# MICHIGAN ROSS

# Joint Inventory and Fulfillment Decisions for Omnichannel Retail Networks

### Aravind Govindarajan
Stephen M. Ross School of Business
University of Michigan

### Amitabh Sinha
Amazon.com

### Joline Uichanco
Stephen M. Ross School of Business
University of Michigan

UNIVERSITY OF MICHIGAN

# Joint Inventory and Fulfillment Decisions for Omnichannel Retail Networks

Aravind Govindarajan

Stephen M. Ross School of Business, University of Michigan, Ann Arbor, MI 48109, arav@umich.edu

Amitabh Sinha

Amazon.com, Seattle, WA 98109, amitabhsi@amazon.com

Joline Uichanco

Stephen M. Ross School of Business, University of Michigan, Ann Arbor, MI 48109, jolineu@umich.edu

With e-commerce growing at a rapid pace compared to traditional retail, many brick-and-mortar firms are supporting their online growth through an integrated omnichannel approach. Such integration can lead to reduction in cost that can be achieved through efficient inventory management. A retailer with a network of physical stores and fulfillment centers facing two demands (online and in-store) has to make important, interlinked decisions – how much inventory to keep at each location and where to fulfill each online order from, as online demand can be fulfilled from any location. We consider order-up-to policies for a general multi-period model with multiple locations and zero lead time, and online orders fulfilled multiple times in each period. We first focus on the case where fulfillment decisions are made at the end of each period, which allows separate focus on the inventory decision. We develop a simple, scalable heuristic for the multi-location problem based on analysis from the two-store case, and prove its asymptotic near-optimality for large number of omnichannel stores under certain conditions. We extend this to the case where fulfillment is done multiple times within a period and combine it with a simple, threshold-based fulfillment policy which reserves inventory at stores for future in-store demand. With the help of a realistic numerical study based on a fictitious retail network embedded in mainland USA, we show that the combined heuristic outperforms a myopic, decentralized planning strategy under a variety of problem parameters, especially when there is an adequate mix of online and in-store demands. Extensions to positive lead times are discussed.

*Key words*: omnichannel; e-commerce; inventory management; fulfillment; heuristic; asymptotic analysis

## 1. Introduction

By the end of 2016, e-commerce sales accounted for around 9% of the total retail sales in the United States (U.S. Census Bureau 2016). Although this is a small portion of the total sales, online sales have been increasing at a rapid growth rate of around 15% each year (Zaroban 2016), and projected to account for 17% of all retail sales within the next five years (Lindner 2017). In comparison, the growth in traditional retail has dwindled to around 2% in recent years. With customers increasingly favoring the online channel, traditional brick-and-mortar (B&M) firms are compelled to develop their e-commerce capabilities to remain competitive against pure play e-commerce firms like Amazon (Leiser

2016), which alone accounted for 53% of the e-commerce sales growth in 2016 (Kim 2017). In order to improve efficiency and flexibility, retailers resort to an omnichannel approach to integrate the online channel with their physical stores.

Omnichannel refers to the seamless integration of a retailer's sales channels, such as in-store and online. Customers can purchase an item in different ways, including placing an order through the online store (websites), through mobile devices (mobile apps), as well as through the traditional practice of walking into physical stores. In addition, customers placing orders online can also choose how they receive the item, which has led to various omnichannel initiatives: they can pick up their items from a nearby physical store (in-store pickup) or from designated self-service kiosks like Amazon Lockers, or simply have the item shipped directly to their homes (ship-to-customer).

Providing an omnichannel customer experience is regarded as a brand differentiator by many retailers, and integrating the online channel with the physical stores increases revenue, reduces shipping costs and improves customer satisfaction (Forrester 2014). Hence, there is an industry-wide shift to omnichannel retailing, with onetime B&M firms like Macy's and Walmart leveraging their existing network of retail stores in their integration of the online channel (Nash 2015). Amazon has also joined these firms through the acquisition of a network of physical stores across the US by means of its purchase of Whole Foods Market. This allows Amazon to not only operate an omnichannel grocery chain, but also absorb the stores into its distribution network to reduce logistic costs.

One of the key aspects of this channel integration is *store fulfillment*, which is the use of physical stores to fulfill online orders. Store fulfillment has now become indispensable for firms like Walmart and Macy's, that rely on a network of physical stores close to population centers to offer same day and next-day delivery options to customers (Giannopoulos 2014). Dedicated floor space and store staff are required to fulfill online orders from stores.

In spite of potential benefits, many firms have struggled in their implementation of channel integration: from 2010 to 2014, even as retail and online sales increased, inventory turnover decreased (Kurt Salmon 2016). One possible cause for this inefficiency could be insufficient planning in inventory management. While firms have traditionally managed inventory levels at stores based on demands in the corresponding locations, such a decentralized approach ceases be optimal in an integrated system.

The optimal inventory decisions depends on the fulfillment policy followed, and there does not seem to be a standard approach to online fulfillment across the industry. Some firms primarily fulfill from online FCs, and resort to store fulfillment in case the online FC runs out of stock. Some firms fulfill online orders from stores, but are agnostic to store inventory levels, while others do not fulfill from stores running low on inventory.

In this paper, we study the problem of an omnichannel firm with a network of physical stores and online FCs facing online (ship-to-customer) and in-store demands, by means of a general multi-period, multi-location model. We consider a dynamic setting, where we allow the online fulfillment decisions are made multiple times within each period. Online orders can be routed to any store or online FC in the network, and items are picked off the shelves, packed, labeled and shipped to the customers' homes. This has several advantages over the dedicated use of online FCs including reduced shipping costs, quicker deliveries and efficient use of store inventory (UPS Compass 2014).

Our goal is to *optimize inventory levels and fulfillment decisions* for a single product. The decisions have to be made based on the network as a whole as opposed to a decentralized approach, in order to take into account demand pooling of online demands across the network, in addition to demand pooling of in-store and online demands in each region.

The firms's problem is described as follows. A retail firm owns a network of stores and online FCs, and has integrated the online channel into the physical stores through store fulfillment. Following a periodic review inventory model, each store orders up to a certain level at the beginning of each review period, to fulfill in-store demand (customers walking into physical stores) and online demand (customers ordering online, expecting items to be shipped directly to them) during the course of the period. The in-store demand at a store is fulfilled as it arrives, until that store runs out of inventory.

Unlike in-store demand, online demand can be fulfilled from any location in the network, and there is typically a delay between the time an order is placed and when items are picked off the shelf. Firms may delay fulfillment decisions due to various reasons:

- for strategic reasons, orders from the same customer or region can be consolidated to lower shipping costs (Xu, Allgor, and Graves 2009, Wei, Jasin, and Kapuscinski 2017),
- or for practical reasons, as the timing of orders fulfilled from stores is affected by store staffing schedules and pick-up times of third-party carriers like UPS.

To model this dynamic, a review period is further divided into $T$ fulfillment epochs, where in-store demands are fulfilled as they arrive, and online fulfillment decisions (assigning online orders to fulfillment locations) are made at the end of each epoch after observing the demands during the epoch, with unmet demands being lost. The inventory and fulfillment decisions are made centrally by the firm to minimize holding, penalty and shipping costs. For the sake of clarity, the two units of time are described below:

- a *review period* is the amount of time between two consecutive inventory replenishments. For stores that are replenished daily, the review period is a single day.

- a *fulfillment epoch* is the time between two fulfillment decisions. Over the course of an epoch, online orders are aggregated, and fulfillment decisions are made at the end of each epoch. For stores replenished daily, the length of an epoch can range from a whole day (e.g. Macy's stores fulfill online orders once a day through UPS (Lewis 2013)) to a few minutes (e.g. firms like Amazon make more frequent fulfillment decisions).

As described, the definition of a fulfillment epoch carries flexibility, and by choosing large enough values for $T$, we can closely approximate the continuous time setting, where firms make fulfillment decisions as online orders arrive.

The online fulfillment decisions are similar to transshipment decisions for online demand, except that instead of items being shipped between stores, they are shipped directly to the customer. The setting can thus be cast as planning of order-up-to levels in a transshipment problem with a replenishment leadtime of $T-1$ periods, with a planning horizon of $T$ periods. This makes the problem hard, as it has been shown that optimal transshipment decisions are intractable, let alone joint optimization of initial inventory levels and transshipment decisions, even for two locations (Tagaras and Cohen 1992).

The general structure of the problem is also subject to complications from other sources - multiple locations, multiple fulfillment epochs, and two non-identical classes of demands. Our main contribution is a combined inventory and fulfillment heuristic for omnichannel retailing, which we derive from a general multi-location, multi-period model shown to be mathematically intractable due to the various generalizations involved. Specifically, the inventory heuristic calculates the stocking levels at each location based on the demands in the network, rather than individually at that location, and the fulfillment heuristic provides location-specific, time-varying inventory thresholds which dictate the rationing between in-store and online demands. The strength of our combined heuristic lies in the

ease of computation and comprehension, and we show by means of a realistic numerical study that our heuristic creates value by planning for virtual pooling of online demands across locations, and diligently reserving inventory at stores for future demands.

The approach we take to address this problem is as follows. We model the general problem in Section 3, and describe the complexities involved. To obtain a heuristic solution to this problem, we first decouple the inventory and fulfillment decisions by considering the case with a single fulfillment epoch ($T = 1$) in Section 4. When there is no leadtime, a myopic fulfillment policy would be optimal in this case - fulfill online demand as much as possible with the available inventory in each review period. Given this fulfillment policy, we discuss the optimal inventory solution for the two-store case, and develop a simple, asymptotically near-optimal inventory heuristic for the multi-location case.

In Section 5, we extend this inventory heuristic to the general problem where online orders are fulfilled multiple times within each review period ($T > 1$), and develop a simple threshold fulfillment policy in each fulfillment epoch, where stores fulfill online orders only when the inventory levels are above a certain threshold.

In Section 6, by means of a realistic numerical study on a network of stores and online FCs embedded in mainland USA, we show that our combined inventory and fulfillment heuristic improves greatly upon a benchmark solution which naively sets inventory levels in a decentralized fashion and fulfills online orders myopically. We test the relative performance of our heuristic over a variety of problem parameters such as shipping costs, online market share, network size, etc. Finally, we conclude with Section 7 by discussing further generalizations including non-identical leadtimes and costs, and areas for future research.

## 2. Literature Review

Omnichannel retailing is a relatively new area in operations management literature, and has been gaining traction in recent years. Readers are referred to Rigby (2011) and Brynjolfsson et al. (2013) for comprehensive reviews of the topic. Existing papers in this area focus on the impact of online channel integration: Gao and Su (2017) study the impact of implementing store pickup on store operations, Bell et al. (2013), Ansari et al. (2008), and Gallino and Moreno (2014) study customer migration due to product information, and Gallino et al. (2017) focus on sales dispersion from implementing store pickup. Gao and Su (2016) analyze the effect of information provided to strategic omnichannel customers on store operations, and Harsha et al. (2016) study the dynamic pricing of omnichannel inventories.

When there is no in-store demand, the problem is analogous to the pure play e-commerce setting, which has enjoyed recent attention in literature: Acimovic and Graves (2017) study the optimal allocation of replenishment to fulfillment centers to reduce shipping costs and mitigate costly spillovers, Lei et al. (2017) consider the joint pricing and fulfillment strategy to maximize the expected profits (revenue minus shipping costs), and Acimovic and Graves (2014) focus on fulfillment strategies to minimize outbound shipping costs.

There have been some studies which discuss integration of online demand to physical stores by means of a separate online fulfillment center, as this was the primary mode of fulfillment in the e-commerce channel in its nascent stages. Seifert et al. (2006) consider the inventory management of a system where an online warehouse handles online orders, and in case of stockouts, stores can fill these orders. Chen et al. (2011) consider a three location system consisting of two stores and an etailer, with a hierarchy to fulfillment - the etailer can fulfill online orders with the least cost, followed by store 1 and then store 2.

We consider a generalized setting representing the current retailing situation wherein physical stores are the primary ports of online fulfillment. To the best of our knowledge, the study closest to ours in emulating the problem setting, where online demand is integrated with the physical stores through store fulfillment is by Jalilipour Alishah et al. (2015). They consider a single store with online and in-store demands, and analyze decisions at three levels — fulfillment structure, inventory optimization and inventory rationing. They show that the optimal rationing policy between in-store and online demands is threshold-based, but their results do not extend to the multi-store case due to the complexity involved in an additional rationing decision - online orders from other regions. This setting is rather important in the context of e-commerce, and falls under the purview of transshipment literature, where it has been shown to be an intractable problem to solve.

The fact that online demands can be fulfilled from any store in the system is analogous to a reactive transshipment setting with zero transshipment lead time, as pointed out by Yang and Qin (2007), who called this 'virtual lateral transshipment'. In addition, our problem has multiple demand classes (online and in-store), where some classes of demand (in-store) cannot be subject to transshipment. For an extensive review of transshipment literature, the readers are referred to Paterson et al. (2011).

The fact the the problem in question can be related to transshipment literature offers little solace. Transshipment problems are infamously hard to solve, and analytical approaches

can be done only for simplified cases with zero replenishment and transshipment leadtimes and two locations (Tagaras 1989) or identical shipping costs across locations (Dong and Rudi 2004). Tagaras and Cohen (1992) show that when there is positive replenishment leadtime, the problem becomes intractable even for two locations, as obtaining the optimal transshipment policy is mathematically complex due to its interdependence on demands during the leadtime, on-hand inventory and in-transit inventory.

Obtaining optimal order-up-to policies are by extension intractable as well, as they need to be calculated based on the optimal transshipment policy. Yao et al. (2016) have recently considered the optimal joint initial stocking and transshipment decisions for the two-store case, where stocking is done once at the beginning of a selling season, and transshipment is done multiple times during the season. Their analysis is limited to two stores, as key mathematical properties like submodularity do not extend to multiple locations.

Due to the complexities involved, one cannot hope to obtain a tight and tractable bound for a problem of this stature, let alone finding the analytical optimal solutions. We will instead develop simple, tractable and scalable heuristics, which perform well compared to naive strategies in most cases, with help from techniques used in literature.

Finally in the zero leadtime case, when online demand is fulfilled only once at the end of each review period, we show that the problem is analogous to a newsvendor network, with virtual lateral transshipment as a 'discretionary policy' (van Mieghem and Rudi 2002). Newsvendor networks have been analyzed in great detail by van Mieghem and Rudi (2002) and van Mieghem (2003), building up from the multi-dimensional newsvendor models proposed by Harrison and van Mieghem (1999). However, as we shall show later, the canonical approach to optimizing inventory levels is difficult even for two stores due to the number of random demands involved, and is intractable for the multi-store case.

## 3. The General Problem - Model and Assumptions

Consider a system composed of a firm which owns $N$ facilities $R_1, R_2, \ldots, R_N$ in different customer regions, selling a single product. Considering multiple products introduces complex combinatorial features to the fulfillment problem as a multi-item order can be fulfilled in different ways (Jasin and Sinha 2015), which we disregard in our analysis to better study the interplay between inventory and fulfillment decisions. There are two classes of demand originating in each region $i$, modeled by non-negative and continuous random variables with well-behaved density functions.
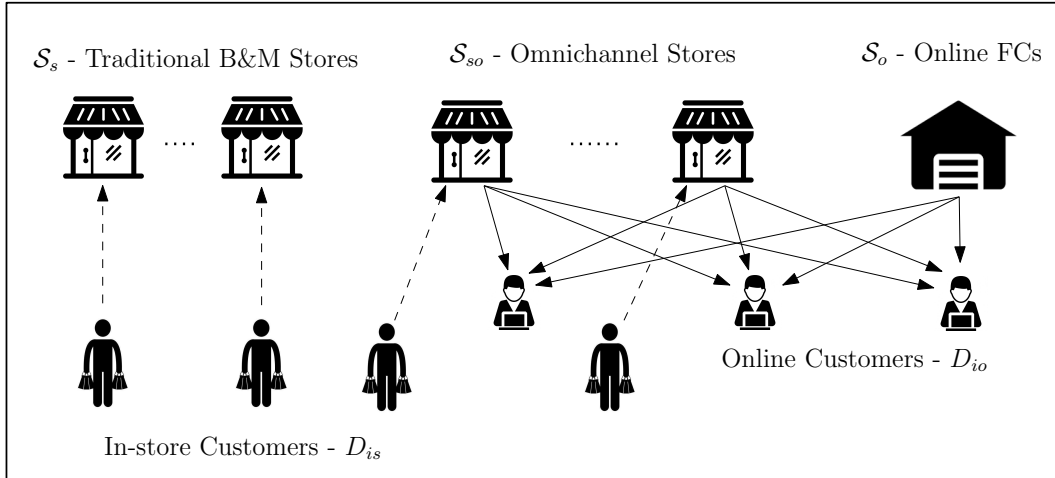
**Figure 1** **Three types of stores in the network a) Traditional brick-and-mortar stores ($\mathcal{S}_s$), b) Omnichannel stores ($\mathcal{S}_{so}$), and c) Online Fulfillment Centers ($\mathcal{S}_o$).**

1. the *in-store demand* ($D_{is}$) consists of customers picking items off the shelves (all the inventory is available on the shelf), with unmet demand lost immediately

2. the *online (ship-to-customer) demand* ($D_{io}$), consisting of customers ordering through the website or mobile app, with items delivered directly to their homes. For orders fulfilled from stores, the store staff pick up the item from the shelf, followed by packing and labeling in the store backroom, and shipping to the customer. A sale is lost when there is no available inventory for fulfillment at any location.

The demands are exogenous and are temporally independent, but can have any general channel or location correlation structure, while we require that the total demands in each region and across the system have continuous and well-defined density functions.

The retail network is shown in Figure 1, where dashed lines represent customers visiting physical stores and solid lines represent items shipped to customers' homes. We consider three different types of facilities described by the following sets:

- $\mathcal{S}_s$ - physical stores which handle only in-store demand.
- $\mathcal{S}_o$ - online fulfillment centers (OFCs) which handle only online orders.
- $\mathcal{S}_{so}$ - omnichannel physical stores which handle both online and in-store demands.

Since traditional B&M stores plan for inventory independent of other facilities in the network, we exclude them from our analysis. We are hence interested in locations involved in online fulfillment, namely the omnichannel stores and online fulfillment centers, denoted by the set of facilities $\mathcal{S} = \mathcal{S}_o \cup \mathcal{S}_{so}$, and the number of such stores is $N = \|\mathcal{S}\|$.

An important feature to be noted in the omnichannel problem is that unfulfilled in-store demand at one region cannot be fulfilled by stores in other regions. Any facility with available inventory can fulfill an online order, and hence there is pooling of online demands across regions in addition to pooling of in-store and online demands within each region.

### 3.1. Periodic Review Setup

We consider a periodic review model, where an order is placed by each facility at the start of each review period, and received with zero replenishment leadtime. The demands are realized during the course of the period based on the facility considered. We are interested in an optimal order-up-to policy, where the order-up-to levels in each period are $y_1, \ldots, y_N$.

Based on conversations with industry executives, there are certain situations in the context of omnichannel stores where the leadtime is effectively negligible: in major cities like New York, store replenishment can only be done at night-time due to traffic restrictions. Such stores handle high volumes of sales, and are usually replenished daily from warehouses in nearby cities. An order placed in the afternoon can often be replenished before the following day. Positive leadtimes can significantly complicate analyses, and we discuss extending our heuristics to the case of non-identical leadtimes across locations in Section 7.

We assume that *online orders are fulfilled in multiple batches in each review period*, which we model by dividing a review period into $T$ fulfillment epochs: in each epoch, in-store demand is fulfilled as it arrives, whereas online fulfillment decisions are made at the end of the epoch after observing demand, and orders are fulfilled with the available inventory.

The assumption reflects practical constraints in store operations: fulfillment activities in stores are usually done by store personnel, who in most cases also share additional store responsibilities. In such situations, it is better to fulfill online orders in batches, as opposed to having store staff picking items every time an online order is received.

### 3.2. Cost Parameters

We consider a per-unit service cost $s_{ij}$ for online demand from region $j$ fulfilled by $R_i$, which encapsulates the cost of picking the item off the shelf, packing and labelling, as well as the shipping cost for delivery. We have $s_{ij} > s_{ii}, \forall j \neq i$, as it is costlier to ship an item over longer distances. We will refer to the service costs $s_{ii}$ (within the same region) as *shipping* costs, and $s_{ij}$ (across regions) as *cross-shipping* costs.

In practice, the handling (pick-pack-and-label) component of the service cost is higher for stores fulfilling online demand, as it involves human labor, than for OFCs where the

process can be automated and streamlined. The shipping component of the service cost can be higher for the OFCs which are usually located farther away from population centers.

We have identical costs at each location, including shipping costs $s_{ii} = s$, $\forall i$. At the end of a fulfillment epoch, each unit of unused inventory incurs an overage cost $h$, and each unit of unfulfilled in-store and online demands incur penalty costs $p_s$ and $p_o$ respectively. We assume that $p_s > p_o - s > 0$, as in-store demand is fulfilled first and costlier to lose, and cross-shipping always leads to a myopic reduction in cost: $s_{ij} (= s_{ji}) < h + p_o$, $\forall i, j$. We ignore the purchasing cost of inventory, but this can be incorporated through linear terms.

### 3.3. Stochastic Programming Formulation

We are now ready to write the total expected per period cost function for the case where online demand is fulfilled over $T$ fulfillment epochs in each review period. We focus on the single period to obtain order-up-to levels, which we show in Section 5 to be optimal in a multi-period setting in the case of negligible replenishment leadtimes.

In each fulfillment epoch $t$, let the starting inventory levels be denoted by $\mathbf{x^t} = (x_i^t)_i$, and $\tilde{D}^t = (D_{is}^t, D_{io}^t)_i$ denotes the demands. From location $R_i$, let $z_i^t$ be the amount of inventory used to fulfill the in-store demand, and $Z_{ij}^t$ be the amount of inventory shipped to fulfill online demand from region $j$, denoted in vector form as $\mathbf{z^t}, \mathbf{Z^t}$ respectively. We have a $T$-stage stochastic program, with the cost-to-go function in epoch $t$, $C_t(\mathbf{x^t}, \tilde{D}^t)$ is given by:

$$C_t(\mathbf{x^t}, \tilde{D}^t) = \min_{\mathbf{z^t}, \mathbf{Z^t} \in \Delta} \left[ P(\mathbf{x^t}, \tilde{D}^t, \mathbf{z^t}, \mathbf{Z^t}) + \mathbb{E}C_{t+1}(x_i^t - z_i^t - \sum_{j=1}^{N} Z_{ij}^t, \tilde{D}^{t+1}) \right] \tag{1}$$

where $P(x^t, \tilde{D}^t, \mathbf{z^t}, \mathbf{Z^t})$ is the total cost in fulfillment epoch $t$, given by:

$$\begin{aligned}
P(x^t, \tilde{D}^t, \mathbf{z^t}, \mathbf{Z^t}) = \sum_{i=1}^{N} h \left( x_i^t - z_i^t - \sum_{j=1}^{N} Z_{ij}^t \right) + \sum_{i=1}^{N} p_s(D_{is}^t - z_i^t) \\
+ \sum_{j=1}^{N} p_o \left( D_{jo}^t - \sum_{i=1}^{N} Z_{ij}^t \right) + \sum_{i=1}^{N} s Z_{ii}^t + \sum_{i=1}^{N} \sum_{j=1, j \neq i}^{N} s_{ij} Z_{ij}^t
\end{aligned} \tag{2}$$

and $\Delta$ is the set of feasible fulfillment decisions, described by the following set of constraints:

$$\begin{aligned}
z_i^t + \sum_{j=1}^{n} Z_{ij}^t &\leq x_i^t, & \forall i \in [N], \forall t \in [T] \\
z_i^t &\leq D_{is}^t, & \forall i \in [N], \forall t \in [T] \\
\sum_{i=1}^{n} Z_{ij}^t &\leq D_{jo}^t, & \forall j \in [N], \forall t \in [T] \\
\mathbf{z^t}, \mathbf{Z^t} &\geq 0, & \forall t \in [T]
\end{aligned} \tag{3}$$

The first ienquality in $\Delta$ represents the supply constraint, and the second and third inequalities model the fulfillment constraints. Note that the online demand in one region can be fulfilled from any facility in the network, as seen in the third inequality in (3).

The goal is to obtain the initial stocking level $\mathbf{y} = (y_i)_i$. The single period, $T$-epoch problem can thus be stated as follows: $\min_{\mathbf{y} \geq \mathbf{0}} \mathbb{E}[C_1(\mathbf{y}, \tilde{D})]$. This is a convex minimization problem, as we will later show in Section 5, but it is intractable to solve. The fulfillment decisions are similar to optimal transshipment decisions with non-negligible lead time, as decisions in any fulfillment epoch depend on future demands in that review period. As pointed out by Tagaras and Cohen (1992) for the two-store case in traditional transshipment, while the optimal fulfillment policy may be threshold-based, the optimization becomes intractable due to the complexity of the decision space in the dynamic programming formulation.

We cannot hope to solve this problem to optimality, and we resort to heuristic solutions that perform well compared to simple, naive strategies and hindsight optimal lower bounds. Note that a heuristic solution specifies both the initial stocking level and fulfillment policy.

We first develop the inventory heuristic in the following way: treat the $T$-epoch problem as a single fulfillment epoch. A similar method was also used by Tagaras and Cohen (1992) to set heuristic inventory levels for the two-location transshipment problem with leadtime, based on numerical evidence that most transshipments took place at or near the end of the planning horizon, when stockouts are more likely to happen.

Our problem is different in two aspects: 1) we have in-store demands which are more costly to lose than online demands and do not have pooling flexibility, and 2) demands follow lost sales. However, we adopt this single fulfillment epoch approximation as it provides a tractable alternative by decoupling inventory and fulfillment decisions, because:

1. a myopic fulfillment policy is optimal, where online demands are fulfilled to the maximum possible extent with the available inventory, and as a result,

2. the inventory problem reduces to a single stage stochastic linear program.

With the help of results obtained through this approximation, we formulate inventory and fulfillment heuristic solutions for the multi-period, multi-location problem in Section 5, and numerically test their performance in Section 6.

## 4.    The Single Fulfillment Epoch Case (T=1) - Model and Analysis

In this setting, items are ordered and received at the beginning of the period with zero lead time, and in-store demand is fulfilled as it arrives. Due to the single fulfillment epoch
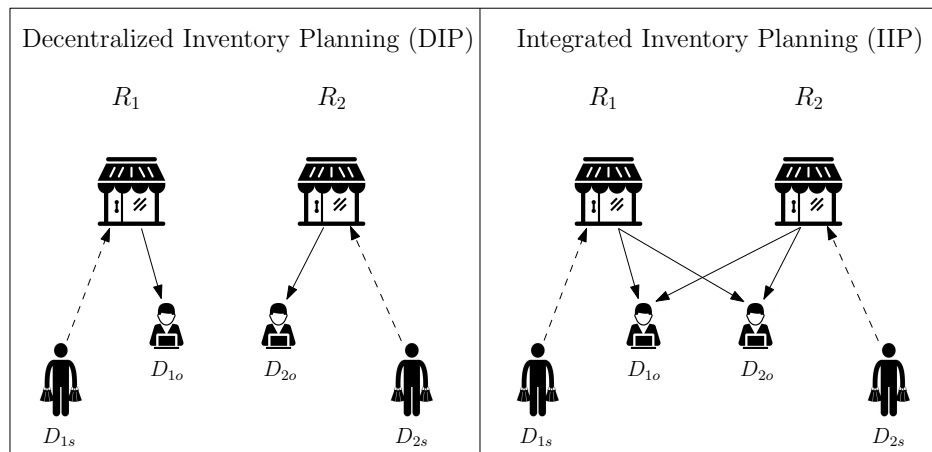
**Figure 2** **Two methods of inventory planning - 1) Decentralized inventory planning (DIP) and 2) Integrated inventory planning (IIP)**

assumption, the fulfillment of online demand is done once at the end of the review period, after in-store demands are fulfilled. There is no benefit to reserving inventory for future demands as replenishments arrive immediately. In such a case, a myopic fulfillment policy is optimal, where online orders are fulfilled to the maximum possible extent in each period.

The case of single fulfillment epoch is quite common in present day omnichannel retailing where stores are replenished daily. Most stores still rely on third party carriers such as UPS and FedEx to ship items to customers. Online orders to be shipped are loaded onto these trucks once a day from the store backroom, usually towards the end of the day. This is especially popular in the context of same-day and next-day deliveries, where stores allow online ordering until a cutoff time, and these orders are ready to be shipped by the end of the day. However with developments in drone technology in the future, one can easily envision stores that fulfill multiple times in a day, which we address through the general case of multiple fulfillment epochs ($T > 1$) in Section 5.

We first consider the two-store setting to exhibit the complicated nature of the decoupled inventory problem alone, given the optimal fulfillment policy is myopic. The insights derived in this case inform our analysis of a generalized multi-location case, which includes a network of omnichannel stores and online FCs.

### 4.1. The Two-store System

A firm owns two retail stores $R_1$ and $R_2$ serving different regions, with two demand streams originating form each region – in-store demand ($D_{1s}$, $D_{2s}$), and online demand ($D_{1o}$, $D_{2o}$). The objective is to set the initial inventory levels $y_1$ and $y_2$ to minimize the total expected

cost. We consider two solutions – decentralized inventory planning (DIP) and integrated inventory planning (IIP), which are represented in Figure 2. The assumptions on cost parameters are recapitulated in the set $\Psi$ in Equation 4.

$$\Psi = \left\{ p_s > p_o - s_i > 0, \ \forall i; \quad h + p_o > s_{ij} > s, \ \forall i, j \neq i \right\} \tag{4}$$

**4.1.1.   The Decentralized Inventory Planning (DIP) Strategy (Pooling within Regions)** We first consider the case where the firm plans for inventory at its stores in a decentralized fashion, without planning in advance for cross-shipping. This serves as a benchmark for any inventory heuristic we may develop for the centralized planning case. The inventory level at store $i$ is set with an objective to minimize the total expected cost incurred in meeting the demands from that region, given by:

$$
\begin{aligned}
C^{DIP}(y_i) = \mathbb{E}\Big[ h\left( (y_i - D_{is})^+ - D_{io}\right)^+ + p_s (D_{is} - y_i)^+ \\
+ p_o \left( D_{io} - (y_i - D_{is})^+ \right)^+ + s \min\left( (y_i - D_{is})^+, D_{io} \right) \Big]
\end{aligned}
\tag{5}
$$

where $x^+ = \max(x, 0)$. The cost function is convex, which can be seen by expressing Equation 5 in terms of the total demands $D_i = D_{is} + D_{io}$ as follows:

$$C^{DIP}(y_i) = s\mu_{io} + \mathbb{E}\Big[ h (y_i - D_i)^+ + (p_o - s)(D_i - y_i)^+ + (p_s - (p_o - s))(D_{is} - y_i)^+ \Big] \tag{6}$$

where $\mu_{io} = \mathbb{E}[D_{io}]$. The simplification is done using the identities $\min(x, y) = y - (y - x)^+$, and $(D_{is} - y_i)^+ + \left( D_{io} - (y_i - D_{is})^+ \right)^+ = (D_i - y_i)^+$, the latter holds when demands are non-negative. The optimal inventory levels $(y_1^{DIP}, y_2^{DIP})$ can obtained from implicit equations:

$$(h + p_o - s) F_i \left( y_i^{DIP} \right) + (p_s - p_o + s) F_{is} \left( y_i^{DIP} \right) = p_s, \qquad \forall i = 1, 2 \tag{7}$$

where $F_i$ is the cumulative distribution function of demand $D_i$. A line search yields unique optimum, as the left hand side is increasing in $y_i^{DIP}$, and the right hand side is constant.

**4.1.2.   The Integrated Inventory Planning (IIP) Strategy (Pooling within and across Regions).** This is similar to the DIP scenario, except that after $R_i$ has fulfilled its own in-store and online demands, unfulfilled online orders from region $j$ ($\neq i$) can be fulfilled using any available inventory at $R_i$. In the two-store problem, the cross-shipped quantity from store $R_i$ to region $j$ can be explicitly calculated as the minimum of the

inventory available at $R_i$ and the unfulfilled online demand at $R_j$, after each store has attempted to fulfill its own demands. The total expected one-period cost function is:

$$
\begin{aligned}
C^{IIP}(y_1, y_2) = \sum_i \mathbb{E}\Big[ & h\left((y_i - D_{is})^+ - D_{io}\right)^+ + p_s(D_{is} - y_i)^+ \\
& + p_o\left(D_{io} - (y_i - D_{is})^+\right)^+ + s \min\left((y_i - D_{is})^+, D_{io}\right) \\
& + (s_{12} - h - p_o)\min\left(\left((y_1 - D_{1s})^+ - D_{1o}\right)^+, \left(D_{2o} - (y_2 - D_{2s})^+\right)^+\right) \\
& + (s_{12} - h - p_o)\min\left(\left((y_2 - D_{2s})^+ - D_{2o}\right)^+, \left(D_{1o} - (y_1 - D_{1s})^+\right)^+\right)\Big]
\end{aligned}
\tag{8}
$$

The additional terms in Equation 8 that are absent in Equation 5 represent the value of cross-shipping: the total savings by cross-shipping a unit from $R_i$ to region $j$, $h + p_o - s_{ij}$, times the total quantity cross-shipped from $R_i$ to region $j$. The total cross-shipped quantity can be expressed as $\sum_i \left(D_{io} - (y_i - D_{is})^+\right)^+ - \left(\sum_i D_{io} - \sum_i (y_i - D_{is})^+\right)^+$. The first term represents the total unfulfilled online demand if there was no cross-shipping allowed, and the second term represents the unfulfilled online demand with cross-shipping. Naturally, the difference yields the cross-shipped quantity. By using this expression, as well as the simplification techniques used in Equation 6, we can simplify Equation 8 as follows:

$$
\begin{aligned}
C^{IIP}(y_1, y_2) = s\sum_i \mu_{io} + \sum_i \mathbb{E}\Big[ & h(y_i - D_i)^+ + (p_s - p_o + s)(D_{is} - y_i)^+ + (p_o - s)(D_i - y_i)^+\Big] \\
& + (s_{12} - h - p_o)\left[\sum_i (D_i - y_i)^+ - \sum_i (D_{is} - y_i)^+ - \left(\sum_i D_{io} - \sum_i (y_i - D_{is})^+\right)^+\right]
\end{aligned}
\tag{9}
$$

We can rearrange the terms to a convex expression, except $\left(\sum_i D_{io} - \sum_i (y_i - D_{is})^+\right)^+$, which is non-convex in $y_i$'s. This is seen by keeping $y_1$ constant and changing $y_2$.

$$
\begin{aligned}
& \left(\sum_i D_{io} - \sum_i (y_i - D_{is})^+\right)^+ \\
& = \begin{cases}
\left(D_{1o} + D_{2o} - (y_1 - D_{1s})^+\right)^+, & \text{if} & y_2 \leq D_{2s} \\
D_2 + D_{1o} - (y_1 - D_{1s})^+ - y_2, & \text{if} & D_{2s} < y_2 < D_2 + D_{1o} - (y_1 - D_{1s})^+ \\
0, & \text{if} & y_2 \geq D_2 + D_{1o} - (y_1 - D_{1s})^+
\end{cases}
\end{aligned}
\tag{10}
$$

In the event that $D_{is} = 0, \forall i$ (similar to traditional transshipment considered by Dong and Rudi 2004), the formulation in Equation 9 would directly yield a convex cost function. Convexity is not obvious in our case, as the nested piecewise linear function in Equation 10 is neither convex nor concave, and this is purely due to the fact that in-store demand

**Govindarajan, Sinha and Uichanco:** *Inventory and Fulfillment Decisions for Omnichannel Retailing*
Stephen M. Ross School of Business, University of Michigan 2018

15

is fulfilled first and cannot be subject to cross-shipment. However, the total cost can be shown to be jointly convex in the inventory levels (Proposition 1):

PROPOSITION 1. *Under the conditions on cost parameters in $\Psi$,*

(a) $C^{IIP}(y_1, y_2)$ *is jointly convex in the order-up-to levels.*

(b) *There exist regions $\Omega_k(y_1, y_2)$ in the demand space, such that in each region the dual-price vector $\lambda^k$ corresponding to the variables $y_1$, $y_2$ remains constant, and the gradient of the IIP cost function can be written as*

$$\nabla C^{IIP}(y_1, y_2) = (h, h)^{\mathsf{T}} - \sum_k \lambda^k \mathbb{P}\left(\Omega_k\left(y_1, y_2\right)\right)) \tag{11}$$

All proofs are relegated to the Appendix. We first observe that under the assumptions in $\Psi$, $C^{IIP}$ can be expressed as the expectation of a linear program, through which joint convexity in inventory levels is established. By noting structural similarities with a newsvendor network (van Mieghem and Rudi 2002), we derive an expression for the gradient based on the dual prices $\lambda = (\lambda_1, \lambda_2)^{\mathsf{T}}$, which are simply the shadow prices of the constraints involving $y_1$ and $y_2$ in the linear program representation (Equation 20, Appendix A).

The demands are shown to be separable into independent regions $\Omega_k$ based on the values of $y_1$ and $y_2$, within which the dual prices $\lambda^k = (\lambda_1^k, \lambda_2^k)$ are constant (refer to Appendix B for a detailed discussion), which enables formulating the gradient as shown in Equation 11. The optimal solution $(y_1^{IIP}, y_2^{IIP})$ can thus be obtained by a gradient descent algorithm. Given values of $(y_1, y_2)$ in each iterative step, the probability of realization of every demand region has to be calculated. As we extend to $N$ stores, we face the following hurdles:

- repeated probability calculations for a $2N$-dimensional multivariate distribution, and
- exponential number of demand regions $\Omega_k$ (in which the dual prices remain constant), whose identification is non-trivial.

The non-triviality arises from the fact that cross-shipment quantities are now set by a transportation linear program, as compared to explicit expressions in the two-store case. Hence we develop a tractable lower bound which yields a heuristic solution for the two-store case, which we later extend to multiple locations.

**4.1.3. Lower Bound and Heuristic for the Two-Location Problem** An important feature which complicates the IIP cost function is that the in-store demands are not pooled across regions, which in turn leads to complex and non-convex coupled terms in the cost

function. We relax this by treating unfulfilled in-store demand as online demand which can be fulfilled by cross-shipping. Specifically, we replace the total unfulfilled demand $\sum_i (D_{is} - y_i)^+ + \left( \sum_i D_{io} - \sum_i (y_i - D_{is})^+ \right)^+$ by its lower bound, which is the total unfulfilled demand when all demands are pooled $\left( \sum_i D_i - \sum_i y_i \right)^+$. Substituting this in Equation 9 and simplifying, we get the following cost function:

$$
\begin{aligned}
C^{LB}(y_1, y_2) = s(\mu_{1o} + \mu_{2o}) + \mathbb{E}\Big[ & h\left( y_1 + y_2 - D \right)^+ + (p_o - s_{12})\left( D - y_1 - y_2 \right)^+ \\
& (p_o - s - (p_o - s_{12}))\left( D_1 - y_1 \right)^+ + (p_o - s - (p_o - s_{12}))\left( D_2 - y_2 \right)^+ \\
& (p_s - (p_o - s))\left( D_{1s} - y_1 \right)^+ + (p_s - (p_o - s))\left( D_{2s} - y_2 \right)^+ \Big]
\end{aligned}
\tag{12}
$$

where $D = D_1 + D_2$, the total demand. Proposition 2 establishes $C^{LB}$ as a lower bound:

PROPOSITION 2. $C^{LB}(y_1, y_2) \le C^{IIP}(y_1, y_2)$, $\forall y_1, y_2 \ge 0$

By removing the nested piecewise linear terms in $C^{IIP}$ from Equation 9, we no longer need the gradient descent approach, as the first order conditions for $C^{LB}$ are greatly simplified:

$$
(h + p_o - s_{12})F_D \left( \sum_{j=1,2} y_j \right) + (s_{12} - s)F_{D_i}(y_i) + (p_s - p_o + s)F_{D_{is}}(y_i) = p_s, \quad \forall i = 1,2 \tag{13}
$$

We have a system of two equations with two variables, which can be solved using numerical methods to yield a heuristic solution $\mathbf{y^{LBH}}$ with expected cost $C^{LBH} = C^{IIP}(\mathbf{y^{LBH}})$. Equation 13 is of a similar structure to the first order conditions obtained by Dong and Rudi (2004) for the case of constant transshipment cost. The difference is that we have an additional term stemming from the presence of in-store demands with a higher underage cost than the online demands.

Note that the relaxation made to formulate the lower bound by replacing $\sum_i (D_{is} - y_i)^+ + \left( \sum_i D_{io} - \sum_i (y_i - D_{is})^+ \right)^+$ with $\left( \sum_i D_i - \sum_i y_i \right)^+$, will be tight when the in-store demand is very small compared to the online demand, as the optimal inventory levels are set based on the total demands. We test this numerically by changing the mix of in-store and online demands in Figure 3. The mean in-store and online demands are calculated as a proportion of a fixed total mean demand ($= 100$) in each region. The demands are normal and identical across regions, with the coefficient of variation fixed at 0.3 for each demand. The cost parameters are: $h = 2$, $p_s = 100$, $p_o = 100$, $s = 8$, $s_{12} = 15$.
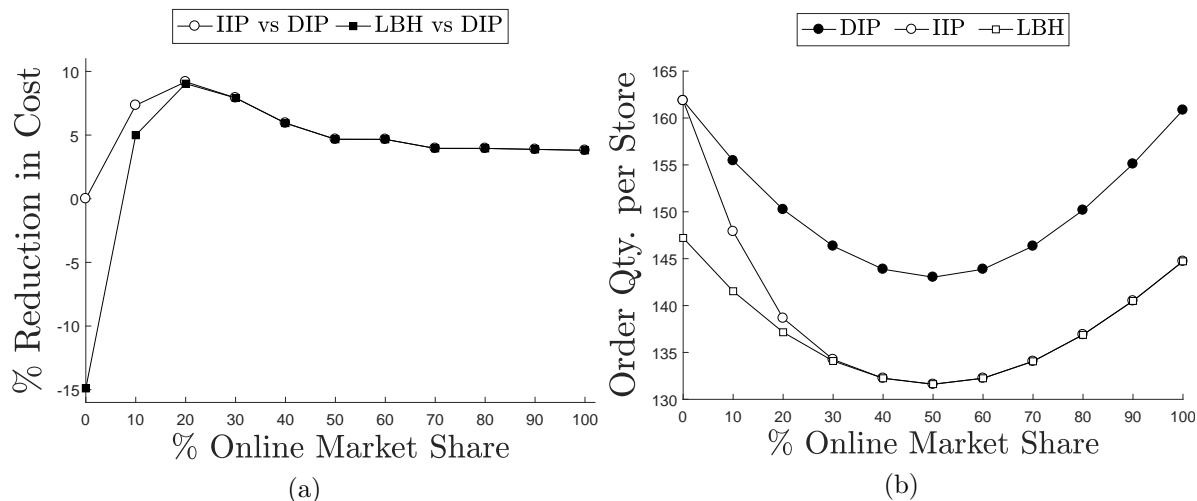
**Figure 3** Shows the effect of online market share on $C^{IIP}$, $C^{LBH}$, and $C^{DIP}$ (left) and the corresponding optimal order quantities per store (right).

From Figure 3a, we see that the heuristic provides savings over the DIP strategy for most cases, except for small values of online market share ($< 10\%$). However, we note that for such small values of online market share, the potential savings from centralized planning is also small, as seen from comparing the IIP and DIP costs. In such cases, one can simply resort to planning for each region separately using the DIP strategy.

Centralized inventory planning is most valuable when there is a moderate mix of online and in-store demands. As online demand grows in comparison to in-store demands, the effect of pooling across regions increases, due to two reasons: 1) more demand is pooled across regions which leads to a bigger reduction in variability of the total online demand, and 2) pooled online demands can better absorb the variability in the in-store demands. Thus, the maximum savings is achieved when there is a good mix of online and in-store demands so that the pooling across channels and locations work in synergy.

As the in-store demand becomes smaller, the probability that there will be unfulfilled in-store demand decreases, and the heuristic solution converges to the optimal IIP solution (Figure 3b). Thus for high values of online market share, in-store demand can effectively be treated as online demand which explains the stable savings achieved by the IIP solution.

The cost savings directly arise from a change in inventory levels in anticipation of pooling across locations. Proposition 3 addresses this observation from Figure 3b that the IIPH solution consistently stocks less than the DIP solution at each store.

18

**Govindarajan, Sinha and Uichanco:** *Inventory and Fulfillment Decisions for Omnichannel Retailing*
Stephen M. Ross School of Business, University of Michigan 2018

PROPOSITION 3. *For identical stores and normally distributed demands, $y^{LB} \leq (\geq) y^{DIP}$ whenever $y^{DIP} \geq (\leq) \mu$, where $\mu$ is the mean total demand at a store. Under perfect positive correlation across locations, $y^{LB} = y^{DIP} = y^{IIP}$.*

Similar to the intuition in newsvendor settings, $y^{DIP} \geq \mu$ would hold when underage costs are greater than overage costs, but this does not translate into an analytical proof due to the structure of the optimality equations in Equation 7, which has a mixture distribution as compared to a simple normal distribution in newsvendor theory. Lastly, positive correlation across locations reduces the pooling benefits achieved by cross-shipping, and under perfect correlation, all locations either have too much or too little inventory without any imbalance.

### 4.2. The Multi-Location Problem

We extend the two-store problem discussed so far to a generalized setting with multiple regions, as described earlier in Section 3 (Figure 1). The cross-shipping costs are taken to be $s_{ij} = s + f(d_{ij})$, where $d_{ij}$ is the distance between location $R_i$ and region $j$, and $f$ is a non-negative, increasing function such that $f(d) \to 0$ as $d \to 0$. Also, $\sup_{d \in \mathcal{D}} f(d) \leq h + p_o - s$, where $\mathcal{D} = \{d_{ij}, \forall i, j\}$, so that the conditions in Equation 4 hold true.

The decentralized solution $\mathbf{y^{DIP}}$ derived from Equation 7 readily extends to the multiple locations as the problem is decoupled by region, whereas the optimal IIP solution cannot be obtained due to the computational infeasibility even of the two-store approach. However, we can extend the heuristic and lower bound developed in the two-store case, by lowering all cross-shipping costs to $s_{\min} = \min_{i \neq j} s_{ij}$, yielding the first order conditions:

$$(h + p_o - s_{\min}) F_D\Big(\sum_{j \in \mathcal{S}} y_j\Big) + (s_{\min} - s) F_{D_i}(y_i) + (p_s - p_o + s) F_{D_{is}}(y_i) = p_s, \quad \forall i \in \mathcal{S} \quad (14)$$

The corresponding cost function yields a lower bound to the multi-location problem, satisfying Propositions 2 and 3 (the proofs are similar to the two-store case, and hence omitted). The optimal solution can be found by iterative root-finding algorithms such as the Newton-Raphson method, but the computational burden of this solution, although reduced from the newsvendor network approach by van Mieghem and Rudi (2002), is still significant for omnichannel networks in practice with thousands of stores due to the number of variables involved. A small change to the parameters: reducing $s_{\min}$ to $s$ yields a weaker lower bound:

$$C^{LBN}(y_1,\ldots,y_N) = s \sum_{i \in \mathcal{S}} \mu_{io} + \left[ \mathbb{E}h \left( \sum_{i \in \mathcal{S}} y_i - D_{\mathcal{S}} \right)^+ + \mathbb{E}(p_o - s) \left( D_{\mathcal{S}} - \sum_{i \in \mathcal{S}} y_i \right)^+ \right.$$
$$\left. + \mathbb{E}(p_s - p_o + s) \sum_{i \in \mathcal{S}_{so}} (D_{is} - y_i)^+ \right] \tag{15}$$

$C^{LBN}$ is convex in the inventory levels, and can be solved to yield a heuristic solution $y^{LBN}$ characterized by the first order conditions:

$$(h + p_o - s)F_{D_{\mathcal{S}}}\left( \sum_{j \in \mathcal{S}} y_j^{LBN} \right) + (p_s - p_o + s)F_{D_{is}}(y_i^{LBN}) = p_s, \quad \forall i \in \mathcal{S} \tag{16}$$

This is similar to the zero transshipment cost case considered by Dong and Rudi (2004), except that the presence of in-store demands allows us to fix inventory levels at each location separately, in contrast to Dong and Rudi (2004) where the setting leads to an optimal system-wide inventory level. As a consequence, the calculation of $y^{LBN}$ is computationally light, established by the following Proposition.

PROPOSITION 4. *The heuristic solution is unique, and when demands follow a multivariate normal distribution, the heuristic inventory levels at stores are at the same critical fractile of their corresponding in-store demands.*

In contrast with Equation 14, we only need to solve for one variable, namely the common critical fractile of the in-store demands. This reduces the computational effort drastically, even for very large networks. There is however, a downside to this computational gain - the optimal solution has zero inventory in the OFCs. This is because all cross-shipping costs are lowered to $s$, a unit of inventory at the OFC can lead to a decrease in total cost if it was instead at a store, as it can also serve to fulfill in-store demands.

We modify $y^{LBN}$ to obtain the heuristic solution $y^{IIPH}$ for multiple locations by calculating order quantities for the OFCs separately, and using them in Equation 16 to compute order quantities for the omnichannel stores. The order-up-to quantities for OFCs are calculated from the pooled total order quantity for OFCs, which is determined using the newsvendor quantity for the combined online demand $D_{\mathcal{S}_o} = \sum_{i \in \mathcal{S}_o} D_{io}$.

$$\sum_{j \in \mathcal{S}_o} y_j^{IIPH} = F_{D_{\mathcal{S}_o}}^{-1}\left( \frac{p_o - s}{h + p_o - s} \right) \tag{17}$$

The actual underage cost for online demands at the OFCs would be less than $p_o - s$ and would depend on inventory information of stores, as stores can fulfill these online orders with available inventory. The calculation of inventory levels at stores and OFCs are dependent on each other, but since we are forced to estimate the inventory at OFCs separately, we inflate the underage cost to $p_o - s$ which yields a higher overall inventory level at the OFCs. This is a limitation that arises out of our heuristic approximation, but it allows us to extend the heuristic to the case where OFCs have a different service cost ($s_o$) compared to the stores ($s$), as the inventory calculation for the OFCs is done separately.

To calculate the individual order quantities $y_i^{IIPH}$, $i \in \mathcal{S}_o$, we use the method of obtaining order-up-to quantities for multiple products with capacity constraints, as described in Chopra and Meindl (2007, p. 367). The total capacity is the total order-up-to quantity calculated from Equation 17, and the order-up-to quantity for each product corresponds to the order-up-to quantity for each OFC. Each unit from $\sum_{j \in \mathcal{S}_o} y_j^{IIPH}$ is allocated incrementally to the OFCs based on the individual expected marginal costs. Once the order-up-to quantities for the OFCs are obtained, they are used in Equation 18 to determine order-up-to levels for other omnichannel stores.

$$(h + p_o - s) F_{D_{\mathcal{S}}} \left( \sum_{j \in \mathcal{S}} y_j^{IIPH} \right) + (p_s - p_o + s) F_{D_{is}} \left( y_i^{IIPH} \right) = p_s, \quad \forall i \in \mathcal{S}_{so} \qquad (18)$$

Calculating the heuristic solution $y^{IIPH}$ is also computationally fast, as Proposition 4 still applies to Equation 18. The cost of the heuristic solution is given by $C^{IIPH} = C^{IIP}(y^{IIPH})$. We capture the effect of virtual pooling among the facilities in this heuristic, and the systematic approach is shown in Algorithm 1.

The performance of the heuristic clearly depends on the structure of the network which directly influences the cross-shipping costs, in addition to the mix of in-store and online demands. However in practice, the range of shipping costs is not too large: for a 5lb package, the ratio $\max_{i,j} s_{ij}/s$ is less than 2 for the UPS Ground option, and less than 3 for the UPS Next Day Air option (UPS 2017) for locations within the mainland US. We test the sensitivity for factors that adversely affect heuristic performance in Section 6 (Figure 5).

As the problem scale increases, and the number of stores grows large within a given area to accommodate the increase in demand, it is highly likely that a store with unfulfilled online demand can find a close-by store with available inventory, and hence, almost all

---

**Algorithm 1** Procedure to calculate the heuristic solution $\mathbf{y}^{IIPH}$

---

1: For physical stores in set $\mathcal{S}_s$, set $y_i^{IIPH} = F_{is}\left(\frac{p_s}{h+p_s}\right), \forall i \in \mathcal{S}_s$.

2: **for** $i \in \mathcal{S}_o$ (OFCs) **do**

3:     Calculate total order quantity: $y^{TOT} = F_{D_{\mathcal{S}_o}}\left(\frac{p_o-s}{h+p_o-s}\right)$, where $D_{\mathcal{S}_o} = \sum\limits_{i \in \mathcal{S}_o} D_{io}$.

4:     Set $y_i^{IIPH} = 0, \forall i \in \mathcal{S}_o$, and $rem = \lfloor y^{TOT} \rfloor$.

5:     Calculate marginal cost $MC_i\left(y_i^{IIPH}\right) = -(p_o - s)(1 - F_{D_{io}}(y_i^{IIPH})) + hF_{D_{io}}\left(y_i^{IIPH}\right)$

6:     Choose $i^* = \min\limits_{i \in \mathcal{S}_o} MC_i(y_i^{IIPH})$.     Set $y_{i^*}^{IIPH} \leftarrow y_{i^*}^{IIPH} + 1$

7:     Set $rem \leftarrow rem - 1$. If $rem > 0$, go to Step 3.

8: **for** $i \in \mathcal{S}_{so}$ **do**

9:     Calculate order quantities implicitly from the optimality equations:
$(h + p_o - s) F_{D_{\mathcal{S}}}\left(\sum\limits_{j \in \mathcal{S}} y_j^{IIPH}\right) + (p_s - p_o + s) F_{D_{is}}\left(y_i^{IIPH}\right) = p_s, \ \forall i \in \mathcal{S}_{so}.$

---

cross-shipping takes place over short distances, at a cost close to $s$. Thus, we can expect the heuristic solution to be close to the optimal solution, and as a consequence of this notion, Proposition 5 shows that the heuristic is near optimal in an asymptotic sense.

PROPOSITION 5. *As the number of omnichannel stores in a given area increases, with demands bounded and i.i.d. across regions, for sufficiently small $h > 0$, the heuristic is near optimal in an asymptotic sense with a constant approximation factor, i.e.*

$$\frac{C^{IIPH}}{C^{LBN}(y^{IIPH})} \leq \frac{h + p_s}{p_s - p_o + s}, \ as \ N \to \infty$$

The proposition holds when all locations have omnichannel stores, and $y^{LBN} = y^{IIPH}$. We first show that reducing cross-shipping costs to $s$ preserves optimality in the asymptotic setting, by considering a simplified setting where the stores are uniformly distributed in the given region, which is in-turn divided into identical sub-regions. As the number of stores grows large, each sub-region has sufficient supply to fulfill its demands, and hence cross-shipping takes place only within the sub-regions with costs converging to $s$.

When in-store demands dominate, the heuristic suffers from its assumption that in-store demands are pooled. However, we bound the heuristic performance by a constant approximation factor dependent on cost parameters. While this bound is not tight, it shows that the heuristic is not critically affected by its assumptions as the problem scale grows.

# 5. Multi-Period, Multi-Location, and Multiple Fulfillment Epochs

So far, we have discussed the single review period setting where online fulfillment is done once, at the end of the period. We now switch back to the general version of the problem described in Section 3, with multiple review periods and online demand fulfilled over $T$ fulfillment epochs in each review period. This is a more realistic representation of practice, as we closely approximate the continuous time case, because the value of $T$ can be flexibly large. We start by proving convexity for the single period problem described in Equation 1.

PROPOSITION 6. *The single-period, $T$-fulfillment epoch expected cost $C(\mathbf{y}) = \mathbb{E}C_1(\mathbf{y}, \tilde{D})$ is jointly convex in the inventory levels $y_i$.*

The proof follows by induction. Let the optimal solution to the single period problem be denoted by $\mathbf{y^{IIP}}$. We extend our analysis to the finite horizon case with multiple periods.

PROPOSITION 7. *For the finite horizon problem with lost sales and zero replenishment leadtime, a stationary base-stock policy is optimal, with order-up-to levels $\mathbf{y^{IIP}}$.*

For the zero replenishment leadtime case with lost sales, the multi-period problem reduces to solving a single-period problem, and the proof is similar to traditional multi-period inventory problems involving lost-sales. As noted earlier, solving for $\mathbf{y^{IIP}}$ is difficult, as optimal fulfillment decisions are intractable. Hence, we resort to heuristic solutions developed from our analysis of the single period problem.

## 5.1. Inventory Levels

To obtain a heuristic solution to set order-up-to levels, we use the procedure described in Algorithm 1, by approximating the problem as a single fulfillment epoch problem. Naturally, the demands used to calculate the heuristic solutions are the total review-period demands at each location. For example, the review-period in-store demand at store $i$ is given by $\mathcal{D}_{is} = \sum_{t=1}^{T} D_{is}^t$. Also, the holding cost parameter used in the algorithm is the review-period holding cost, which is given by $\bar{h} = h^* T$.

We compare this heuristic solution with the naive strategy which plans for inventory in a decentralized fashion. We extend the DIP solution derived in Equation 7, by using the total review-period demands for each location and holding cost $\bar{h}$. We will continue to denote the heuristic solution derived in this fashion by IIPH and the decentralized solution as DIP for the numerical studies in the following sections.

### 5.2.    Fulfillment Policies

We consider two fulfillment policies, which dictate how online orders are fulfilled:

1. the *myopic fulfillment (MF)* policy, where online demands in the current fulfillment epoch are fulfilled to the maximum possible extent with the available inventory, without consideration for demands in the future, and

2. the *threshold fulfillment (TF)* policy, which reserves inventory at each location for future in-store demands, by halting online fulfillment from a location when the inventory level falls below a certain threshold in each fulfillment epoch.

As future in-store demands are costlier to lose and do not have the additional flexibility of cross-shipping, it is intuitive that the TF policy can lead to reduction in costs compared to the MF policy when implemented well. Rationing inventory between high-priority and low-priority demands has been studied in literature (for a review, refer to Kleijn and Dekker 1999), and along similar lines, Jalilipour Alishah et al. (2015) prove the existence of an optimal threshold rationing policy between in-store and online demands at a single store.

In our case it is not straightforward to estimate the underage cost for the low-priority (online) demand, as it is endogenized by the fulfillment policy followed and depends on where an order is fulfilled from. The optimal thresholds depend on in-store and online demands in a complicated, network-based fashion, as online demands are pooled across locations, and their calculation is akin to obtaining optimal transshipment decisions based on such a threshold structure. We propose simple newsvendor-based thresholds which only take into account future in-store demands. In any fulfillment epoch $t$, an amount $w_i^t$ is reserved at store $i$ for future in-store demands in that review period, where

$$w_i^t = F_{\mathcal{D}_{is}^t}\left(\frac{p_s}{h(T-t+1)+p_s}\right), \quad \text{where} \quad \mathcal{D}_{is}^t = \sum_{\hat{t}=t+1}^{T} D_{is}^{\hat{t}} \tag{19}$$

We have developed a static fulfillment policy, as these thresholds can be evaluated at the start of the review period based on the demand forecasts. We formalize the TF policy in Algorithm 2. The MF policy places no such restriction on fulfillment, and can simply be recovered from Algorithm 2 by setting the thresholds $w_i^t$ to be zero in step 1.

Note that the fulfillment heuristic is agnostic to current inventory levels and online demands. While including such information would be valuable, we show that such a simple policy, when combined with a good inventory heuristic which positions inventory in a calculated fashion, can provide considerable savings compared to naive strategies.

---

**Algorithm 2** Implementation of the Threshold Fulfillment (TF) Policy

---

1: At the start of the review period, evaluate thresholds $w_i^t$, $\forall i, t$ using Equation 19.

2: In each fulfillment epoch $t$, each location first fulfills its own in-store demand to the maximum possible extent, and the leftover inventory at location $i$ is $\hat{x}_i^t$.

3: Calculate fulfillment capacities for each location $i$ as $K_i^t = (\hat{x}_i^t - w_i^t)^+$.

4: Online fulfillment decisions $Z_{ij}^t$ are obtained from the transportation linear program:
$$\left\{ \min \ \sum_{i,j}(s_{ij} - h - p_o)Z_{ij}^t, \quad \text{subject to: } \sum_k Z_{kj}^t \le D_{jo}, \ \sum_k Z_{ik}^t \le K_i^t, \ Z_{ij}^t \ge 0, \ \forall i,j \right\}$$

---

To evaluate the performance of the fulfillment policies, we compare them with the so-called hindsight-optimal policy. The cost of this policy can be evaluated through a linear program which minimizes the total cost in the review period, given that all uncertainty is realized at the beginning of the period. Given inventory levels, the cost of such a policy is a natural lower bound for the cost of any fulfillment policy, and we numerically show that the simple TF policy performs very well compared to this lower bound in Section 6.

## 6. Numerical Analysis

We employ a realistic setting to test the performance of the inventory and fulfillment heuristic solutions, based on a fictitious network embedded in mainland US. We shall mainly focus on the case with zero lead time and multiple fulfillment epochs.

We evaluate the total expected costs through a Monte-Carlo simulation with a sample size of $10^4$, for two inventory heuristics - IIPH (integrated planning heuristic) and DIP (decentralized planning), and two fulfillment heuristics - MF (myopic) and TF (threshold-based). We mostly focus on comparing our combined heuristic, the ⟨IIPH,TF⟩ strategy, to the benchmark ⟨DIP,MF⟩ strategy, which represents a naive solution.

### 6.1. Network Setup

We take the locations of the stores to be at the most populous cities in mainland US (Wikipedia 2016) and the OFCs are located according to the list of most efficient warehouses in the US, in terms of possible transit lead-times (Chicago Consulting 2013). The shipping costs are calculated using the cost equation estimated by Jasin and Sinha (2015) based on UPS Ground shipping rates for an item weighing one pound: $s_{ij} = 9.182 + 0.000541d_{ij}$, where $d_{ij}$ is the distance in miles from region $i$ to region $j$. We also perform sensitivity analysis for the slope of the shipping cost with respect to distance, to study the

effect of shipping costs on the relative performance of our combined heuristic. Other cost parameters used are: $\bar{h} = 2$, $p_s = p_o = 100$, $s = 9.182$.

The review-period demands are taken to be independent and normally distributed with mean and standard deviations calculated based on the population of the cities. To study the effect of online market share ($\alpha$) on the performance of the heuristic solutions, we take that the sum of the mean in-store and online demands in each region to be a fixed proportion of the cities' populations. This represents the average market size of the region, and the review-period mean in-store and online demands are calculated as $1 - \alpha$ and $\alpha$ proportions respectively of this mean market size in each region. The coefficient of variation of the review-period demands are fixed at 0.2. Demands are identical and independently distributed across fulfillment epochs, with parameters calculated from the review-period demands. In the base case, $\alpha = 0.5$ and $T = 5$, and we perform sensitivity analyses with respect to these parameters.

Let $n_s$ be the number of physical stores and $n_o$ be the number of OFCs. 90% of the physical stores are assumed to be omnichannel stores, and the rest are traditional B&M stores. Further details on the numerical setup and a brief overview of the simulation process can be found in Appendix C.

### 6.2. Results

We tabulate the results obtained. We mainly focus on comparing the cost of the combined heuristic ⟨IIPH,TF⟩ to that of the naive strategy ⟨DIP,MF⟩. We ignore the costs associated with traditional B&M stores to focus on the effect of the heuristic on online fulfillment.

**6.2.1. Network Size.** As the network size increases, centralized inventory planning and strategic fulfillment can be valuable, as there is more flexibility in terms of options available in fulfillment. Figure 4a shows that increasing network size have a positive and marginally decreasing effect on the relative performance of the combined heuristic.

We also compare the strategies based on two important metrics, inventory imbalance and inventory efficiency, and the results are shown in Figure 4b for $n_o = 2$. Higher imbalance can lead to costly spillovers and local stockouts (Acimovic and Graves 2017), which in turn can cause markdowns in stores. We measure imbalance by recording the variance of ending inventory positions across locations at the end of each epoch, and taking the average value over the review period. Although this is different from the metric used by Acimovic and
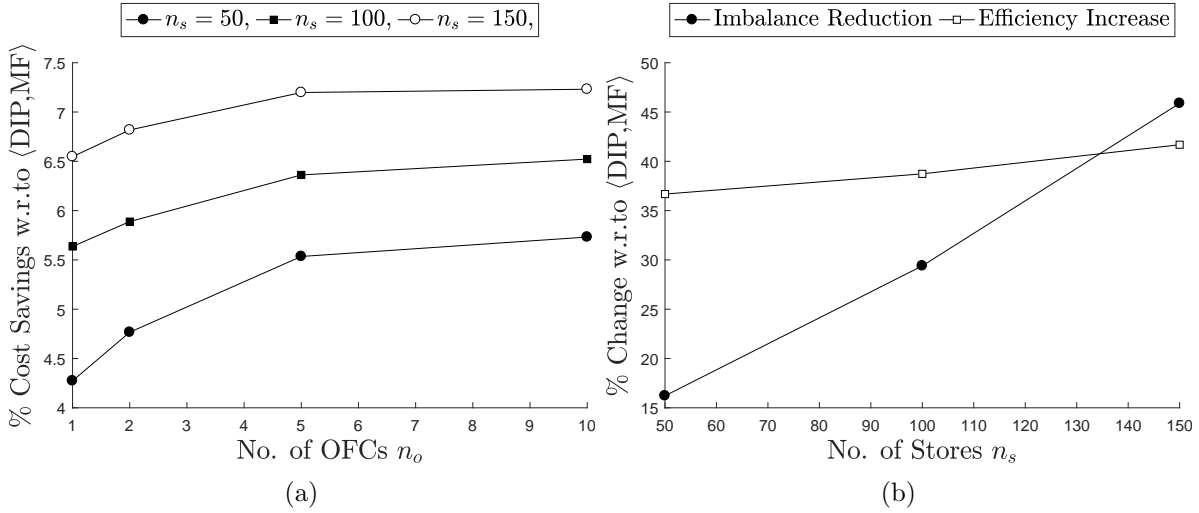
26

**Govindarajan, Sinha and Uichanco:** *Inventory and Fulfillment Decisions for Omnichannel Retailing*
Stephen M. Ross School of Business, University of Michigan 2018

**Figure 4**   **Shows the effect of network size on the performance of** $\langle$IIPH,TF$\rangle$ **compared to** $\langle$DIP,MF$\rangle$**, in terms of cost (left), inventory imbalance and inventory efficiency (right)**

Graves (2017), it captures the essence of imbalance among locations in an omnichannel network. We see that our combined heuristic achieves a lower imbalance across locations as compared to the $\langle$DIP,MF$\rangle$ strategy, and this effect is more pronounced for larger networks.

We define another metric, inventory efficiency, as an equivalent measure for inventory turnover, calculated as the ratio of the total fulfilled demand to the average inventory level of the system in a review period (calculated as the mean of the starting inventory level and expected ending inventory at the end of the review period). Higher efficiency achieved by the heuristic stems from a reduction in inventory levels without a considerable decrease in service levels, due to planning in advance for cross-shipping. This offers a potential solution to decreasing trend in turnovers in the retail industry in recent years(Kurt Salmon 2016).

**6.2.2.   Cross-shipping Costs and Online Market Share.**   As discussed in Section 4.2, two major factors affect the inventory heuristic performance – shipping cost structure and online market share. For fixed fulfillment policy TF, we compare the $\langle$IIPH,TF$\rangle$ and $\langle$DIP,TF$\rangle$ strategies to understand the effect of these parameters on the inventory heuristic. We found similar results when comparing to the $\langle$DIP,TF$\rangle$ strategy.

We first vary the slope of shipping costs with respect to distance, thereby increasing the ratio $s_{\max}/s$ (value of 1.2 corresponds to the base case setting). As expected, the relative performance of the heuristic decreases as shipping costs become more sensitive to distance (Figure 5a). For a perspective, the costliest shipping option, the UPS Next Day Air, has
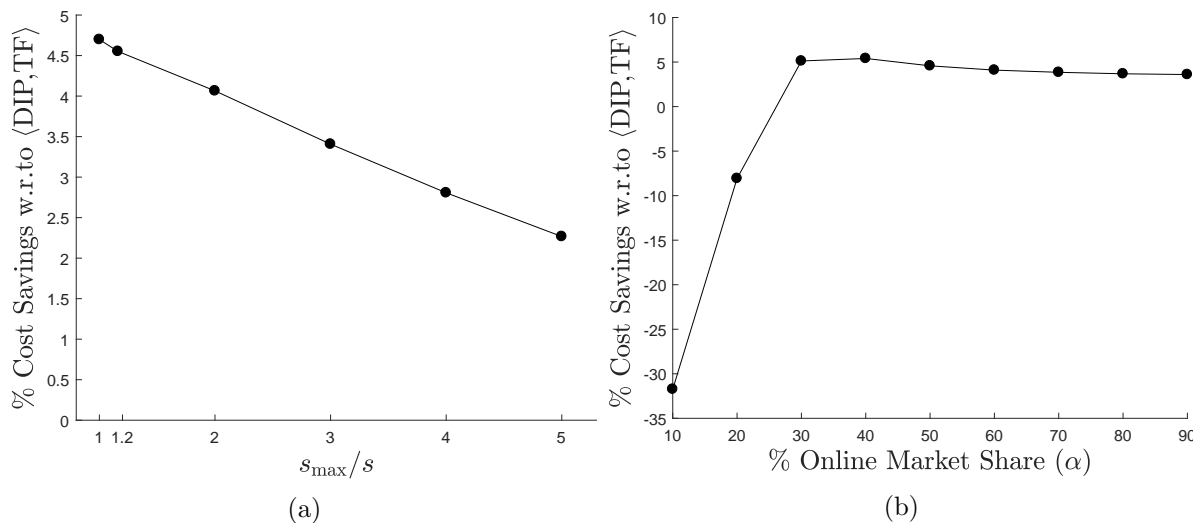
**Figure 5** Shows the effect of the slope of cross-shipping costs with distance by varying the ratio $s_{\max}/s$ **(left)** and online market share **(right)** on the performance of $\langle\text{IIPH,TF}\rangle$ **compared to** $\langle\text{DIP,MF}\rangle$.

a ratio of $s_{\max}/s$ less than 3 for shipping a 5lb package within mainland US. Hence the heuristic provides significant savings for most existing shipping cost structures.

Figure 5b shows the effect of online market share. We see that the heuristic performs worse than the decentralized solution when the online market share is low. This reflects the deficiency noted in the two-store case, as the heuristic assumes that in-store demands are pooled across locations. However, the heuristic provides a valuable alternative to the decentralized solution for products that have adequate online market shares – e.g. books, computers and consumer electronics have an online market share of about 50% (FTI Consulting 2015). Additionally, with rapidly increasing online sales, firms can obtain considerable savings through centralized inventory strategies, and for most cases, our heuristic serves as a viable proxy for inaccessible optimal decisions.

**6.2.3.   Number of Fulfillment Epochs** ($T$)   By increasing the number of times online fulfillment decisions are made, we can closely model the continuous time case. We keep the total review-period demand parameters constant, and keep demands across fulfillment epochs independent and identically distributed. To reduce the computational burden associated with higher values of $T$, we use a smaller network with $n_s = 10$, $n_o = 2$.

The results are shown in Figure 6. In Figure 6a, we compare the MF and TF fulfillment strategies with IIPH inventory levels, against the hindsight optimal strategy HF, which makes fulfillment decisions with all uncertainty realized at the start of the review period.
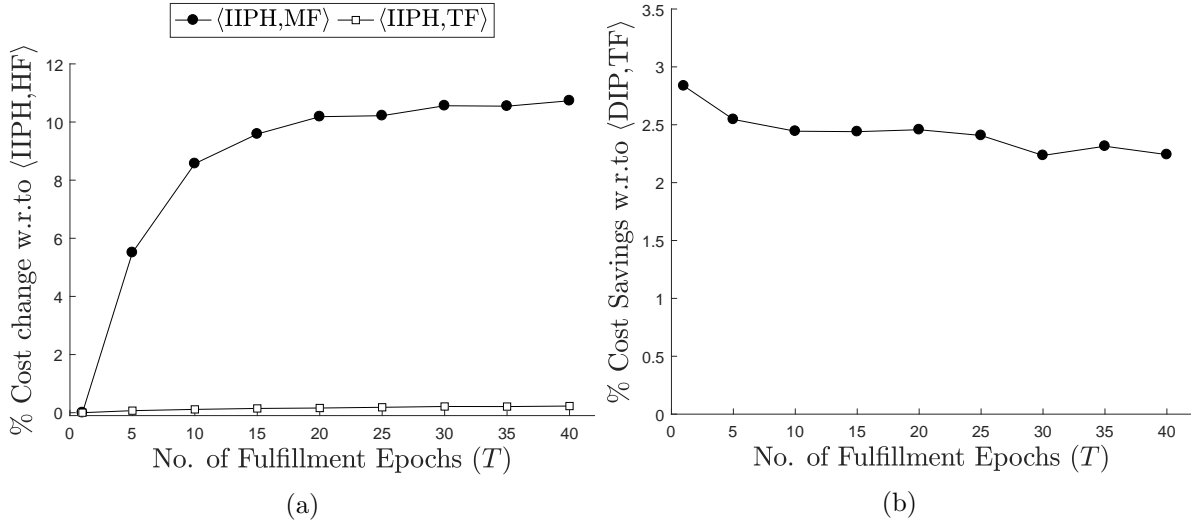
**Figure 6** Shows the effect of increasing the number of fulfillment epochs in a single review period on the cost of fulfillment policies with respect to the hindsight optimal policy (left) and the performance of $\langle$IIPH,TF$\rangle$ compared to $\langle$DIP,TF$\rangle$ (right).

As $T$ increases, the MF policy is punished for failing to reserve inventory for future in-store demands (Figure 6a). The TF policy on the other hand proves to be a simple but effective fulfillment strategy, achieving costs within 0.5% of the HF lower bound.

For a fixed fulfillment policy TF, we compare the $\langle$IIPH,TF$\rangle$ and $\langle$DIP,TF$\rangle$ strategies in Figure 6b, and see that the effect of increasing $T$ has a decreasing effect on the relative performance of the inventory heuristic.

Finally, we note that our heuristics are extremely scalable with respect to network size - for a network with $n_s = 150$, $n_o = 10$ and $T = 5$, calculating the inventory levels using the heuristic takes only around 10 seconds, and the calculation of fulfillment thresholds takes around 2 minutes. Real-life retail networks are often much bigger in size – for instance, Target ships online orders from more than 1000 stores (Lindner 2016), and our heuristic can provide considerable improvements compared to traditional strategies in most cases.

## 7. Conclusion

Despite numerous retailers struggling with the operational problems posed by omnichannel retailing, the area has received comparatively less attention in literature. Our research addresses an important facet of omnichannel retailing — inventory management, by demonstrating the value in utilizing the pooling benefits offered by omnichannel retailing, through a combined inventory and fulfillment policy.

Our heuristic policies, though derived from a complicated multi-location and multi-period model, are quite generalizable. We can extend our analysis to demands originating from abstract regions, by treating them as OFCs that carry zero inventory. Disparity in service costs at OFCs and stores can also be taken into account by using $s_o$, the service cost from OFCs, instead of $s$ in Equation 17, as inventory planning is done separately for OFCs. We still need to make the assumption that demands from a region with an omnichannel store can be fulfilled from that store with the least cost. Otherwise, the demand at this store will be assigned to be fulfilled from the online FC with the least fulfillment cost, which can lead to different first order conditions in inventory planning for the online FC.

We can also extend the heuristic solutions to the case of positive leadtimes as follows: assuming each location $i$ has a replenishment leadtime of $L_i$ review periods, the total planning horizon for order-up-to policies is $(L_i + 1)$ review periods, or equivalently, $(L_i + 1)T$ fulfillment epochs for each location. Using the total demands during the planning period for each location instead of review period demands, we can directly extend our inventory heuristic to set order-up-to levels for each location.

For the fulfillment heuristic, an additional threshold for inventory position needs to be calculated based on future in-store demands in the remainder of the current planning horizon, which can also be computed based on a simple newsvendor formula. Online fulfillment from a location is temporarily stopped in an epoch when either threshold is violated.

An important direction for future research is to include multiple classes of online demand, especially in-store pickups, which is a popular mode of omnichannel fulfillment. A heuristic control for managing multiple products is also an interesting and important extension to be considered. Future research may also focus on further extensions such as capacities and stochastic leadtimes. We believe that our framework provides a platform to build further complexities on, which can yield important decision support tools for the industry.

## Acknowledgments

## References

Acimovic, Jason, Stephen C Graves. 2014. Making better fulfillment decisions on the fly in an online retail environment. *Manufacturing & Service Oper. Management* **17**(1) 34–51.

Acimovic, Jason, Stephen C Graves. 2017. Mitigating spillover in online retailing via replenishment. *Manufacturing & Service Oper. Management* **19**(3) 419–436.

Ansari, Asim, Carl F Mela, Scott A Neslin. 2008. Customer channel migration. *J. of Marketing Res.* **45**(1) 60–76.

Bell, David, Santiago Gallino, Antonio Moreno. 2013. Inventory showrooms and customer migration in omni-channel retail: The effect of product information. Working Paper.

Bijvank, Marco, Iris FA Vis. 2011. Lost-sales inventory theory: A review. *Eur. J. of Oper. Res.* **215**(1) 1–13.

Brynjolfsson, Erik, Yu Jeffrey Hu, Mohammad S Rahman. 2013. Competing in the age of omnichannel retailing. *MIT Sloan Management Review* **54**(4) 23.

Chen, Jian, Youhua Chen, Mahmut Parlar, Yongbo Xiao. 2011. Optimal inventory and admission policies for drop-shipping retailers serving in-store and online customers. *IIE Transactions* **43**(5) 332–347.

Chen, Xin, Peng Hu, Simai He. 2013. Preservation of supermodularity in parametric optimization problems with nonlattice structures. *Operations Res.* **61**(5) 1166–1173.

Chicago Consulting. 2013. Ten best warehouse networks. URL `http://www.chicago-consulting.com/ten-best-warehouse-networks/`.

Chopra, Sunil, Peter Meindl. 2007. Supply chain management. strategy, planning & oper. *Das Summa Summarum des Management*. Springer, 265–275.

Dong, Lingxiu, Nils Rudi. 2004. Who benefits from transshipment? Exogenous vs. endogenous wholesale prices. *Management Sci.* **50**(5) 645–657.

Forrester. 2014. Customer desires vs. retailer capabilities: Minding the omnichannel commerce gap. URL `http://global.sap.com/asia/campaigns/2014-07/hybris-accenture-forrester-tlp-omni-channel.pdf`.

FTI Consulting. 2015. The U.S. online retail forecast. URL `http://www.fticonsulting.com/~/media/Files/us-files/insights/featured-perspectives/fti-2015onlineretailforecast.pdf`.

Gallino, Santiago, Antonio Moreno. 2014. Integration of online and offline channels in retail: The impact of sharing reliable inventory availability information. *Management Sci.* **60**(6) 1434–1451.

Gallino, Santiago, Antonio Moreno, Ioannis Stamatopoulos. 2017. Channel integration, sales dispersion, and inventory management. *Management Sci.* **63**(9) 2813–2831.

Gao, Fei, Xuanming Su. 2016. Online and offline information for omnichannel retailing. *Manufacturing & Service Oper. Management* **19**(1) 84–98.

Gao, Fei, Xuanming Su. 2017. Omnichannel retail operations with buy-online-and-pick-up-in-store. *Management Sci.* **63**(8) 2478–2492.

Giannopoulos, Nicole. 2014. The macy's advantage: Secrets to same-day delivery, omnichannel. *RIS News.* Retrieved from: `http://risnews.edgl.com/retail-news/The-Macy-s-Advantage--Secrets-to-Same-Day-Delivery,-Omnichannel97157`.

**Govindarajan, Sinha and Uichanco:** *Inventory and Fulfillment Decisions for Omnichannel Retailing*
Stephen M. Ross School of Business, University of Michigan 2018

31

Harrison, J Michael, Jan A van Mieghem. 1999. Multi-resource investment strategies: Operational hedging under demand uncertainty. *Eur. J. Oper. Res.* **113**(1) 17–29.

Harsha, Pavithra, Shivaram Subramanian, Joline Uichanco. 2016. Dynamic pricing of omnichannel inventories. Available at `http://web.mit.edu/people/pavithra/papers/HarshaSubramanianUichanco2017_OCPX.pdf`.

Heyman, Daniel P, Matthew J Sobel. 2003. *Stochastic models in operations research: Stochastic optimization*, vol. 2. Courier Corporation.

Hoeffding, Wassily. 1963. Probability inequalities for sums of bounded random variables. *J. of the American Statistical Association* **58**(301) 13–30.

Jalilipour Alishah, Elnaz, Kamran Moinzadeh, Yong-Pin Zhou. 2015. Inventory fulfillment strategies for an omni-channel retailer. Available at SSRN: `https://ssrn.com/abstract=2659671`.

Jasin, Stefanus, Amitabh Sinha. 2015. An LP-based correlated rounding scheme for multi-item ecommerce order fulfillment. *Oper. Res.* **63**(6) 1336–1351.

Kim, Eugene. 2017. More than half of online sales growth in the US came from Amazon last year. *Business Insider.* Retrieved from: `http://www.businessinsider.com/amazon-drives-more-than-half-us-ecommerce-growth-2016-2017-2`.

Kleijn, Marcel J, Rommert Dekker. 1999. An overview of inventory systems with several demand classes. *Lecture Notes in Economics and Mathematical Systems* 253–266.

Kurt Salmon. 2016. Building a solid omnichannel foundation with effective inventory management. URL `http://www.kurtsalmon.com/en-us/Retail/vertical-insight/1478/Building-a-Solid-Omnichannel-Foundation-with-Effective-Inventory-Management`.

Lei, Yanzhe, Stefanus Jasin, Amitabh Sinha. 2017. Dynamic joint pricing and order fulfillment for e-commerce retailers. *Manufacturing & Service Oper. Management.* Forthcoming.

Leiser, Jordy. 2016. Think tank: Why an omnichannel approach can help retail escape the Amazon. *WWD.* Retrieved from: `http://wwd.com/retail-news/forecasts-analysis/blazing-the-amazon-and-why-an-omnichannel-approach-can-help-retail-10312172/`.

Lewis, Sam. 2013. Macy's grows order fulfillment centers to 500. *Retail Solutions Online.* Retrieved from: `https://www.retailitinsights.com/doc/macy-s-grows-order-fulfillment-centers-to-0001`.

Lindner, Matt. 2016. Target now ships online orders from more than 1,000 stores. *InternetRetailer.* Retrieved from: `https://www.internetretailer.com/2016/11/16/target-now-ships-online-orders-more-1000-stores`.

Lindner, Matt. 2017. Ecommerce is expected to grow to 17% of US retail sales by 2022. *InternetRetailer.* Retrieved from: `https://www.digitalcommerce360.com/2017/08/09/e-commerce-grow-17-us-retail-sales-2022/`.

Nash, Kim. 2015. Walmart build supply chain to meet e-commerce demands. *The Wall Street Journal.* Retrieved from: `http://www.wsj.com/articles/wal-mart-builds-supply-chain-to-meet-e-commerce-demands-1431016708`.

Paterson, Colin, Gudrun Kiesmüller, Ruud Teunter, Kevin Glazebrook. 2011. Inventory models with lateral transshipments: A review. *Eur. J. Oper. Res.* **210**(2) 125–136.

Rigby, Darrell. 2011. The future of shopping. *Harvard Business Review* **89**(12) 65–76.

Seifert, Ralf W, Ulrich W Thonemann, Marcel A Sieke. 2006. Relaxing channel separation: Integrating a virtual store into the supply chain via transshipments. *IIE Trans.* **38**(11) 917–931.

Şen, Alper, Alex X Zhang. 1999. The newsboy problem with multiple demand classes. *IIE Trans.* **31**(5) 431–444.

Tagaras, George. 1989. Effects of pooling on the optimization and service levels of two-location inventory systems. *IIE Trans.* **21**(3) 250–257.

Tagaras, George, Morris A Cohen. 1992. Pooling in two-location inventory systems with non-negligible replenishment lead times. *Management Sci.* **38**(8) 1067–1083.

UPS. 2017. 2017 UPS rate and service guide. URL `https://www.ups.com/media/en/daily_rates.pdf`. [Online; accessed 10-Oct-2017].

UPS Compass. 2014. Ship from store: A smart competitive strategy for retailers. URL `https://compass.ups.com/ship-from-store-benefits/`.

U.S. Census Bureau. 2016. Quarterly retail e-commerce sales. Retrieved from: `https://www.census.gov/retail/mrts/www/data/pdf/ec_current.pdf`.

van Mieghem, Jan A. 2003. Commissioned paper: Capacity management, investment, and hedging: Review and recent developments. *Manufacturing Service Oper. Management* **5**(4) 269–302.

van Mieghem, Jan A, Nils Rudi. 2002. Newsvendor networks: Inventory management and capacity investment with discretionary activities. *Manufacturing Service Oper. Management* **4**(4) 313–335.

Wei, Lai, Stefanus Jasin, Roman Kapuscinski. 2017. Shipping consolidation with delivery deadline and expedited shipment options Available at SSRN: `https://ssrn.com/abstract=2920899orhttp://dx.doi.org/10.2139/ssrn.2920899`.

Wikipedia. 2016. List of United States cities by population. URL `https://en.wikipedia.org/wiki/List_of_United_States_cities_by_population`. [Online; accessed 12-Dec-2016].

Xu, Ping Josephine, Russell Allgor, Stephen C Graves. 2009. Benefits of reevaluating real-time order fulfillment decisions. *Manufacturing Service Oper. Management* **11**(2) 340–355.

Yang, Jian, Zhaoqiong Qin. 2007. Capacitated production control with virtual lateral transshipments. *Oper. Res.* **55**(6) 1104–1119.

**Govindarajan, Sinha and Uichanco:** *Inventory and Fulfillment Decisions for Omnichannel Retailing*
Stephen M. Ross School of Business, University of Michigan 2018

33

Yao, David D, Sean X Zhou, Weifen Zhuang. 2016. Joint initial stocking and transshipment–Asymptotics and bounds. *Production and Oper. Management* **25**(2) 273–289.

Zaroban, Stefany. 2016. U.S. e-commerce grows 14.6% in 2015. *InternetRetailer.* Retrieved from: `https://www.internetretailer.com/2016/02/17/us-e-commerce-grows-146-2015`.

# Appendices

## Appendix A: Proofs of Propositions

### A.1. Proof of Proposition 1

*Proof:* We first observe that given a realization of the demands, the optimal cost can be obtained using a linear program. The proof follows in similar fashion to Seifert et al. (2006, Proposition 1). Consider the linear program $P(y_1, y_2, \tilde{D})$, where $z_i$ represents the amount of inventory at $R_i$ used to fulfill its in-store demand, and $z_{ij}$ represents the amount of inventory of $R_i$ used to fulfill online demand from region $j$.

$$
\begin{aligned}
P(y_1, y_2, \tilde{D}) = \min_{z_i, z_{ii}, z_{ij}} & \sum_i h(y_i - z_i - \sum_j z_{ij}) + \sum_i p_s(D_{is} - z_i) \\
& + \sum_i p_o(D_{io} - \sum_j z_{ji}) + \sum_i s z_{ii} + \sum_i \sum_{j \neq i} s_{ij} z_{ij} \\
\text{subject to} \quad & z_i + \sum_j z_{ij} \leq y_i, & \forall i \\
& z_i \leq D_{is}, & \forall i \\
& \sum_j z_{ji} \leq D_{io}, & \forall i \\
& z_i, z_{ij} \geq 0, & \forall i, j
\end{aligned}
\tag{20}
$$

To show that the function $P$ represents $C^{IIP}$ for a given demand $\tilde{D}$, notice that the coefficients of the decision variables $z_i$, $z_{ii}$, $z_{ij,(j\neq i)}$ in the objective function follow $(-h - p_s) < (s - h - p_o) < (s_{ij} - h - p_o)$, under the conditions in $\Psi$ in Equation 4. The linear program can be solved greedily, and it is easy to see that the optimal solution is given by $z_i = \min(y_i, D_{is})$, $z_{ii} = \min\left((y_i - D_{is})^+, D_{io}\right)$, $z_{ij} = \min\left(\left((y_i - D_{is})^+ - D_{io}\right)^+, \left(D_{jo} - (y_j - D_{js})^+\right)^+\right)$.

The sequence of fulfillment is clear: in-store demand is fulfilled first, followed by online demand from the same region, and finally cross-shipment to other regions. Hence, we have $C^{IIP}(y_1, y_2) = \mathbb{E}_{\tilde{D}}\left(P\left(y_1, y_2, \tilde{D}\right)\right)$. The objective function is linear and the constraint set in (20) is a polyhedral convex set with linear constraints, and hence by Heyman and Sobel (2003, Proposition B-4), P is jointly convex in $y_1$, $y_2$, $\tilde{D}$. As the expectation of a convex function is convex, it follows that $C^{IIP}(y_1, y_2)$ is jointly convex in $y_1$ and $y_2$.

The structure of $C^{IIP}$ as an expectation of a linear program draws direct comparison with the value function in newsvendor networks (van Mieghem and Rudi 2002). Similar to Proposition 2 in Harrison and van Mieghem (1999), the gradient of the function $P(y_1, y_2, \tilde{D})$ with respect to $(y_1, y_2)$ can be written as:

$$
\nabla_{y_1, y_2} P\left(y_1, y_2, \tilde{D}\right) = (h, h)^T - \lambda\left(y_1, y_2, \tilde{D}\right)
\tag{21}
$$

where $\lambda(y_1, y_2, \tilde{D})$ is the dual-price vector corresponding to the constraints with $y_1$ and $y_2$ in (20). The 4-dimensional demand space can be divided into domains $\Omega_i(y_1, y_2)$ such that in each domain, the optimal values of the decision variables $z_i$, $z_{ii}$ and $z_{ij}$ are linear in $y_1$ and $y_2$, and hence the dual-price vector $\lambda(y_1, y_2, \tilde{D})$ is constant (refer to Appendix B for a discussion). The first-order conditions are:

$$
\nabla_{y_1, y_2} C^{IIP}(y_1, y_2) = 0 = \nabla_{y_1, y_2} \mathbb{E}_{\tilde{D}}\left(P\left(y_1, y_2, \tilde{D}\right)\right)
\tag{22}
$$

We can interchange the gradient and expectation on the right hand side of Equation 22 (see Harrison and van Mieghem (1999) for a proof), and thus Equation 22 becomes

$$\nabla_{y_1,y_2} C^{IIP}(y_1, y_2) = 0 = \mathbb{E}_{\tilde{D}} \nabla_{y_1,y_2} P\left(y_1, y_2, \tilde{D}\right) = (h, h)^T - \mathbb{E}_{\tilde{D}} \lambda\left(y_1, y_2, \tilde{D}\right)$$
$$= (h, h)^T - \sum_i \lambda^i \mathbb{P}\left(\Omega_i(y_1, y_2)\right) \tag{23}$$

where $\lambda^i$ is the constant $\lambda\left(y_1, y_2, \tilde{D}\right)$ for $\tilde{D} \in \Omega_i(y_1, y_2)$. $\qquad \square$

### A.2. Proof of Proposition 2

Based on the approximation used to formulate $C^{LB}$, the difference in costs between $C^{IIP}$ and $C^{LB}$ is:

$$C^{IIP}(\mathbf{y}) - C^{LB}(\mathbf{y}) = (h + p_o - s_{12}) \mathbb{E}\left[\left(\sum_i D_{io} - \sum_i (y_i - D_{is})^+\right)^+ + \sum_i (D_{is} - y_i)^+ - \left(D - \sum_i y_i\right)^+\right]$$
$$\geq (h + p_o - s_{12}) \mathbb{E}\left[\left(\sum_i D_{io} - \sum_i (y_i - D_{is})^+ + \sum_i (D_{is} - y_i)^+\right)^+ - \left(D - \sum_i y_i\right)^+\right]$$
$$= 0$$

The first inequality follows from : $a^+ + b^+ \geq (a+b)^+$, and further simplification uses $x^+ - (-x)^+ = x$. $\qquad \square$

The proof follows for any number of stores, as long as the cross-shipping cost is a constant and $s_{12} < h + p_o$. The proof also follows when $s_{12}$ is reduced to $s$, as done in Equation 15.

### A.3. Proof of Proposition 3

A similar result is proved in Dong and Rudi (2004, Lemma 1), who consider the case of traditional trans-shipment. Substituting $\mathbf{y^{DIP}}$ into the first order condition for $C^{LB}$ in Equation 13, we have:

$$(h + po - s_{12}) F_D\left(\sum_j y_j^{DIP}\right) + (s_{12} - s) F_{D_i}(y_i^{DIP}) + (p_s - p_o + s) F_{D_{is}}((y_i^{DIP}) - p_s$$
$$= (h + p_o - s_{12})\left(\Phi\left(z^{DIP} \sum_i \sigma_i / \sigma\right) - \Phi\left(z^{DIP}\right)\right)$$

The equality follows from the fact that $\mathbf{y^{DIP}}$ satisfies Equation 7, and the normality of demands, as we can write $y_i^{DIP} = \mu_i + z^{DIP}\sigma_i$, where $D_i \sim \mathcal{N}(\mu_i, \sigma_i)$, and $D \sim \mathcal{N}(\mu, \sigma)$. As $\sum_i \sigma_i / \sigma \geq 1$, it follow that the gradient of $C^{LB}$ at $\mathbf{y^{DIP}}$ is $\geq 0 (\leq 0)$ whenever $z^{DIP} \leq (\geq)\mu_i$. Also, writing $\sigma = \sqrt{\sum_i \sigma_i^2 + \sum_j 2\rho_l \sigma_i \sigma_j}$, where $\rho_l$ is the correlation coefficient between locations, $y^{DIP}$ is optimal to $C^{LB}$ and $C^{IIP}$ when $\rho_l = 1$. $\qquad \square$

### A.4. Proof of Proposition 4

Due to similarities to Dong and Rudi (2004), we have a similar solution where the optimal inventory at each location is at the same critical fractile of the location's demands. Equation 16 can be written as:

$$y_i^{LBN} = F_{D_{is}}^{-1}\left(\frac{m}{p_s - p_o + s}\right), \quad \forall i \in \mathcal{S}_{so} \tag{24}$$

where $m = p_s - (h + p_o - s) F_{D_{\mathcal{S}}}(\sum_{j \in \mathcal{S}} y_j^{LBN})$. Substituting Equation 24 into the definition of $m$, we have:

$$\sum_{j \in \mathcal{S}} F_{D_{is}}^{-1}\left(\frac{m}{p_s - p_o + s}\right) = F_{D_{\mathcal{S}}}^{-1}\left(\frac{p_s - m}{h + p_o - s}\right) \tag{25}$$

Solving this yields a unique solution for $m$, which in turn yields a unique solution $\mathbf{y^{LBN}}$, where each stores stocks at the same critical fractile of their in-store demand, as seen from Equation 24. $\qquad \square$

For OFCs $(i \in \mathcal{S}_o)$, $y_i^{LBN} = 0$, as otherwise, the value of $m$ is forced to be $p_s - p_o + s$, which renders Equation 24 to infinity.

### A.5.  Proof of Proposition 5

Consider a square of unit area in which $N$ stores are uniformly distributed. Let the square be divided into $\sqrt{N}$ identical cells, such that each cell contains $\sqrt{N}$ stores. The dimensions of each cell are thus $\frac{1}{N^{\frac{1}{4}}} \times \frac{1}{N^{\frac{1}{4}}}$. The superscript $l$ for a demand variable (e.g. $D_{is}^l$) denotes that the demand belongs to a store in cell $l$.

Let $C^{LB'}$ be the cost function obtained from $C^{IIP}$ by lowering all cross-shipping costs to the within-region shipping cost $s$. Let $C^{IIP_c}$ and $C^{LB'_c}$ be the functions obtained by restricting $C^{IIP}$ and $C^{LB'}$ respectively, so that cross-shipments can only be made between two stores belonging to the same cell. Clearly, $C^{IIP}(y) \le C^{IIP_c}(y)$ and $C^{LB'}(y) \le C^{LB'_c}(y)$ for any $y \ge 0$. Let $g(y, N)$ denote the cost incurred by $N$ stores starting with inventory $y$ each, without the option of cross-shipping:

$$g(y, N) = \sum_{i=1}^{N} \left[ h(y - D_i)^+ + p_s(D_{is} - y)^+ + p_o\left(D_{io} - (y - D_{is})^+\right)^+ + s\min\left(D_{io}, (y - D_{is})^+\right) \right]$$

Note that $g(y, N)$ represents the sum of costs incurred by individual stores, and hence, $\mathbb{E}g(y, N) = \mathbb{E}\sum_{l=1}^{\sqrt{N}} g(y, \sqrt{N}) = \sqrt{N}g(y, \sqrt{N})$. Let $CS_{ij}(y, N)$ denote the cross-shipped quantity between stores $i$ and $j$, when there are $N$ stores with order-up-to quantity $y$ each ($CS_{ij}^l$ when defined within a cell). Note that both the functions $g$ and $CS_{ij}$ also depend on the demand vector, but the dependency is ignored for notational convenience. As the cells are identical in terms of demands and costs, we have:

$$
\begin{aligned}
C^{IIP_c}(y^{IIPH}) = \quad & \mathbb{E}\left( \sum_{l=1}^{\sqrt{N}} \left[ g(y^{IIPH}, \sqrt{N}) + \sum_{i=1}^{\sqrt{N}} \sum_{j=1, j\neq i}^{\sqrt{N}} (s_{ij}^l - h - p_o) CS_{ij}^l(y^{IIPH}, \sqrt{N}) \right] \right) \\
= \quad & \mathbb{E}g(y^{IIPH}, N) + \mathbb{E}\left( \sum_{l=1}^{\sqrt{N}} \left( \sum_{i=1}^{\sqrt{N}} \sum_{j=1, j\neq i}^{\sqrt{N}} (s_{ij}^l - h - p_o) CS_{ij}^l(y^{IIPH}, \sqrt{N}) \right) \right)
\end{aligned}
$$

$$
\begin{aligned}
C^{LB'}(y^{IIPH}) = \quad & C^{LB'_c}(y^{IIPH}) \\
& + (s - h - p_o)\mathbb{E}\left[ \sum_{l=1}^{\sqrt{N}}\left( \sum_{i=1}^{\sqrt{N}} D_{io}^l - (y^{IIPH} - D_{is}^l)^+ \right)^+ - \left( \sum_{i=1}^{N} D_{io} - (y^{IIPH} - D_{is})^+ \right)^+ \right] \\
= \quad & \mathbb{E}g(y^{IIPH}, N) + \mathbb{E}\left( \sum_{l=1}^{\sqrt{N}} \left( \sum_{i=1}^{\sqrt{N}} \sum_{j=1, j\neq i}^{\sqrt{N}} (s - h - p_o) CS_{ij}^l(y^{IIPH}, \sqrt{N}) \right) \right) \\
& + (s - h - p_o)\left[ \sqrt{N}\mathbb{E}\left( \sum_{i=1}^{\sqrt{N}} D_{io}^l - (y^{IIPH} - D_{is}^l)^+ \right)^+ - \mathbb{E}\left( \sum_{i=1}^{N} D_{io} - (y^{IIPH} - D_{is})^+ \right)^+ \right]
\end{aligned}
$$

The expression for $C^{LB'}$ is written as the sum of $C^{LB'_c}$ which restricts cross-shipping to within each cell, and the cost of the additional cross-shipped units with this restriction removed. We know that $C^{LB}(y^{IIPH}) \le C^{LB'}(y^{IIPH}) \le C^{IIP}(y^{IIPH}) \le C^{IIP_c}(y)$, where the first inequality follows from Proposition 5. We first show that $\frac{C^{IIP_c}(y^{IIPH})}{C^{LB'}(y^{IIPH})} \to 1$ as $N \to \infty$. We have:

$$
\begin{aligned}
\frac{C^{IIP_c}(y^{IIPH})}{C^{LB'}(y^{IIPH})} - 1 = \; & \frac{\mathbb{E}\left( \sum_{l=1}^{\sqrt{N}} \left( \sum_{i=1}^{\sqrt{N}} \sum_{j=1, j\neq i}^{\sqrt{N}} (s_{ij}^l - s) CS_{ij}^l(y^{IIPH}, \sqrt{N}) \right) \right)}{C^{LB'}(y^{IIPH})} \\
& + \frac{(h + p_o - s)\left[ \sqrt{N}\mathbb{E}\left( \sum_{i=1}^{\sqrt{N}} D_{io}^l - (y^{IIPH} - D_{is}^l)^+ \right)^+ - \mathbb{E}\left( \sum_{i=1}^{N} D_{io} - (y^{IIPH} - D_{is})^+ \right)^+ \right]}{C^{LB'}(y^{IIPH})}
\end{aligned}
$$

We have $s_{ij}^l - s = f(d_{ij}^l) \leq f\left(\frac{\sqrt{2}}{N^{\frac{1}{4}}}\right)$, as the maximum distance within a cell is $\frac{\sqrt{2}}{N^{\frac{1}{4}}}$. Thus, using $C^{LB'}(y^{IIPH}) \geq$
$\mathbb{E}\left(\sum_{l=1}^{\sqrt{N}}\left(\sum_{i=1}^{\sqrt{N}}\sum_{j=1,j\neq i}^{\sqrt{N}}(s)CS_{ij}^l(y^{IIPH},\sqrt{N})\right)\right)$ for the first term, and $C^{LB'}(y^{IIPH}) \geq s\mu_o N$ for the second term,
we have

$$\frac{C^{IIP_c}(y^{IIPH})}{C^{LB'}(y^{IIPH})} - 1 \leq \frac{f\left(\frac{\sqrt{2}}{N^{\frac{1}{4}}}\right)}{s} + \left(\frac{h+p_o-s}{s\mu_o\sqrt{N}}\right)\mathbb{E}\left(\sum_{i=1}^{\sqrt{N}}D_{io} - \left(y^{IIPH}-D_{is}\right)^+\right)^+ \tag{26}$$

The first term on the right hand side vanishes to zero as $N \to \infty$, as $f(d) \to 0$ as $d \to 0$. To simplify the second term, we need the following lemmas.

LEMMA 1. *If $h < p_o - s$, then $y^{IIPH} > \mu$ where $\mu = \mu_s + \mu_o$, and if additionally $h < (p_s - p_o + s)F_s(\mu)$,*

$$y^{IIPH} \to F_s^{-1}\left(\frac{p_s - p_o + s - h}{p_s - p_o + s}\right) \in (0, \infty), \text{ as } N \to \infty \tag{27}$$

Proof: Lemma 1 is proved from the optimality equations of $C^{LBN}$ (Equation 16) for identical stores:

$$(h + p_o - s)\mathbb{P}\left(\sum_{i=1}^N D_i \leq Ny^{IIPH}\right) + (p_s - p_o + s)F_{D_{1s}}(y^{IIPH}) = p_s$$

From the above equation, when $h < p_o - s$, we have $p_s < 2(p_o - s)\mathbb{P}\left(\sum_{i=1}^N D_i \leq Ny^{IIPH}\right) + (p_s - p_o + s)$. This simplifies to yield $y^{IIPH} > \mu$. Now, by applying the central limit theorem as $N \to \infty$ and $y^{IIPH} > \mu$, $\mathbb{P}\left(\sum_{i=1}^N D_i/N \leq y^{IIPH}\right) \to 1$, and the result follows. Note that the asymptotic solution should also satisfy $y^{IIPH} > \mu$, which translates to the condition $h < (p_s - p_o + s)F_s(\mu)$. □

LEMMA 2. *When $h < p_o - s$ and $h < p_s - (p_o - s)$, and the demands are bounded above as $D_{is} \leq M_s$ and $D_{io} \leq M_o$ for all $i$,*

$$\mathbb{P}\left(\sum_{i=1}^{\sqrt{N}}D_{io} > \sum_{i=1}^{\sqrt{N}}\left(y^{IIPH}-D_{is}\right)^+\right) \leq \exp\left\{\frac{-2\sqrt{N}(y^{IIPH}-\mu)^2}{M_o+M_s}\right\} \tag{28}$$

Proof:
$$\mathbb{P}\left(\sum_{i=1}^{\sqrt{N}}D_{io} > \sum_{i=1}^{\sqrt{N}}\left(y^{IIPH}-D_{is}\right)^+\right) = \mathbb{P}\left(\sum_{i=1}^{\sqrt{N}}\left(D_i - \left(D_{is}-y^{IIPH}\right)^+\right) > \sqrt{N}y^{IIPH}\right) \leq \mathbb{P}\left(\sum_{i=1}^{\sqrt{N}}D_i > \sqrt{N}y^{IIPH}\right)$$

$$\leq \exp\left\{\frac{-2\sqrt{N}(y^{IIPH}-\mu)^2}{M_o+M_s}\right\} \to 0, \text{ as } N \to \infty$$

The final inequality follows from the Hoeffding bound for tail probabilities Hoeffding (1963), as $y^{IIPH} > \mu$ and demands are bounded, and the limit exists as $y^{IIPH}$ approaches a finite positive quantity as $N \to \infty$ by Lemma 1. The expectation in the second term of Equation 26 can be bounded as follows:

$$\mathbb{E}\left(\sum_{i=1}^{\sqrt{N}}\left(D_{io} - \left(y^{IIPH}-D_{is}\right)^+\right)\right)^+$$

$$= \mathbb{E}\left[\left(\sum_{i=1}^{\sqrt{N}}\left(D_{io} - \left(y^{IIPH}-D_{is}\right)^+\right)\right)^+ \middle| \sum_{i=1}^{\sqrt{N}}D_{io} > \sum_{i=1}^{\sqrt{N}}\left(y^{IIPH}-D_{is}\right)^+\right]\mathbb{P}\left(\sum_{i=1}^{\sqrt{N}}D_{io} > \sum_{i=1}^{\sqrt{N}}\left(y^{IIPH}-D_{is}\right)^+\right)$$

$$\leq \mathbb{E}\left[\sum_{i=1}^{\sqrt{N}}D_{io}\middle|\sum_{i=1}^{\sqrt{N}}D_{io} > \sum_{i=1}^{\sqrt{N}}\left(y^{IIPH}-D_{is}\right)^+\right]\mathbb{P}\left(\sum_{i=1}^{\sqrt{N}}D_{io} > \sum_{i=1}^{\sqrt{N}}\left(y^{IIPH}-D_{is}\right)^+\right)$$

$$\leq M_o\sqrt{N}\exp\left\{\frac{-2\sqrt{N}(y^{IIPH}-\mu)^2}{M_o+M_s}\right\}$$

The last inequality follows from Lemma 2 and the boundedness of the demands as $D_{is} \leq M_s$, and $D_{io} \leq M_o$ for all $i$ with $0 < M_s, M_o < \infty$. □

Thus, we have:

$$\frac{C^{IIP_c}(y^{IIPH})}{C^{LB'}(y^{IIPH})} \leq 1 + \frac{f\left(\frac{\sqrt{2}}{N^{\frac{1}{4}}}\right)}{s} + \left(\frac{h + p_o - s}{s\mu_o}\right)\left(M_o\sqrt{N}\exp\left\{\frac{-2\sqrt{N}(y^{IIPH} - \mu)^2}{M_o + M_s}\right\}\right) \tag{29}$$

$$\rightarrow 1, \; as \; N \rightarrow \infty$$

The next step is to show the $C^{LBN}$ is off by a constant factor from the $C^{LB'}$. From the proof of Proposition 2, the difference simplifies to:

$$C^{LB'}(y^{IIPH}) - C^{LBN}(y^{IIPH})$$

$$= (h + p_o - s)\mathbb{E}\left[\left(\sum_{i=1}^{N} D_{io} - \left(y^{IIPH} - D_{is}\right)^+\right)^+ + \sum_{i=1}^{N}\left(D_{is} - y^{IIPH}\right)^+ - \left(D - \sum_{i=1}^{N} y^{IIPH}\right)^+\right]$$

where $D = \sum_{i=1}^{N} D_{is} + D_{io}$.

Similar to what was done to bound the second term in Equation 26, we can show that whenever the conditions in Lemma 2 are satisfied, $\mathbb{E}\left(\sum_{i=1}^{N} D_{io} - (y^{IIPH} - D_{is})^+\right)^+ \leq M_o N \exp\left\{\frac{-2N(y^{IIPH} - \mu)^2}{M_o + M_s}\right\}$. Thus, we have:

$$C^{LB'}(y^{IIPH}) - C^{LBN}(y^{IIPH}) \leq (h + p_o - s)\left[M_o N \exp\left\{\frac{-2N(y^{IIPH} - \mu)^2}{M_o + M_s}\right\} + \sum_{i=1}^{N}\left(D_{is} - y^{IIPH}\right)^+\right]$$

Using $C^{LBN}(y^{IIPH}) \geq s\mu_o N$ and $C^{LBN}(y^{IIPH}) \geq (p_s - p_o + s)\sum_{i=1}^{N}(D_{is} - y^{IIPH})^+$, we have:

$$\frac{C^{LB'}(y^{IIPH})}{C^{LBN}(y^{IIPH})} - 1 \leq \left(\frac{h + p_o - s}{s\mu_o}\right)\left(M_o \exp\left\{\frac{-2N(y^{IIPH} - \mu)^2}{M_o + M_s}\right\}\right) + \left(\frac{h + p_o - s}{p_s - p_o + s}\right) \tag{30}$$

Thus, from Equations 29 and 30, as $N \rightarrow \infty$, we have

$$\frac{C^{IIP_c}(y^{IIPH})}{C^{LBN}(y^{IIPH})} \leq 1 + \frac{h + p_o - s}{p_s - p_o + s}$$

$$\Rightarrow \frac{C^{IIPH}}{C^{LBN}(y^{IIPH})} \leq \frac{h + p_s}{p_s - p_o + s}$$

The final step follows from $C^{IIP_c}(y^{IIPH}) \geq C^{IIP}(y^{IIPH}) = C^{IIPH}$. □

The result may hold subject to some generalizations, such as the unit square can be replaced with any finite area, and non-identical cells as long as the number of stores in each cell grows to infinity as $N \rightarrow \infty$. The resulting cases may call for a more complicated proof, and is outside the scope of this study.

### A.6. Proof of Proposition 6

Consider the single period case, where items are ordered at the start of the period, and online demands are fulfilled over $T$ fulfillment epochs. Assume that $C_{T+1}(\mathbf{x^{T+1}}, \tilde{D}^{T+1}) = 0$ without loss of generality. Thus, $C_T(\mathbf{x^T}, \tilde{D}^T)$ is given by a simple linear program which is jointly convex in $(\mathbf{x^T}, \tilde{D}^T)$. This leads to the base case result that $C_T(\mathbf{x^T}, \tilde{D}^T)$ is convex in $x^T$ given any $\tilde{D}^T$. By backward induction, we need to show that $C_t(\mathbf{x^t}, \tilde{D}^t)$ is convex in $\mathbf{x^t}$ for any given $\tilde{D}^t$, with the assumption that $C_{t+1}(\mathbf{x^{t+1}}, \tilde{D}^{t+1})$ is convex in $\mathbf{x^{t+1}}$ given any $\tilde{D}^{t+1}$. The cost-to-go function can be represented by $C_t(\mathbf{x^t}, \tilde{D}^t) = \min_{\mathbf{z^t}, \mathbf{Z^t} \in \Delta} \mathcal{G}(\mathbf{x^t}, \tilde{D}^t, \mathbf{z^t}, \mathbf{Z^t})$, where

$$\mathcal{G}(\mathbf{x^t}, \tilde{D}^t, \mathbf{z^t}, \mathbf{Z^t}) = \left[P(\mathbf{x^t}, \tilde{D}^t, \mathbf{z^t}, \mathbf{Z^t}) + \mathbb{E}C_{t+1}(x_i^t - z_i^t - \sum_{j=1}^{N} Z_{ij}^t, \tilde{D}^{t+1})\right] \tag{31}$$

Consider any $\mu \geq 0$, and $\mathbf{x_1^t}, \mathbf{x_2^t} \geq 0$. Let $(\mathbf{z_i^t}, \mathbf{Z_i^t}) = \underset{\mathbf{z^t}, \mathbf{Z^t} \in \Delta}{\arg \min} \mathcal{G}(\mathbf{x_i^t}, \tilde{D}^t, \mathbf{z^t}, \mathbf{Z^t})$. Note that $P$ is a linear function in its variables (Equation 2), and $\mathbb{E}C_{t+1}(\mathbf{x^{t+1}}, \tilde{D}^{t+1})$ is convex in $\mathbf{x^{t+1}}$, as expectation preserves convexity. Let $\bar{\mathbf{x}}^t = \mu \mathbf{x_1^t} + (1-\mu)\mathbf{x_2^t}$, $\bar{\mathbf{z}}^t = \mu \mathbf{z_1^t} + (1-\mu)\mathbf{z_2^t}$ and $\bar{\mathbf{Z}}^t = \mu \mathbf{Z_1^t} + (1-\mu)\mathbf{Z_2^t}$. We have:

$$
\begin{aligned}
C_t(\bar{\mathbf{x}}^t, \tilde{D}^t) &= \min_{\mathbf{z^t}, \mathbf{Z^t} \in \Delta} \left[ P(\bar{\mathbf{x}}^t, \tilde{D}^t, \mathbf{z^t}, \mathbf{Z^t}) + \mathbb{E}C_{t+1}(\bar{x}_i^t - z_i^t - \sum_{j=1}^N Z_{ij}^t, \tilde{D}^{t+1}) \right] \\
&\leq P(\bar{\mathbf{x}}^t, \tilde{D}^t, \bar{\mathbf{z}}^t, \bar{\mathbf{Z}}^t) + \mathbb{E}C_{t+1}(\bar{x}_i^t - \bar{z}_i^t - \sum_{j=1}^N \bar{Z}_{ij}^t, \tilde{D}^{t+1}) \\
&\leq \mu P(\mathbf{x_1^t}, \tilde{D}^t, \mathbf{z_1^t}, \mathbf{Z_1^t}) + (1-\mu)P(\mathbf{x_2^t}, \tilde{D}^t, \mathbf{z_2^t}, \mathbf{Z_2^t}) + \mathbb{E}C_{t+1}(\bar{x}_i^t - \bar{z}_i^t - \sum_{j=1}^N \bar{Z}_{ij}^t, \tilde{D}^{t+1})
\end{aligned}
\tag{32}
$$

The first inequality follows from the feasibility of $\bar{\mathbf{z}}^t, \bar{\mathbf{Z}}^t$ in $\Delta$, as $(\mathbf{z_1^t}, \mathbf{Z_1^t})$ and $\mathbf{z_2^t}, \mathbf{Z_2^t})$ are feasible in $\Delta$. The second inequality follows from the convexity of $P$. As $\mathbb{E}C_{t+1}(\mathbf{x^{t+1}}, \tilde{D}^{t+1})$ is convex in $\mathbf{x^{t+1}}$, we have:

$$
\begin{aligned}
\mathbb{E}C_{t+1}\left(\bar{x}_i^t - \bar{z}_i^t - \sum_{j=1}^N \bar{Z}_{ij}^t, \tilde{D}^{t+1}\right) &= \mathbb{E}C_{t+1}\left[ \mu \left( x_{1,i}^t - z_{1,i}^t - \sum_{j=1}^N Z_{1,ij}^t \right) + (1-\mu) \left( x_{2,i}^t - z_{2,i}^t - \sum_{j=1}^N Z_{2,i}^t \right), \tilde{D}^{t+1} \right] \\
&\leq \mu \mathbb{E}C_{t+1}\left[ x_{1,i}^t - z_{1,i}^t - \sum_{j=1}^N Z_{1,ij}^t, \tilde{D}^{t+1} \right] + (1-\mu)\mathbb{E}C_{t+1}\left[ x_{2,i}^t - z_{2,i}^t - \sum_{j=1}^N Z_{2,i}^t, \tilde{D}^{t+1} \right]
\end{aligned}
\tag{33}
$$

Thus, from Equation 31, we have:

$$
\begin{aligned}
C_t(\bar{\mathbf{x}}^t, \tilde{D}^t) &\leq \mu \mathcal{G}(\mathbf{x_1^t}, \tilde{D}^t, \mathbf{z_1^t}, \mathbf{Z_1^t}) + (1-\mu)\mathcal{G}(\mathbf{x_2^t}, \tilde{D}^t, \mathbf{z_2^t}, \mathbf{Z_2^t}) \\
&= \mu C_t(\mathbf{x_1^t}, \tilde{D}^t) + (1-\mu)C_t(\mathbf{x_2^t}, \tilde{D}^t)
\end{aligned}
\tag{34}
$$

The equality follows from the definitions of $(\mathbf{z_1^t}, \mathbf{Z_1^t})$ and $(\mathbf{z_2^t}, \mathbf{Z_2^t})$. $\qquad\square$

### A.7. Proof of Proposition 7

Let the single period cost function be given by $C^{IIP}(\mathbf{y}) = \mathbb{E}C_1(\mathbf{y}, \tilde{D})$, and let $\mathbf{y}^{IIP}$ be the optimal solution. When the initial level of inventory $x_i$ at region $i$ before ordering, the cost function is as follows:

$$
V^{IIP}(\mathbf{x}) = \min_{\mathbf{y} \geq \mathbf{x}} C^{IIP}(\mathbf{y}) = C^{IIP}(\mathbf{y}^{IIP})
\tag{35}
$$

As $\mathbf{y}^{IIP}$ minimizes the function $C^{IIP}$, for any $\{\mathbf{x} : \mathbf{x} \leq \mathbf{y}^{IIP}\}$, it is optimal to order up to $\mathbf{y}^{IIP}$. We ignore cases where $x_i > y_i^{IIP}$ for some $i$, as eventually the system is brought to the state $\mathbf{x} \leq \mathbf{y}^{IIP}$.

For the multiple period case, we have M time periods: $m = 1, 2, .., M$. The in-store demands $\{D_{is}^m, m > 0\}$ and online demands $\{D_{io}^m, m > 0\}$ are assumed to be i.i.d. The available inventory at the end of a review period serves as the initial inventory for the next review period, and we assume zero purchasing costs. The discount factor is $\delta \in (0, 1]$.

The proof is by induction, and similar to the proof of Proposition 4 in van Mieghem and Rudi (2002). If we show that $V_m^{IIP}(\mathbf{x}^m)$, the expected cost-to-go function evaluated in review period $m$ with the initial inventory $\mathbf{x}^m$, is convex and affine, a stationary base stock policy would be optimal. For the $M + 1^{th}$ review period, the cost function is $V_{M+1}^{IIP}(\mathbf{x}^{M+1}) = 0$ (assuming zero purchasing costs) which is trivially convex and affine in $\mathbf{x}^{M+1}$. Let $V_{m+1}^{IIP}$ be convex and affine in $\mathbf{x}^{m+1}$. The cost function for review period $m$ is:

$$
V_m^{IIP}(\mathbf{x}) = \min_{\mathbf{y} \geq \mathbf{x}} \left[ C^{IIP}(\mathbf{y}) + \delta \mathbb{E}V_{m+1}^{IIP}\left( f\left(\mathbf{y}, \tilde{D}\right) \right) \right] = \min_{\mathbf{y} \geq \mathbf{x}} U_m^{IIP}(\mathbf{y})
\tag{36}
$$

where $f(\mathbf{y}, D)$ is the vector of ending inventories. $\tilde{D}$ is the demand vector constituting the in-store and online demands for both the regions. As taking expectation preserves convexity, and the sum of convex functions is convex, $U_m^{IIP}(\mathbf{y})$ is convex in $\mathbf{y}$. It only remains to be shown that $V_m^{IIP}$ is affine in $\mathbf{x}$. To show this, consider any $\mathbf{y} \le \mathbf{y}^{IIP}$, so that $f\left(\mathbf{y}, \tilde{D}\right) \le \mathbf{y} \le \mathbf{y}^{IIP}$. We have

$$U_m^{IIP}(\mathbf{y}) = C^{IIP}(\mathbf{y}) + \delta \mathbb{E} V_{m+1}^{IIP}\left(f\left(\mathbf{y}, \tilde{D}\right)\right) = C^{IIP}(\mathbf{y}) + \delta \mathbb{E} V_{m+1}^{IIP}\left(\mathbf{y}^{IIP}\right) \tag{37}$$

as $V_{m+1}^{IIP}$ is affine in $\mathbf{x}^{m+1}$ and the purchasing cost is zero. Clearly, $\mathbf{y} = \mathbf{y}^{IIP}$ minimizes $U_m^{IIP}$ for $\mathbf{y} \le \mathbf{y}^{IIP}$. Thus, $V_m^{IIP}(\mathbf{x}) = \max_{\mathbf{y} \ge \mathbf{x}} U_m^{IIP}(\mathbf{y})$ is affine (constant) in $\mathbf{x}$ for all $\mathbf{x} \le \mathbf{y}^{IIP}$, and hence a stationary base-stock policy $\mathbf{y}^{IIP}$ is optimal if $\mathbf{x} \le \mathbf{y}^{IIP}$. If there is some $i$ for which $x_i > y_i^{IIP}$, the optimal policy will be more complicated, but eventually, the system comes back to $\mathbf{x} \le \mathbf{y}^{IIP}$. $\qquad \square$

## Appendix B:   Demand Regions for the IIP Solution

We illustrate the identification of demand regions in which the dual vector $\lambda$ is constant (as discussed in Section 3.1.3) and the calculation of the corresponding probabilities. For any given $(y_1, y_2)$, the demand space $(D_{1s}, D_{1o}, D_{2s}, D_{2o})$ can be divided into a number of independent regions. Based on the values taken by the variables in the optimal solution in (20), Table 1 shows the different cases that are possible given $y_1$ and $y_2$. From these cases, the independent demand regions are listed in Table 2 along with the constant dual prices in those regions. The underlined cases are redundant, and can be discarded while calculating the probability for each region. The dual prices $\lambda_1, \lambda_2$ are the shadow prices of the constraints which contain $y_1$

**Table 1**    Table showing the various demand cases based on the values of $y_1, y_2$

|   | A | B | C | D |
|---|---|---|---|---|
| 1 | $y_1 < D_{1s}$ | $D_{1s} \le y_1 < D_1$ | $D_1 \le y_1 < D_1 + D_{2o}$ | $y_1 \ge D_1 + D_{2o}$ |
| 2 | $y_2 < D_{2s}$ | $D_{2s} \le y_2 < D_2$ | $D_2 \le y_2 < D_2 + D_{1o}$ | $y_2 \ge D_2 + D_{1o}$ |
| 3 | $y_1 + y_2 < D_1 + D_2$ | $y_1 + y_2 \ge D_1 + D_2$ | | |

**Table 2**    Table showing the various demand regions and the corresponding constant dual-prices.

| Region | Case | $\lambda_1$ | $\lambda_2$ | Region | Case | $\lambda_1$ | $\lambda_2$ |
|---|---|---|---|---|---|---|---|
| $\Omega_1$ | A1,A2,$\underline{A3}$ | $h+p_s$ | $h+p_s$ | $\Omega_{11}$ | C1,A2,$\underline{A3}$ | $h+p_o-s_{12}$ | $h+p_s$ |
| $\Omega_2$ | A1,B2,$\underline{A3}$ | $h+p_s$ | $h+p_o-s$ | $\Omega_{12}$ | C1,B2,A3 | $h+p_o-s_{12}$ | $h+p_o-s$ |
| $\Omega_3$ | A1,C2,$\underline{A3}$ | $h+p_s$ | $h+p_o-s_{12}$ | $\Omega_{13}$ | C1,B2,B3 | $0$ | $s_{12}-s$ |
| $\Omega_4$ | $\underline{A1}$,D2,A3 | $h+p_s$ | $0$ | $\Omega_{14}$ | C1,C2,$\underline{B3}$ | $0$ | $0$ |
| $\Omega_5$ | A1,$\underline{D2}$,B3 | $h+p_s$ | $0$ | $\Omega_{15}$ | C1,D2,$\underline{B3}$ | $0$ | $0$ |
| $\Omega_6$ | B1,A2,$\underline{A3}$ | $h+p_o-s$ | $h+p_s$ | $\Omega_{16}$ | D1,$\underline{A2}$,A3 | $0$ | $h+p_s$ |
| $\Omega_7$ | B1,B2,$\underline{A3}$ | $h+p_o-s$ | $h+p_o-s$ | $\Omega_{17}$ | $\underline{D1}$,A2,B3 | $0$ | $h+p_s$ |
| $\Omega_8$ | B1,C2,A3 | $h+p_o-s$ | $h+p_o-s_{12}$ | $\Omega_{18}$ | D1,B2,$\underline{B3}$ | $0$ | $s_{12}-s$ |
| $\Omega_9$ | B1,C2,B3 | $s_{12}-s$ | $0$ | $\Omega_{19}$ | D1,C2,$\underline{B3}$ | $0$ | $0$ |
| $\Omega_{10}$ | B1,D2,$\underline{B3}$ | $s_{12}-s$ | $0$ | $\Omega_{20}$ | D1,D2,$\underline{B3}$ | $0$ | $0$ |

and $y_2$ respectively, namely the first set of constraints $z_i + \sum_{j=1}^{2} z_{ij} \leq y_i, \forall i$ in the linear program in (20), and can be obtain in a standard fashion from linear programming theory. For example, for the demand regions with the case D1, that is, $y_1 \geq D_1 + D_{2o}$, irrespective of the value of $y_2$, there will be inventory left over at retail store 1 at the end of the period. Thus the constraint $z_1 + \sum_{j=1}^{2} z_{1j} \leq y_1$ will not bind, and hence $\lambda_1 = 0$.

The probability for each region is calculated as follows, when demands follow normal distributions. The region is expressed as an inequality of the form $R_k \tilde{D} <= S_k Y$, where $\tilde{D} = [D_{1s}, D_{1o}, D_{2s}, D_{2o}]^{\intercal}$ and $Y = [y_1, y_2]^{\intercal}$. For example, $\Omega_3 = (A1, C2) = \{y_1 < D_{1s}, D_2 \leq y_2 < D_2 + D_{1o}\}$. This can be expressed as:

$$
\begin{bmatrix} -1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 1 \\ 0 & -1 & -1 & -1 \end{bmatrix} \begin{bmatrix} D_{1s} \\ D_{1o} \\ D_{2s} \\ D_{2o} \end{bmatrix} \leq \begin{bmatrix} -1 & 0 \\ 0 & 1 \\ 0 & -1 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix}
$$

$R_k \tilde{D}$ is multivariate normal with mean $R_k \mu$ and covariance matrix $R_k \Sigma \Sigma^{\intercal} R_k^{\intercal}$, where $\mu$ and $\Sigma$ are the mean and covariance matrices of $\tilde{D}$. The probability of region $k$ reduces to evaluating the cumulative distribution function of $A_k \tilde{D}$ at $B_k Y$. For general demand distributions, numerical methods have to be employed.

## Appendix C:    Additional Details for Numerical Analyses

All numerical analyses were done on a desktop computer (i7-3770 CPU @3.7GHz, 16GB RAM). The total market is assumed to be the top 300 most populous cities in mainland US. The demands for the OFCs are calculated based on the population not covered by omnichannel stores. This online demand is allocated to each OFC based on the optimal throughput rates estimated by Chicago Consulting (2013).

### C.1.    Simulation Procedure

A brief overview of the simulation is listed below:

1. The parameters for demands in each fulfillment epoch are calculated based on review-period demands estimated from population data. The starting inventory level vectors $\mathbf{y^{DIP}}$ and $\mathbf{y^{IIPH}}$ are calculated using the demand information based on Equation 7 and Algorithm 1 respectively.

2. We generate a sample of size $10^4$, where each sample is a realization of demands in a review period, although fulfillment decisions in each fulfillment epoch are made without knowing future demands. For each sample, we iterate over steps 3-8, and take the sample averages as approximations for expectations.

3. The fulfillment thresholds for the TF policy are calculated based on Equation 19. For the MF policy, these thresholds are set to zero.

4. For $t = 1, \ldots, T$, iterate over steps 6-7. The starting inventory levels are set based on the inventory policy followed (IIPH or DIP).

5. Implement Algorithm 2 based on the fulfillment policy followed (MF or TF) and the corresponding thresholds calculated in Step 3.

6. At the end of each fulfillment epoch, the holding, penalty and fulfillment costs are calculated. The ending inventory at a location becomes the starting inventory for the next epoch.

7. The total cost in a review period is the sum of the costs in each fulfillment epoch in that period.