This assumption is not necessary for our analysis. It is only used in Theorem 3.3.6 below to show that the optimal base-stock level is non-negative and finite. This avoids anomalous behavior such as ordering an infinite amount in some period or not placing orders despite the presence of backorders (a special case that satisfies this assumption is $h_t = h, c_t = c$ and $b_t = b$ for all t).

The following describes the sequence of events in each period $t, t = 0, \dots, T$: first, an order of size $q_t \geq 0$ is placed with the outside supplier, then possibly a delivery of an outstanding order would be made, finally demand for the current period in the amount of d_t is realized. Holding and backorder costs are incurred at the end of the period. Let IN_t denote the net inventory at the beginning of period t , where $IN_t = I_t - B_t$, I_t is the level of on-hand inventory, and B_t is the number of backorders. Then $I_t = IN_t^+ \equiv \max(0, IN_t)$ and $B_t = IN_t^- \equiv \max(0, -IN_t)$. The holding cost incurred in period t is $h_t IN_{t+1}^+$ and the backorder cost is $b_t IN_{t+1}^-$ (we assume that costs are charged based on the ending net inventory in period t which is the same as the starting net inventory in period $t + 1$).

The above problem could be formulated as a stochastic dynamic program, where in each

period t the decision variable is q_t , and the state of the system is determined by the net inventory level, the set of outstanding orders, and the history of the process up to time t . As mentioned previously, such an approach would lead to a large multi-dimensional dynamic program that is difficult to solve and analyze. In this paper, we take a different approach, the customer-item decomposition approach. We view demand as a stream of individual customers with a batch of such customers arriving in each period and where each customer requests one unit of the product. To differentiate between the customers, we assign at the beginning of the planning horizon a unique index $i, i = 1, 2, \dots$, to each customer in increasing order of the customers' arrival times. This means that customer i refers to the i th customer that will arrive to the system, with ties broken arbitrarily. At any time, the arrival times of future customers are of course not known, and we can define a random variable X_i to describe the arrival time of customer i , $X_i = \inf\{s : D_0 + \dots + D_s \geq i\}$ with $X_i = T$ if no such s exists. However, we can use knowledge of the demand distributions at period t to assign a probability $p_{t,s}^i$ to customer i if the customer has not arrived yet, where $p_{t,s}^i$ denotes the probability evaluated at the beginning of period t that the i th customer will arrive in period s . More specifically, given that k customers have already arrived prior to period t , and given the history H_{t-1} observed so far, we have:

$$p_{t,s}^i = \sum_{n=0}^{i-k-1} P(D_t + \dots + D_{s-1} = n, D_s \geq i - n - k | H_{t-1}), s \geq t \quad (3.1)$$

for $i > k$. If $\sum_{s=t}^T p_{t,s}^i < 1$, we can simply redefine $p_{t,T}^i \equiv 1 - \sum_{s=t}^{T-1} p_{t,s}^i$, since the customer arrivals at period T will not affect system costs. Thus, we have $\sum_{s=t}^T p_{t,s}^i = 1$ (recall that we assume $h_T = b_T = 0$)

Similarly, we view supply as a stream of physical items with a batch of such items delivered in each period. We assign index $j', j' = 1', 2', \dots$, to each item in increasing order of the items' delivery times, with ties again broken arbitrarily. We refer the j th item as item j' . The delivery times of items are obviously determined by when orders are placed with

the outside supplier. In each period the observed information \hat{H}_{t-1} is used to update the distributions of future leadtimes, namely the probability $\hat{p}_{l,l+m}^t$ evaluated at the beginning of period t that an order placed in a future period l will be received in period $l + m$.

At the beginning of a period t , some customers would have already arrived to the system and their orders fulfilled; some customers would have arrived but their orders would have not been fulfilled (they are currently backlogged), while the remaining customers would have not arrived yet. This situation is graphically depicted in Figure 3.1. Similarly, at the beginning of the same period, some items would have been delivered and used to fulfill customer orders, some items would be in transit, some items would be in stock, while the remaining items (potentially an infinite number of them) would have not been ordered yet; see Figure 3.1 for an illustration.

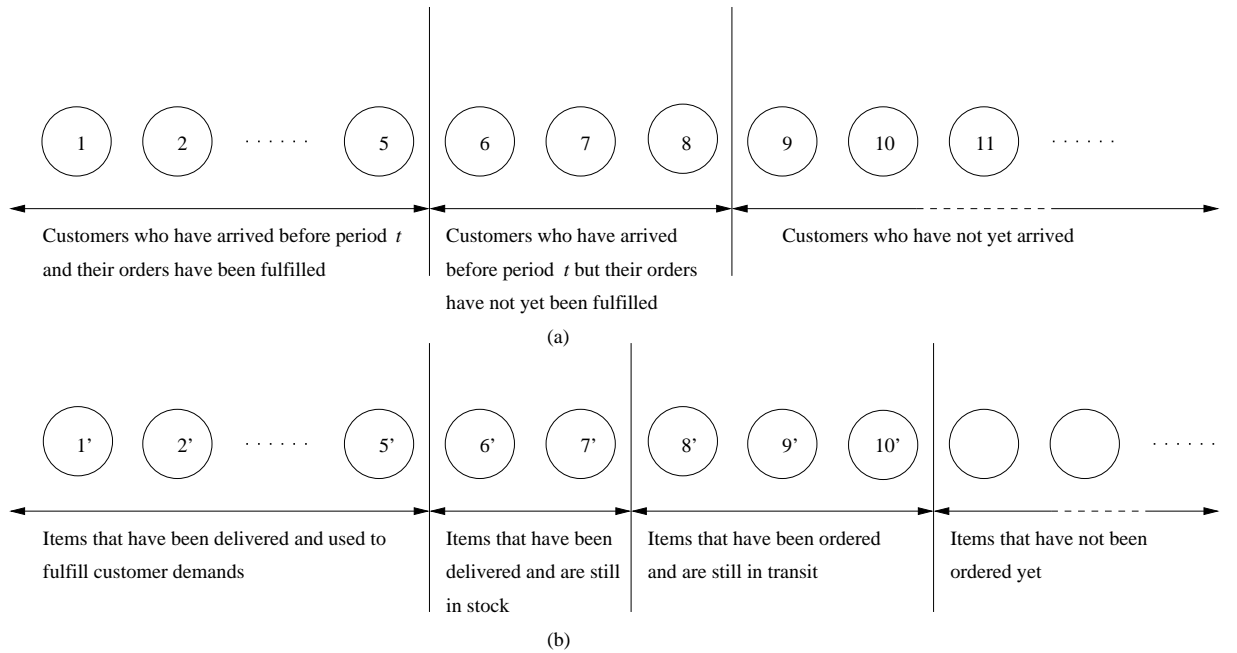


Figure 3.1: An illustration of the demand and supply processes

The inventory control problem can now be restated for each period in terms of (1)

determining which items to order, among those who have not been ordered yet and (2) determining how to allocate available stock to customers who have arrived but whose orders have not been fulfilled yet. More specifically, at the beginning of each period, we first decide on how many additional items to order from the outside supplier (this is equivalent to determining the order quantity), we then observe the delivery of items and the arrival of new customers, and finally we decide on how to allocate items in stock to customers who have arrived but whose orders have not been fulfilled yet. This includes customers who have arrived in the current period as well as any backlogged from previous periods.

Proposition 3.3.1 *If there is inventory on-hand, it is optimal to use this inventory to satisfy customers that are currently backlogged. Moreover, it is optimal to allocate this inventory on a first-come, first-served basis, so that the item with the smallest index is assigned to the customer, among those backlogged, with the smallest index.*

Proof: It is trivial to show that if there are items in stock and there are backlogged customers, then it is optimal to use the items in stock to satisfy as many customers as possible. Since we assume that both customers and items are homogeneous, then any feasible allocation of items to customers that maximizes the number of customers that are satisfied is optimal. This includes the allocation that assigns items with the smallest index to customers with the smallest index. ■

Proposition 3.3.2 *If there is inventory in transit, it is optimal to allocate this inventory once it arrives in stock to backlogged customers who have not been satisfied from on-hand inventory, also on a first-come, first-served basis.*

The proof of Proposition 3.3.2 is similar to that of Proposition 3.3.1. Note that because leadtimes are sequential, items that are ordered first arrive first in stock. Hence, once an item is ordered it can be immediately allocated to the customer with the earliest arrival

time among those that have not been assigned items yet.

Propositions 3.3.1 and 3.3.2 show that once an item has been ordered, it is optimal to immediately commit it to the customer with the earliest arrival time among those who have not been assigned items yet. Consider now an extension of this allocation policy where we commit all items to specific customers whether or not they have been ordered. We will show that such a policy is indeed optimal. In order to do so, we will show that it leads to assigning items once they have been ordered to customer on a first-come, first-served basis.

Suppose that at the beginning of period 0, there are k items that have already been ordered and are either in transit or in stock, then by the virtue of Proposition 3.3.1 and Proposition 3.3.2, we know that it is optimal to allocate these k items to the first k customers on a first-come, first-served basis. For each of the other customers, let us also assign an item (yet to be ordered) to which we will refer as item $a(i)$ for $i = k+1, k+2, \dots$. This is possible since the outside supplier can be viewed as maintaining an infinite stock of items. These items are, by assumption, homogeneous and undifferentiated prior to ordering. Hence, we can couple each item with the supplier to a specific customer. By coupling a specific item to a specific customer, we are committed to allocating the item to the customer regardless of when the item would eventually be ordered and regardless of when the customer would eventually arrive. We refer to this allocation policy as the committed allocation policy. Since under a committed allocation policy there is a one-to-one correspondence between customers and items, we can use a single index, say i , to refer to a coupled customer-item pair i , so that item $a(i)$ now refers to the item that is allocated to customer i .

Given a committed allocation policy, the inventory control problem reduces to making decisions in each period on whether or not we should order an item $a(i)$ (destined for customer i) if the item has not been ordered yet to minimize the expected cost incurred by each customer-item pair. Note that the expected cost associated with an ordering decision

depends only on the distribution of arrival time X_i of customer i . To see why this is the true, let us consider the case where $X_i = x_i$, where x_i is a particular realization of the arrival time random variable X_i . Also, let $\tau_{a(i)}$ be the time at which item $a(i)$ is ordered. Consider a particular realization $L_{\tau_{a(i)}}(\omega) = m$. If $\tau_{a(i)} + m < x_i$ (the item is delivered before its corresponding customer arrives), then the system incurs a holding cost

$$g(\tau_{a(i)}, m, x_i) = \sum_{k=\tau_{a(i)}+m}^{x_i-1} \beta^k h_k; \quad (3.2)$$

if $\tau_{a(i)} + m = x_i$, (the item is delivered in the same period the customer arrives), no holding nor a backorder cost is incurred; while if $\tau_{a(i)} + m > x_i$, (the customer arrives before its item is delivered), the system incurs a backorder cost

$$g(\tau_{a(i)}, m, x_i) = \sum_{k=x_i}^{\tau_{a(i)}+m-1} \beta^k b_k. \quad (3.3)$$

As we can see, although the dynamics of different customers are correlated (i.e., knowing the arrival of customer j will give us new information concerning customer $j + 1$), the *marginal inventory cost* due to a customer-item pair, depends only on when the item is ordered and when the customer arrives and is independent of ordering decisions made for other customer-item pairs under a committed allocation policy.

If customer i has already arrived in period $x_i < t$ and item $a(i)$ has not been ordered yet, the expected marginal cost evaluated in period t of ordering item $a(i)$ in period $l = t, \dots, T$ is given by $\beta^l c_l + \hat{p}_{l,l+m}^t \sum_{s=x_i}^{l+m} \beta^s b_s$. It is easy to verify that, in this case, it is optimal to order item $a(i)$ in period t . In other words for customers who have arrived, it is optimal to order their corresponding items as soon as possible. On the other hand, if customer i has not arrived yet, then the expected marginal cost, evaluated in period t , of ordering item $a(i)$ in period $l = t, \dots, T$ conditioned on $F_{t-1} = H_{t-1} \cup \hat{H}_{t-1}$ is given by:

$$G_l(i|F_{t-1}) = \sum_{m=0}^{r_{max}} \hat{p}_{l,l+m}^t (\beta^l c_l + \sum_{s=t}^{l+m} \hat{p}_{t,s}^i \sum_{k=s}^{l+m-1} \beta^k b_k + \sum_{s=l+m+1}^T \hat{p}_{t,s}^i \sum_{k=l+m}^{s-1} \beta^k h_k). \quad (3.4)$$

Note that the probabilities $\hat{p}_{i,l+m}^t$ and $p_{t,s}^i$ depend on the history of the process F_{t-1} observed at time t . For a given a leadtime realization $L_l(\omega) = m$, the first term inside the bracket on the right-hand side of the above equality is the discounted ordering cost associated with ordering item $a(i)$ in period l ; the second term is the expected backorder cost; and the third term is the expected holding cost. A backorder cost is incurred if customer i arrives in a period s such that $t \leq s < l + m$. The corresponding cost is then $\sum_{k=s}^{l+m-1} \beta^k b_k$, and the expected cost over all possible values of s is $\sum_{s=t}^{l+m-1} p_{t,s}^i \sum_{k=s}^{l+m-1} \beta^k b_k$. A holding cost is incurred if customer i arrives in a period s such $l + m < s$. The corresponding cost is then $\sum_{k=l+m}^{s-1} \beta^k h_k$ and the expected cost is $\sum_{s=l+m+1}^T p_{t,s}^i \sum_{k=l+m}^{s-1} \beta^k h_k$.

In every period t , we must decide whether or not to place an order for item $a(i)$, if it has not been ordered already. In making this decision, we compare the expected cost of ordering the item in period t to the cost of postponing the ordering decision until a later period. Hence, the ordering decision can be formulated as an *optimal stopping* problem (in an optimal stopping problem, the decision is when to terminate a process by taking a certain action, which in our case corresponds to placing an order). The optimality equation in period t is given by

$$V_t(i|F_{t-1}) = \min\{G_t(i|F_{t-1}), E[V_{t+1}(i|F_t)|F_{t-1}]\}, \quad (3.5)$$

for $t = 0, \dots, T$, so that it is optimal to order in period t if $G_t(i|F_{t-1}) \leq E[V_{t+1}(i|F_t)|F_{t-1}]$ and to postpone the decision otherwise. Note that the expectation in (3.5) is taken with respect to demand in period t , D_t , \mathbf{E}_{t+1} , $\hat{\mathbf{E}}_{t+1}$ and L_t conditional on F_{t-1} .

Since under a committed allocation policy, the cost incurred by a particular customer-pair depends only on when the item is ordered and when the customer arrives and is independent of ordering decisions made for other customer-item pairs, the decision for each item can be treated independently of decisions for other items. This immediately leads to the following result.

Proposition 3.3.3 *Under a committed allocation policy, the ordering decision for each customer-item pair can be made independently of decisions for other customer-item pairs.*

The following proposition is needed for the proof of Theorem 3.3.5 .

Proposition 3.3.4 *The function $G_l(i|F_{t-1})$ is submodular in (l, i) for $t = 0, \dots, T$, $l \geq t$ and all customers $i \geq 1$ whose items have not already been ordered at the beginning of period t . That is, $G_l(i|F_{t-1}) - G_{l+1}(i|F_{t-1}) \leq G_l(i+1|F_{t-1}) - G_{l+1}(i+1|F_{t-1})$ for all $F_{t-1} \in \mathcal{F}_{t-1}$.*

Proof. First, let us note that for every possible realization x_i and x_{i+1} of the random variables X_i and X_{i+1} respectively, we have $x_i \leq x_{i+1}$. Therefore, we also have $Pr(X_i > a) \leq Pr(X_{i+1} > a)$ for any $a > 0$. In other words, X_i is *stochastically smaller* than X_{i+1} . Consider a sample path realization ω of the leadtimes $L_0(\omega), \dots, L_T(\omega)$. The expected marginal cost associated with ordering item $a(i)$ in period l under realization ω is given by:

$$G_l(i|H_{t-1}, \omega) = \beta^l c_l + \sum_{s=t}^{l+L_l(\omega)} p_{t,s}^i \sum_{k=s}^{l+L_l(\omega)-1} \beta^k b_k + \sum_{s=l+L_l(\omega)+1}^T p_{t,s}^i \sum_{k=l+L_l(\omega)}^{s-1} \beta^k h_k$$

and

$$G_{l+1}(i|H_{t-1}, \omega) - G_l(i|H_{t-1}, \omega) = \beta^{l+1} c_{l+1} - \beta^l c_l + \sum_{m=L_l(\omega)}^{L_{l+1}(\omega)} \left(\beta^{l+m} b_{l+m} \sum_{s=t}^{l+m} p_{t,s}^i - \beta^{l+m} h_{l+m} \sum_{l+m+1}^T p_{t,s}^i \right).$$

Since X_i is stochastically smaller than X_{i+1} , i.e., $\sum_{l+m+1}^T p_{t,s}^i \leq \sum_{l+m+1}^T p_{t,s}^{i+1}$ for all l , and m , we have $G_{l+1}(i|H_{t-1}, \omega) - G_l(i|H_{t-1}, \omega) \geq G_{l+1}(i+1|H_{t-1}, \omega) - G_l(i+1|H_{t-1}, \omega)$ for all sample paths ω . Consequently, we also have $G_{l+1}(i|F_{t-1}) - G_l(i|F_{t-1}) \geq G_{l+1}(i+1|F_{t-1}) - G_l(i+1|F_{t-1})$ for all F_{t-1} . That is, the function, $G_l(i|F_{t-1})$ is submodular in (i, l) . ■

The submodularity of the function $G_l(i|F_{t-1})$ leads to the following important result.

Theorem 3.3.5 *It is always optimal to order item $a(i)$ for customer i no later than item $a(i+1)$ for customer $i+1$. That is, $G_t(i|F_{t-1}) - E[V_{t+1}(i|F_t)|F_{t-1}] \leq G_t(i+1|F_{t-1}) - E[V_{t+1}(i+1|F_t)|F_{t-1}]$ for all t and all i and all $F_t \in \mathcal{F}_t$. Consequently, the committed allocation policy is optimal.*

Proof. We prove the result by induction. First note that in period $T - 1$, the optimal decision is either to order in period $T - 1$ or to order in period T (which is equivalent to not ordering). Hence,

$$V_{T-1}(i|F_{T-2}) = \min\{G_{T-1}(i|F_{T-2}), E[V_T(i|F_{T-1})]\} = \min\{G_{T-1}(i|F_{T-2}), G_T(i|F_{T-2})\}.$$

Next, suppose that

$$G_{k+1}(i|F_k) - E[V_{k+2}(i|F_{k+1})|F_k] \leq G_{k+1}(i+1|F_k) - E[V_{k+2}(i+1|F_{k+1})|F_k],$$

where $k < T - 2$. It follows that

$$\begin{aligned} & G_k(i|F_{k-1}) - E[V_{k+1}(i|F_k)|F_{k-1}] \\ = & G_k(i|F_{k-1}) - E[\min\{G_{k+1}(i|F_k), E[V_{k+2}(i|F_{k+1})|F_k]\}|F_{k-1}] \\ = & G_k(i|F_{k-1}) - E[G_{k+1}(i|F_k) + \min\{0, E[V_{k+2}(i|F_{k+1})|F_k] - G_{k+1}(i|F_k)\}|F_{k-1}] \\ = & G_k(i|F_{k-1}) - G_{k+1}(i|F_{k-1}) - E[\min\{0, E[V_{k+2}(i|F_{k+1})|F_k] - G_{k+1}(i|F_k)\}|F_{k-1}] \\ \leq & G_k(i+1|F_{k-1}) - G_{k+1}(i+1|F_{k-1}) - E[\min\{0, E[V_{k+2}(i+1|F_{k+1})|F_k] - G_{k+1}(i+1|F_k)\}|F_{k-1}] \\ = & G_k(i+1|F_{k-1}) - E[V_{k+1}(i+1|F_k)|F_{k-1}], \end{aligned}$$

the inequality is due to Proposition 3.3.4 and the inductive assumption.

The optimality of the committed allocation policy follows from the fact that the i th item to be ordered is always assigned to the i th customer to arrive ($a(i) = i'$), and from Propositions 3.3.1 and 3.3.2. ■

To the best of our knowledge, the results below are the first to characterize the structure of the optimal policy under a demand and leadtime correlation model as general as ours.

Theorem 3.3.6 *The optimal ordering policy in periods $t = 0, \dots, T - 1$ is a base-stock policy with state-dependent base-stock level $s_t(F_{t-1})$, such that if $IP_t < s_t(F_{t-1})$, we order $s_t(F_{t-1}) - IP_t$ items and if $IP_t \geq s_t(F_{t-1})$, we order nothing, where IP_t is the inventory*

position at the beginning of period t defined as $IP_t \equiv I_t - B_t + IO_t$, where IO_t is the number of outstanding orders with the outside supplier at time t . Moreover, $0 \leq s_t(F_{t-1}) < \infty$ for $t = 0, \dots, T-1$ and $s_t(F_{t-1})$ is independent of all previous ordering decisions.

Proof: First, we show that it is optimal to order items for any customers who are currently backlogged and whose items have not been ordered yet. Suppose customer i arrived at some period $x_i < t$ (period t is the current period). Consider a particular sample path for leadtimes $L_0(\omega), \dots, L_T(\omega)$. Consider the marginal cost $g(l, x_i) = \beta^l c_l + \sum_{k=x_i}^{l+L_l(\omega)-1} \beta^k b_k$ incurred by the system if item $a(i)$ is ordered in some period $l \geq t$. It is easy to show that

$$g(l+1, x_i) - g(l, x_i) > 0,$$

where the inequality follows from Assumption 1 and the fact that $L_l(\omega) + l < L_{l+1}(\omega) + l + 1$, which implies that it is indeed optimal to order item $a(i)$ in period t . Hence, the only non-trivial decision is for how many future customers to place orders at period t . Suppose that there are k customers have arrived at the beginning of period t . Let $k + j^*(F_{t-1})$ be the index of the largest customer for which it is optimal to order in period t . By virtue of Theorem 3.3.5, if it is optimal to order for customer $k + j^*(F_{t-1})$, then it is optimal to order items for customers $k + 1, k + 2, \dots, k + j^*(F_{t-1}) - 1$ if these items have not been ordered yet in previous periods. This would bring the inventory position to $j^*(F_{t-1})$. On the other hand, if the current inventory position is greater than or equal to $j^*(F_{t-1})$, then we order nothing. Hence, the optimal policy is a base-stock policy with base-stock level $s_t(F_{t-1}) = j^*(F_{t-1})$.

The fact that we always order enough to bring the inventory position to zero (we always order items for customers that have arrived) implies that $s_t(F_{t-1}) \geq 0$. To show that $s_t(F_{t-1}) < \infty$, we first note that $\sum_{s=t}^{l+L} p_{t,s}^{k+j} \rightarrow 0$ monotonically as $j \rightarrow \infty$ for $t \leq l < T$.

Given a particular sample path of $L_0(\omega), \dots, L_T(\omega)$, by Assumption 1, we have

$$\begin{aligned} & \lim_{j \rightarrow \infty} (G_{t+1}(j|H_{t-1}, \omega) - G_t(j|H_{t-1}, \omega)) \\ &= \lim_{j \rightarrow \infty} \left(\beta^{t+1} c_{t+1} - \beta^t c_t + \sum_{m=L_t(\omega)}^{L_{t+1}(\omega)} (\beta^{t+m} b_{t+m} \sum_{s=t}^{t+m} p_{t,s}^j - \beta^{t+m} h_{t+m} \sum_{t+m+1}^T p_{t,s}^i) \right) < 0. \end{aligned}$$

Therefore, there always exists a large enough finite n_{max} for which $G_{t+1}(k+n|F_{t-1}) - G_t(k+n|F_{t-1}) < 0$ for all $n \geq n_{max}$. This means that for any customer with index $i \geq k+n, n \geq n_{max}$, it is not optimal to order the corresponding item in period t . Thus, the base-stock level in period t is finite. The fact that $s_t(F_{t-1})$ is independent of all previous ordering decisions is due to that different subsystems are independent. \blacksquare

In general, because of the multi-dimensional nature of the state of the system in the presence of correlation, computing the optimal policy using the customer-item decomposition approach cannot be carried out in polynomial time (this is of course also true for traditional dynamic programming formulations that use aggregate demand and aggregate order quantities). However, as we show in the following proposition, this is possible for the important special case where demand is strictly positive in each period (or equivalently when the arrival times for future customers are bounded) and the number of all possible realizations of the random vectors \mathbf{E}_t and $\hat{\mathbf{E}}_t$ is upper-bounded by a nonnegative integer K , which can be arbitrarily large. The strictly positive demand assumption is easily satisfied in many practical applications, when the demand in each period tends to be large or the periods are relatively long.

Proposition 3.3.7 *If $D_t \geq \delta > 0$ for all t , the optimal policy can be computed using an algorithm with complexity $O((r_{max} + 1)^{\lfloor \frac{s_{max}}{\delta} \rfloor + 1} K^{\lfloor \frac{s_{max}}{\delta} \rfloor + 1} f(s_{max}, \delta) s_{max}^2 T)$, where $s_{max} = \max\{s_0(F_{-1}), \dots, s_{T-1}(F_{T-2})\}$ and $f(s_{max}, \delta) = \max_{\{k=1, \dots, \lfloor \frac{s_{max}}{\delta} \rfloor + 1\}} \binom{s_{max} - (\delta - 1)k}{k}$.*

Proof. Consider a period t where k customers have already arrived. To compute $V_i(i|F_{t-1})$, we need to consider each realization d_t, \dots, d_{l-1} and l_0, \dots, l_{t-1} such that $0 \leq d_t + \dots +$

$d_{l-1} \leq i - k \leq s_{max}$ for each realization of $\{\mathbf{e}_{t+1}, \dots, \mathbf{e}_l, \hat{\mathbf{e}}_{t+1}, \dots, \hat{\mathbf{e}}_l\}$ (If $d_t + \dots + d_{l-1} > i - k$, then customer i arrives at or before period $l - 1$). Since $\delta > 0$, the number of all possible such demand realizations d_t, \dots, d_{l-1} is $\binom{s_{max} - \delta(l-t) + (l-t)}{l-t}$ given a specific realization of $\{\mathbf{e}_{t+1}, \dots, \mathbf{e}_l, \hat{\mathbf{e}}_{t+1}, \dots, \hat{\mathbf{e}}_l\}$ and the number of all possible such leadtime realization is $(r_{max} + 1)^{l-t}$. We only need to compute the value functions of customer i from period t to period $t + \lfloor \frac{s_{max}}{\delta} \rfloor + 1$ since customer i must arrive before period $t + \lfloor \frac{i-k}{\delta} \rfloor$ and it is not optimal to order customer i such that $i > k + s_{max}$. We need to consider at most s_{max} customers in each period where each of those customers will arrive in at most s_{max} periods. Hence, the total number of computations over the entire planning horizon is of order $((r_{max} + 1)^{\lfloor \frac{s_{max}}{\delta} \rfloor + 1} K^{\lfloor \frac{s_{max}}{\delta} \rfloor + 1} f(s_{max}, \delta) s_{max}^2 T)$. \blacksquare

Note that s_{max} is upper-bounded by the base-stock level of a myopic policy, as described in Proposition 3.5.1 of Section 3.5. As shown in Section 3.5, see equation (3.16), this upper bound is independent of the number of periods T .

To our knowledge, the result in Proposition 3.3.7 is the first to provide a polynomial time algorithm (in T) for systems with demand correlation as general as ours. Using a traditional dynamic programming approach, the complexity is exponential in T as there is a need to compute in each period the optimal base-stock level for every sample path and the number of sample paths grows exponentially in T . The computational effectiveness of our approach follows from the fact that we associate an arrival time with each customer. When demand in each period is strictly positive, the arrival time for each customer is bounded. Therefore, the number of periods in which it may be optimal to order for a particular customer is always finite. In the appendix, we provide numerical results for a simple example which illustrate the differences in computational effort between the traditional dynamic programming approach and the decomposition method and highlight the magnitude of savings achievable with the decomposition method.

We conclude this section by noting that, when demand is strictly positive, the decomposition approach can be extended to systems with infinite horizon. This is possible even when the problem parameters, including demand, are time-dependent. In particular, using the decomposition approach, it is possible in each period to identify the optimal ordering decision for that period (since it continues to be possible to formulate the problem for each customer whose item has not been ordered yet as an optimal stopping problem over a finite horizon). In other words, the decomposition approach can be used to determine optimal decisions in an *online* fashion, with optimal decisions made for each period at the time when that period occurs. Moreover, these optimal ordering decisions can be obtained using a polynomial time algorithm with complexity $O((r_{max} + 1)^{\lfloor \frac{s_{max}}{\delta} \rfloor + 1} K^{\lfloor \frac{s_{max}}{\delta} \rfloor + 1} f(s_{max}, \delta) s_{max}^2)$. The proof of this is similar to that of the finite horizon case and is omitted for brevity.

3.4 Systems with Multiple Stages in Series

In this section, we extend our analysis to inventory systems with multiple stages in series. For ease of exposition, we consider systems with fixed leadtimes first, then we show how the analysis can be extended to systems with stochastic and correlated leadtimes. Consider a system consisting of M stages, with stage $m = 1, \dots, M$ replenished from items in stage $m + 1$. Stage M is replenished from an outside supplier, assumed to have unlimited stock. We refer to the outside supplier as stage $M + 1$. An item shipped in period t from stage $m, m = 2, \dots, M + 1$ arrives at stage $m - 1$ in period $t + L_{m-1}$, where L_m corresponds to the leadtime at stage m . Demand is satisfied from on-hand inventory at stage 1, if any is available; otherwise it is backlogged. The characteristics of the demand process remain unchanged from those described for systems in a single stage. Inventory can be held at any stage but incurs a holding cost $h_{m,t}$ per unit held in stage m in period t . Items in transit from stage $m + 1$ to stage m also incur the per unit cost $h_{m,t}$. There is a per-unit shipping

cost $c_{m,t} \geq 0$ for initiating a shipment of an item from stage $m+1$ to stage m . The per unit backorder cost b_t is incurred for each unit of demand that is backordered at time t . Finally, without loss of generality, we assume that $h_{1,T} = \dots = h_{M,T} = 0$ and $b_T = 0$.

The following describes the sequence of events in each period $t, t = 0, \dots, T$. First, items that were shipped from stage $m+1$ in period $t - L_m$ arrive at stage m , for $m = 1, \dots, M$. Then, a decision is made as to how many items to ship from stage m to stage $m-1$ for $m = 2, \dots, M+1$. Finally, demand for the current period in the amount of d_t is realized and is either fulfilled from on-hand inventory at stage 1 or backordered. All costs are incurred at the end of the period. We use the following notation. We denote by B_t the number of backorders in period t , $I'_{m,t}$ the local inventory at stage m in period t , $IT_{m,t}$ the inventory in transit to stage m , $I_{m,t}$ the echelon inventory at stage m , where $I_{m,t} = I'_{m,t} + \sum_{k=1}^{m-1} (IT_{k,t} + I'_{k,t})$, $IN_{m,t}$ the echelon net inventory at stage m $IN_{m,t} = I_{m,t} - B_t$, and by $IP_{m,t}$ the echelon inventory position at stage m , where $IP_{m,t} = IN_{m,t} + IT_{m,t}$.

We assume that there is no speculative motivation to holding inventory or backordering. Therefore, we assume that the following conditions on the cost parameters are satisfied. These conditions are needed only to ensure that the policy parameters in each stage and each period are non-negative and finite. All others results presented in this section remain valid even if these inequalities do not hold.

Assumption 2 : $c_{m,t} - \beta c_{m,t+1} > h_{m+1,t} - h_{m,t}$, $t = 1, \dots, T$, $m = 1, \dots, M$, $e_{m,t+1} - e_{m,t} + \beta^t h_{m,t} - \beta^{t+\bar{L}_m} h_{1,t+\bar{L}_m} < 0$, and $e_{m,t+1} - e_{m,t} + \beta^t h_{m,t} + \beta^{t+\bar{L}_m} b_{t+\bar{L}_m} > 0$, for all t and $m = 2, \dots, M+1$, where

$$e_{m,t} = \sum_{n=1}^{m-1} \beta^{t+\sum_{l=n+1}^{m-1} L_l} c_{n,t+\sum_{l=n+1}^{m-1} L_l} + \sum_{n=2}^m \sum_{k=\sum_{l=n}^{m-1} L_l}^{\sum_{l=n-1}^{m-1} L_l-1} \beta^{t+k} h_{n-1,t+k}, \quad (3.6)$$

and $\bar{L}_m = L_1 + \dots + L_m$.

We refer to the parameter $e_{m,t}$ as the *echelon leadtime cost* at stage $m \in \{2, \dots, M+1\}$.

Echelon leadtime cost corresponds to the cost that would be incurred if an item is shipped from stage m at period t and pushed through to stage 1, without ever being held in inventory at any of the intermediate stages. Hence, the associated costs are the holding and shipping costs incurred by the item as it travels from stage m to stage 1 during the time interval t to $t + L_{m-1} + \dots + L_1$.

As we did in the single stage case, we associate an index i with the i -th arriving customer. For items, we associate an index j' with the j -th item to be shipped from the outside supplier. At the beginning of a period t , some items would have been shipped from the outside supplier, moved through the entire supply system and used to fulfill customer orders, some items would have been shipped from the outside supplier but are currently either in inventory at one of the stages or in transit between stages, and the remaining items (potentially an infinite number of them) would have not been shipped yet from the outside supplier. The inventory control problem can be stated for each period in terms of (1) determining which items, among those currently in stock at a particular stage (including the outside supplier, stage $M + 1$) to ship to the next stage and (2) determining how to allocate available stock in the last stage, stage 1, to customers who have arrived but whose orders have not been fulfilled yet.

Proposition 3.4.1 *If there is on-hand inventory in stage 1 or inventory in transit to stage 1, then it is optimal to use this inventory to satisfy customers that are currently backlogged. Moreover, it is optimal to allocate this inventory on a first-come, first-served basis so that the item with the smallest index is assigned to the customer, among those backlogged, with the smallest index.*

The proof for the above proposition is similar to the proofs of Propositions 3.3.1 and 3.3.2 and, therefore, we omit it. We can extend the definition of the committed allocation policy to serial systems as follows. Suppose at the beginning of period 0, there are k items in stock

in stage 1 or in transit to stage 1, we allocate these k items to the first k customers on a first-come, first-served basis. In addition, suppose that there are j items that have been shipped from the outside supplier (stage $M + 1$) but have not been shipped to stage 1 yet (these include items in stock at, or in transit to, stage 2, \dots , M). Under the committed allocation policy, we allocate item $(k + i)'$ to customer $k + i$ for $i = 1, \dots, j$. For each of the other customers, we assign to customer i , for $i = k + j + 1, k + j + 2, \dots$, an item $a(i)$, yet to be released from stage $M + 1$. Given a committed allocation policy, the inventory control problem reduces to making decisions in each period on whether or not to ship an item $a(i)$ if it is currently in inventory at stage $m, m = 2, \dots, M + 1$ to the next stage.

Let us first consider decisions associated with items currently in stock in stage 2. The expected marginal cost, evaluated in period t , of shipping an item $a(i)$ from stage 2 in period $l, l = t, \dots, T - L_1 - 1$ is given by

$$G_{2,l}(i|H_{t-1}) = \beta^l c_{1,l} + \sum_{k=l}^{l+L_1-1} \beta^k h_{1,k} + \sum_{r=t}^{l+L_1} p_{t,r}^i \sum_{k=r}^{l+L_1-1} \beta^k b_k + \sum_{r=l+L_1+1}^T p_{t,r}^i \sum_{k=l+L_1}^{r-1} \beta^k h_{1,k} \quad (3.7)$$

In making the decision of whether or not to ship item $a(i)$ in the current period, we compare the expected cost of shipping the item in period t to the cost of postponing the shipping decision until a later period. Hence the shipping decision at stage 2 can be formulated as an optimal stopping problem. The optimality equation in period t is given by

$$V_{2,t}(i|H_{t-1}) = \min\{G_{2,t}(i|H_{t-1}), \beta^t h_{2,t} + E[V_{2,t+1}(i|H_t)|H_{t-1}]\} \quad (3.8)$$

for $t = 0, \dots, T - \bar{L}_2 - 1$.

For items in stages $m = 3, \dots, M + 1$, the formulation is more complicated since the process does not stop once the item is shipped from its current stage. Once the item reaches the next stage, a similar decision must be made then as to whether or not to release it to the subsequent stage. This process continues until the item reaches stage 2, where

a decision to ship does terminate the process. Hence, we have a *nested* optimal stopping problem. The optimality equation, evaluated at period t , associated with an item $a(i)$ that is currently in stock in stage $m = 3, \dots, M + 1$ is given by

$$V_{m,t}(i|H_{t-1}) = \min\{r_{m,t} + E[V_{m-1,t+L_{m-1}}(i|H_{t+L_{m-1}-1})|H_{t-1}], \\ \beta^t h_{m,t} + E[V_{m,t+1}(i|H_t)|H_{t-1}]\} \quad (3.9)$$

where $r_{m,t} = \beta^t c_{m-1,t} + \sum_{k=t}^{t+L_{m-1}-1} \beta^k h_{m-1,k}$.

Since under a committed policy the expected cost incurred by a particular customer-item pair depends only on when the item is shipped from stage to stage and when the corresponding customer arrives to the system, decisions for each customer-item pair can be made again independently of decisions made for other items.

The following proposition provides a submodularity result for stage 2 similar to the result in Proposition 3.3.4. This result is needed for the proof of Theorem 3.4.3. The proof is similar to that of Proposition 3.3.4 and is, therefore, omitted.

Proposition 3.4.2 *For all $l \geq t$ and for all t , and for all customers i whose items are currently in stage 2 at the beginning of period t , The function $G_{2,l}(i|H_{t-1})$ is submodular in (l, i) for $t = 0, \dots, T - \bar{L}_2 - 1$ and all customer i such that $a(i)$ is at stage 2 at the beginning of period t . That is, $G_{2,l}(i|H_{t-1}) - G_{2,l+1}(i|H_{t-1}) \leq G_{2,l}(i+1|H_{t-1}) - G_{2,l+1}(i+1|H_{t-1})$ for all $H_{t-1} \in \mathcal{F}_{t-1}$.*

Theorem 3.4.3 *It is always optimal to ship item $a(i)$ destined for customer i no later than item $a(i+1)$ destined for customer $i+1$ at each stage $m, m = 2, \dots, M + 1$. That is,*

$$G_{2,t}(i|H_{t-1}) - E[V_{2,t+1}(i|H_t)|H_{t-1}] \leq G_{2,t}(i+1|H_{t-1}) - E[V_{2,t+1}(i+1|H_t)|H_{t-1}] \quad (3.10)$$

for all $H_{t-1} \in \mathcal{F}_{t-1}$ and all items $a(i)$ at stage 2, and

$$\begin{aligned} & E[V_{m-1,t+L_{m-1}}(i|H_{t+L_{m-1}-1})|H_{t-1}] - E[V_{m,t+1}(i|H_t)|H_{t-1}] \\ & \leq E[V_{m-1,t+L_{m-1}}(i+1|H_{t+L_{m-1}-1})|H_{t-1}] - E[V_{m,t+1}(i+1|H_t)|H_{t-1}], \end{aligned} \quad (3.11)$$

for all t , all $H_{t-1} \in \mathcal{F}_{t-1}$ and all items $a(i)$ at stage $m = 3, \dots, M+1$. Consequently, the committed allocation policy is optimal.

A proof of the above theorem, which uses an inductive argument, can be found in the Appendix. As stated in the following theorem, the above results can be used to show that the optimal policy at each stage and in each period is a state-dependent *echelon base-stock* policy. A proof can be found in the Appendix.

Theorem 3.4.4 *At each stage $m, m = 1, \dots, M$, the optimal policy in periods $t = 0, 1, \dots, T - \bar{L}_m - 1$ is a state-dependent echelon base-stock policy with an echelon base-stock level $s_{m,t}(H_{t-1})$, such that if $IP_{m,t} < s_{m,t}(H_{t-1})$, we ship $s_{m,t}(H_{t-1}) - IP_{m,t}$ items if there are at least $s_{m,t}(H_{t-1}) - IP_{m,t}$ available at stage $m+1$, otherwise we ship all items at stage $m+1$, and if $IP_{m,t} \geq s_{m,t}(H_{t-1})$, we ship nothing. Moreover, $0 \leq s_{m,t}(H_{t-1}) < \infty$ and $s_{m,t}(H_{t-1})$ is independent of all previous ordering decisions for $t = 0, \dots, T$ and $m = 1, \dots, M$.*

Although the optimal stopping formulation described in this section does not lead to a polynomial time solution algorithm in general, it does so again for the case where demand is strictly positive. In particular, it can be shown (the proof is similar to that of Proposition 3.3.7) that the optimal policy can be computed using an algorithm with complexity $O(MK^{\lfloor \frac{s_{max}}{\delta} \rfloor + 1} f(s_{max}, \delta) s_{max}^2 T)$, where s_{max} is the upper bound on all possible echelon base-stock levels.

We can extend the analysis to systems with stochastic and correlated leadtimes (the approach is similar to the one used for the single stage case). Let $L_{m,t}$ be a random variable

that describes the leadtime for stage m in period t . As we did for the single stage case, we assume that leadtimes are sequential so that $t + L_{m,t}(\omega) < s + L_{m,s}(\omega)$ for all $s > t$ and for all sample paths ω of the random variables $L_{m,0}, \dots, L_{m,T}$. We also assume that leadtimes are bounded so that there exists $r_{max} > 0$ such that $L_{m,t}(\omega) < r_{max}$ for all sample paths ω , all period t , and all stages $m = 1, \dots, M$. For simplicity, we assume that the leadtime $L_{m,t}$ is independent of leadtimes $L_{n,s}$ for all $n \neq m$ and all $s = 0, \dots, T$ for all t . Finally, we assume that leadtimes are independent of inventory policies and of inventory status, including the current inventory level or the number of orders outstanding.

Let $\mathbf{L}_m = (L_{m,0}, \dots, L_{m,T})$ and $\hat{\mathbf{E}} = (\hat{\mathbf{E}}_0, \dots, \hat{\mathbf{E}}_T)$. In each period t , in addition to observing the information set H_{t-1} for demand, we observe an information set $\hat{H}_{t-1} = \{\hat{\mathbf{e}}_0, \dots, \hat{\mathbf{e}}_t, l_{m,0}, \dots, l_{m,t-1} \text{ for } m = 1, \dots, M\} \in \hat{\mathcal{F}}_{t-1}$ for lead time, where $\hat{\mathbf{e}}_t$ is the realized value of the vector of random variables $\hat{\mathbf{E}}_t$ associated with the leadtime indicators in period t and $l_{m,t}$ is the realized value of leadtime in period t at stage m and $\hat{\mathcal{F}}_{t-1}$ is the set of all possible such realizations.

In each period the observed information set \hat{H}_{t-1} is used to update the distributions of future leadtimes, namely the probabilities $p_{l,l+n,m}^t$ evaluated at the beginning of period t that an item released in a future period l to stage m will be received in period $l + n$. All of our analysis and results can be shown to continue to hold, except that the expected costs associated with ordering a particular item in a particular period have now to take into account the distribution of leadtimes. For stage 2, the optimality equation for the customer-item pair i is given by

$$V_{2,t}(i|H_{t-1}, \hat{H}_{t-1}) = \min\{G_{2,t}(i|H_{t-1}, \hat{H}_{t-1}), \beta^t h_{2,t} + E[V_{2,t+1}(i|H_t, \hat{H}_t)|H_{t-1}, \hat{H}_{t-1}]\}, \quad (3.12)$$

where

$$G_{2,i}(i|H_{t-1}, \hat{H}_{t-1}) = \beta^l c_l + \sum_{n=0}^{r_{max}} p_{i,l+n,2}^t \left(\sum_{s=l}^{l+n-1} \beta^s h_{1,s} + \sum_{s=t}^{l+n} p_{t,s}^i \sum_{k=s}^{l+n-1} \beta^k b_k + \sum_{s=l+n+1}^T p_{t,s}^i \sum_{k=l+n}^{s-1} \beta^k h_{1,k} \right). \quad (3.13)$$

For stage $m > 2$, the optimality equation is given by

$$V_{m,t}(i|H_{t-1}, \hat{H}_{t-1}) = \min\{E[r_{m,t}(L_{m-1,t}) + V_{m-1,t+L_{m-1,t}}(i|H_{t+L_{m-1,t}-1}, \hat{H}_{t+L_{m-1,t}-1}|H_{t-1}, \hat{H}_{t-1})], \\ \beta^t h_{m,t} + E[V_{m,t+1}(i|H_t, \hat{H}_t)|H_{t-1}, \hat{H}_{t-1}]\}, \quad (3.14)$$

where $r_{m,t}(L_{m-1,t}) = \beta^t c_{m-1,t} + \sum_{k=t}^{t+L_{m-1,t}-1} \beta^k h_{m-1,k}$ (note that here $L_{m-1,t}$ is a random variable).

We can use the above formulation to show that the optimal policy is a state-dependent echelon base-stock policy with an echelon base-stock level $s_{m,t}(H_{t-1}, \hat{H}_{t-1})$ at stage $m, m = 1, \dots, M$ in period $t = 0, \dots, T$. The optimal solution can again be obtained in polynomial time when demand in each period is strictly positive.

3.5 Myopic Policies

In this section, we study myopic heuristics as alternatives to optimal policies. Such heuristics can be easier to implement and compute than the optimal policy. Such heuristics are also common in practice. We restrict ourselves to one-period look-ahead policies that involve comparing the cost of ordering in a particular period to the cost of ordering in the next period and, in doing so, use only information available in that period. We illustrate how the decomposition approach leads to a simple characterization of myopic policies and conditions under which these policies are optimal. In some cases, we also obtain explicit expressions for the optimal parameters of myopic policies in the form of critical fractiles. We focus mostly on the case with fixed leadtimes and then discuss extensions to the case of stochastic leadtimes.

3.5.1 Systems with a Single Stage

Consider a heuristic policy that makes decisions about whether or not to place an order for customer i in period t by comparing only $G_{t+1}(i|H_{t-1})$ and $G_t(i|H_{t-1})$ (i.e., by comparing only the expected marginal cost of ordering item $a(i)$ in period t to ordering it in period $t + 1$ based on the available information in period t). More specifically, item $a(i)$ is ordered in period t if the following inequality holds

$$G_{t+1}(i|H_{t-1}) - G_t(i|H_{t-1}) > 0,$$

or equivalently if

$$\sum_{s=t}^{t+L} p_{t,s}^i > \gamma_t = \frac{\frac{1}{\beta^L}(c_t - \beta c_{t+1}) + h_{t+L}}{h_{t+L} + b_{t+L}}, \quad (3.15)$$

otherwise, the ordering is postponed until the next period where a similar evaluation is carried out. The above policy is in general sub-optimal. However, it is intuitively appealing. It states that if the probability that customer i will arrive during the leadtime interval $[t, t + L]$ exceeds the threshold γ_t , then we should go ahead and order item $a(i)$. In other words, the ordering policy is determined by a critical fractile similar to the way ordering decisions are made for a newsvendor problem, except that the critical fractile now applies to an individual customer.

In fact, we can recover a newsvendor-type solution by noting that if Inequality (3.15) above specifies that an item for customer i should be ordered, then the policy would also specify that an item for customer $i - 1$ should also be ordered, if customer $i - 1$ has not arrived yet (by Proposition 3.3.4). Let k denote the number of customers that arrived prior to period t and let $k + \hat{s}_t(H_{t-1})$ be the index of the largest customer for which Inequality (3.15) is satisfied. Then, the condition $\sum_{s=t}^{t+L} p_{t,s}^{k+\hat{s}_t(H_{t-1})} > \gamma_t$ can be equivalently restated as the probability that at least $\hat{s}_t(H_{t-1})$ customers arrive in the interval $[t, t + L]$ is greater than γ_t . Hence, the heuristic policy is a base-stock policy with base-stock level $\hat{s}_t(H_{t-1})$ in

period t . This base-stock level is the largest integer s that satisfies the inequality

$$P(D_t + \dots + D_{t+L} \geq s | H_{t-1}) > \gamma_t. \quad (3.16)$$

The following proposition states that the base-stock level under the myopic heuristic is always larger than (or equal to) the base-stock level under the optimal policy.

Proposition 3.5.1 $\hat{s}_t(H_{t-1}) \geq s_t(H_{t-1})$ for $t = 0, \dots, T - L - 1$.

The result follows from the fact that $G_{t+1}(i|H_{t-1}) - G_t(i|H_{t-1}) \geq E[V_{t+1}(i|H_t)] - G_t(i|H_{t-1})$, which is itself due to $G_t(i|H_{t-1}) \geq V_t(i|H_{t-1})$ for $t = 0, \dots, T - L - 1$.

In the next proposition, we provide a condition under which the heuristic is optimal. The condition is based on the *monotone condition* for optimal stopping problems (see for example Chow *et al.* 1971). We reinterpret the condition here for our context and using our notation. Define the function $A_i(t) \equiv G_t(i|H_{t-1}) - G_{t+1}(i|H_{t-1})$. We say that the optimal stopping problem is monotone if $A_i(t) < 0$ implies that $A_i(t+1) < 0$ for $t = 0, \dots, T - L - 1$ regardless of the demand realization in period t and for all customers i whose items have not been ordered yet at time t . If the monotone condition holds, then if the heuristic calls for ordering item $a(i)$ in period t , then it would also call for ordering $a(i)$ in period $t + 1$ no matter what realization of D_t takes place.

Proposition 3.5.2 *The myopic policy is optimal if $A_i(t)$ is monotone for $t = 0, \dots, T - L - 1$ and for $i \geq 0$.*

The proof, which uses arguments similar to those of Chow *et al.* (1971), is omitted for the sake of brevity. The following proposition describes a setting where the monotone condition is satisfied.

Corollary 3.5.3 *If $\beta c_{t+1} - c_t$ and $b_t + h_t$ are both non-decreasing in t , h_t is non-increasing in t , and D_0, D_1, \dots, D_T are stochastically non-decreasing, i.e., $D_0 \leq_{st} D_1 \leq_{st} \dots \leq_{st} D_T$, then the myopic policy is optimal.*

Proof. Suppose that at the beginning of period t , there are k customers who have already arrived. It suffices to show that if $G_{t+1}(i|H_{t-1}) - G_t(i|H_{t-1}) > 0$ for all $i > k$, then we must have $G_{t+2}(i|H_t) - G_{t+1}(i|H_t) > 0$. Note that

$$\begin{aligned} & G_{t+1}(i|H_{t-1}) - G_t(i|H_{t-1}) \\ &= \beta^{t+1}c_{t+1} - \beta^t c_t + \beta^{t+L}b_{t+L} \sum_{s=t}^{t+L} p_{t,s}^i - \beta^{t+L}h_{t+L} \sum_{s=t+L+1}^T p_{t,s}^i \\ &= \beta^t \left(\beta c_{t+1} - c_t + (\beta^L b_{t+L} + \beta^L h_{t+L}) \sum_{s=t}^{t+L} p_{t,s}^i - \beta^L h_{t+L} \right), \end{aligned}$$

and

$$\begin{aligned} & G_{t+2}(i|H_t) - G_{t+1}(i|H_t) \\ &= \beta^{t+2}c_{t+2} - \beta^{t+1}c_{t+1} + \beta^{t+L+1}b_{t+L+1} \sum_{s=t+1}^{t+L+1} p_{t+1,s}^i - \beta^{t+L+1}h_{t+L+1} \sum_{s=t+L+2}^T p_{t+1,s}^i \\ &= \beta^{t+1} \left(\beta c_{t+2} - c_{t+1} + (\beta^L b_{t+L+1} + \beta^L h_{t+L+1}) \sum_{s=t+1}^{t+L+1} p_{t+1,s}^i - \beta^L h_{t+L+1} \right). \end{aligned}$$

Also note that

$$\sum_{s=t+1}^{t+L+1} p_{t+1,s}^i \geq P(D_{t+1} + \dots + D_{t+L+1} \geq i - k),$$

where $P(D_{t+1} + \dots + D_{t+L+1} \geq i - k|H_t)$ corresponds to the probability that customer i arrives between periods $t + 1$ and $t + L + 1$ given H_t subject to the condition that there is no customer arrival in period t and $\sum_{s=t+1}^{t+L+1} p_{t+1,s}^i$ corresponds to the probability that customer i arrives between periods $t + 1$ and $t + L + 1$. Because the D_t 's are stochastically non-decreasing, we have

$$P(D_{t+1} + \dots + D_{t+L+1} \geq i - k|H_t) \geq P(D_t + \dots + D_{t+L} \geq i - k|H_{t-1}) = \sum_{s=t}^{t+L} p_{t,s}^i.$$

Hence, we have $\sum_{s=t+1}^{t+L+1} p_{t+1,s}^i \geq \sum_{s=t}^{t+L} p_{t,s}^i$. Furthermore, because $\beta c_{t+1} - c_t$ and $b_t + h_t$ are both non-decreasing in t and h_t is non-increasing in t , we can conclude that if $G_{t+1}(i|H_{t-1}) - G_t(i|H_{t-1}) > 0$, then we must have $\beta c_{t+1} - c_t + (\beta^L b_{t+L} + \beta^L h_{t+L}) \sum_{s=t}^{t+L} p_{t,s}^i - \beta^L h_{t+L} > 0$.

This in turn implies that $\beta c_{t+2} - c_{t+1} + (\beta^L b_{t+L+1} + \beta^L h_{t+L+1}) \sum_{s=t+1}^{t+L+1} p_{t+1,s}^i - \beta^L h_{t+L+1} > 0$. Thus, $G_{t+2}(i|H_t) - G_{t+1}(i|H_t) > 0$ and the monotone condition is satisfied. Consequently the one-period look-ahead policy is optimal. ■

The monotone condition of Corollary 3.5.3 is sufficient but not necessary. In the following proposition (the proof of which can be found in the Appendix), we provide another condition under which the heuristic policy is optimal.

Proposition 3.5.4 *Suppose that $D_t(H_{t-1}) \geq \eta_t$ for some nonnegative integer η_t and for all H_{t-1} , then the myopic policy in period t is optimal if $s_t(H_{t-1}) \leq s_{t+1}(H_t) + \eta_t$.*

To the best of our knowledge, the results of Corollary 3.5.3 and Proposition 3.5.4 are new. Veinott (1965a) proves that if demands are independent and stochastically increasing and the leadtime is 0, then the myopic policy is optimal, which is a special case of Corollary 3.5.3. Johnson and Thompson (1975) extend the result in Veinott (1965a) to a stationary autoregressive demand process again with no leadtime. Veinott (1965b) proves that the myopic policy is optimal in period t if $s_t \leq s_{t+1}$ under independent demands, hence it is also a special case of Proposition 3.5.4, although the proof in his case is more involved. Lu et al. (2006) show a similar result as in Proposition 3.5.4 under the MMFE demand forecast model. Finally, we should note that the use of the monotone condition from the theory of optimal stopping for proving the optimality of myopic policies for inventory systems also appears to be new.

All of our analysis and results can be extended to the case with stochastic leadtimes with stationary leadtime distributions (that is, when $\hat{p}_{t,t+m} = \hat{p}_m$). In this case, the decision (under the myopic policy) is to order if

$$\begin{aligned} & G_{t+1}(i|H_{t-1}) - G_t(i|H_{t-1}) \\ &= \beta^{t+1} c_{t+1} - \beta^t c_t + \sum_{m=0}^{r_{max}} \hat{p}_m \left(\beta^{t+m} b_{t+m} \sum_{s=t}^{t+m} p_{t,s}^i - \beta^{t+m} h_{t+m} \sum_{s=t+m+1}^T p_{t,s}^i \right) > 0. \end{aligned}$$

Using similar arguments to the ones we used for the fixed leadtime case, we can show that both Corollary 3.5.3 and Proposition 3.5.4 continue to hold. These results are of course also new. Extending the results further to the case where leadtimes are also correlated appears however difficult without additional assumptions on the leadtime distributions and the form of the correlation.

3.5.2 Systems with Multiple Stages

For systems with multiple stages, a myopic heuristic can be constructed as follows. A decision to ship item $a(i)$ from stage m , $m = 2, \dots, M + 1$ in period t (if item $a(i)$ at stage m in period t) is made by comparing the *optimal marginal cost* of shipping the item in period t versus shipping it in period $t + 1$.

Let us consider first the case with fixed leadtimes. At stage 2, the decision is to ship item $a(i)$ to stage 1 (the end stage) if and only if

$$G_{2,t+1}(i|H_{t-1}) + \beta^t h_{2,t} - G_{2,t}(i|H_{t-1}) > 0,$$

while at stage $m = 3, \dots, M + 1$, the decision is to ship item $a(i)$ to stage $m - 1$ if and only if

$$\begin{aligned} & \beta^t h_{m,t} + r_{m,t+1} + E[V_{m-1,t+L_{m-1}+1}(i|H_{t+L_{m-1}})|H_{t-1}] \\ & - r_{m,t} - E[V_{m-1,t+L_{m-1}-1}(i|H_{t+L_{m-1}-1})|H_{t-1}] > 0. \end{aligned}$$

As we can see, other than for stage 2, this policy is not amenable to a simple characterization (in the form of an easy to compute critical fractile) as it is the case for single stage systems. This is due to the *nested* nature of decisions for stages $m = 3, \dots, M + 1$, where the cost at stage m depends on the optimal decision at stages $m - 1, \dots, 2$. Because of this nestedness, the heuristic does not also offer a significant computational advantage over the optimal policy.

An alternative myopic heuristic can be constructed as follows. Let $\tilde{G}_{m,l}(i|H_{t-1})$ denote the expected cost of shipping item $a(i)$ from stage m at time l (given the available information at time t) assuming, once it is shipped from stage m , it is not held in inventory at any subsequent stages. Hence,

$$\tilde{G}_{m,l}(i|H_{t-1}) = \beta^l c_{m-1,l} + e_{m,l} + \sum_{r=t}^{l+\bar{L}_{m-1}} p_{t,r}^i \sum_{k=r}^{l+\bar{L}_{m-1}-1} \beta^k b_k + \sum_{r=l+\bar{L}_{m-1}+1}^T p_{t,r}^i \sum_{k=l+\bar{L}_{m-1}}^{r-1} \beta^k h_{1,k}, \quad (3.17)$$

for $l, l = t, \dots, T - \bar{L}_{m-1} - 1$. Under this policy, we ship item $a(i)$ from stage m in period t if and only if the following holds

$$\beta^t h_{m,t} + \tilde{G}_{m,t+1}(i|H_{t-1}) - \tilde{G}_{m,t}(i|H_{t-1}) > 0,$$

or equivalently

$$\sum_{s=t}^{t+\bar{L}_{m-1}} p_{t,s}^i > \gamma_{m,t} = \frac{\beta^{\bar{L}_{m-1}} h_{1,t+\bar{L}_{m-1}} - h_{m,t} - \frac{1}{\beta^t} (e_{m,t+1} - e_{m,t})}{\beta^{\bar{L}_{m-1}} b_{t+\bar{L}_{m-1}} + \beta^{\bar{L}_{m-1}-1} h_{1,t+\bar{L}_{m-1}}}. \quad (3.18)$$

Otherwise, the decision to ship the item is postponed until the next period where a similar comparison is carried out.

As we did with the single stage case, we can show that the above heuristic leads to a newsvendor-type solution by re-interpreting the inequality in terms of the demand distributions in each period over the leadtime. In particular, we can show that the heuristic is equivalent to a state-dependent echelon base-stock policy with echelon base-stock level $\tilde{s}_{m,t}(H_{t-1})$ in period t , where $\tilde{s}_{m,t}(H_{t-1})$ corresponds to the the largest integer s that satisfies the inequality

$$P(D_t + \dots + D_{t+\bar{L}_{m-1}} \geq s | H_{t-1}) > \gamma_{m,t}. \quad (3.19)$$

In the case of stationary problem parameters ($c_{m,t} = c_m$, $h_{m,t} = h_m$, and $b_t = b$) and no discounting ($\beta = 1$), the parameter $\gamma_{m,t}$ simplifies to $\gamma_{m,t} = \frac{h_1 - h_m}{b + h_1}$.

The following proposition states that the echelon base stock level under the above heuristic is always larger than (or equal to) the echelon base stock level under the optimal policy. A proof can be found in the Appendix.

Proposition 3.5.5 $\infty > \tilde{s}_{m,t}(H_{t-1}) \geq s_{m,t}(H_{t-1})$ for $t = 0, \dots, T - \bar{L}_{m-1} - 1$, all $H_{t-1} \in \mathcal{F}_{t-1}$.

The above proposition generalizes the result in Shang and Song (2003) for systems with independent and identically distribution demand in to systems with non-stationary parameters and correlated demand and leadtimes. It also complements the results obtained by Dong and Lee (2003) who provide a lower bound for the optimal base-stock levels in the special case of systems with an MMFE forecast model.

The following proposition provides conditions under which $\tilde{s}_{m,t}(H_{t-1})$ is optimal.

Proposition 3.5.6 *The myopic policy is optimal at each stage m in each period t , for all m and t , if $s_{m,t}(H_{t-1}) \leq s_{m-1,t+L_{m-1}}(H_{t+L_{m-1}-1})$, and $s_{m,t}(H_{t-1}) \leq s_{m,t+1}(H_t)$ for all t , all $H_T \in \mathcal{F}_T$, and all $m = 3, \dots, M + 1$.*

Proposition 3.5.6 complements the result for single stage systems in Proposition 3.5.5. Although it provides a useful insight into what would make a myopic policy optimal, Proposition 3.5.6 does not provide explicit conditions in terms of the primitives (input parameters) of the problem. The following theorem provides such conditions for the important case of time-dependent demand but stationary cost parameters.

Theorem 3.5.7 *The myopic policy is optimal at each stage if (1) $h_{m,t} = h_m, c_{m,t} = c_m$ and $b_t = b$ and (2) $D_0 \leq_{st} D_1 \leq_{st} \dots \leq_{st} D_T$ and $P(D_t + \dots + D_{t+\bar{L}_m} \geq s | H_{t-1}) < \gamma_{m+1}$ implies $P(D_{t+L_m} + \dots + D_{t+\bar{L}_m} \geq s | H_{t+L_m-1}) < \gamma_m$ for all m and all t and H_{t-1} .*

The result in Theorem 3.5.7 appears to be the first in the literature to provide explicit conditions for the optimality of myopic policies in serial systems and we view this as an

important contribution of the paper. The condition in Theorem 3.5.7 is straightforward to verify for specific demand distributions and cost parameters (a simple example of when the condition holds is a two-stage system where demands are i.i.d and the upstream stage has zero leadtime).

All the above results can be extended to the case of stochastic and leadtimes with stationary distributions. For the reasons mentioned earlier, the results however cannot be extended to systems with correlated leadtimes without additional assumptions on the form of the correlation.

3.6 Extensions

In this section, we briefly discuss how the analysis can be extended to additional systems. We consider two systems whose analysis can be challenging using a traditional dynamic programming approach, namely systems with advance demand information and systems with batch ordering.

3.6.1 Inventory Systems with Advanced Demand Information

Consider a system where demand is announced N periods ahead of its due date, where N is referred to as the demand leadtime. Demand announced in period t can be fulfilled, without a backorder penalty, any time between period t and period $t + N$ (that is, demand can be fulfilled early). Demand fulfilled past the due date $t + N$ incurs a backorder penalty in each period the demand remains unfulfilled. A version of this problem was recently studied by Wang and Toktay (2007) for a system with a single stage, independent demands, and fixed leadtimes. A related problem that has been widely studied is one where demand announced in period t can be fulfilled only in period $t + N$ or later (that is, demand cannot be fulfilled early); see for example Gallego and Ozer (2001), and more recently Gayon et al. (2009),

and the references therein.

The analysis is largely similar to the case without advanced demand information except for how the marginal costs are computed. To illustrate, consider a single stage system with independent demands and independent but stochastic sequential leadtimes. Demand can be still viewed as a stream of customers that arrive over time. A customer is indexed by the time it is announced. To see how the costs associated with each customer-item pair (and the corresponding optimal ordering decision) can be formulated, consider the following two scenarios associated with a customer i that has not become due yet at time t . In the first scenario, customer i that has been announced in some period $t - N < x_i < t$, but its corresponding item has not been ordered yet (if $x_i \leq t - N$, then it is obviously optimal to order for customer i in period t). If we order for customer i in period $l \geq t$, then the corresponding marginal cost is given by $g(l, x_i) = \beta^l c_l + \sum_{m=0}^{r_{max}} \hat{p}_{l,l+m}^t (\sum_{k=x_i+N}^{l+m} \beta^k b_k)$. Hence, we order for customer i in period t if and only if $t \in \arg \min_{\{l=t, \dots, x_i+N-1\}} g(l, x_i)$. Note that the optimal ordering period for customer i depends on when the customer has been announced. Second, consider the case where customer i has not been announced yet at time t . In this case, the marginal cost of ordering for the customer in period l is given by

$$G_l(i|H_{t-1}) = \beta^l c_l + \sum_{m=0}^{r_{max}} \hat{p}_{l,l+m}^t \left(\sum_{s=t}^{l+m} p_{t,s}^i \sum_{k=s+N}^{l+m-1} \beta^k b_k + \sum_{s=l+m+1}^T p_{t,s}^i \sum_{k=l+m}^{s-1} \beta^k h_k \right). \quad (3.20)$$

Using a sub-modularity argument, we can show that, in both cases, it is optimal to order for customer i before we order for customer $i + 1$. Consequently, the decomposition approach continues to be optimal. We can also show that the optimal policy in period t is a base-stock policy but now the base-stock level is state-dependent. The state in period t is specified by the vector of announced demands $(d_{t-N,t}, \dots, d_{t-1,t+N-1})$, where $d_{t',t'+N}$ denotes demand announced in period t' due in period $t' + N$. The state dependency is due to the fact that the ordering decision for a customer depends on whether or not the

customer has been announced and in what period the announcement occurred.

Using a traditional dynamic programming formulation, the multi-dimensionality of the state-space leads to computational intractability. In particular, computational effort grows exponentially in N , the demand leadtime, making it difficult to study systems where demand leadtime is more than few periods. In contrast, the decomposition approach leads to a polynomial time algorithm in both N and T . To see why, note that in each period the only additional effort over the case of no advance demand information involves decisions for customers who have been announced but are not due yet. Since the number of periods in which there may be announced demand is N (namely periods $t - N, \dots, t - 1$) and the number of periods in which it may be optimal to order for this demand is also N (periods $t, \dots, t + N - 1$), the total number of computations in period t is N^2 . This means that the overall computational complexity is of order $O((s_{max}^2 + N^2)T^2)$.

As we did for systems without advance demand information, the analysis can be extended to systems with correlation and multiple stages. For these systems, we can show that the optimal policy retains the same structure, namely state-dependent base-stock policies (state-dependent echelon base-stock policies for systems with multiple stages), except that the state now includes the vector of announced demand. We can also show that optimal solutions can be obtained in polynomial time, in both T and N , when demand is strictly positive. Similar results regarding the myopic policy can also be derived.

3.6.2 Inventory systems with batch ordering

Consider a system where items must be ordered in multiples of batches of size Q , where $Q \geq 1$ (the system is otherwise the same as the one described in Sections 3 and 4). Versions of this problem have been studied previously using a traditional dynamic programming approach for systems with independent demands and fixed leadtimes; see for example Veinott (1965c)

for an early reference and Chen (2000) for a more recent one. To briefly illustrate how the decomposition approach can be extended to systems with batch ordering, let us consider a system with a single stage, fixed leadtimes, and correlated leadtimes. The analysis can be further extended to the more complicated settings with stochastic leadtimes and multiple stages.

Suppose that at time t , we are deciding whether or not to order for the set of customers $Q(k, j) = \{k + jQ, k + jQ + 1, \dots, k + (j + 1)Q - 1\}$ where $j \geq 1$ and k is the number of customers that have already arrived. That is, we are deciding whether or not to order for the next j th batch. This decision can be formulated as the following optimal stopping problem:

$$\hat{V}_t(k, j|H_{t-1}) = \min\left\{ \sum_{i=k+jQ}^{k+(j+1)Q-1} G_t(i|H_{t-1}), E[\hat{V}_{t+1}(k, j|H_t)|H_{t-1}] \right\},$$

$G_t(i|H_{t-1})$ corresponding again to the expected marginal cost incurred by the customer-item pair i if we order item $a(i)$ at period t and can be computed as in equation (3.4). Note that the expected cost associated with the j -th batch is independent of ordering decisions we make for other batches.

Using again a submodularity argument, we can show that it is optimal to order for batch j before ordering for batch $j + 1$, which can then be used to show the optimality of the batch-based decomposition approach. Moreover, we can show that the optimal policy follows a so-called (R, nQ) policy. More specifically, there is a state-dependent order up to level, $R(H_{t-1})$, such that we order enough batches in period t to bring inventory position to a level between $(R(H_{t-1}) - Q$ and $R(H_{t-1})]$. Similarly, we can define a myopic policy by ordering items for the set of customers $Q(i, j)$ at period t if and only if $\sum_{k=i+jQ}^{i+(j+1)Q-1} G_t(k|H_{t-1}) - \sum_{k=i+jQ}^{i+(j+1)Q-1} G_{t+1}(k|H_{t-1}) < 0$. It is possible to derive conditions, similar to those for the single item case, for which this myopic policy is optimal.

Similarly, we can define myopic policy for the case with batch ordering: we order

items for the set of customers $Q(i, j)$ at period t if and only if $\sum_{k=i+jQ}^{i+(j+1)Q-1} G_t(k|H_{t-1}) - \sum_{k=i+jQ}^{i+(j+1)Q-1} G_{t+1}(k|H_{t-1}) < 0$. Once, we get the optimal $j^*(i)$ under myopic policy for different i s, we can figure out the optimal reordering level under myopic policy. Although here we do not have a critical-fractile solution for the optimal reordering level $R_t(H_{t-1})$. Similarly, by the monotone condition, we can prove the following corollary

Corollary 3.6.1 *Under stationary sequential leadtimes, if $\beta = 1$ and $c_{t+1} - c_t$, $b_t + h_t$ and h_t are all non-decreasing in t , h_t is non-increasing in t , and D_0, D_1, \dots, D_T are stochastically nondecreasing, i.e., $D_0 \leq_{st} D_1 \leq_{st} \dots \leq_{st} D_T$, then the myopic policy is optimal.*

The idea of the proof is as follows: if we can show that $\sum_{k=i+jQ}^{i+(j+1)Q-1} G_t(k|H_{t-1}) - \sum_{k=i+jQ}^{i+(j+1)Q-1} G_{t+1}(k|H_{t-1})$ is always no more than $\sum_{k=i+jQ}^{i+(j+1)Q-1} G_{t+1}(k|H_t) - \sum_{k=i+jQ}^{i+(j+1)Q-1} G_{t+2}(k|H_t)$, then we are done. If we follows this, the rest of the proof is very similar to the case with unit ordering. However if $0 < \beta < 1$, unlike the unit ordering case, the myopic policy may not be optimal. Therefore, if there is a discount factor, the myopic policy may not be optimal even if for inventory systems with stationary demand and cost parameters and fixed leadtimes under batch ordering. To our best knowlege, that is the first myopic policy for inventory systems with batch ordering and first optimality condition for such myopic policy.

Next, we consider the serial system with fixed leadtimes. Specifically, we consider the same inventory model as the correlated demand and fixed leadtime in serial system except now that the release quantity in each period from stage $m, m = 2, \dots, M + 1$ must be a multiple of the lot size Q_m such that $Q_{m+1} = n_m Q_m$ as in Chen (2000). Let $Q_m(i, j) = \{i + jQ_m, \dots, i + (j + 1)Q_m - 1\}$. The releasing decision for the set of customer-item pairs $Q_2(i, j)$ at stage 2 can be formulated as an *optimal stopping* problem. The optimality

equation in period t is given by

$$\hat{V}_{2,t}(i, j|H_{t-1}) = \min\left\{ \sum_{k=i+jQ_2}^{i+(j+1)Q_2-1} G_{2,t}(k|H_{t-1}), \beta^t Q_2 h_{2,t} + E[\hat{V}_{2,t+1}(i, j|H_t)|H_{t-1}] \right\} \quad (3.21)$$

for $t = 0, \dots, T - \bar{L}_2 - 1$, where $\hat{V}_{2,t}(i, j|H_{t-1})$ is the optimal expected cost accounted from period t to period T given $H_{t-1} \in \mathcal{F}_{t-1}$ for customers $Q_2(i, j)$. Similarly, the releasing decision of the set of customers $Q_m(i, j)$ at stage m can be formulated as a *nested optimal stopping* problem. The optimality equation is given by

$$\begin{aligned} \hat{V}_{m,t}(i, j|H_{t-1}) = \min\{ & Q_m r_{m,t} + \sum_{k=jn_{m-1}}^{(j+1)n_{m-1}-1} E[\hat{V}_{m-1,t+L_{m-1}}(i, k|H_{t+L_{m-1}-1})|H_{t-1}], \\ & Q_m \beta^t h_{m,t} + E[\hat{V}_{m,t+1}(i, j|H_t)|H_{t-1}] \} \end{aligned} \quad (3.22)$$

where $\hat{V}_{m,t}(i, j|H_{t-1})$ is the optimal expected cost accounted from period t to period T given $H_{t-1} \in \mathcal{F}_{t-1}$ for $t = 0, \dots, T - \bar{L}_m - 1$ and customers $Q_m(i, j)$ is at stage $m, m = 3, \dots, M+1$ at the beginning of period t . Similarly, we can prove that

$$\begin{aligned} & \sum_{k=i+jQ_2}^{i+(j+1)Q_2-1} G_{2,t}(k|H_{t-1}), \beta^t Q_2 h_{2,t} + E[\hat{V}_{2,t+1}(i, j|H_t)|H_{t-1}] \\ \leq & \sum_{k=i+(j+1)Q_2}^{i+(j+2)Q_2-1} G_{2,t}(k|H_{t-1}), \beta^t Q_2 h_{2,t} + E[\hat{V}_{2,t+1}(i, j+1|H_t)|H_{t-1}] \end{aligned}$$

and

$$\begin{aligned} & \sum_{k=jn_{m-1}}^{(j+1)n_{m-1}-1} E[\hat{V}_{m-1,t+L_{m-1}}(i, k|H_{t+L_{m-1}-1})|H_{t-1}] - E[\hat{V}_{m,t+1}(i, j|H_t)|H_{t-1}] \\ \leq & \sum_{k=(j+1)n_{m-1}}^{(j+2)n_{m-1}-1} E[\hat{V}_{m-1,t+L_{m-1}}(i, k|H_{t+L_{m-1}-1})|H_{t-1}] - E[\hat{V}_{m,t+1}(i, j+1|H_t)|H_{t-1}] \end{aligned}$$

for all i and $j = 0, 1, \dots$ and all H_{t-1} .

Theorem 3.6.2 *If it is optimal to release the items for the set of customers $Q_m(i, j)$ no later than the set of customers $Q_m(i, j+1)$, hence the committed policy is optimal. The optimal releasing policy is a state-dependent echelon $(R_m(H_{t-1}), nQ_m)$ policy, where*

$R_m(H_{t-1}) > -Q_m$ is the echelon reorder point, i.e., we order nQ_m items in period t such that the echelon inventory position immediately after ordering is between $(R_m(H_{t-1}) - Q_m, R_m(H_{t-1})]$.

The proof again follows from that if it is optimal to order the set of customers $Q_m(i, j)$ no later than the set of customers $Q_m(i, j + 1)$ at each stage $m = 2, \dots, M + 1$, and the optimality of FCFS allocation policy at stage 1. Hence our result generalizes Chen (2000)'s result which focuses on serial systems with i.i.d demand and fixed leadtimes to the cases with generalized demand correlation and sequential stochastic leadtimes with correlation.

Note that if $Q_m = Q$ for all $m = 2, \dots, M + 1$. Then we can carry out similar heuristic analysis in the serial systems. Specifically, if we release the items for the set of customers $Q_m(i, j)$ if and only if

$$Q\beta^t h_{m,t} + \sum_{k=i+jQ}^{i+(j+1)Q-1} (G_{m,t+1}(k|H_{t-1}) - G_{m,t}(k|H_{t-1})) > 0.$$

Let $\tilde{R}_m(H_{t-1})$ be the reorder point calculated by this heuristic policy, then we can also prove that $\tilde{R}_m(H_{t-1}) \geq R_m(H_{t-1})$.

By mimicking the proof in the time-correlated sequential leadtime case for serial systems, we can extend our result to that case under batch ordering as well. Hence our result generalizes Chen (2000)'s result which focuses on serial systems with i.i.d demand and fixed leadtimes to the cases with generalized demand correlation and sequential stochastic leadtimes with correlation. Similarly if we can obtain polynomial time algorithms for both Markov modulated demand case and positive demand case. It is easy to see that the complexity for batching ordering is upperbounded by the complexity without batch ordering.

3.7 Conclusion

In this paper, we considered a stochastic inventory system with correlation in both demands and leadtimes. The demand and leadtime correlation models we considered are general and incorporate most models previously treated in the literature. In doing so, we treated systems with both single and multiple stages and discussed extensions to more general cases, involving advance demand information and batch ordering. In each case, we showed how a customer-item decomposition approach can be used to formulate the inventory control problem as a set of independent optimal stopping problems, each involving a single customer-item pair. We then used properties that arise from this formulation to characterize the structure of the optimal policy. We showed how our formulation leads in the important case of positive demand to polynomial time algorithms. We also used the formulation to construct myopic heuristics. We described how the myopic heuristics can lead to explicit solutions for the optimal policy in the form of critical fractiles. We also identified conditions under which the myopic heuristic is optimal and described settings under which these conditions hold.

As we have seen, the potential applicability of the approach is broad. We expect additional problems, such as those involving multiple stages with more complex structures, to be also amenable to analysis. However, we must caution that the approach is not without its limitations. An important assumption is the sequentiality of leadtimes. This means that the approach would not be applicable when there is the possibility of order cross over (e.g., when leadtimes are stochastic and independent). A second important assumption is the optimality of first-come, first served allocation of inventory. This assumption would not hold when there are reasons to withhold inventory and to allocate it based on a criterion other than first-come, first-served (e.g., this would be the case when demand arises from multiple customer classes, with backorder penalties varying from class to class). The approach is also

not applicable to systems with lost sales since a first-come first served allocation cannot be guaranteed. A third important assumption is that inventory decisions cannot affect supply leadtimes. This assumption is not valid when leadtimes are sensitive to the number of orders already placed (e.g., production facilities at the supplier are subject to congestion). Please refer to Chapter 4 of this thesis for possible extensions of this approach.

Appendix

Proof of Proposition 3.5.4: Suppose that there are k customers who have already arrived. We know that it is optimal to order customer j in period t if $j \leq k + s_t$. Since $0 \leq s_t(H_{t-1}) \leq s_{t+1}(H_t) + \eta_t$, it is equivalent to say that if it is optimal to order for customer j in period t , then it must be optimal to order for customer j in period $t + 1$ (since we know in period t there will be at least η_t customer arrivals). Specifically,

$$\begin{aligned} & G_t(j|H_{t-1}) - E[V_{t+1}(j|H_t)|H_{t-1}] \\ &= G_t(j|H_{t-1}) - G_{t+1}(j|H_{t-1}) - E[\min\{0, E[V_{t+2}(j|H_{t+1})|H_t] - G_{t+1}(j|H_t)\}]. \end{aligned}$$

Hence, for $j \leq k + s_t(H_{t-1})$, we have $E[V_{t+2}(j|H_{t+1})|H_t] - G_{t+1}(j|H_t) \geq 0$. It follows that the one-period lookahead policy is optimal in period t for customers $j \leq k + s_t(H_{t-1})$. ■

Proof of Theorem 3.4.3: From proposition 3.4.2, and using arguments as in the proof of Theorem 3.3.5, we have

$$G_{2,t}(i|H_{t-1}) - E[V_{2,t+1}(i|H_t)|H_{t-1}] \leq G_{2,t}(i+1|H_{t-1}) - E[V_{2,t+1}(i+1|H_t)|H_{t-1}].$$

Next, we show that

$$V_{2,t}(i|H_{t-1}) - E[V_{2,t+1}(i|H_t)|H_{t-1}] \leq V_{2,t}(i+1|H_{t-1}) - E[V_{2,t+1}(i+1|H_t)|H_{t-1}].$$

To do so, note that

$$V_{2,t}(i|H_{t-1}) - E[V_{2,t+1}(i|H_t)|H_{t-1}] = \min\{G_{2,t}(i|H_{t-1}) - E[V_{2,t+1}(i|H_t)|H_{t-1}], \beta^t h_{2,t}\}.$$

Hence,

$$V_{2,t}(i|H_{t-1}) - E[V_{2,t+1}(i|H_t)|H_{t-1}] \leq V_{2,t}(i+1|H_{t-1}) - E[V_{2,t+1}(i+1|H_t)|H_{t-1}]$$

for all $H_{t-1} \in \mathcal{F}_{t-1}$. To prove the result for $m > 2$, we use induction. Suppose that for some $m > 2$, we have the result in (1):

$$V_{m,t}(i|H_{t-1}) - E[V_{m,t+1}(i|H_t)|H_{t-1}] \leq V_{m,t}(i+1|H_{t-1}) - E[V_{m,t+1}(i+1|H_t)|H_{t-1}],$$

and (2):

$$\begin{aligned} & E[V_{m,t+L_m}(i|H_{t+L_m-1})|H_{t-1}] - E[V_{m+1,t+1}(i|H_t)|H_{t-1}] \\ & \leq E[V_{m,t+L_m}(i+1|H_{t+L_m-1})|H_{t-1}] - E[V_{m+1,t+1}(i+1|H_t)|H_{t-1}]. \end{aligned}$$

Then, for $m+1$, we have

$$\begin{aligned} & V_{m+1,t}(i|H_{t-1}) - E[V_{m+1,t+1}(i|H_t)|H_{t-1}] \\ & = \min\{r_{m+1,t} + E[V_{m,t+L_m}(i|H_{t+L_m-1})|H_{t-1}] - E[V_{m+1,t+1}(i|H_t)|H_{t-1}], \beta^{t+1}h_{m+1,t}\} \\ & \leq \min\{r_{m+1,t} + E[V_{m,t+L_m}(i+1|H_{t+L_m-1})|H_{t-1}] - E[V_{m+1,t+1}(i+1|H_t)|H_{t-1}], \beta^{t+1}h_{m+1,t}\} \\ & = V_{m+1,t}(i+1|H_{t-1}) - E[V_{m+1,t+1}(i+1|H_t)|H_{t-1}]. \end{aligned}$$

The inequality is due to the inductive assumption (2). To prove the second part, we again use an inductive argument. Note that for period $T - \bar{L}_{m+1} - 1$, we have

$$\begin{aligned} & V_{m+2,T-\bar{L}_{m+2}-1}(i|H_{T-\bar{L}_{m+2}-2}) \\ & = \min\{r_{m+2,T-\bar{L}_{m+2}-1} + V_{m+1,T-\bar{L}_{m+1}-1}(i|H_{T-\bar{L}_{m+1}-2}), \\ & \quad \beta^t h_{m+2,T-\bar{L}_{m+2}-1} + E[V_{m+2,T-\bar{L}_{m+2}}(i|H_{T-\bar{L}_{m+2}-1})]\} \\ & = \min\{r_{m+2,T-\bar{L}_{m+2}-1} + E[V_{m+1,T-\bar{L}_{m+1}-1}(i|H_{T-\bar{L}_{m+1}-2})|H_{T-\bar{L}_{m+1}-2}], \\ & \quad \beta^t h_{m+2,T-\bar{L}_{m+2}-1} + r_{m+1,T-\bar{L}_{m+2}} + E[V_{m+1,T-\bar{L}_{m+1}}(i|H_{T-\bar{L}_{m+1}-1})|H_{T-\bar{L}_{m+2}-2}]\} \end{aligned}$$

By the inductive assumption (1), we have

$$\begin{aligned} & E[V_{m+1,T-\bar{L}_{m+1}-1}(i|H_{T-\bar{L}_{m+1}-2}) - V_{m+2,T-\bar{L}_{m+2}}(i|H_{T-\bar{L}_{m+2}-1})|H_{T-\bar{L}_{m+2}-2}] \\ & \leq E[V_{m+1,T-\bar{L}_{m+1}-1}(i+1|H_{T-\bar{L}_{m+1}-2}) - V_{m+2,T-\bar{L}_{m+2}}(i+1|H_{T-\bar{L}_{m+2}-1})|H_{T-\bar{L}_{m+2}-2}]. \end{aligned}$$

Suppose now that

$$\begin{aligned} & E[V_{m+1,n+1+L_{m+1}}(i|H_{n+L_{m+1}})|H_n] - E[V_{m+2,n+2}(i|H_{n+1})|H_n] \\ & \leq E[V_{m+1,n+1+L_{m+1}}(i+1|H_{n+L_{m+1}})|H_n] - E[V_{m+2,n+2}(i+1|H_{n+1})|H_n], \end{aligned}$$

for some period $n + 1 < T - \bar{L}_{m+2} - 1$. It then follows that for period n

$$\begin{aligned}
& E[V_{m+1,n+L_{m+1}}(i|H_{n+L_{m+1}-1})|H_{n-1}] - E[V_{m+2,n+1}(i|H_n)|H_{n-1}] \\
&= E[V_{m+1,n+L_{m+1}}(i|H_{n+L_{m+1}-1})|H_{n-1}] \\
&\quad - \min\{r_{m+2,n+1} + E[V_{m+1,n+L_{m+1}+1}(i|H_{n+L_{m+1}})|H_{n-1}], \\
&\quad \beta^{n+1}h_{m+2,n+1} + E[V_{m+2,n+2}(i|H_{n+1})|H_{n-1}]\} \\
&= E[V_{m+1,n+L_{m+1}}(i|H_{n+L_{m+1}-1}) - V_{m+1,n+L_{m+1}+1}(i|H_{n+L_{m+1}})|H_{n-1}] \\
&\quad - \min\{r_{m+2,n+1} + E[V_{m+2,n+2}(i|H_{n+1}) - V_{m+1,n+L_{m+1}+1}(i|H_{n+L_{m+1}})|H_{n-1}], \beta^{n+1}h_{m+2,n+1}\} \\
&\leq E[V_{m+1,n+L_{m+1}}(i + 1|H_{n+L_{m+1}-1})|H_{n-1}] - E[V_{m+1,n+L_{m+1}+1}(i + 1|H_{n+L_{m+1}})|H_{n-1}] \\
&\quad - \min\{r_{m+2,n+1} + E[V_{m+2,n+2}(i + 1|H_{n+1}) - V_{m+1,n+L_{m+1}+1}(i + 1|H_{n+L_{m+1}})|H_{n-1}], \\
&\quad \beta^{n+1}h_{m+2,n+1}\} \\
&= E[V_{m+1,n+L_{m+1}}(i + 1|H_{n+L_{m+1}-1})|H_{n-1}] - E[V_{m+2,n+1}(i + 1|H_n)|H_{n-1}].
\end{aligned}$$

The optimality of committed allocation policies follows from the fact that the i th item to be released is assigned to the i th customer to arrive at every stage $m = 2, \dots, M + 1$, i.e., $a(i) = i'$ for all i and from Proposition 3.4.1. \blacksquare

Proof of Theorem 3.4.4: First, given Assumption 2, it is straightforward to show that it is optimal to ship items to the next stage if their corresponding customers are backlogged. Hence, the only non-trivial decision in each period is for how many future customers to ship items from stage their current stage to subsequent stages. Let $k + j_{m-1}^*(H_{t-1})$ be the index of the largest customers for whom it is optimal to ship an item from stage m to stage $m - 1$ in period t , assuming that item is available at stage m . Then, by virtue of Theorem 3.4.3, it is also optimal to ship items for customers $k + 1, k + 2, \dots, k + j_{m-1}^*(H_{t-1}) - 1$ if these items are in stage m . This would bring the echelon inventory position in stage $m - 1$ to $j_{m-1}^*(H_{t-1})$. If the current echelon inventory position in stage $m - 1$ is greater than or equal to $j_{m-1}^*(t)$, then no items are shipped. Thus, the optimal policy is an echelon base-stock

policy with echelon base-stock level $s_{m,t}(H_{t-1}) = j_m^*(H_{t-1})$ for $m = 1, \dots, M$. The fact that we always release enough to bring the echelon inventory position at stage $m-1$ to zero implies that $s_{m,t} \geq 0$. We will show in the Proof of Proposition 3.5.5 that $s_{m,t}(H_{t-1}) < \infty$ by showing that an upper bound on $s_{m,t}(H_{t-1})$ is finite. \blacksquare

Proof of Proposition 3.5.5: We prove the result by induction. For stage 2, we have

$$G_{2,t}(i|H_{t-1}) - G_{2,t+1}(i|H_{t-1}) \leq V_{2,t}(i|H_{t-1}) - E[V_{2,t+1}(i|H_t)|H_{t-1}].$$

Thus, $G_{2,t}(i|H_{t-1}) - G_{2,t+1}(i|H_{t-1}) \geq 0 \Rightarrow V_{2,t}(i|H_{t-1}) - E[V_{2,t+1}(i|H_t)|H_{t-1}] \geq 0$. Next suppose that for stage $m > 2$,

$$\begin{aligned} \tilde{G}_{m,t}(i|H_{t-1}) - \tilde{G}_{m,t+1}(i|H_{t-1}) - \beta^t h_{m,t} &\geq 0 \\ \Rightarrow V_{m,t}(i|H_{t-1}) - E[V_{m,t+1}(i|H_t)|H_{t-1}] - \beta^t h_{m,t} &\geq 0. \end{aligned}$$

Then for stage $m+1$, we would have

$$\begin{aligned} &V_{m+1,t}(i|H_{t-1}) - E[V_{m+1,t+1}(i|H_t)|H_{t-1}] - \beta^t h_{m+1,t} \\ &= \min\{r_{m+1,t} + E[V_{m,t+L_m}(i|H_{t+L_m-1})|H_{t-1}] - E[V_{m+1,t+1}(i|H_t)|H_{t-1}] - \beta^{t+1} h_{m+1,t}, 0\} \\ &\geq \min\{r_{m+1,t} + E[V_{m,t+L_m}(i|H_{t+L_m-1})|H_{t-1}] \\ &\quad - r_{m+1,t+1} - E[V_{m,t+L_m+1}(i|H_{t+L_m})|H_{t-1}] - \beta^{t+1} h_{m+1,t}, 0\} \\ &\geq \min\{r_{m+1,t} + \tilde{G}_{m,t+L_m-1}(i|H_{t-1}) - r_{m+1,t+1} - \tilde{G}_{m,t+L_m+1}(i|H_{t-1}) - \beta^{t+1} h_{m+1,t}, 0\} \\ &= \min\{\tilde{G}_{m+1,t}(i|H_{t-1}) - \tilde{G}_{m+1,t+1}(i|H_{t-1}) - \beta^{t+1} h_{m+1,t}, 0\}, \end{aligned}$$

which leads to

$$\begin{aligned} \tilde{G}_{m+1,t}(i|H_{t-1}) - \tilde{G}_{m+1,t+1}(i|H_{t-1}) - \beta^t h_{m+1,t} &\geq 0 \\ \Rightarrow V_{m+1,t}(i|H_{t-1}) - E[V_{m+1,t+1}(i|H_t)|H_{t-1}] - \beta^t h_{m+1,t} &\geq 0. \end{aligned}$$

Hence, $\tilde{s}_{m,t}(H_{t-1}) \geq s_{m,t}(H_{t-1})$. To show that $\tilde{s}_{m,t}(H_{t-1}) < \infty$, we first note that $\sum_{r=t}^{l+\bar{L}_m} p_{t,r}^{k+j} \rightarrow 0$ monotonically as $j \rightarrow \infty$ for $t \leq l < T - \bar{L}_m$ for all $H_{t-1} \in \mathcal{F}_{t-1}$.

Furthermore, by Assumption 2, we have

$$e_{m,t+1} - e_{m,t} + \beta^t h_{m,t} - \beta^{t+\bar{L}_m} h_{t+\bar{L}_m} < 0.$$

Therefore, for $m = 2, \dots, M + 1$,

$$\begin{aligned} & \lim_{j \rightarrow \infty} \left(\tilde{G}_{m,t+1}(k+j|H_{t-1}) - \tilde{G}_{m,t}(k+j|H_{t-1}) \right) \\ &= \lim_{j \rightarrow \infty} \left(e_{m,t+1} - e_{m,t} + \beta^t h_{m,t} + \beta^{t+\bar{L}_m} b_{t+\bar{L}_m} \sum_{r=t}^{t+\bar{L}_m} p_{t,r}^j \right. \\ & \quad \left. - \beta^{t+\bar{L}_m} h_{1,t+\bar{L}_m} \sum_{r=t+\bar{L}_m+1}^T p_{t,r}^j \right) < 0. \end{aligned}$$

Thus, there always exists a large enough finite n_{max} for which $G_{m,t+1}(k+n|H_{t-1}) - G_{m,t}(k+n|H_{t-1}) < 0$ for all $n > n_{max}$. This means that it is not optimal to ship an item to the subsequent stage in period t if it is destined for a customer with index $i > k+n, n > n_{max}$. Since $\tilde{s}_{m,t}(H_{t-1}) \geq s_{m,t}(H_{t-1})$, this also means that the echelon base-stock level in period t is finite for each stage $m, m = 2, \dots, M + 1$. \blacksquare

Proof of Proposition 3.5.6: Note that it is sufficient (by Proposition 3.5.5) to prove that

$$\begin{aligned} & \beta^t h_{m,t} + \tilde{G}_{m,t+1}(i|H_{t-1}) - \tilde{G}_{m,t}(i|H_{t-1}) > 0 \\ \Rightarrow & \beta^t h_{m,t} + r_{m,t+1} + E[V_{m-1,t+L_{m-1}+1}(i|H_{t+L_{m-1}})|H_{t-1}] \\ & - r_{m,t} - E[V_{m-1,t+L_{m-1}}(i|H_{t+L_{m-1}-1})|H_{t-1}] > 0. \end{aligned}$$

for all H_{t-1} and $m = 2, \dots, M + 1$. We prove the result by induction. For stage 2, Proposition 3.5.4, the result is true as it is in the single stage system case. Suppose the result is true for some stage $m > 2$. Since we have $s_{m+1,t}(H_{t-1}) \leq s_{m+1,t+1}(H_t)$, then the one-period lookahead policy is optimal by a similar argument as in Proposition 3.5.4. Now, suppose that there are k customers that have already arrived prior to period t , the optimal

echelon base stock level for stage $m + 1$ at period t is the largest index $i, i > k$ such that

$$\begin{aligned} & \beta^t h_{m+1,t} + r_{m+1,t+1} + E[V_{m,t+L_m+1}(i|H_{t+L_m})|H_{t-1}] \\ & - r_{m+1,t} - E[V_{m,t+L_m}(i|H_{t+L_m-1})|H_{t-1}] > 0. \end{aligned}$$

However, note that

$$\begin{aligned} & \beta^{t+L_m} h_{m,t+L_m} + E[V_{m,t+L_m+1}(i|H_{t+L_m})|H_{t-1}] - E[V_{m,t+L_m}(i|H_{t+L_m-1})|H_{t-1}] \\ = & \beta^{t+L_m} h_{m,t+L_m} - E[\min\{r_{m,t+L_m} + E[V_{m-1,t+L_m+L_{m-1}}(i|H_{t+L_m+L_{m-1}-1})|H_{t+L_m-1}], \\ & \beta^{t+L_m} h_{m,t+L_m} + E[V_{m,t+L_m+1}(i|H_{t+L_m})|H_{t+L_m-1}]\}|H_{t-1}] + E[V_{m,t+L_m+1}(i|H_{t+L_m})|H_{t-1}] \\ = & \beta^{t+L_m} h_{m,t+L_m} - E[\min\{r_{m,t+L_m} + E[V_{m-1,t+L_m+L_{m-1}}(i|H_{t+L_m+L_{m-1}-1})|H_{t+L_m-1}] \\ & - E[V_{m,t+L_m+1}(i|H_{t+L_m})|H_{t+L_m-1}], \beta^{t+L_m} h_{m,t+L_m}\}|H_{t-1}] + E[V_{m,t+L_m+1}(i|H_{t+L_m})|H_{t-1}]. \end{aligned} \tag{3.23}$$

Since we have $s_{m+1,t}(H_{t-1}) \leq s_{m,t+L_m}(H_{t+L_m-1})$, then if it is optimal to ship item $a(i)$ from stage $m + 1$ in period t , it must also be optimal to ship item $a(i)$ from stage m in period $t + L_m$. By the inductive assumption we must have

$$\begin{aligned} & \tilde{G}_{m,t+L_m}(i|H_{t+L_m-1}) - \tilde{G}_{m,t+L_m+1}(i|H_{t+L_m-1}) - \beta^{t+L_m} h_{m,t+L_m} \\ = & r_{m,t+L_m} + \tilde{G}_{m-1,t+L_m+L_{m-1}}(i|H_{t+L_m-1}) \\ & - r_{m,t+L_m+1} - \tilde{G}_{m-1,t+L_m+L_{m-1}+1}(i|H_{t+L_m-1}) - \beta^{t+L_m} h_{m,t+L_m} < 0. \end{aligned}$$

Hence, for such an item $a(i)$, we have

$$\begin{aligned}
& \min\{r_{m,t+L_m} + E[V_{m-1,t+L_m+L_{m-1}}(i|H_{t+L_m+L_{m-1}-1})|H_{t+L_m-1}] \\
& \quad - E[V_{m,t+L_m+1}(i|H_{t+L_m})|H_{t+L_m-1}], \beta^{t+L_m} h_{m,t+L_m}\} \\
& = \beta^{t+L_m} h_{m,t+L_m} + \min\{r_{m,t+L_m} + \tilde{G}_{m-1,t+L_m+L_{m-1}}(i|H_{t+L_m-1}) \\
& \quad - r_{m,t+L_m+1} - \tilde{G}_{m-1,t+L_m+L_{m-1}+1}(i|H_{t+L_m-1}) - \beta^{t+L_m} h_{m,t+L_m}, 0\} \\
& = r_{m,t+L_m} + \tilde{G}_{m-1,t+L_m+L_{m-1}}(i|H_{t+L_m-1}) - r_{m,t+L_m+1} - \tilde{G}_{m-1,t+L_m+L_{m-1}+1}(i|H_{t+L_m-1}).
\end{aligned} \tag{3.24}$$

The last equality is due to the inductive assumption and the fact that it is optimal to ship item $a(i)$ from stage m in period $t + L_m$. Combining results in Equations (3.23) and (3.24), leads to

$$\begin{aligned}
& \beta^t h_{m+1,t} + \tilde{G}_{m+1,t+1}(i|H_{t-1}) - \tilde{G}_{m+1,t}(i|H_{t-1}) > 0 \\
& \Rightarrow \beta^t h_{m+1,t} + r_{m+1,t+1} + E[V_{m,t+L_m+1}(i|H_{t+L_m})|H_{t-1}] \\
& \quad - r_{m+1,t} - E[V_{m,t+L_m}(i|H_{t+L_m-1})|H_{t-1}] > 0,
\end{aligned}$$

which means that $\tilde{s}_{m+1,t}(H_{t-1}) \leq s_{m+1,t}(H_{t-1})$. However, we also know that $\tilde{s}_{m+1,t}(H_{t-1}) \geq s_{m+1,t}(H_{t-1})$. Thus, $\tilde{s}_{m+1,t}(H_{t-1}) = s_{m+1,t}(H_{t-1})$. \blacksquare

Proof of Theorem 3.5.7 First, note that the optimal base-stock level $s_{2,t+L_2}(H_{t+L_2-1})$ under the myopic policy for the stage 2 problem is given by the largest index s such that

$$P(D_{t+L_2} + \cdots + D_{t+\bar{L}_2} \geq s | H_{t+L_2-1}) > \gamma_2. \tag{3.25}$$

We also know that $\tilde{s}_{t+L_2}(H_{t+L_2-1}) = s_{t+L_2}(H_{t+L_2-1})$ by Corollary 3.5.3. Hence, we have

$$P(D_{t+L_2} + \cdots + D_{t+\bar{L}_2} \geq s_{2,t+L_2}(H_{t+L_2-1}) | H_{t-1}) > \gamma_2. \tag{3.26}$$

The optimal echelon base stock level for the stage 3 problem under the myopic policy is given by the largest s such that

$$P(D_t + \cdots + D_{t+\bar{L}_2} \geq s | H_{t-1}) > \gamma_3. \tag{3.27}$$

Let $\tilde{s}_{3,t}(H_{t-1})$ be such index. But based on our assumption, we must have

$$P(D_{t+L_2} + \cdots + D_{t+\bar{L}_2} \geq s | H_{t+L_2-1}) > \gamma_2, \quad (3.28)$$

for every s satisfying (3.27). Hence we must have $s_{3,t}(H_{t-1}) \leq \tilde{s}_{3,t}(H_{t-1}) \leq \tilde{s}_{2,t+L_2}(H_{t+L_2-1})$.

We know that under our assumption, we have $\tilde{s}_{2,t+L_2}(H_{t+L_2-1}) = s_{2,t+L_2}(H_{t+L_2-1})$. Hence,

$s_{3,t}(H_{t-1}) \leq \tilde{s}_{3,t}(H_{t-1}) \leq s_{2,t+L_1}(H_{t+L_1-1})$. Then by Proposition 3.5.6, $\tilde{s}_{3,t}(H_{t-1})$ is optimal.

Next, we prove our result by induction. Suppose that at stage m , we have $\tilde{s}_{m,t}(H_{t-1}) =$

$s_{m,t}(H_{t-1})$. Since $P(D_t + \cdots + D_{t+\bar{L}_m} \geq s | H_{t-1}) < \gamma_{m+1}$ implies $P(D_{t+L_m} + \cdots + D_{t+\bar{L}_m} \geq$

$s | H_{t+L_m-1}) < \gamma_m$, we must have $s_{m+1,t}(H_{t-1}) \leq \tilde{s}_{m+1,t}(H_{t-1}) \leq \tilde{s}_{m,t+L_m}(H_{t+L_m-1}) =$

$s_{m,t+L_m}(H_{t+L_m-1})$. Then by Proposition 3.5.6, $\tilde{s}_{m+1,t}(H_{t-1})$ is optimal.

Numerical Results for Section 3.3

In Tables 3.7 and 3.7, we illustrate for a simple example with a single stage the differences in computational effort between a traditional dynamic programming approach and the decomposition method. The results are based on a correlation model where the demand in period t depends on the demand in the three previous periods as follows:

$D_t = \lfloor \rho_1 d_{t-3} + \rho_2 d_{t-2} + \rho_3 d_{t-1} \rfloor + \epsilon_t$, d_t is the realized demand in period t , $\rho_1 = .3, \rho_2 =$

$0.15, \rho_3 = 0.1$, and ϵ_t is an independent random variable with a discrete uniform distribution.

For the results shown in Table 3.7, ϵ_t is uniformly distributed over the values

$\{10, 15, 20, 25, 30\}$. For the results shown in Table 3.7, the planning horizon is fixed at

$T = 10$ but the distribution of ϵ_t is varied by considering uniform distributions over the

values $\{10, 15, 20, 25, 30\}, \{11, 16, 21, 26, 31\}, \dots, \{22, 27, 32, 37, 42, 47\}$. These distributions

which we index from 10 to 22 based on their minimum support represent distributions where

demand values are higher. Therefore, the associated state-space under the traditional dynamic

programming approach is larger. The results shown were obtained using a work

station with the following specifications: 4 x (real) Intel(R) Xeon(TM) CPU 3.73GHz, 16

Gb memory, and Linux operating system.

As we can see, significant savings in computational effort can be achieved using the decomposition method. In all cases, the computational effort remains modest under the decomposition method. Moreover, while it continues to be possible to obtain solutions for larger problems (problems with either a larger number of periods or higher demand values) using the decomposition methods, this is not possible (due to memory limitations) using the traditional approach (those cases are indicated by N/A). The difference in performance between the two approaches increases with correlation models involving correlation over more periods or correlation over the entire planning horizon. Solution of such problems is, in most cases, impossible under the traditional approach because of excessive memory requirements (for brevity corresponding numerical results are not shown but can be obtained from the authors upon request). However, these problems continue to be solvable with modest computational times using the decomposition methods.

Table 3.1: Solution times (in seconds) for varying planning horizons

Number of Periods	10	15	20	25	30	35	40
Traditional DP approach	719	1170	1533	1898	2256	N/A	N/A
Decomposition method	50	111	162	226	282	339	397

Table 3.2: Solution times (in seconds) for varying demand distributions

Minimum Support	10	11	12	13	14	15	16	17	18	19	20	21	22
Traditional DP approach	719	1333	1499	1678	2080	2312	2571	2843	3134	3453	3793	N/A	N/A
Decomposition method	50	51	52	53	54	54	55	56	57	58	60	62	63

Chapter 4

Future Research and Other Completed and Ongoing Research Projects

In this chapter, we briefly discuss possible future research directions related to Chapters 2 and 3. We also discuss briefly other completed and ongoing research.

4.1 Possible Extensions to Capacity Pooling among Independent Firms

There are numerous possible avenues for future research related to the topic described in Chapter 2. We plan to study systems with alternative assumptions about customer workloads and customer service priorities. Although exact analysis might be difficult, we may be able to take advantage of several effective approximations. We also intend to study more complex queueing systems, such as systems with parallel servers, servers in series, or servers

in a general network configuration, and also systems with alternative cost structures. Most such systems have no closed form expressions for performance measures of interest, but it may still be possible to obtain ordinal results leading to non-emptiness of the core in some cases and corresponding allocation rules. In other cases, we may be able to identify conditions under which pooling is not beneficial, which may serve to guide the formation of stable partial pooling arrangements.

In our current treatment, we assumed that cooperation among independent firms, in the form of capacity sharing, affects only their costs. In practice, cooperation often takes place between competing firms (e.g., cooperation among airlines, hospitals, or original equipment manufacturers in the same industry). In that case, cooperation, by reducing the firms' costs, could affect the intensity of the competition. More importantly, if cooperation is more beneficial to one firm than to others, it may affect their market share. In turn, this could affect the desirability of cooperation. Clearly, there is a need for models that capture both competition and cooperation and incorporate the effect of competition in the design of allocation rules.

We are also interested in the applications of other notions of cooperation and coalition stability, including stability over time. Concepts such as farsighted stability (see for example Chwe 1994) might be useful to investigate for our capacity sharing setting. It would be of interest to investigate settings in which one firm is responsible for choosing the capacity level and examined the impact of the fraction of capacity cost incurred by this lead firm on overall capacity investment. In practice, the cost portion paid by the lead firm is likely to be affected by the bargaining power of that firm. Therefore, it would be useful to construct models that explicitly capture the relative bargaining power of the different firms and the role bargaining power plays in cost allocation. This might also be useful for studying settings where all firms outsource to a common (third party) supplier. The supplier still decides on

how much capacity to invest and on the fees to charge each firm.

Finally, we are also interested in the applications of other notions of cooperation and coalition stability, including stability over time. Concepts such as farsighted stability (see for example Chwe 1994) might be useful to investigate for our capacity sharing setting.

4.2 Possible Extensions of the Customer-Item Decomposition Approach to Other Stochastic Inventory Systems

There are various inventory problems for which the decomposition approach might be possible. This includes complex systems, such as assembly and distribution systems, which have remained largely intractable or whose optimal solution is unknown. For these systems, it would be particularly of interest to identify conditions under which the myopic policy, or other simple policies, are optimal. It would also be of interest to investigate the applicability of the decomposition approach to systems with capacity limits. Here it may be possible to decompose the problem into simpler subproblems involving the notion of *unit capacity*, the amount of capacity involved in the production or supply of a single unit. The approach might also have applicability beyond traditional inventory systems, to systems where the amount of inventory is a fixed amount and must be allocated to random future demand arising from multiple customer classes, as in revenue management.

Moreover, it would be of interest to investigate how the customer-item decomposition approach can be used to construct effective heuristics for complex systems where either characterizing or implementing the optimal policy is difficult. This may include such as inventory systems with lost sales, multiple demand classes, etc. The purpose of such research is to find efficient heuristic algorithms with certain performance guarantees.

4.3 Other Completed and Ongoing Research Projects

In addition to the research presented in Chapter 2 and Chapter 3, we have completed, or nearly completed, several other research projects. The following is a brief summary of some of these projects.

4.3.1 Inventory Systems with Concave Ordering Costs

In this research project (Yu and Benjaafar 2008a) we analyze inventory systems with concave ordering costs. Concave ordering costs are common in practice because of economies of scale effects, the possibility of quantity discounts and the potential availability of multiple suppliers. However, the existing research on inventory systems with concave ordering costs is fairly limited. In this paper, we introduce a local monotone property for the optimal policy and a new generalized convexity property to show that the structure of the optimal policy for systems with piecewise linear concave ordering costs is a generalized (s, S) policy. The policy in each period is specified by parameters $s_n \leq \dots \leq s_1 \leq S_1 \leq \dots \leq S_n$ such that it is optimal to order up to S_i if the starting inventory level is less than s_i and more than s_{i+1} for $i = 1, \dots, n-1$, and it is optimal to order up to S_n if the starting inventory level is less than s_n , otherwise, it is optimal to order nothing. We carry out the analysis under quite general assumptions on ordering costs and demand distributions. This is in contrast with the limited existing literature where known results are obtained only for certain classes of demand distributions. In fact, it is an open problem whether this generalized (s, S) is optimal for arbitrary demand distribution. Possible future research avenues include problems with joint pricing and inventory control and problems with advanced demand information. I am also interested in analyzing other complicated inventory control models.

4.3.2 Optimal Incentives in Supply Chain with Asymmetric Information

In this research project (Yu and Benjaafar 2008c) we analyze the mechanism design problem in a principal-agent setting where agents have semi-independent utilities (the utility of an agent is not affected by the decisions of other agents). The principal has its own utility that depends on the decision strategies of all agents, and the principal wishes to maximize the sum of all utilities. We allow the principal and agents to have private information. Our objective is to design a mechanism such that when individuals interact through the mechanism, they have incentives to report their true private information that leads to a socially desirable result. We show that if we allow the transfer function (or contract) for an agent to depend not only the messages the principal receives but also on the decision of the agent itself, then we can design a class of dominant strategy incentive-compatible schemes. Our mechanism is an extension of the mechanism provided by Groves (1973), where he analyzes the case where the principal does not have its own utility and the agents' utilities are semi-independent. We apply the results to the optimal procurement for supply chains under asymmetric information, where retailers are replenished by a procurement center and retailers have private information concerning local market conditions or demand. Retailers may have their own pricing decisions on the products received from the procurement center. One of such examples is the capacity allocation problem among a car manufacturer and its subsidized car dealerships. Car dealerships usually have private information concerning the market conditions in their covering regions and have the ability to price cars. In order to optimally allocate the limited capacity to different car dealerships, the car manufacturer can use subsidized contract to induce car dealerships to give their true market conditions.

4.3.3 A Marginal Revenue Approach to Airline Seats Allocation Problems

In this research project (Yu and Benjaafar 2008b) we consider a new marginal revenue approach to the airline seat allocation problem with nested multiple fare classes, where our objective is to set the optimal protection levels for different fare classes to maximize the expected total revenue. By using this approach, we are able to obtain the explicit solution of the optimal protection level for each fare class under more general assumptions than the existing literature. It might be possible to extend this approach to analyze more complex systems, such as *network* revenue management problems.

4.3.4 Optimal Control of Make-to-Stock Queues

Make-to-stock queues are important operational models used in the design, planning and analysis of production and inventory systems. We are currently using such models to study assemble-to-order (ATO) systems with multiple products and varying components. We have been so far been able to characterize several properties of the optimal policy. These results are among the first to describe the structure of the optimal policy for such complex class of problems. We continue to investigate additional properties and to consider several special cases for which we may be able to fully characterize the optimal policy. .

Bibliography

O. Z. Aksin, F. Karaesmen, and E. L. Ormeci, “A review of workforce cross-training in call centers from an operations management perspective,” in *Workforce Cross Training Handbook*, D. Nembhard (editors), CRC Press, 211-240, 2007.

O.Z. Aksin, F. de Vericourt, and F. Karaesmen, “Call center outsourcing contract design and choice,” *Management Science*, **54**, 354-368 , 2008.

S. L. Albin, “Approximating a point process by a renewal process, II: superposition arrival processes to queues,” *Operations Research*, **32**, 1133-1162, 1984.

G. Allon and A. Federgruen, “Competition in service industries,” *Operations Research*, **55**, 37-55, 2007.

G. Allon and A. Federgruen, “ Service competition with general queueing facilities,” *Operations Research*, **56**, 827-849, 2008.

S. Anily and M. Haviv, “Cooperation in service systems,” working paper, Tel Aviv University, 2008.

- S. Anily and M. Haviv, “Cost-allocation for the first order interaction joint replenishment model,” *Operations Research*, **55**, 292-302, 2007.
- S. Axsäter, “Simple solution procedures for a class of two echelon inventory problem,” *Operations Research*, **38**, 64-69, 1990.
- S. Axsäter, “Exact and approximate evaluation of batch-ordering policies for two level inventory systems,” *Operations Research*, **41**, 777-785, 1993.
- A. Bassombo, R. S. Randhawa, and J. A. van Mieghem, “A little flexibility is all you need: optimality of tailored chaining and pairing,” working paper, Northwestern University, 2008.
- S. Benjaafar, “Performance bounds for the effectiveness of pooling in multi-processing systems,” *European Journal of Operational Research*, **87**, 375-388, 1995.
- S. Benjaafar, W. L. Cooper and J. S. Kim, “On the benefits of pooling in production-inventory systems,” *Management Science*, **51**, 548-565, 2005.
- S. Benjaafar, E. Elahi and K. Donohue, “Outsourcing via service quality competition,” *Management Science*, **53**, 241-259, 2007.
- N. Ben-Zvi and Y. Gerchak, “Inventory centralization when shortage costs differ: priorities and costs allocation,” working paper, Tel-Aviv University, 2005.

- J. A. Buzacott, "Commonalities in reengineered business processes: models and issues," *Management Science*, **42**, 768-782, 1996.
- G. Cachon and P. Harker, "Competition and outsourcing with scale economies," *Management Science*, **48**, 1314-1333, 2002.
- G. Cachon and F. Zhang, "Obtaining fast service in a queuing system via performance-based allocation of demand," *Management Science*, **53**, 408-420, 2007.
- F. Chen, "Optimal policies for multi-echelon inventory problems with batch ordering," *Operations Research*, **48**, 376-389, 2000.
- F. Chen and J.-S. Song, "Optimal policies for multi-echelon inventory problems with Markov modulated demand," *Operations Research*, **49**, 226-234, 2001.
- X. Chen and J. Zhang, "Duality approaches to economic lot sizing games," working paper, New York University, 2006.
- X. Chen and J. Zhang, "A stochastic programming duality approach to inventory centralization games," working paper, New York University, 2007.
- F. Chen and Y.-S. Zheng, "Lower bounds for multi-echelon stochastic inventory systems," *Management Science*, **40**, 1426-1443, 1994.

Y.S. Chow, H. Robbins and D. Siegmund, *Great expectations: the theory of optimal stopping*, Houghton Mifflin Company, Boston, 1971.

M.S. Chwe, “Farsighted Coalitional Stability,” *Journal of Economic Theory*, **63**, 299-325, 1994.

de Albeniz and Lago, “Myopic inventory policies using individual customer arrival information”, working paper, IESE Business School of the University of Navarra, 2007.

S. Dewan and H. Mendelson, “User delay costs and internal pricing for a service facility,” *Management Science*, **36**, 1502-1517, 1990.

L. Dong and H. Lee, “Optimal policies and approximations for a serial multiechelon inventory system with time-correlated demand,” *Operations Research*, **51**, 969-980, 2003.

M. Dror and B. Hartman, “Shipment consolidation: who pays for it and how much?” *Management Science*, **53**, 7887, 2007.

N. Erkip, W.H. Hausman and S. Nahmias, “Optimal centralized ordering policies in multi-echelon inventory systems with correlated demands,” *Management Science*, **36**, 381-392, 1990.

G. Gallego and O. Ozer, “Integrating replenishment decision with advance demand information,” *Management Science*, **47**, 1344-1360, 2001.

G. Gallego and P. Zipkin, "Stocking positioning and performance estimation in serial production-transportation systems," *Manufacturing & Service Operations Management*, **1**, 77-88, 1999.

N. Gans and Y-P. Zhou, "A call-routing problem with service-level constraints," *Operations Research*, **51**, 255-271, 2003.

N. Gans and Y-P. Zhou. "Call-routing schemes for call-center outsourcing," to appear in *Manufacturing and Service Operations Management*, **9**, 33-50, 2007.

M. D. Garcia-Sanz, F. R. Fernandez, M. G. Fiestras-Janeiro, I. Garcia-Jurado and J. Puerto, "Cooperation in markovian queueing models," to appear in *European Journal of Operational Research*, 2007.

J. P. Gayon, S. Benjaafar and F. de Vericourt, "Using imperfect demand information in production-inventory systems with multiple demand classes," *Manufacturing and Service Operations Management*, **11**, 128-143, 2009.

P. Gonzalez and C. Herrero, "Optimal sharing of surgical costs in the presence of queues," *Mathematical Methods of Operations Research*, **59**, 435-446, 2004.

S.C. Graves, "A single-item inventory model for a nonstationary demand process," *Manufacturing & Service Operations Management*, **1**, 50-61, 1999.

S.C. Graves, H. C. Meal, S. Dasu, and Y. Qiu, "Two-stage production planning in a dynamic environment," S. Axsäter, C. Schneeweiss, E. Silver, eds. *Multi-Stage Production Planning and Control. Lecture Notes in Economics and Mathematical Systems*, Springer-Verlag, Berlin, Germany, **266**, 9-43, 1986.

S.C. Graves, D. B. Kletter, and W. B. Hetzel, "A dynamic model for requirements planning with application to supply chain optimization," *Operations Research*, **46**, S35-S49, 1998.

T. Groves, "Incentives in teams," *Econometrica*, **41**, 617-633, 1973.

T. Groves, "Incentive compatible control of decentralized organizations," in *Directions in large-scale systems many-person optimization and decentralized control*, ed. by Y. C. Ho and S. K. Mitter, New York, 149-185, 1976.

T. Groves and M. Loeb, "Incentives in a divisionalized firm," *Management Science*, **25**, 221-230, 1979.

S. Gurusurthi and S. Benjaafar, "Modeling and analysis of flexible queueing systems," *Naval Research Logistics*, **51**, 755-782, 2004.

A. Y. Ha, "Optimal pricing that coordinates queues with customer-chosen service requirements," *Management Science*, **47**, 915-930, 2001.

A. Y. Ha, "Incentive-compatible pricing for a service facility with joint production and congestion externalities," *Management Science*, **44**, 1623-1636, 1998.

E. Hanany and Y. Gerchak, "Nash Bargaining over inventory and pooling contracts," forthcoming in *Naval Research Logistics*, 2008.

J. M. Harrison, *Brownian Motion and Stochastic Flow Systems*, John Wiley, New York, New York, 1985.

R. Hassin and M. Haviv, *To Queue or not to Queue*, Kluwer, Boston, 2003.

D.C. Heath, P. L. Jackson, "Modeling the evolution of demand forecasts with application to safety stock analysis in production/distribution systems," *IIE Transactions*, **26**, 17-30, 1994.

W. J. Hopp, E. Tekin and M. P. Van Oyen, "Benefits of skill chaining in production lines with cross-trained workers," *Management Science*, **50**, 83-98, 2004.

E. Ignall, A. Veinott, "Optimality of myopic inventory policies for several substitute products," *Management Science*, **15**, 284-304, 1969.

T. Iida, P. Zipkin, "Approximation solutions of a dynamic forecast-inventory model," *Manufacturing & Service Operations Management*, **8**, 407-425, 2006.

S. M. Iravani, M. P. Van Oyen and K.T. Sims (2005). "Structural flexibility: a new perspective on the design of manufacturing and service operations," *Management Science*, **51**, 151-166.

- G. Janakiraman, J. Muckstadt , “A decomposition approach for a class of capacitated serial systems,” working paper, New York University, 2005.
- G. D. Johnson and H. E. Thompson, “Optimality of myopic inventory policies for certain dependent demand processes,” *Management Science*, **21**, 1303-1307, 1975.
- W. C. Jordan, R.R. Inman and D. E. Blumenfeld, “Chained cross-training of workers for robust performance,” *IIE Transactions*, **36**, 953-967, 2004.
- O. Jouini, Y. Dallery and R. Nait-Abdallah, “Analysis of the impact of team-based organizations in call center management,” *Management Science*, 400-414, 2008.
- J. Kahn, “Inventory and the volatility of production,” *American Economic Review*, **77**, 667-679, 1987.
- E. Kalai, M. I. Kamien and M. Rubinovitch, “Optimal service speeds in a competitive environment,” *Management Science*, **38**, 1154-1163, 1992.
- S. Karlin, “Dynamic inventory policy with varying stochastic demands,” *Management Science*, **6**, 231-258, 1960.
- K. Katircioglu and D. Atkins, “New optimal policies for a unit demand inventory problem,” working paper, University of British Columbia, 1998.

- A. Katta and J. Sethuraman, "Cooperation in queues," working paper, Columbia University, 2006.
- E. Kemahlioglu-Ziya, "Formal methods of value sharing in supply chains", Ph.D Thesis, Georgia Institute of Technology, 2004.
- L. Kleinrock, *Queueing Systems, Computer Applications, Volume 2*, John Wiley & Sons, 1975.
- G. Koole and A. Pot, "An overview of routing and staffing in multi-skill customer contact centers," working paper, Vrije Universiteit, 2005.
- H. Lee, V. Padmanabhan, and S. Whang, "Information distortion in a supply chain: The bullwhip effect," *Management Science*, **43**, 546-558, 1997.
- H. Lee, K.C. So, and C. S. Tang, "The value of information sharing in a two-level supply chain," *Management Science*, **46**, 626-643, 1999.
- H. Lee and S. Whang, "Value of postponement, " *Product Variety Management*, edited by T. Ho and C. Tang, Kluwer Academic Publishers, 65-84, 1998.
- R. Levi, M. Pal, R. Roundy, and D. Shmoys, "Approximation algorithms for stochastic inventory control models," working paper, School of Operations Research and Industrial Engineering, Cornell University, 2005.

- X. Lu, J. Song, and A. Regan, "Inventory planning with forecast updates: approximate solutions and cost error bounds," *Operations Research*, **54**, 1079-1097, 2006.
- A. Mandelbaum, and M. I. Reiman, "On pooling in queueing networks," *Management Science*, **44**, 971-981, 1998.
- F. Maniquet, "A characterization of the Shapley Value in Queueing Problems", *Journal of Economic Theory*, **109**, 90-103, 2003.
- H. Mendelson and S. Whang, "Optimal incentive-compatible priority pricing for the M/M/1 queue," *Operations Research*, **38**, 870-883, 1990.
- H. Moulin, *Cooperative Microeconomics: a Game-Theoretic Introduction*, Princeton University Press, 1995.
- H. Moulin and R. Strong, "Fair queueing and other probabilistic allocation methods," *Mathematics of Operations Research*, **27**, 1-30, 2002.
- A. Muharremoglu and J. Tsitsiklis, "A single-unit decomposition approach to multi-echelon inventory systems," working paper, Graduate School of Business, Columbia University, 2008.
- A. Muller, M. Scarsini and M. Shaked, "The Newsvendor game has a nonempty core," *Games and Economic Behavior*, **38**, 118-126, 2002.

M. Nagarajan and G. Sošić, “Game-theoretic analysis of cooperation among supply chain agents: review and extensions,” to appear in *European Journal of Operational Research*, 2007.

H. Scarf, “Bayes solutions of the statistical inventory problem,” *Annals of Mathematical Statistics*, **30**, 490-508, 1959.

H. Scarf, “Some remarks on Bayes solution to inventory problem,” *Naval Research Logistics*, **7**, 591-596, 1960.

K. Shang and J.-S. Song, “Newsvendor bounds and heuristic for optimal policies in serial supply chains,” *Management Science*, **49**, 618-638, 2003.

M. Sheikhzadeh, S. Benjaafar and D. Gupta, “Machine sharing in manufacturing systems: flexibility versus chaining,” *International Journal of Flexible Manufacturing Systems*, **10**, 351-378, 1998.

D. R. Smith and W. Whitt, “Resource sharing for efficiency in traffic systems,” *The Bell System Technical Journal*, **60**, 1981.

J.-S. Song and P. Zipkin, “Inventory control in a fluctuating demand environment,” *Operations Research*, **41**, 351-370, 1993.

J.-S. Song and P. Zipkin, “Evaluation of base-stock policies in multiechelon inventory systems with state-dependent demands: Part I: State-independent policies,” *Naval Research*

Logistics, **39**, 715-728, 1996a.

J.-S. Song and P. Zipkin, "Evaluation of base-stock policies in multiechelon inventory systems with state-dependent demands: Part II: state-dependent depot policies," *Naval Research Logistics*, **43**, 381-396, 1996b.

S. Stidham, "On the optimality of single-server queueing systems," *Operations Research*, **18**, 708-732, 1970.

S. Stidham, "Pricing and capacity decisions for a service facility: stability and multiple local optima," *Management Science*, **38**, 1121-1139, 1992.

E. Tekin, W. Hopp and M. V. Oyen, "Pooling strategies for call center agent crosstraining," working paper, Northwestern University, 2004.

W. van den Heuvel, P.E.M. Borm and H. Hamers, "Economic lot-sizing games," *European Journal of Operational Research*, **176**, 1117-1130, 2007.

A. Veinott, "Optimal stockage policies with non-stationary stochastic demands," In H. Scarf, D. Gilford, and M. Shelly (eds.), *Multistage Inventory Models and Techniques*, Stanford University Press, 85-115, 1963.

A. Veinott, "Optimal policy for a single product, non-stationary inventory model with several demand classes," *Operations Research*, **13**, 761-778, 1965a.

A. Veinott, "Optimal policy for a multi-product, dynamic, non-stationary inventory problem," *Management Science*, **12**, 206-222, 1965b.

A. Veinott, "Optimal inventory policy for batch ordering," *Operations Research*, **13**, 424-432, 1965c.

R. B. Wallace and W. Whitt, "A staffing algorithm for call centers with skill-based routing," *Manufacturing and Service Operations Management*, **7**, 276-294, 2005.

T. Wang and B. Toktay, "Inventory management with advance demand information and flexible delivery," forthcoming in *Management Science*, 2007.

W. Whitt, "Approximating a point process by a renewal process, I: two basic methods," *Operations Research*, **30**, 125-147, 1982.

Y. Yu and S. Benjaafar, "A customer-item decomposition approach to stochastic inventory systems with correlation," working paper, University of Minnesota, 2006.

Y. Yu and S. Benjaafar, "An analysis of inventory systems with concave ordering costs," working paper, University of Minnesota, 2008a.

Y. Yu and S. Benjaafar, "A marginal revenue approach to an airline seat allocation problem with nested multiple fare classes," working paper, University of Minnesota, 2008b.

Y. Yu and S. Benjaafar, “Optimal incentives in supply chains with asymmetric information,” working paper, University of Minnesota, 2008c.

P. H. Zipkin, *Foundation of Inventory Management*, McGraw-Hill, 2000.