

Tilburg University

On-demand last-mile distribution network design with omnichannel inventory

Snoeck, André; Winkenbach, Matthias; Fransoo, Jan C.

Published in:

Transportation Research Part E: Logistics and Transportation Review

DOI:

[10.1016/j.tre.2023.103324](https://doi.org/10.1016/j.tre.2023.103324)

Publication date:

2023

Document Version

Peer reviewed version

[Link to publication in Tilburg University Research Portal](#)

Citation for published version (APA):

Snoeck, A., Winkenbach, M., & Fransoo, J. C. (2023). On-demand last-mile distribution network design with omnichannel inventory. *Transportation Research Part E: Logistics and Transportation Review*, 180, Article 103324. <https://doi.org/10.1016/j.tre.2023.103324>

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

On-Demand Last-Mile Distribution Network Design with Omnichannel Inventory

André Snoeck^a, Matthias Winkenbach^a, Jan C. Fransoo^b

^aMassachusetts Institute of Technology, Center for Transportation and Logistics, 77 Massachusetts Avenue, Cambridge, MA 02139, USA

^bTilburg University, Tilburg School of Economics and Management, Warandelaan 2, NL-5000 LE Tilburg, The Netherlands

Abstract

E-commerce delivery deadlines are getting increasingly tight, driven by a growing ‘I-want-it-now’ instant gratification mindset of consumers and the desire of online and omnichannel retailers to capitalize on the growth of on-demand e-commerce. On-demand deliveries with delivery deadlines as tight as one or two hours force companies to rethink their last-mile distribution network, since tight delivery deadlines require decentralization of order picking and inventory holding to ensure close proximity to consumers. This fundamentally changes the strategic design process of last-mile distribution networks. We study the impact of incorporating inventory order-up-to level decisions into the strategic design process of last-mile distribution networks with tight delivery deadlines. We develop an approximate inventory model by including an estimate of the cost of late delivery and additional transportation due to local stock-outs in a newsvendor formulation. Such local stock-outs require an order to be delivered from a more distant facility, which may lead to late delivery and additional transportation cost. We integrate our approximate inventory model and a location-allocation mixed-integer program that determines optimal facility locations, associated order-up-to inventory levels, and fleet composition, into a metamodel simulation-based optimization approach. Our numerical analyses demonstrate that pooling the additional online inventory with brick-and-mortar (B&M) inventories leads to cannibalization by the B&M network and higher B&M service levels. However, the pooling benefits to the online network outweigh the cost of inventory cannibalization. Furthermore, we show under which circumstances omnichannel retailers may have an incentive to consolidate online inventory in specific B&M facilities.

Keywords: last-mile logistics; strategic network design; omnichannel; simulation-based optimization; newsvendor model

1. Introduction

E-commerce is growing rapidly on a global scale. Worldwide retail e-commerce sales may increase to \$6.54 trillion by the end of 2023, up from \$3.54 trillion in 2019 (eMarketer 2019). To capitalize on this opportunity, online and omnichannel retailers recognize high-service delivery as an important factor to boost brand perception and gain

Email addresses: asnock@mit.edu (André Snoeck), mwinkenb@mit.edu (Matthias Winkenbach), Jan.Fransoo@tilburguniversity.edu (Jan C. Fransoo)

market share (Zebra Technologies 2018). For example, Nike stopped selling shoes directly through Amazon, to retake control over the entire online purchasing process, including their interaction with consumers at the ordering and delivery stages (Novy-Williams and Soper 2019). In addition, driven by the rise of the on-demand economy and the associated ‘I-want-it-now’ instant-gratification mindset, 78% of logistics companies expect to provide same-day delivery by 2023, while 39% even anticipate delivery within two hours by 2028 (Colby and Bell 2016, Zebra Technologies 2018). Indeed, Amazon Prime Now, JD Express, and Instacart Express are examples of e-commerce companies promising one-hour delivery already today, while Mediamarkt and Gucci are examples of traditional retail companies offering similarly tight delivery deadlines (Farfetch 2017, MediaMarkt 2020).

Such tight delivery deadlines force companies to rethink their last-mile distribution network. In particular, tight delivery deadlines require decentralization of order picking and inventory holding to ensure close proximity to consumers expecting high-service delivery. As a result, companies increasingly move their logistics facilities into city centers. For example, Amazon is investing heavily in so-called Prime Now hubs to support one- and two-hour delivery (MVPL 2020). Alternatively, companies leverage new or existing stores to support highly responsive distribution of online orders. For example, Target uses its new small-format store concept for order pick-up and same-day delivery (Redman 2019), while on-demand grocery service Gorillas collaborates with grocery retailer Tesco to offer 10-minute delivery from supermarkets and “dark stores” in affluent residential areas of large cities. Indeed, with the rise of omnichannel strategies, brick-and-mortar (B&M) and online retail are increasingly integrated to accommodate consumer shopping behavior (Bell et al. 2014). Enabling stores for e-commerce fulfillment provides an opportunity to (1) fulfill orders from existing facilities and reduce the investment into network assets, e.g., new warehouse space, and (2) pool safety stock between the online and offline channel and reduce the investment in additional inventory (Melacini et al. 2018). However, pooling inventory for online and offline demand can lead to on-shelf availability and priority issues when a product is ordered online while a walk-in consumer is in the process of purchasing the product (Ferne and Grant 2008). Furthermore, the retail store set-up is typically not optimized for online order fulfillment, leading to additional operational costs due to picking and shipping inefficiencies (Hübner et al. 2015). Nonetheless, the potential of leveraging B&M networks increases for last-mile distribution networks with tight delivery deadlines since such deadlines are often promised in dense urban areas. These are often the same areas where traditional B&M stores are located, and required investments in additional network assets are high due to competition for space.

However, distributing inventory and order processing capacity over decentralized satellite facilities (SFs) leads to local inventory effects, which should be accounted for in the design and associated inventory policy of distribution networks. If delivery deadlines become increasingly tight, aggregate inventory availability across the last-mile distribution network alone might not suffice to satisfy consumer promises. Particularly when the time available for delivery is short or the city is large or congested, assigning an order to a facility with the required inventory on hand at the other end of town may preclude on-time delivery of the order. Decentralized inventory in close proximity to consumer demand, leveraging both dedicated facilities and existing B&M stores, may reduce delivery distances and times, and thus increase the time available for order processing, increase the potential for order consolidation, and reduce the

likelihood of late delivery.

A local-stock out at a single facility may negatively impact the performance at other locations in the network in two ways. First, if all facilities that would be able to ship before the promised delivery deadline run out of inventory, the order must be shipped from a more distant facility. This leads to late delivery and associated late delivery costs. These could involve a discount to the consumer or a loss in consumer satisfaction, often reflected in metrics such as the net promoter score (NPS). Second, shipping orders from a more distant facility impacts transportation cost through distance increases, throughput reductions, and consolidation limitations. Transportation agents need to travel further distances, thus increasing the distance-based component of delivery cost. The increase in required delivery time also leads to a reduction in the throughput of transportation agents, as well as to a reduction in consolidation potential, leading to higher transportation capacity requirements. Traditionally, the inventory trade-off considers global stock-out effects, balancing the cost of lost sales and excess inventory. Hence, in last-mile distribution networks with tight delivery deadlines and decentralized inventory, the trade-off should be extended with local stock-out considerations. This fundamentally changes the network design perspective.

The impact of a local inventory stock-out is directly affected by the number and location of facilities in the network (*facility activation*) and by the composition of the fleet. If more facilities are available in the network, and another facility is close to the facility stocking out, this alternative facility can be used to satisfy the order. Obviously, this may lead to additional transportation cost due to increased distance, reduced throughput, or reduced consolidation. If more transport capacity is available in the network, such use of alternative facilities can be served using the upfront fixed investment, i.e., through scheduled employees and proprietary or contracted vehicles. Hence, the occurrence of local inventory effects structurally changes the strategic network design decision problem. Solving such complex network design problems with extremely high service requirements, requires new conceptual models with advanced solution strategies. In this paper, we develop a two-stage analytical network design model to support the design of last-mile distribution networks while accounting for such local-inventory effects. In the first stage, we build upon the newsvendor model by incorporating local inventory effects when determining the optimal order-up-to inventory level. While the traditional newsvendor model is based on global inventory effects, i.e., lost-sales and excess inventory cost, we additionally account for late delivery and additional transportation cost. In the second stage, we integrate our newsvendor formulation with a location-allocation mixed-integer linear program (MILP) that determines the optimal facility activation decisions, facility basestock levels, and scheduled transportation capacity. We integrate this two-stage analytical network design model into a metamodel-based simulation-based optimization (SO) approach to support the strategic design of last-mile distribution networks with tight delivery deadlines.

The methodological contribution of our paper is two-fold and centered around a novel model formulation that addresses the new problem introduced above. First, building on established newsvendor considerations, we incorporate tactical inventory decisions for both the B&M and the online channel and their inter-dependencies with strategic facility activation decisions in a last-mile network design problem. Specifically, our work captures the complex trade-offs in serving two concurrent retail channels from potentially integrated (omnichannel) inventory positions. This novel

model formulation allows us to provide important managerial insights by showing analytically that there are conditions under which the integration of channel inventories leads to a situation where the orders of one channel are served at the expense of the service level of the other channel (see Section 5). This is particularly relevant for an increasing number of retailers seeking to leverage B&M stores as shared fulfillment locations across channels. Such shared assets are often a critical enabler of their on-demand delivery strategies in major metropolitan markets in which suitable real-estate is scarce and costly. Second, unlike most inventory-location models in the literature, our proposed model allows for orders to be fulfilled from any facility in the network as online orders arrive and operational order allocation decisions are made. This is a particularly relevant feature in the context of last-mile distribution with very short delivery deadlines (e.g., for on-demand sub-same-day delivery) from a network of stores or store-like facilities with limited inventory positions due to tight space constraints. As retailers seek to make larger parts of their assortments available to consumers on-demand, the inherently tight space constraints of urban fulfillment locations such as B&M stores are becoming an increasingly important limitation to the delivery experience these retailers can cost-efficiently and reliably provide to their consumers.

Our methodological contributions are complemented with extensive numerical experiments inspired by real-world data that underline their real-world impact and relevance. Specifically, our numerical analyses (see Section 6) indicate that the two-stage metamodel SO approach we propose provides value in terms of an improved cost performance of the final network design, a faster speed of algorithmic convergence, and a greater inter-restart consistency of the model results. Moreover, we demonstrate the value of our proposed modeling and solution approach in capturing multiple real-world operational complexities, while maintaining the core structure of the newsvendor model. Lastly, we show how our proposed modeling and solution approach can help to assess the impact of service level cannibalization effects between online and B&M channels in case of shared inventories.

The remainder of this paper is structured as follows. In Section 2, we review the relevant literature on last-mile distribution network design and inventory planning in omnichannel fulfillment, as well as SO approaches. In Section 3, we formulate a two-stage analytical metamodel that builds on a newsvendor model formulation while accounting for local stock-out effects with a location-allocation MILP. In Section 4, we then propose a SO solution approach for our model. Since we explicitly incorporate the inventory pooling benefit of leveraging B&M stores in our newsvendor formulation, we can analytically derive conditions for when it is cost-optimal to integrate online and B&M inventories, which we present in Section 5. In Section 6, we present and discuss the results of our extensive numerical analyses inspired by a real-world case study in the fashion retail industry. We conclude in Section 7.

2. Related Literature

In the following, we position our work in the extant literature on last-mile network design and inventory optimization in omnichannel distribution, as well as the relevant optimization methods.

2.1. Last-Mile Distribution Network Design

The growth of e-commerce and growing concerns about the impact of freight on inner-city livability in recent years sparked an increased attention to the strategic design of last-mile distribution networks in the literature (see, e.g., [Crainic et al. 2004](#), [Boccia et al. 2011](#), [Winkenbach et al. 2016a,b](#), [Janjevic et al. 2019](#), [Snoeck and Winkenbach 2020](#)). Contrary to distribution networks with tight delivery deadlines, traditional urban distribution networks are typically operated as periodic order fulfillment systems. Here, the order collection period and the order delivery period are separated by means of an order cut-off time. After the cut-off time, an operational plan is constructed to deliver the accrued orders, which is referred to as the day-before planning problem ([Crainic et al. 2009](#)). This design assumption limits the applicability of existing strategic design models to networks with tight delivery deadlines. However, we observe a growing prevalence of last-mile distribution networks with tight delivery deadlines, both in industry and academic discourse ([Savelsbergh and Van Woensel 2016](#)). In such networks, the required order processing and delivery lead-times are relatively long compared to the time available until the delivery deadline. This sense of urgency in order fulfillment reduces the potential for order consolidation, both during order processing at the facility and during last-mile delivery. Moreover, it increases the susceptibility of such networks to order processing delays due to facility congestion. Consequently, the design and planning problems that emerge in the context of last-mile distribution with tight delivery deadlines are typically of highly stochastic and dynamic nature. While a growing body of literature studies these problems (see, e.g., [Voccia et al. 2019](#), [Klapp et al. 2018](#), and references therein), there are, to the best of our knowledge, only two studies that focus on the strategic design of such networks. First, [Ulmer \(2017\)](#) studies the effect of the network design on the performance of last-mile distribution networks with tight delivery deadlines. The author conducts a simulation study to explore the effect of the tightness of delivery deadlines on the cost and layout of the distribution network. Second, [Snoeck and Winkenbach \(2022\)](#) propose a metamodel SO approach for the strategic design of last-mile distribution networks with tight delivery deadlines. The authors explicitly incorporate facility congestion effects in order processing at the distribution facilities when determining their optimal location and processing capacity and the optimal transportation fleet composition. None of these authors incorporates inventory considerations in the strategic design of such networks. In this paper, we expand on the work of [Snoeck and Winkenbach \(2022\)](#) to address this gap in the literature.

[Savelsbergh and Van Woensel \(2016\)](#) argue that inventory decisions are typically not explicitly considered in traditional urban network design problems, as the nature of the day-before planning problem ensures that there is enough time to make goods available at urban distribution facilities. However, [Amiri-Aref et al. \(2018\)](#), [Daskin et al. \(2002\)](#), [Shen et al. \(2003\)](#), and [Shen and Qi \(2007\)](#), amongst others, show that the cost optimal design of a network of distribution facilities and the optimal allocation of inventory to these facilities are tightly interrelated problems that need to be tackled in an integrated manner. [Ozsen et al. \(2009\)](#) further show that the benefits from integrating network design and inventory decisions are even more pronounced if customer demands can be fulfilled from multiple distribution facilities, depending on inventory availability. This typically applies to e-commerce fulfillment networks offering tight delivery deadlines, especially if they are leveraging an exiting network of B&M facilities. As pointed out

by [Amiri-Aref et al. \(2018\)](#), in such a setting, the available inventory levels and thus the observed demand at a given facility at a given time depends on the inventory allocation decisions. This gives rise to computational challenges due to non-linear effects when optimizing for both inventory and facility location decisions simultaneously.

In this paper, we address this issue by proposing a two-stage analytical model and a metamodel-based SO solution approach that (i) integrates strategic facility activation and fleet composition decisions with tactical decisions on local inventory order-up-to levels, and (ii) is able to account for non-linear interdependencies between these decisions. Moreover, the proposed model allows for the integration of existing B&M facilities and inventories in the strategic design of last-mile distribution networks with tight delivery deadlines.

2.2. Network Design and Inventory Decisions in Omnichannel Fulfillment

Driven by changing customer preferences towards instant gratification and on-demand consumption, many retailers are pursuing omnichannel retail strategies that require them to re-configure their supply and distribution networks ([Lim and Winkenbach 2019](#)). Due to the rise of e-commerce and mobile internet, consumers today are using a variety of online and offline channels throughout their interaction with a company and its products - from information discovery, to placing an order, to choosing a product fulfillment option. In omnichannel retailing, an increasing variety of sales channels - from traditional B&M stores to online stores to mobile apps or social media - are supported by the same supply chain and distribution infrastructure. Here, a key challenge consists in the fact that for any given order, the chosen fulfillment option by the consumer is not necessarily aligned with the sales channel that generated the order ([Bell et al. 2014](#)). Specifically, retailers are now facing complex operational trade-offs when deciding whether to fulfill e-commerce orders from dedicated e-commerce facilities, from the same warehouses that support conventional B&M retail stores, from the actual B&M stores themselves, or from a combination of these sites ([Millstein and Campbell 2018](#)). These operational trade-offs also have immediate repercussions on the optimal strategic design and tactical planning of distribution networks that are intended to optimally serve demand across multiple channels. For instance, [Janjevic et al. \(2021\)](#) propose a continuum approximation-based MILP formulation to solve a three-echelon location routing problem (LRP) for the strategic design of multi-modal last-mile delivery networks for omnichannel retailers with time-differentiated delivery services and a variety of in-person and unattended product-exchange options. Their model does not consider any tactical inventory positioning decisions throughout the network. Similarly, [Snoeck and Winkenbach \(2022\)](#) consider the optimal strategic design of a ship-from-store network to serve time-critical online orders of an omnichannel retailer, considering facility processing capacity and transportation decisions, but abstracting from inventory decisions. However, the overall literature on the (strategic) design of omnichannel distribution networks is still sparse ([Melacini et al. 2018](#)), and to the best of our knowledge, none of the existing works in the literature jointly consider strategic network design decisions, operational order fulfillment, and tactical inventory decisions for omnichannel retailers with tight delivery deadlines in the online channel. A potential reason for the lack of academic contributions addressing the integrated design of omnichannel distribution networks lies in the significant additional complexities that arise in the modeling and optimization of such networks. A major driver of these com-

plexities is the need to integrate location-specific inventory decisions into the network design problem. Outside of the omnichannel context, [Miranda and Garrido \(2004\)](#) present a non-linear MILP model and heuristic solution approach to integrate inventory control decisions in the strategic design of a distribution network. However, their model does not capture the complexities of a multi-tiered network of distribution centers (DCs) and stores and the competition between offline and online demand faced in an omnichannel setting. [Acimovic and Graves \(2017\)](#) point out that - unlike many conventional single-channel distribution problems - the decision which facility to ship an order from in an omnichannel distribution network can no longer be decoupled from the decision how much inventory to hold in each facility of the network.

Since the early days of online retailing, a growing body of literature has emerged that focuses on how to fulfill online orders from a given set of distribution facilities (see, e.g., [Acimovic and Graves 2015](#), [Jasin and Sinha 2015](#)). Several authors also study replenishment policies tailored to online distribution networks. For instance, [Acimovic and Graves \(2017\)](#) study demand spillover between fulfillment centers in an e-commerce setting. Based on real data, they show how local stock-outs influence the demand distribution at other facilities and consequently the performance of the inventory policy. [Govindarajan et al. \(2021\)](#) indicate that online fulfillment decisions from multiple facilities with flexible demand allocation are similar to transshipment decisions, except that items are shipped directly to the customer. [Yang and Qin \(2007\)](#) refer to such reactive transshipments with zero lead time as ‘virtual lateral transshipments’, and a growing body of literature studies their effect on inventory decisions (see, e.g., [Paterson et al. 2011](#), for a review).

As offline and online channels are becoming increasingly interconnected, several authors began to study how retailers’ operational decisions, such as pricing and inventory decisions, are affected by the introduction of an omnichannel strategy. For instance, [Bendoly et al. \(2007\)](#) study if inventory to support online sales should be stored in a central warehouse, or decentralized in B&M stores. They conclude that the attractiveness of decentralized inventory reduces with increasing online demand. [Guo and Keskin \(2018\)](#) study the optimal supply chain strategy for a dual-channel retailer that operates a physical and a web-based store. They determine the optimal degree of channel integration using a two-stage stochastic programming model, and the associated optimal order-up-to-level of the physical store based on a newsvendor-approach. [Wei and Li \(2020\)](#) study the pricing and inventory decisions of retail brands in the luxury segment as they consider moving from a physical retail-only to an omnichannel strategy. They analyze these operational decisions analytically, using the concept of rational expectation equilibrium, and assuming a fixed underlying store and logistics network. Other authors investigate the benefits of buy-online, ship-from-store policies in which manufacturers leverage the warehouses inventories of offline retailers to fulfill online orders (see, e.g., [He et al. 2020](#)), and examine optimal inventory policies for these emerging dual-channel warehouses (see, e.g., [Alawneh and Zhang 2018](#)). In contrast, [He et al. \(2021\)](#) consider the case in which an online-first e-retailer opens an offline channel with physical stores that can serve walk-in customers as well as time-sensitive online customers via a ship-from-store approach. The authors focus on determining the value of the ship-from-store option to the retailer and its effects on pricing decisions using a game-theoretic analytical model. They explicitly do not incorporate inventory

or network design considerations. [Govindarajan et al. \(2021\)](#) explicitly consider the effect of virtual transshipments in a problem with online and offline demand, where offline demand is not subject to transshipments. The authors develop a scalable heuristic to determine order-up-to policies for a multi-location network. [Atamtürk et al. \(2012\)](#) present a joint location-inventory problem with uncertain and correlated demands across facilities and propose a conic integer programming solution approach. The authors allow for transshipments between facilities. Similar to [Govindarajan et al. \(2021\)](#) and [Yang and Qin \(2007\)](#), they assume that online orders that can be delivered from any location in the network can be considered as ‘virtual transshipments’ with zero lead time. While [Atamtürk et al. \(2012\)](#) consider various generalized formulations of location-inventory models, the paper doesn’t consider some of the real-world dependencies of the transportation cost of an order, which typically depend on where the order was placed, and the cost of late orders. Using a newsvendor approach, [Gao and Su \(2017\)](#) study how a buy-online, pick-up-from-store option influences consumer choice. They conclude that it is rarely optimal to assign the revenue from such an integrated channel offering only to a single channel, but that revenue should be shared across channels to optimize the commercial success of such an offering. Similarly, [Gallino et al. \(2017\)](#) show that providing a buy-online, pick-up-from-store service increases sales dispersion. Using a newsvendor model, they conclude that this leads to an increase in safety stocks and thus inventory cost. Recently, [Jiu \(2022\)](#) studies the optimal inventory replenishment and allocation across a network of DCs and stores of an omnichannel retailer. In a multi-period stochastic optimization model, the author assumes that offline demand occurs throughout each time period and can result in lost demand due to store-level stock-outs, while online demand is realized at the end of each period and can be fulfilled from any store or DC in the network, resulting in lost demand only in case of a network-level stock-out. While this study does consider the optimal allocation of inventory across a given network, it does not consider strategic decisions on the design of the network itself. Further, it does not capture the possibility of online orders with tight delivery deadlines that would have to be met within the course of a given time period. While many of these examples from the literature show that it is not straightforward to determine whether integrating online and B&M distribution networks and their respective inventories has a positive effect on overall economic performance, [Millstein and Campbell \(2018\)](#) conclude that the ability to leverage in-store inventory for e-commerce sales will be increasingly economically beneficial with the growing trend towards tight delivery deadlines. [Difrancesco et al. \(2021\)](#) use a combination of simulation and exploratory modeling to prescribe optimal in-store fulfillment policies for online orders of an omnichannel retailer. [Bayram and Cesaret \(2021\)](#) propose a heuristic approach for omnichannel retailers to make near-optimal dynamic fulfillment decisions (i.e., whether to ship an incoming order from a nearby store or from an online fulfillment center) under demand and shipment cost uncertainty. [Arslan et al. \(2021\)](#) investigate two alternative omnichannel deployment strategies, considering order allocation, inventory positioning, delivery, and inbound flow related decisions under uncertainty, based on a two-stage stochastic program with mixed-integer recourse.

Apart from the work of [Atamtürk et al. \(2012\)](#), the works reviewed above, regardless of whether they are studying independent single-channel or integrated multi-channel / omnichannel inventory and fulfillment decisions, do so assuming a fixed and given set of distribution network facilities. In the context of highly responsive (i.e., on-demand

/ sub-same-day) delivery services, research on integrated facility and inventory decisions is still lacking. Our work therefore expands the literature by integrating tactical inventory decisions in the strategic design of distribution networks with such tight delivery deadlines. This enables us to explicitly incorporate the effect of local inventory effects and the resulting cost of virtual transshipments into the strategic design process. In addition, it allows us to evaluate the impact of leveraging existing retail stores for the fulfillment of online orders and of integrating inventory positions across both channels.

2.3. Simulation-Based Optimization

The dynamic nature and complex interdependencies in distribution networks with tight delivery deadlines make the problem sufficiently complex to render state-of-the-art stochastic programming and robust optimization methods intractable (Powell 2016). However, state-of-the-art simulation models are able to capture arbitrarily complex agent behavior, interactions with the distribution network, and demand patterns, allowing us to obtain detailed, disaggregate performance indications of the network (Osorio and Bierlaire 2013). Simulations can be incorporated in SO approaches to search for the objective optimizing set of input parameter settings (Amaran et al. 2016). While such approaches are uncommon in supply chain design, a large body of literature from other domains explores SO, (see, e.g., Andradóttir (1998), Fu et al. (2005), and Amaran et al. (2016) for reviews on SO, including methodological advancements and applications).

While the majority of SO research focuses on problems with continuous decision variables, the reviews of Nelson (2014) and Hong et al. (2015) focus particularly on discrete SO algorithms. Examples of discrete algorithms include Convergent Optimization via Most-Promising-Area Stochastic Search (COMPASS) and the Adaptive Hyperbox Algorithm (AHA), both guaranteeing local convergence (Hong and Nelson 2006, Xu et al. 2013). By focusing on finding a local optimum, these algorithms deliver good finite-time performance because they only explore a small fraction of the feasible solution space (Hong and Nelson 2006, Xu et al. 2013). Xu et al. (2010) propose a framework integrating COMPASS into a global search algorithm. The global search phase explores the entire solution space to identify promising areas for intensive local search, which are then explored using COMPASS. Xu et al. (2013) develop a similar algorithm based on AHA.

Inspired by urban transportation problems, Osorio and Bierlaire (2013) introduce a so-called metamodel SO approach. In contrast to the prevalent general-purpose SO algorithms, metamodel SO leverages knowledge of the underlying system by incorporating an analytical approximation of the objective function in the solution approach. Metamodel SO follows two main steps. First, the metamodel is fitted based on a set of simulation observations. Second, the metamodel is optimized to derive a new trial point that is evaluated by the simulator, leading to an updated set of observations. By iterating these two steps, the accuracy of the metamodel improves, leading to gradually improving trial points. Osorio and Bierlaire (2013) combine a physical and a functional metamodel. Physical metamodels are problem-specific functions that attempt to capture the structure of the underlying decision problem. Functional metamodels are general-purpose functions chosen based on their mathematical properties (Søndergaard 2003). The most

common form of metamodels used to perform SO are functional metamodels, since they can be used to approximate any objective function (Osorio and Bierlaire 2013). However, only with a physical metamodel component we can fully exploit our knowledge about the problem structure.

Zhou et al. (2019) and Snoeck and Winkenbach (2022) extend the work on continuous metamodel SO by proposing metamodel SO approaches for discrete problems, such as car sharing and last-mile distribution network design problems. In this paper, we build on the previous work by extending the AHA with a two-stage metamodel approach that explicitly incorporates inventory decisions into the last-mile distribution network design problem by means of a newsvendor approximation. While the newsvendor fractile provides for an exact solution to the single-period problem (e.g., Porteus 1990), it has been established by de Kok (2018) that the newsvendor fractile structure characterizes the optimal solution in almost any multi-period inventory control policy, such as for instance the base stock level in an infinite horizon (s, S) system.

3. Model Formulation

Model overview. Our proposed model minimizes total last-mile distribution network cost, including facility fixed cost, as well as variable transportation cost, inventory cost, lost sales cost, excess inventory disposal cost, and late delivery penalties. These model components are introduced in Section 3.1 below. The model optimizes over a number of strategic, tactical, and operational decision variables, which are introduced in Sections 3.2 and 3.3. On a strategic level, the model decides on the activation of network facilities, as well as the advance contracting of scheduled transportation capacities. Following these strategic decisions, the model determines the tactical order-up-to inventory level at each active facility. On an operational level, each customer order is allocated to an active facility for fulfillment and to a transportation agent for delivery. The agent can either part of the strategically contracted scheduled transportation capacity or procured on-demand. Each order also gets assigned to a specific delivery trip of the transportation agent. For each trip, the model decides on its specific dispatch time. All of these decisions are made subject to a number of constraints, some of which are imposed by other decisions made on the strategic and tactical level. Specifically, we consider capacity-constrained facilities, carrying capacity-constrained transportation agents, inventory availability constraints at the active facilities, scheduled transportation agent availability constraints, and required adherence to promised delivery deadlines.

Modeling approach. The complexity of the strategic design decision problem for last-mile distribution networks with tight delivery is such that approximations need to be made to develop a tractable model formulation. The interrelatedness of the strategic design decisions with the tactical and operational decisions causes the decision problem to suffer from the ‘curses of dimensionality’ associated with the size of the decision space, state space, and action space. In addition, our problem is defined by various complex sources of uncertainty, e.g., the demand and traffic conditions may be defined by potentially non-trivial and non-stationary distributions. Therefore, we formulate a two-stage approximate analytical model to tractably incorporate the key factors and relationships that govern the performance of

the strategic network design in Sections 3.4 and 3.5. To account for inaccuracies in our approximation, we include several correction factors, represented by the vector β , to capture any real-life complexities excluded in the analytical model. While the correction factors can be determined in various ways, e.g., through real-life piloting and experimentation, we introduce a metamodel SO approach to estimate these factors in Section 4. Consequently, we refer to the model introduced in this chapter as the metamodel. A high-level overview of the methodological approach presented in the following sections is presented in Figure 1.

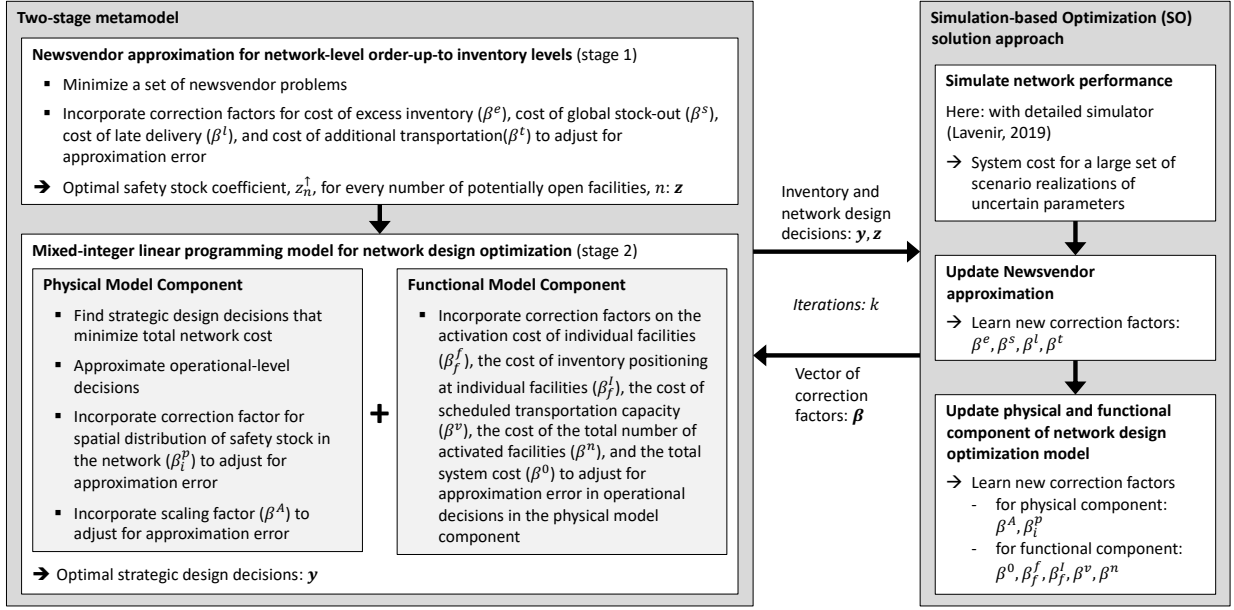


Figure 1: High-level overview of proposed methodological approach

3.1. Distribution Network

We define the last-mile distribution network as a set of capacity-constrained facilities and transportation agents. In our model, we define $\mathcal{F} = \{1, 2, \dots, F\}$ as the set of candidate facilities with associated inventory capacity I_f^{\max} and fixed activation cost, K_f^f . Subset $\mathcal{B} \subseteq \mathcal{F}$ contains the existing B&M facilities. We consider the cost of activating a B&M facility for the distribution of online demand as negligible. We define $\mathcal{V} = \{s, o\}$ as the set of capacity-constrained transportation types, where s refers to scheduled transportation capacity which is paid by the hour, e.g., through scheduled employees and proprietary vehicles, or through contracted external capacity, and o refers to on-demand transportation capacity which is paid per delivery, e.g., through summoning (crowd-sourced) couriers. Transportation agents of type v have a carrying capacity of ξ_v^c customers per trip. In this paper, we consider a distribution network designed for a single stock keeping unit (SKU). Considering multiple SKUs introduces complex combinatorial features to the fulfillment problem as multi-item orders can be fulfilled in different ways (see, e.g., Jasin and Sinha 2015), which we disregard to focus our study on the effect of tactical inventory decisions on the strategic design and the associated performance of last-mile distribution networks. Furthermore, in line with Govindarajan et al. (2021), we

consider a periodic review inventory model, where an order is placed by each activated facility at the start of each review period to satisfy a basestock level. The order is received with zero replenishment lead time. This resembles the often observed operational practice where an order is placed at the end of a business day and delivered overnight.

Online customer requests arrive throughout the service period. To satisfy a customer request c , which is characterized by a location θ_c and time τ_c , the order needs to be delivered within a promised service time l , i.e., the delivery deadline is at time $\tau_c + l$. Delivery after time $\tau_c + l$ is considered late, causing a late delivery penalty cost, c^l , which accounts for a negative impact on the NPS due to a negative customer experience. c^l is independent of the degree of lateness, since as soon as an order is late, the promised premium service is not delivered, impacting the brand perception of the consumer. If a customer request cannot be satisfied due to a global stock-out, the company incurs a lost-sales cost of c^s , similarly accounting for a negative impact on the NPS. At the end of each service period, e.g., at the end of the day, left-over inventory incurs a per-unit cost of c^e to account for opportunity cost and depreciation on the value of inventory. The total online demand at the network level during the service period is normally distributed and defined with mean μ and standard deviation σ , in line with a large body of literature on inventory modeling (see, e.g., [Acimovic and Graves 2017](#)).

In addition, in-store retail demand at each B&M facility is independent of the demand at other B&M facilities or the online demand and is defined with mean μ^b and standard deviation $\chi^b \sigma$, where χ^b scales the standard deviation of the network-level online demand. The inventory order-up-to level at each B&M facility is defined by the B&M safety stock coefficient z^b . While benefits of demand pooling exist between the online inventory and the B&M inventory, there is no transshipment option, and therefore no potential pooling benefit, between the B&M facilities. We refer the reader to Tables [A.2](#) through [A.6](#) in [Appendix A](#) for an overview of notation.

3.2. Strategic and Tactical Network Design Decisions

We consider three types of strategic and tactical decisions that are made while future demand realizations are still unknown. First, the strategic decision to activate facility location $f \in \mathcal{F}$ is represented by the binary variable a_f . Note that $a_f = 1$ for $f \in \mathcal{B}$, i.e., B&M facilities are always activated, since the opening or closing B&M stores is outside the scope of this study, and the fixed cost of activating B&M facilities for online distribution is assumed to be negligible. Second, the tactical order-up-to inventory level at activated facility f is represented by $I_f \in \mathbb{Z}_0$. Third, the strategic contracting of a certain number of scheduled transportation agents is represented by q^s . We define \mathbf{a} and \mathbf{I} to be the vectors containing a particular realization for all facility activation and inventory decisions, and $\mathbf{y} = \mathbf{a} \cap \mathbf{I} \cap \{q^s\}$ to be the combined set of strategic and tactical decisions.

3.3. Operational Fulfillment Decisions

The strategic and tactical decisions \mathbf{y} limit the operational decisions of planners in allocating customer requests to particular facility and transportation agent combinations. Operational planners aim to deliver demand at the lowest cost by making three decisions. First, they allocate each customer request c to a particular combination of facility f

and individual transportation agent of type v , either scheduled or on-demand, within the constraints imposed by the strategic design and tactical planning decisions. When allocating an order to an agent, the planner decides whether to add the order to an existing planned trip, i.e., to consolidate the order, or to plan a new trip. Second, if the order is allocated to a planned trip, the planner can decide the sequence of delivery on that particular trip. Third, the operational planner decides when to dispatch the transportation agent based on the trade-off between the likelihood of delivering late and the potential for future order consolidation, both of which increase as dispatching is delayed further. Order picking at active facilities, and consequently the allocation of inventory to specific orders, happens on a first-come, first-serve basis, independent from the type of order, i.e., online or B&M. These operational decisions are explicitly captured by a detailed simulator as part of our proposed SO solution approach (see Section 4.1, Step 2).

3.4. Stage 1: Determining Network-Level Order-Up-To Inventory Level

We define the first stage of our analytical metamodel by the minimization of a set of newsvendor problems. The model finds the optimal safety stock coefficient z_n^\dagger for every number of potentially activated facilities $n \in \{1, \dots, F\}$, by minimizing the total global and local inventory cost, while accounting for the opportunity to leverage safety stock of B&M facilities. Based on z_n^\dagger , we compute the network-wide optimal basestock level required to serve the online demand as

$$I_n = \mu + \sigma z_n^\dagger, \quad n \in \{1, \dots, F\}. \quad (1)$$

We define $z_n^\dagger = z_n^* \gamma_n$ as the product of the optimal safety stock coefficient without leveraging B&M facilities, z_n^* , and a correction factor accounting for the opportunity of leveraging B&M facilities, γ_n .

Optimal dedicated safety stock coefficient. We determine the vector $\mathbf{z}^* = \{z_1^*, \dots, z_F^*\}$ by minimizing the sum of the cost of global stock-out, $U^s(\cdot)$, excess inventory $U^e(\cdot)$, late delivery, $U^l(\cdot)$, and additional transportation $U^t(\cdot)$, as

$$\min_{\mathbf{z}} C(\mathbf{z}, \boldsymbol{\beta}) = c^e U^e(z_n, \beta^e) + c^s U^s(z_n, \beta^s) + c^l U^l(z_n, \beta_n^l) + c^t U^t(z_n, \beta_n^t), \quad n \in \{1, \dots, F\}. \quad (2)$$

Both the expected global stock-out quantity and the expected global level of excess inventory are independent of the number of activated facilities and are defined in the traditional newsvendor literature. We generalize these definitions by including correction factors β^e and β^s , for the computation of expected excess inventory and stock-out respectively, as

$$U^e(z_n, \beta^e) = \beta^e \int_0^{I_n} (I_n - x) f(x) dx = \beta^e (I_n - (\mu - U^s(z_n))), \quad (3)$$

$$U^s(z_n, \beta^s) = \beta^s \int_{I_n}^{\infty} (x - I_n) f(x) dx = \beta^s \sigma G(z_n), \quad (4)$$

where $f(\cdot)$ is the probability density function (PDF) of the normal distribution, and $G(\cdot)$ the unit normal loss function. Pooling benefits can be overestimated if a normal distribution is used to approximate demand that follows a heavy-tailed distribution (Bimpikis and Markakis 2016). Including the correction factors β^e and β^s accounts for potential arbitrary, non-stationary, geographic and temporal distributions that may describe the arrival of customer requests.

Two factors complicate the determination of the expected number of late deliveries and the expected units that incur additional transportation. First, the demand distribution at a facility is dependent on the inventory levels and demand distributions of the other facilities in the network due to spillover effects that might occur (Acimovic and Graves 2017). Second, the demand distribution depends on the location of the facilities, which is part of the optimization problem addressed in the second stage of the metamodel.

We rely on four approximations to address these factors in providing a tractable formulation of $U^l(\cdot)$ and $U^t(\cdot)$. First, we consider the facility closest to a customer request as the preferred facility to serve the customer. In practice, the preferred facility depends on the state of the system, e.g., location and availability of transportation agents, and delivering from the closest facility could lead to additional transportation costs compared to delivering from the preferred facility. Second, the demand distribution at a facility is independent of its inventory level and the demand distributions and inventory levels at other facilities. Consequently, when determining the probability of a local stock-out, we neglect any spillover effects from other facilities. Since it is common in practice that the demand in arbitrary geographic subareas of the service area, and consequently the inventory level at individual facilities, is correlated, our approximation overestimates the expected units of late delivery and the resulting additional transportation cost. However, it is more important that our approximation captures the nature rather than the level of the relationship between \mathbf{z}^* and $U^l(\cdot)$ and $U^t(\cdot)$, since the correction factors β_n^l and β_n^t can ensure proper scaling. Third, demand is identically normally distributed at every facility with mean $\frac{\mu}{n}$ and variance $\frac{\sigma^2}{n}$ and, similarly, the inventory I_n is uniformly distributed over the facilities. Fourth, we consider the location of each facility to be governed by a uniform distribution throughout the service area.

We define the expected units that incur additional transportation by

$$U^t(z_n, \beta_n^t) = \beta_n^t \left(n \frac{\sigma}{\sqrt{n}} G\left(\frac{z_n}{\sqrt{n}}\right) - U^s(z_n) \right), \quad (5)$$

where $\frac{\sigma}{\sqrt{n}} G\left(\frac{z_n}{\sqrt{n}}\right)$ is the expected local units short at each facility, and β_n^t is the correction factor on the approximation. The demand distributions at facilities are independent and identical, so we find the total expected units short by multiplying the expected units short at one facility by the number of facilities, while controlling for the absence of additional transportation if there is no inventory on the network-level.

To determine the expected units late, we define αA as the subarea of the service area A that lies within range of an arbitrary customer request, with $0 \leq \alpha \leq 1$. An order is delivered late due to a local stock-out if there is no facility within this area with sufficient inventory while inventory is available in the network. Additionally, let i be the number

of facilities within αA . The expected units short within that area is

$$U^l(z_n, n, i) = \begin{cases} \alpha\mu & \text{if } i = 0 \\ \int_{\frac{il_n}{n}}^{\infty} (x - \frac{il_n}{n}) f_i(x) dx & \text{if } i > 0, \end{cases} \quad (6)$$

where $f_i(\cdot)$ is the PDF of a normally distributed random variable defined by mean $\frac{il_n}{n}$ and variance $\frac{i\sigma^2}{n}$. As facilities are uniformly distributed throughout the service area, the probability of finding i facilities within range, given that n facilities are activated follows a binomial distribution with parameters α and n , i.e., $i|n \sim \text{Binom}(n, \alpha)$. To find the expected units late, we combine the probability distribution of $i|n$ and Equation (6) using the law of total expectation as

$$U^l(z_n, \beta'_n) = \beta'_n \left(\frac{1}{\alpha} \sum_{i=0}^F U^l(z_n, n, i) P(i|n) - U^s(z_n) \right), \quad (7)$$

where $\frac{1}{\alpha}$ ensures that we scale our result to capture the entire service area, $-U^s(z_n)$ controls for the absence of late deliveries if there is no inventory on the network-level, and β'_n is the correction factor on the approximation.

Safety stock correction factor accounting for B&M inventory. We determine γ_n , i.e., the safety stock coefficient correction factor, to account for the inventory pooling benefits causing a difference in required safety stock for a dedicated online network, and the required additional safety stock for a combined network supplementing the existing B&M safety stock. Since the absence of transshipment between B&M facilities renders the joint distribution over all stores and the online demand non-trivial, we assume that the joint demand distribution over all stores is normally distributed. The mean of the distribution is equivalent to the sum of the means of all B B&M facilities, $B\mu^b$, and the standard deviation of the distribution would lead to the current network-level safety stock in B&M stores given the safety stock coefficient of B&M stores, z^b , $B\chi^b\sigma$. Consequently, the joint distribution of online and B&M demand is normally distributed with mean $\mu + B\mu_b$ and variance $\sigma^2 + (B\chi^b\sigma)^2$. We find the safety stock coefficient correction factor as

$$\gamma_n = \sqrt{1 + (B\chi^b)^2} - \frac{z_b}{z_n^*} B\chi^b, \quad (8)$$

Consequently, if $\gamma_n < 1$, we get $z_n^\dagger \leq z_n^*$ and there exist inventory pooling benefits to the online network of leveraging B&M facilities when designing a last-mile distribution network for online demand with tight delivery deadlines.

3.5. Stage 2: Designing the Network

We define the second stage of the metamodel as a MILP given by

$$\min_{\mathbf{y}} m(\mathbf{y}; \boldsymbol{\beta}) = \beta_{AG_A}(\mathbf{y}) + \phi(\mathbf{y}; \boldsymbol{\beta}). \quad (9)$$

The model combines a physical component, $g_A(\mathbf{y})$, which attempts to capture the structure of the underlying last-mile urban distribution network, and a functional component, which parametrically corrects the physical component through a scaling term β_A and an additive error term, $\phi(\mathbf{y}; \boldsymbol{\beta})$ (Zhou et al. 2019).

Physical component of the second stage of the metamodel. We define the physical component as a MILP that explicitly models the strategic and tactical decisions, \mathbf{y} , while approximating the operational decisions introduced in Section 3.2. The formulation of the model extends the functional component of the metamodel proposed by Snoeck and Winkenbach (2022) by incorporating three approximations and aggregations to ensure tractability of the model: (1) we develop an expected value based deterministic analytical approximation of the stochastic decision problem; (2) we aggregate demand temporally in time periods, and spatially in discrete geographical pixels; and (3) we aggregate transportation capacity by computing the expected number of transportation agents required per pixel and time period, rather than modeling individual agents making individual trips. Additionally, since safety stock is worthless in a deterministic world, we introduce a mechanism to optimally distribute safety stock based on a pixel-specific correction parameter β_i^p . We formally define the physical component of the metamodel in Appendix C.

Functional component of the second stage of the analytical metamodel. We define the functional component of the analytical model as a linear function of the facility activation and transportation capacity variables. The function is defined by the correction parameter vector $\boldsymbol{\beta}$ as

$$m(\mathbf{y}; \boldsymbol{\beta}) = \beta^A g_A(\mathbf{y}) + \beta^0 + \sum_{f \in \mathcal{F}} \beta_f^f a_f + \sum_{f \in \mathcal{F}} \beta_f^l I_f + \beta^v q^t + \beta^n \sum_{f \in \mathcal{F}} a_f. \quad (10)$$

The correction terms β^0 , β_f^f , β_f^l , β^v , and β^n in the correction factor vector $\boldsymbol{\beta}$ capture relationships that are not explicitly modeled between the total cost of the network and (i) the activation of individual facilities, (ii) the order-up-to level of individual facilities, (iii) the scheduled transportation capacity, and (iv) the total number of facilities.

At this point, all elements of the correction factor vector $\boldsymbol{\beta}$ have been defined, with

$$\boldsymbol{\beta} = (\beta^e, \beta^s, \beta^l, \beta^t, \beta_i^p, \beta^A, \beta^0, \beta_f^f, \beta_f^l, \beta^v, \beta^n)$$

4. Solution Approach

The metamodel introduced in Section 3 allows us to approximate the decision problem that is subject to the well-known ‘curses of dimensionality’ associated with the size of the decision, state, and action spaces (e.g., Powell 2016). The metamodel includes parameter vector $\boldsymbol{\beta}$ to correct for the approximation error. In the following, we estimate $\boldsymbol{\beta}$ using an iterative SO approach (Osorio and Bierlaire 2013, Zhou et al. 2019).

4.1. Solution by Simulation-Based Optimization

Figure 2 summarizes the solution approach. The formal algorithm is presented in detail in [Appendix D](#). We discuss here a conceptual overview of the five steps of the solution approach.

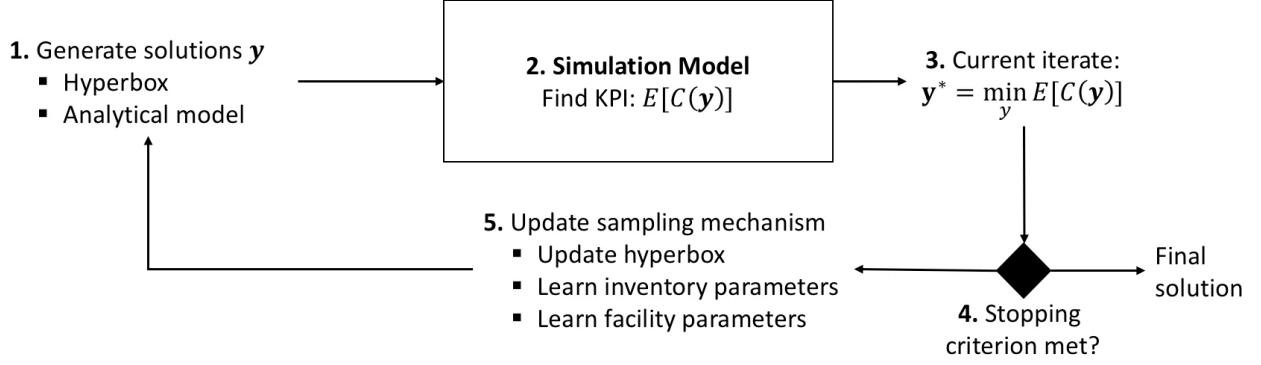


Figure 2: Overview of SO solution approach

Step 1: Solution generation via optimization. We generate feasible solutions based on a sampling mechanism with two components. First, following AHA, we sample random solutions from the so-called hyperbox, i.e., the most promising area in the solution space ([Xu et al. 2013](#)). To formally define the hyperbox, let $\tilde{\mathbf{y}}$ be the abbreviated decision vector, including order-up-to level inventory and transportation capacity decisions, $\tilde{\mathbf{y}} = \mathbf{I} \cap \{q^s\}$. Furthermore, let $\tilde{y}^{(d)}$ be the d^{th} coordinate of decision vector $\tilde{\mathbf{y}}$, which consists of D elements, and $l_k^{(d)}$ and $u_k^{(d)}$ be the lower and upper bound of the hyperbox for coordinate d at iteration k of the algorithm. We can then formally define the hyperbox at iteration k as

$$\mathcal{H}(k) = \{\tilde{\mathbf{y}} : l_k^{(d)} \leq \tilde{y}^{(d)} \leq u_k^{(d)}, 1 \leq d \leq D\}. \quad (11)$$

We translate the abbreviated vector $\tilde{\mathbf{y}}$ to a full decision vector \mathbf{y} per the procedure in [Appendix D.1](#). Second, we solve the two-stage metamodel that approximates the non-linear, probabilistic, and non-convex last-mile distribution network design problem of networks with tight delivery deadlines and locally distributed inventory positions, which we introduced in [Section 3](#).

Step 2: Performance evaluation via simulation. We evaluate the performance of generated solutions leveraging an disaggregate, in-depth simulator. We can use simulation to obtain good approximations for the performance of our network for an individual realization of the uncertain parameters, i.e., for a given scenario ω . By performing a set of r simulation runs for different realizations of the uncertain parameters, i.e., for different scenarios, we approximate the expected performance of network, $\hat{G}(\mathbf{y}) = \mathbb{E}[G(\mathbf{y})]$, by its sample average. The particular simulator used for the purpose of our analyses (documented in [Lavenir 2019](#)) accurately captures the dynamic and stochastic arrival of orders, and the repeated execution of the operational decisions, i.e. the operational allocation, dispatching, and routing decisions.

Step 3: Current iterate determination. We update the current iterate, i.e., the best-found solution hitherto, based on the simulation-based performance evaluation of the solutions generated in both the current and previous iterations of the algorithm.

Step 4: Stopping criterion evaluation. We evaluate the current iterate using the stopping criterion provided by [Xu et al. \(2013\)](#). This criterion is a combination of a statistical guarantee on the local optimality of the current iterate and a computational budget.

Step 5: Sampling mechanism update. We update the parameters governing the sampling mechanism based on previously evaluated solutions. This includes updating the hyperbox (see [Appendix D](#) for details), and the correction parameters, β , of the analytical metamodel. We introduce the updates to both stages of our metamodel in [Section 4.2](#).

4.2. Update of Sampling Mechanism Parameters

To update the parameters governing the sampling mechanism, we need to update the hyperbox, as well as the correction factors β in the metamodel.

The hyperbox is bounded from above (below) in the d^{th} -dimension by the solution with the lowest (highest) value for $y^{(d)}$ that is higher (lower) than the value of that particular element of the current iterate. Throughout the algorithm, the hyperbox changes in size (it typically shrinks) and position in two ways: (1) due to the exploration of new solutions, and (2) by increasing the number of simulations for already explored solutions. We formally describe the updating of the hyperbox in [Appendix D](#).

For the updating of the vector of correction factors β we have developed a two-stage procedure, since we need to update both the factors in the newsvendor model and in the MILP. Note that the values for β are iteration dependent, so throughout the remainder of this paper we include an iteration-dependent index k when referring to (any subcomponent of) β .

Stage 1: Update of Newsvendor Approximation. The iteration-dependent variant of Equation (2) is formulated as

$$\min_{z_k} C(z_k, \beta_k) = c^e U^e(z_{nk}, \beta_k^e) + c^s U^s(z_{nk}, \beta_k^s) + c^l U^l(z_{nk}, \beta_{nk}^l) + c^t U^t(z_{nk}, \beta_{nk}^t), \quad n \in \{1, \dots, F\} \quad (12)$$

Let $\lambda \in \mathcal{L}$ be the set of evaluated solutions, i.e., evaluated strategic design and tactical planning decisions defined by a particular combination of values for \mathbf{y} . To update the correction factors, we rely on the observed cost performance for a particular solution λ of the individual terms of Equation (12) based on the simulator, i.e., \hat{C}_λ^e , \hat{C}_λ^s , \hat{C}_λ^l , and \hat{C}_λ^t . Furthermore, for each solution λ , we can determine the safety stock coefficient z_n^λ based on the number activated B&M and dedicated online facilities and the aggregate order-up-to inventory level in the network.

Consequently, we can learn β_k^e at iteration k by minimizing the weighted least squares problem

$$\min_{\beta_k^e} \sum_{\lambda \in \mathcal{L}} w_k(\lambda) (c^e U^e(\beta_k^e, z_n^\lambda) - \hat{C}_\lambda^e)^2, \quad (13)$$

where the weight function $w_k(\lambda)$, which is defined in [Appendix D.3](#), provides additional weight to solutions closer to the current iterate, which is defined as the best solution found until iteration k . Consequently, the least-squares problem minimizes a weighted distance between the simulated cost \hat{C}_λ^e and the analytical predictions $c^e \mathbb{E}[U^e(\beta_k^e, z^\lambda)]$, where each observation is weighted based on their proximity to the current iterate \mathbf{y}_k^* to improve the local fit of the analytical metamodel around the current iterate.

Similarly, we learn β_k^s at iteration k by minimizing

$$\min_{\beta_k^s} \sum_{\lambda \in \mathcal{L}} w(\lambda) (c^s e U^s(\beta_k^s, z^\lambda) - \hat{C}_\lambda^s)^2. \quad (14)$$

Since β_{nk}^l and β_{nk}^t are dependent on the number of active facilities n , we define \mathcal{L}_n the set of evaluated solutions with n facilities activated. Consequently, we learn β_{nk}^l and β_{nk}^t by finding the values that minimize

$$\min_{\beta_{nk}^l} \sum_{\lambda \in \mathcal{L}_n} w(\lambda) (c^l U^l(\beta_{nk}^l, z^\lambda) - \hat{C}_\lambda^l)^2, \quad \min_{\beta_{nk}^t} \sum_{\lambda \in \mathcal{L}_n} w(\lambda) (c^t U^t(\beta_{nk}^t, z^\lambda) - \hat{C}_\lambda^t)^2, \quad (15)$$

for $n \in \{1, \dots, F\}$.

Stage 2: Update of the MILP. We provide two types of feedback to the second stage of the metamodel. First, we update the physical component of the MILP by updating the pixel specific safety stock correction factor. Second, we update the functional component of the MILP.

The safety stock correction factor β_{ik}^p is defined for every pixel i , based on the fraction of customer requests in pixel i for which inventory is available at the preferred facility, p_{ik} , in the current iterate k . Note that $1 - p_{ik}$ indicates the fraction of customer requests that could not be served from the preferred facility. We determine β_{ik}^p for $i \in \mathcal{I}$ by solving the following set of equations,

$$p_{ik} \beta_{ik}^p = p_{i'k} \beta_{i'k}^p \quad i, i' \in \mathcal{I}, \quad (16)$$

$$\sum_{i \in \mathcal{I}} \beta_{ik}^p \mu_i = \mu. \quad (17)$$

We update the functional component of the second stage of the analytical metamodel defined by Equation (9) at every iteration by solving a weighted least squares problem to find the parameters β that minimize the weighted distance function between the metamodel, $m(\lambda)$, and the observed performance, $\hat{G}(\lambda)$, for every solution λ . Formally,

$$\begin{aligned} \min_{\beta_k^A, \beta_k^0, \beta_k^f, \beta_k^l, \beta_k^v, \beta_k^n} \sum_{\lambda \in \mathcal{L}} [w_k(\lambda) (\hat{G}(\lambda) - m_k(\lambda; \beta_k^A, \beta_k^0, \beta_k^f, \beta_k^l, \beta_k^v, \beta_k^n))]^2 \\ + w_0 (\beta_k^A - 1)^2 + w_0 \beta_k^0 + \sum_{j=0}^F (w_0 \beta_{jk}^f) + \sum_{j=0}^F (w_0 \beta_{jk}^l) + w_0 \beta_k^v + w_0 \beta_k^n. \end{aligned} \quad (18)$$

Consequently, the least squares problem minimizes a weighted distance between the simulated cost estimates \hat{G} and the metamodel predictions m_k , where each point is weighted based on their proximity to the current optimal solution \mathbf{y}_k^* . The additional terms and associated weights w_0 ensure a full rank least squares matrix when the number of observations is smaller than the number of parameters to be fitted. This estimation approach is formulated and discussed in greater detail in [Osorio and Bierlaire \(2013\)](#).

5. Formal Analysis of Benefits of Integrating Online and Brick-and-Mortar Inventory

Based on the formulation of γ_n in Equation (8), we can derive several properties about the interaction between B&M inventory and the newly added online inventory based on their respective safety stock coefficients. Note that $B\chi^b \geq 0$ by definition, since both the number of B&M facilities as well as the standard deviation of B&M demand have to be positive. Furthermore, it is easily verifiable that γ_n is convex in $B\chi^b$. Based on these observations, we can derive five propositions. The propositions show when it is beneficial to integrate online and B&M inventory, i.e., when inventory for online and B&M should be stored together and customer requests can be fulfilled first-come, first-served, without consideration of the origin of the request. We refer to [Appendix B](#) for the proofs.

Proposition 1. *If $B\chi^b = 0$, $z_n^\dagger = z_n^*$.*

Proposition 1 implies that no inventory pooling effects exist when B&M demand is deterministic since there is no safety stock in the B&M network.

Proposition 2. *If $z_n^* = z^b$ and $z_n^* > 0$, $0 \leq z_n^\dagger \leq z_n^*$ and z_n^\dagger is decreasing in $B\chi^b$.*

Proposition 2 implies that pooling benefits always exist when the safety stock coefficients are equal for online and B&M demand. Furthermore, the pooling benefits to the online network increase if the standard deviation of B&M demand relative to the standard deviation of online demand increases. Consequently, it is desirable to integrate online and B&M inventory when $z_n^* = z^b$.

Proposition 3. *If $z_n^* > z^b > 0$, $z_n^* \sqrt{1 - \frac{z^b}{z_n^*}^2} \leq z_n^\dagger \leq z_n^*$ when $B\chi^b \leq \frac{-2\frac{z^b}{z_n^*}}{\frac{z^b}{z_n^*} - 1}$, and $z_n^\dagger > z_n^*$ when $B\chi^b > \frac{-2\frac{z^b}{z_n^*}}{\frac{z^b}{z_n^*} - 1}$.*

If $z_n^* > z^b$, the target safety stock factor for the online demand is larger than the target safety stock factor for B&M demand. This implies that, in the joint network, additional inventory is required to increase the B&M service level to satisfy the online service level targets, i.e., the additional inventory subsidizes the B&M service level. In this case, it may not be desirable to integrate online and B&M inventory, and separate inventory positions should be kept for both networks. However, Proposition 3 implies that if the standard deviation of the B&M relative to the online demand is sufficiently small, i.e., $B\chi^b < \frac{-2\frac{z^b}{z_n^*}}{\frac{z^b}{z_n^*} - 1}$, the pooling benefits outweigh the subsidization of the B&M service level, i.e., $z_n^\dagger \leq z_n^*$, and it is desirable to integrate online and B&M inventory.

Proposition 4. *If $z_n^* > z^b$ while $z_n^* > 0$ and $z_b < 0$, $z_n^\dagger \geq z_n^*$.*

Proposition 4 implies that it is not desirable to integrate online and B&M inventory when the target safety stock factor for the B&M is negative. In this case, there are no pooling benefits while additional online inventory is required to subsidize the B&M service level.

Proposition 5. *If $z_n^* < z^b$ and $z_n^* > 0$, $z_n^\dagger \leq z_n^*$ and $\lim_{B\chi^b \rightarrow \infty} z_n^\dagger = -\infty$.*

Proposition 5 implies that, if $z_n^* < z^b$, i.e., the target safety stock level of online demand is smaller than for B&M demand, the online network can cannibalize on the B&M safety stock, thus reducing the need for additional inventory. Consequently, integrating online and B&M inventory by definition leads to a reduction in B&M service level. While for lower levels of relative standard deviation of B&M demand, the reduction in B&M service level might be outweighed by the benefits of requiring less inventory in the online network, this cannot be determined without including the B&M network into the optimization explicitly, which is outside the scope of this paper and remains an important topic for future research.

Note that we presented the propositions solely for $z_n^* > 0$. It is trivial to derive the propositions when $z_n^* < 0$, and by definition there is no safety stock when $z_n^* = 0$.

6. Numerical Analyses

In the following, we summarize how the results of our numerical study demonstrate (i) the value of our proposed modeling framework, (ii) the benefit of the enhanced structure of the proposed approximate newsvendor model, and (iii) the potential role of B&M stores in online fulfillment.

6.1. Experimental Design

To validate our methodology for various demand scenarios, we develop six stylized problem instances based on a real-world case study in corporation with a global fashion retailer (GFR) in Manhattan, NY. Inspired by the geography of Manhattan, we define the stylized service area as a rectangular area of 100km² (5km by 20km). The parameter values characterizing the transportation agents and distribution facilities are identical across all instances, and are chosen to reflect real-world conditions at the GFR. In particular, both scheduled and on-demand transportation agents are modeled as bike couriers. Scheduled agents are paid by the hour, while on-demand couriers are paid a fixed price per trip, in addition to a distance-based cost component. Existing B&M stores and dedicated online facilities differ in terms of their fixed activation cost and available B&M inventory, but are otherwise identical. We generate 10 potential locations for dedicated online facilities using Algorithm 3 in Appendix F.2. In addition, we generate 10 locations with existing B&M facilities based on a p-median problem on 100 points uniformly located throughout the service area, of which a subset is activated depending on the problem instance. The problem instances differ in two dimensions. First, we define four types of geographic demand distributions: (i) *uniform* (U), i.e., within the service area, orders occur following a homogeneous two-dimensional point process, and their spatial distribution is static throughout the day; (ii) *concentrated* (C), i.e., the majority of order occurrences is concentrated in one

geographic area, and their spatial distribution is static throughout the day; (iii) *evolving* (E), i.e., order occurrences are also concentrated, however their spatial distribution changes dynamically, i.e., the centroid of the concentration moves throughout the day; and (iv) *independent* (I), i.e., the service area is split into twenty equally-sized sub-areas, and orders in these sub-area occur following independent and identical homogeneous two-dimensional point processes. Second, the number of existing B&M facilities can take values of 0, 5, or 10. Four out of our six stylized problem instances do not contain any existing B&M facilities and only differ in their respective demand distribution, i.e., (U-0), (C-0), (E-0), and (I-0). The remaining two instances assume a uniform demand distribution and differ in the number of existing B&M facilities, i.e., (U-5) and (U-10).

6.2. Value of a Two-Stage Metamodel

We analyze the value of each stage of the two-stage analytical metamodel introduced in Section 3, *MetaAHA+*, by proposing two alternative algorithms based on only one of the stages of the metamodel: First, in *MetaAHA+SS*, the analytical metamodel is defined solely by its second stage, i.e., the MILP. Essentially, we reduce the metamodel to a one-stage metamodel, similar to the approaches of Zhou et al. (2019) and Snoeck and Winkenbach (2022). Formally, we omit Constraints (C.3), set the parameter vectors \mathbf{z}^* and \mathbf{z}^\dagger to 0, and constrain every element of the safety stock allocation decision variable vector \mathbf{u} to be equal to 0. Further, we omit Step 4.2 of Algorithm 1. Second, in *MetaAHA+FS*, the correction factors in the second-stage of the analytical metamodel are not updated and all learning takes place in the first stage of the metamodel. Formally, Steps 4.3 and 4.4 of Algorithm 1 are omitted. Essentially, at every iteration of the algorithm, the physical component of the MILP is solved, with the caveat that the global order-up-to level, I , might differ in every iteration, based on our proposed newsvendor approximation.

Figure 3 provides an overview of the algorithmic performance and solution quality of *MetaAHA+*, *MetaAHA+FS*, and *MetaAHA+SS* for the six problem instances defined in Section 6.1. Note that the solution quality of the first iteration of each approach is the solution to analytical model with default correction factors, i.e., no learning has taken place. This is equivalent to solving a traditional deterministic optimization model, which is still the predominant solution approach in industry. On average across instances, *MetaAHA+* outperforms the deterministic model by 48.5%, confirm the need for a SO-approach.

Figure 3 illustrates that the value of each stage of the metamodel is contingent on the root cause of the complexity of the problem. In the absence of B&M facilities, the approximations and aggregations required to render the underlying location-allocation problem tractable have a limited impact on the validity of the MILP. Therefore, the algorithmic performance, i.e., speed of convergence and inter-restart consistency, as well as the cost performance of the optimal network predominantly depend on the correct determination of the network-wide order-up-to inventory level. Consequently, there is a negligible performance gap between *MetaAHA+* and *MetaAHA+FS*. Since *MetaAHA+SS* does not directly learn parameters associated with the cause of the problem complexity, i.e., the network-wide inventory level, it performs significantly worse along all three performance dimensions considered. However, in the presence of B&M stores, the approximations and aggregations required to make the problem tractable have a more pronounced effect,

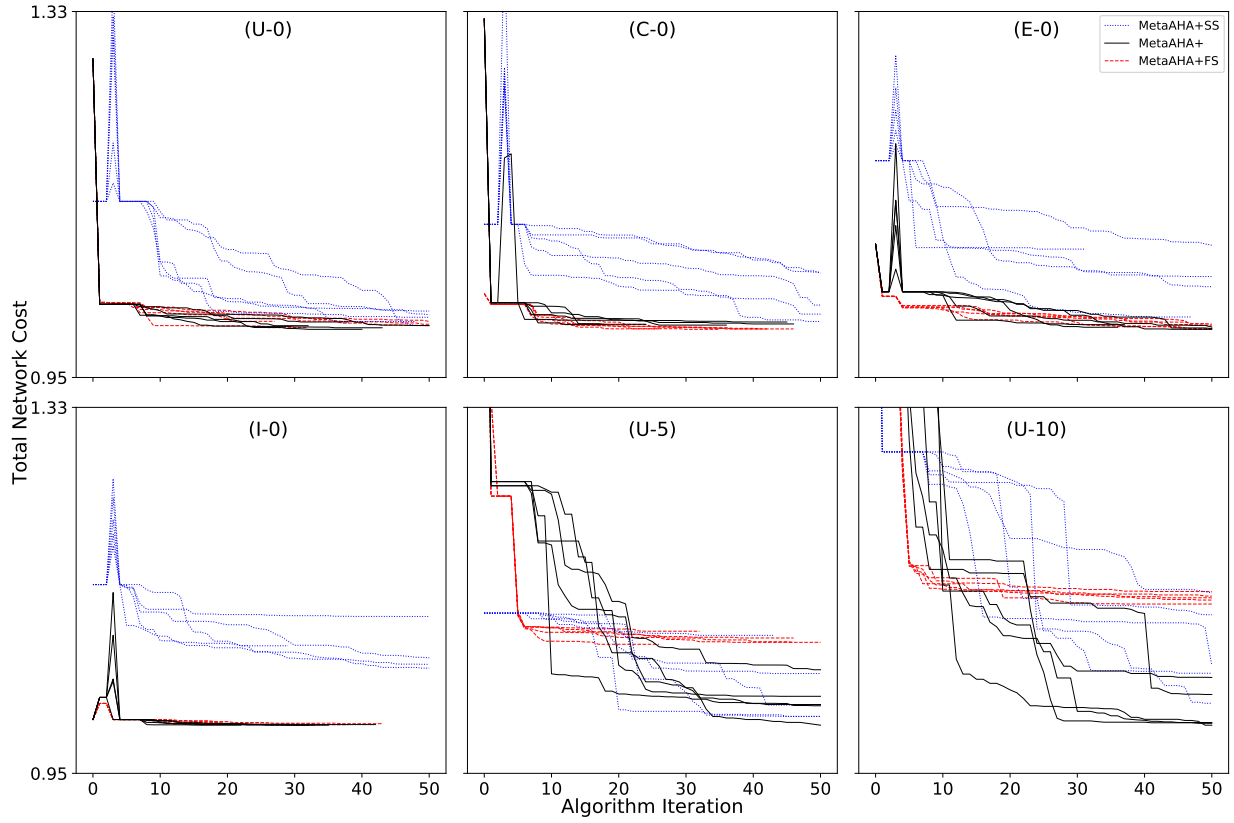


Figure 3: Total network cost evolution (normalized against the best found design) for each solution method and algorithm restart for each of the problem instances.

thus increasing the value of learning the correction factors in the functional component of the MILP. Consequently, if we were to solely rely on our newsvendor approximation, as in *MetaAHA+FS*, we would fail to appropriately capture the location-allocation trade-offs inherent to the design problem, leading to a poor cost performance of the suggested design. Nonetheless, determining the network-wide inventory level remains non-trivial. Since *MetaAHA+SS* includes no explicit consideration of inventory trade-offs, *MetaAHA+* outperforms *MetaAHA+SS* in terms of network cost of the final design, inter-restart consistency, and speed of convergence. Additionally, this becomes more pronounced as the problem becomes more complex with the addition of additional B&M stores. Especially in scenario (U-10), the speed of convergence of *MetaAHA+SS* lags compared to *MetaAHA+*. Speed of convergence is critically important for real-world network design decisions. Decision makers often evaluate several different market scenarios, each of which requires a different network design. Therefore, it is impractical to employ algorithms that yield good final solutions but might run for an arbitrarily long time to obtain these solutions and it is of high practical relevance to assess the performance of our algorithm based on speed of convergence. This indicates the importance of integrating a newsvendor component in the metamodel.

These results indicate that the two-stage metamodel SO approach we propose with *MetaAHA+* provides a robust

approach to solve the strategic design problem for highly responsive distribution networks. Hence, for most problems of realistic size and topology, both solving and iteratively updating the first- and second stage models provides value in terms of improved cost performance of the final design, higher speed of convergence, and greater inter-restart consistency.

6.3. Value of Learning-Based Correction Factors to Enhance the Newsvendor Model

The correction factors in the first stage of the metamodel allow us to enhance the traditional newsvendor model by accounting for three types of real-life complexities in highly responsive distribution networks, namely (i) spatial demand correlation, (ii) the spatial configuration of the underlying facility network, and (iii) distance-dependent additional transportation cost. It should be noted that – unlike the strict newsvendor assumptions – we do have inventory carryover in our system. However, it is well-known from inventory theory (e.g., [Zipkin 2000](#)), that optimal basestock levels have a newsvendor-type formulation. Also note that in our problem setting, replenishments are done overnight with orders placed after store closure and deliveries before store opening, so effectively with zero leadtime. Finally, note that we estimate the newsvendor parameters with a simulation optimization model that incorporates the full real-life complexity.

Spatial demand correlation. The relative impact of local inventory stock-outs on the cost of a highly responsive last-mile distribution system depends on the level of spatial correlation of demand throughout the service area. A local stock-out impacts the network cost if an order arrives for which the inventory at the nearest facility has run out, while inventory is available in another facility. If demand is perfectly positively correlated across facilities, they stock out simultaneously and no virtual transshipments occur. Consequently, there are no late delivery and additional transportation cost due to local inventory effects, and the first stage metamodel can be reduced to a traditional newsvendor model accounting for global stock-out and excess inventory cost. In a similar vein, if demand is perfectly negatively correlated, our model underestimates the impact of local inventory effects. As the spatial correlation of demand decreases, the impact of local stock-outs on cost increases. Since the computation of local inventory effects in our model is based on the assumption of spatially uncorrelated demand, the cost estimation is most accurate for a spatially uncorrelated demand distribution. In real-world operations, however, demand is typically neither perfectly positively or negatively correlated, nor uncorrelated. Rather, it is characterized by some intermediate level of correlation. For instance, the demand correlations observed in our case study with a GFR vary from 0.25 for (I-0) to 0.83 - 0.87 for (U-0), (C-0) and (E-0). Our proposed two-stage metamodel SO approach allows us to iteratively learn a set of correction factors that enhance the newsvendor model by correcting for local stock-out effects and associated transportation costs. It thus provides the versatility to provide solutions for arbitrary demand distributions and arbitrary spatial demand correlations. For example, per [Table 1](#), both (U-0) and (I-0) have four activated facilities in the optimal design. However, the safety stock coefficient is 22% higher for (I-0) compared to (U-0), because the correction factors for late delivery and additional transportation more than tripled due to limited spatial demand correlation.

Problem instance	(U-0)	(I-0)
Activated facilities	4	4
Spatial demand correlation	0.84	0.25
Safety stock coefficient	0.69	0.85
Correction factor late deliveries	0.05	0.17
Correction factor additional transportation	0.14	0.44

Table 1: Impact of spatial demand correlation on correction factors

Facility network. The relevance and impact of the above-mentioned correction factors is amplified by the fact that real-world distribution operations rely on a discrete set of network facilities with fixed and deterministic locations. This implies that the optimized facility network resulting from the first stage of our metamodel violates the implicit assumption of the newsvendor model that facility locations are probabilistic and homogeneously distributed across the service area (cf., Section 3.4). For instance, unlike the newsvendor assumption, the optimal network design seeks to ensure that the entirety of the service area can be reached from at least one facility location within the promised delivery lead time. To illustrate this further, Figure 4 shows the fraction of late delivery cost (captured by the correction factor β_n^l) as a function of the safety stock factors, for different numbers of facilities. The fraction of late delivery costs relative to the safety stock increases when the number of facilities decreases. If there are many stock locations, then the value of network optimization is relatively limited and the locations are more homogeneously distributed. In these cases, the costs do not need to be corrected much as the standard newsvendor equation provides a good estimate of the overall costs as delivery can be done easily. However, if the number of stock locations is small, the newsvendor model needs to be corrected accordingly, due to the higher resulting costs of delivery. It is interesting to note that when the network optimization has more degrees of freedom (with fewer locations opened) the basic newsvendor equation needs to be corrected most. In practice, optimized distribution networks avoid redundant activation of facilities, leading to an increased importance of the correction factors.

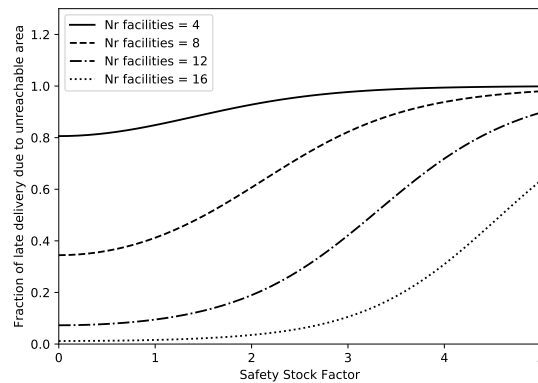


Figure 4: Overview of fraction of late delivery cost in the newsvendor approximation due to areas of the service area being unreachable within the promised delivery lead time.

Distance-dependent additional transportation cost. As outlined above, the correction factors provide our proposed two-stage metamodel SO approach with the versatility to account for non-trivial additional transportation cost effects due to local stock-outs. In our newsvendor model, the additional transportation cost are approximated by a constant parameter, independent of the distance between an order affected by a local stock-out and the alternative network facility serving the order. This provides accurate results if the additional transportation cost are constant and independent of distance, e.g., based on a contracted courier rate per shipment. However, in practice, the additional transportation cost for a particular order often depend on a variety of factors, e.g., on the location of facilities and on the availability and utilization of transportation agents at the time of the customer request. In such cases, the correction factors improve the fit of the newsvendor model to real-life complexity.

6.4. Channel Service Level Effects of Inventory Cannibalization

As very tight delivery deadlines are becoming increasingly common in online fulfillment, deliveries need to rely on inventory being stored more locally in the urban distribution network. This implies that the cost of a network-level global stock-out is higher than in the case of less tight delivery deadlines, since there is less opportunity to deliver from a facility outside of the last-mile network, such as a remote distribution center or a suburban B&M store. Furthermore, tight delivery deadlines go hand-in-hand with increased service expectations of the customer, which suggests a higher cost of lost sales for such orders. Missing a tight delivery deadline for a high margin product like a rare pair of sneakers could imply losing the customer for good, thus valuing a lost sale up to the customer lifetime value (Gupta et al. 2006).

Many retailers have responded to this by making B&M store inventories available for online order fulfillment. Benefits from such inventory pooling have been suggested by Millstein and Campbell (2018). In Propositions 2 and 3 (see Section 5), we show that such benefits to integrating inventory across channels do exist without harming the service level of individual channels if safety stock factors are aligned or if the standard deviation of the demand at the channel with the lower safety stock coefficient is sufficiently small.

However, once the online distribution network is optimized taking into account the tight time constraints and high delivery cost associated to online orders, this pooling benefit may get skewed, such that the online service benefits disproportionately at the expense of a deterioration of the in-store service level. This entails that the additional inventory required to support the online network is positioned in such a way that it allows the online channel to fully leverage the existing B&M inventory, while it limits the accessibility of the additional online inventory to be used by the B&M channel. The key enabler for this phenomenon is the ability to ship online orders from a subset of alternative facilities in case of a local stock-out at the preferred facility (i.e., a virtual transshipment occurs), while B&M demand can only be served from the facility where it is observed, in which case a local stock-out results in a lost sale. The magnitude of this phenomenon is exacerbated by the higher cost of lost sales in the online distribution network, which increasingly outweigh the cost of additional transportation that result from inventory not being positioned in the facility closest to an online order.

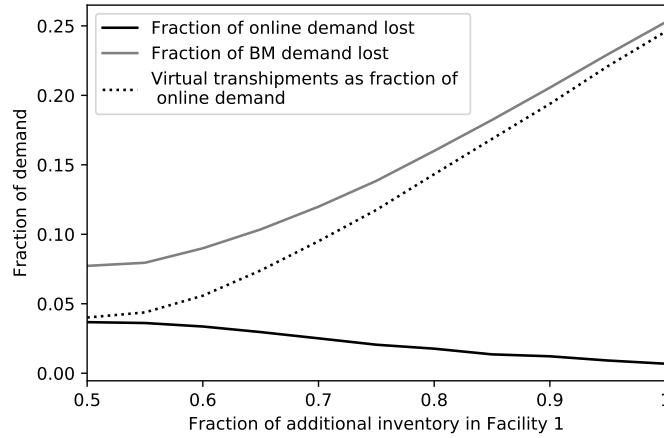


Figure 5: Relationship between service level and inventory distribution for a simulation with two demand regions, two facilities, online and B&M demand, and the option to virtually transship online demand.

We illustrate this phenomenon in Figure 5, which displays the result of a stylized simulation with two demand regions and two facilities, that observe independently and identically distributed online and B&M demand in both regions. If inventory is distributed equally across both facilities, the item fill-rate of the B&M channel is lower than that of the online channel, since online demand can be fulfilled by transshipping from the other facility. As the allocation of additional online inventory gets increasingly concentrated in one facility (here, Facility 1), the item fill-rate of the B&M channel decreases, since the available inventory in Facility 2 reduces. However, the potential for virtual transshipments from Facility 1 to Facility 2 increases as well, thus the item-fill rate of the online channel increases at the expense of the B&M fill-rate. Figure 6 provides numerical evidence for this phenomenon on the network level. The figure reveals a mismatch between the additional inventory to support the online network allocated to each B&M facility and the fraction of total demand that is present in the service sub-area associated to each facility in the optimal network designs obtained for problem instances (U-5) and (U-10), i.e., for 5 and 10 B&M stores, respectively. In particular, additional inventory is pooled in a subset of B&M facilities.

7. Conclusion

In response to the observed market trend towards increasingly tight delivery deadlines in e-commerce, we propose a metamodel-based SO approach to support the strategic design of last-mile distribution networks with tight delivery deadlines. We develop a two-stage analytical metamodel to inform facility location, inventory order-up-to level, and fleet composition decisions. The first stage of the model determines the optimal system-wide order-up-to inventory level based on an approximated newsvendor model that incorporates the cost of late delivery and additional transportation due to local stock-outs. The second stage of the model determines the optimal network configuration and inventory allocation based on a deterministic MILP. We integrate the analytical metamodel in the SO approach

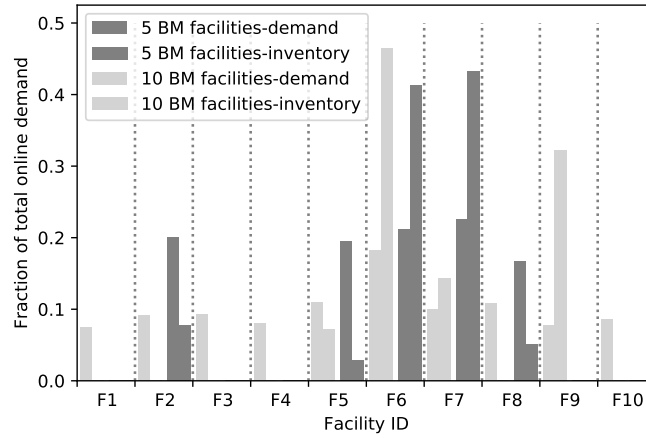


Figure 6: Percentage of total inventory allocated to a facility and the percentage of total demand in the influence area of a facility for 5 and 10 active B&M facilities.

by solving it to generate solutions in the sampling stage of the SO algorithm. We update the metamodel at every iteration of the algorithm by learning a set of correction factors based on extensive simulations of the performance of the last-mile distribution network designs proposed by our metamodel. The evolution of these correction factors throughout the iterations of the proposed algorithm reflect our gradually improving understanding of the performance of the system. In particular, the complex operational interplay between local stockouts (one store being out of stock while another store in the network still has available inventory) and last-mile delivery requires adjustments in the inventory deployment and in the operational cost estimates that inform our network design decisions. These adjustments are operationalized by adding and estimating correction factors in the newsvendor model. This maintains the structure of the newsvendor model, but enhances it with the above-mentioned operational complexities. We believe that this approach is not only computationally effective to resolve such complex problems at scale, but it also allows us to keep a formal structure to the problem that helps managers understand the underlying trade-offs. In particular for omnichannel retailers in segments where consumer baskets consist of expensive single SKUs, such as fashion sneakers or watches shipped from flagship stores, our model and insights are valuable. For retailers where the basket size is larger, the assortment decision would need to be linked to the strategic inventory decision in our problem. This is beyond the scope of our study, but is an interesting avenue for future work.

We analytically define conditions for when integrating the online network with an existing B&M network and its inventory positions leads to a cost reduction. Our numerical study inspired by a highly responsive omnichannel fashion retailer in Manhattan, NY, shows that integrating our model in the SO approach improves the cost performance of the suggested network design. This cost improvement is primarily enabled by enhancing the newsvendor model with correction factors for the operational complexities of last-mile delivery and local stockouts.

Furthermore, our results also indicate that the online network does not distribute additional inventory proportion-

ally to demand over the active B&M facilities, but consolidates inventory in a subset of B&M facilities as an implicit mechanism to avoid cannibalization of online channel safety stock by the B&M channel. However, this negatively impacts the performance of the B&M network, as it actually causes cannibalization of B&M inventory by the online channel. These results demonstrate that while more abstract models may suggest that there are significant inventory pooling benefits from integrating online fulfillment and B&M store networks, these effects might quickly diminish in actual networks. Real-life networks are often not evenly spread over the service area; the resulting operational costs to mitigate these imperfections of the network structure are then such that it may be better to keep channel-specific inventories separate.

There are several fruitful avenues for future research. For instance, integrating pricing decisions on the product and the delivery service-level into the strategic network design and tactical inventory allocation would pose additional methodological challenges and yield relevant insights for omnichannel retailers rolling out buy-online, ship-from-store services with potentially multiple concurrent delivery speeds. Moreover, while we assume overnight replenishment of store inventories, extending our model to incorporate parts of the upstream supply chain could yield interesting insights on optimal replenishment policies and their repercussions on inventory levels and network design decisions.

References

- Acimovic, J., Graves, S.C., 2015. Making better fulfillment decisions on the fly in an online retail environment. *Manufacturing & Service Operations Management* 17, 34–51.
- Acimovic, J., Graves, S.C., 2017. Mitigating spillover in online retailing via replenishment. *Manufacturing & Service Operations Management* 19, 419–436.
- Alawneh, F., Zhang, G., 2018. Dual-channel warehouse and inventory management with stochastic demand. *Transportation Research Part E: Logistics and Transportation Review* 112, 84–106.
- Amaran, S., Sahinidis, N.V., Sharda, B., Bury, S.J., 2016. Simulation optimization: a review of algorithms and applications. *Annals of Operations Research* 240, 351–380.
- Amiri-Aref, M., Klibi, W., Babai, M.Z., 2018. The multi-sourcing location inventory problem with stochastic demand. *European Journal of Operational Research* 266, 72–87.
- Andradóttir, S., 1998. Simulation optimization, in: Banks, J. (Ed.), *Handbook of simulation: Principles, methodology, advances, applications, and practice*. John Wiley & Sons, New York. chapter 9, pp. 307–333.
- Arslan, A.N., Klibi, W., Montreuil, B., 2021. Distribution network deployment for omnichannel retailing. *European Journal of Operational Research* 294, 1042–1058. doi:<https://doi.org/10.1016/j.ejor.2020.04.016>.
- Atamtürk, A., Berenguer, G., Shen, Z.J., 2012. A conic integer programming approach to stochastic joint location-inventory problems. *Operations Research* 60, 366–381.
- Bayram, A., Cesaret, B., 2021. Order fulfillment policies for ship-from-store implementation in omni-channel retailing. *European Journal of Operational Research* 294, 987–1002. doi:<https://doi.org/10.1016/j.ejor.2020.01.011>.
- Bell, D.R., Gallino, S., Moreno, A., 2014. How to win in an omnichannel world. *MIT Sloan Management Review* 56, 45–53.

- Bendoly, E., Blocher, D., Bretthauer, K.M., Venkataramanan, M., 2007. Service and cost benefits through clicks-and-mortar integration: Implications for the centralization/decentralization debate. *European Journal of Operational Research* 180, 426–442.
- Bimpikis, K., Markakis, M.G., 2016. Inventory pooling under heavy-tailed demand. *Management Science* 62, 1800–1813.
- Boccia, M., Crainic, T.G., Sforza, A., Sterle, C., 2011. Location-routing models for designing a two-echelon freight distribution system. Technical Report, CIRRELT-2011-06, Université de Montréal .
- Colby, C., Bell, K., 2016. The on-demand economy is growing, and not just for the young and wealthy. *Harvard Business Review* URL: <https://hbr.org/2016/04/the-on-demand-economy-is-growing-and-not-just-for-the-young-and-wealthy>.
- Crainic, T.G., Ricciardi, N., Storchi, G., 2004. Advanced freight transportation systems for congested urban areas. *Transportation Research Part C: Emerging Technologies* 12, 119–137.
- Crainic, T.G., Ricciardi, N., Storchi, G., 2009. Models for evaluating and planning city logistics systems. *Transportation Science* 43, 432–454.
- Daskin, M.S., Coullard, C.R., Shen, Z.M., 2002. An inventory-location model: Formulation, solution algorithm and computational results. *Annals of operations research* 110, 83–106.
- Difrancesco, R.M., van Schilt, I.M., Winkenbach, M., 2021. Optimal in-store fulfillment policies for online orders in an omnichannel retail environment. *European Journal of Operational Research* 293, 1058–1076.
- eMarketer, 2019. Global ecommerce 2019. URL: <https://www.emarketer.com/content/global-ecommerce-2019>. accessed April 21, 2020.
- Farfetch, 2017. Gucci in 90 minutes. URL: <https://www.farfetech.com/editorial/gucci-in-90-minutes.aspx>. accessed September 30, 2019.
- Fernie, J., Grant, D.B., 2008. On-shelf availability: the case of a uk grocery retailer. *The International Journal of Logistics Management* 19, 293–308.
- Fu, M.C., Glover, F.W., April, J., 2005. Simulation optimization: a review, new developments, and applications, in: Kuhl, M.E., Steiger, N.M., Armstrong, F.B., Joines, J.A. (Eds.), *Proceedings of the 2005 Winter Simulation Conference*, IEEE. pp. 83–95.
- Gallino, S., Moreno, A., Stamatopoulos, I., 2017. Channel integration, sales dispersion, and inventory management. *Management Science* 63, 2813–2831.
- Gao, F., Su, X., 2017. Omnichannel retail operations with buy-online-and-pick-up-in-store. *Management Science* 63, 2478–2492.
- Govindarajan, A., Sinha, A., Uichanco, J., 2021. Joint inventory and fulfillment decisions for omnichannel retail networks. *Naval Research Logistics (NRL)* 68, 779–794.
- Guo, J., Keskin, B.B., 2018. Designing a centralized distribution system for omni-channel retailing. Available at SSRN: <https://ssrn.com/abstract=3317160> .
- Gupta, S., Hanssens, D., Hardie, B., Kahn, W., Kumar, V., Lin, N., Ravishanker, N., Sriram, S., 2006. Modeling customer lifetime value. *Journal of service research* 9, 139–155.

- He, P., He, Y., Xu, H., 2020. Buy-online-and-deliver-from-store strategy for a dual-channel supply chain considering retailer's location advantage. *Transportation Research Part E: Logistics and Transportation Review* 144, 102127.
- He, Y., Xu, Q., Shao, Z., 2021. "Ship-from-store" strategy in platform retailing. *Transportation Research Part E: Logistics and Transportation Review* 145, 102153.
- Hong, L.J., Nelson, B.L., 2006. Discrete optimization via simulation using compass. *Operations Research* 54, 115–129.
- Hong, L.J., Nelson, B.L., Xu, J., 2015. Discrete optimization via simulation, in: Fu, M. (Ed.), *Handbook of simulation optimization*. International series in Operations Research & Management Science. Springer, New York, NY. volume 216, pp. 9–44.
- Hübner, A., Holzapfel, A., Kuhn, H., 2015. Operations management in multi-channel retailing: an exploratory study. *Operations Management Research* 8, 84–100.
- Janjevic, M., Merchán, D., Winkenbach, M., 2021. Designing multi-tier, multi-service-level, and multi-modal last-mile distribution networks for omni-channel operations. *European Journal of Operational Research* 294, 1059–1077.
- Janjevic, M., Winkenbach, M., Merchán, D., 2019. Integrating collection-and-delivery points in the strategic design of urban last-mile e-commerce distribution networks. *Transportation Research Part E: Logistics and Transportation Review* 131, 37–67.
- Jasin, S., Sinha, A., 2015. An lp-based correlated rounding scheme for multi-item ecommerce order fulfillment. *Operations Research* 63, 1336–1351.
- Jiu, S., 2022. Robust omnichannel retail operations with the implementation of ship-from-store. *Transportation Research Part E: Logistics and Transportation Review* 157, 102550.
- Klapp, M.A., Erera, A.L., Toriello, A., 2018. The dynamic dispatch waves problem for same-day delivery. *European Journal of Operational Research* 271, 519–534.
- de Kok, T., 2018. Inventory management: Modeling real-life supply chains and empirical validity. *Foundations and Trends in Technology, Information and Operations Management* 11, 343–437.
- Lavenir, X., 2019. The Strategic Design and Environmental Footprint of Highly Responsive Urban Distribution Networks. Master's thesis. Massachusetts Institute of Technology, Cambridge, MA.
- Lim, S.F.W.T., Winkenbach, M., 2019. Configuring the last-mile in business-to-consumer e-retailing. *California Management Review* 61, 132–154.
- MediaMarkt, 2020. Entrega inmediata en 2 horas. URL: <https://specials.mediamarkt.es/entrega-inmediata-2-horas>. accessed January 15, 2020.
- Melacini, M., Perotti, S., Rasini, M., Tappia, E., 2018. E-fulfilment and distribution in omni-channel retailing: a systematic literature review. *International Journal of Physical Distribution & Logistics Management* 48, 391–414.
- Millstein, M.A., Campbell, J.F., 2018. Total hockey optimizes omnichannel facility locations. *INFORMS Journal on Applied Analytics* 48, 340–356.
- Miranda, P.A., Garrido, R.A., 2004. Incorporating inventory control decisions into a strategic distribution network design model with stochastic demand. *Transportation Research Part E: Logistics and Transportation Review* 40, 183–207.
- MVPL, 2020. Amazon global fulfillment center network. URL: <http://www.mwpvl.com/html/amazon-com.html>.

- Nelson, B.L., 2014. Optimization via simulation over discrete decision variables, in: *INFORMS TutORials in Operations Research*, pp. 193–207.
- Novy-Williams, E., Soper, S., 2019. Nike pulling its products from amazon in e-commerce pivot. Bloomberg URL: <https://www.bloomberg.com/news/articles/2019-11-13/nike-will-end-its-pilot-project-selling-products-on-amazon-site>.
- Osorio, C., Bierlaire, M., 2013. A simulation-based optimization framework for urban transportation problems. *Operations Research* 61, 1333–1345.
- Ozsen, L., Daskin, M.S., Coullard, C.R., 2009. Facility location modeling and inventory management with multisourcing. *Transportation Science* 43, 455–472.
- Paterson, C., Kiesmüller, G., Teunter, R., Glazebrook, K., 2011. Inventory models with lateral transshipments: A review. *European Journal of Operational Research* 210, 125–136.
- Porteus, E.L., 1990. Stochastic inventory theory. *Handbooks in Operations Research and Management Science* 2, 605–652.
- Powell, W.B., 2016. A unified framework for optimization under uncertainty, in: *INFORMS TutORials in Operations Research*. INFORMS, pp. 45–83.
- Redman, R., 2019. Target energizes store base with small formats, remodels. *Supermarket News* URL: <https://www.supermarketnews.com/store-design-construction/target-energizes-store-base-small-formats-remodels>.
- Savelsbergh, M., Van Woensel, T., 2016. 50th anniversary invited article—city logistics: Challenges and opportunities. *Transportation Science* 50, 579–590.
- Shen, Z.M., Coullard, C., Daskin, M.S., 2003. A joint location-inventory model. *Transportation science* 37, 40–55.
- Shen, Z.M., Qi, L., 2007. Incorporating inventory and routing costs in strategic location models. *European journal of operational research* 179, 372–389.
- Snoeck, A., Winkenbach, M., 2020. The value of physical distribution flexibility in serving dense and uncertain urban markets. *Transportation Research Part A: Policy and Practice* 136, 151–177.
- Snoeck, A., Winkenbach, M., 2022. A discrete simulation-based optimization algorithm for the design of highly responsive last-mile distribution networks. *Transportation Science* 56, 201–222.
- Søndergaard, J., 2003. Optimization using surrogate models - by the space mapping technique. Ph.D. thesis. Technical University of Denmark, Kgs. Lyngby, Denmark.
- Ulmer, M., 2017. Delivery deadlines in same-day delivery. *Logistics Research* 10, 1–15.
- Voccia, S.A., Campbell, A.M., Thomas, B.W., 2019. The same-day delivery problem for online purchases. *Transportation Science* 53, 167–184.
- Wei, Y., Li, F., 2020. Omnichannel supply chain operations for luxury products with conspicuous consumers. *Transportation Research Part E: Logistics and Transportation Review* 137, 101918.
- Winkenbach, M., Kleindorfer, P.R., Spinler, S., 2016a. Enabling urban logistics services at La Poste through multi-echelon location-routing. *Transportation Science* 50, 520–540.

- Winkenbach, M., Roset, A., Spinler, S., 2016b. Strategic redesign of urban mail and parcel networks at La Poste. *Interfaces* 46, 445–458.
- Xu, J., Nelson, B.L., Hong, L.J., 2010. Industrial strength compass: A comprehensive algorithm and software for optimization via simulation. *ACM Transactions on Modeling and Computer Simulation (TOMACS)* 20, 1–29.
- Xu, J., Nelson, B.L., Hong, L.J., 2013. An adaptive hyperbox algorithm for high-dimensional discrete optimization via simulation problems. *INFORMS Journal on Computing* 25, 133–146.
- Yang, J., Qin, Z., 2007. Capacitated production control with virtual lateral transshipments. *Operations Research* 55, 1104–1119.
- Zebra Technologies, 2018. Reinventing the supply chain: the future of fulfillment vision study. Technical Report.
- Zhou, T., Osorio, C., Fields, E., 2019. Large-scale data-driven simulation-based car-sharing network design. MIT working paper.
- Zipkin, P., 2000. *Foundations of Inventory Management*. McGraw-Hill.

Appendix A. Notation

\mathcal{B}	set of B&M facilities, $\mathcal{B} \subseteq \mathcal{F}$
\mathcal{I}	set of pixels
\mathcal{I}_f	set of pixels within reach of facility f within l using any vehicle in any time period
\mathcal{I}_{fvt}	set of pixels within reach of facility f within l using vehicle type v in time period t
\mathcal{F}	set of candidate facilities
$\mathcal{N}(i)$	set of pixels that lie within the neighborhood of i , defined for consolidation purposes
\mathcal{V}	set of transportation agent types
\mathcal{T}	set of discrete time periods

Table A.2: Notation: Sets defining the distribution network

a_f	binary variable indicating whether a facility at location f is activated
I_f	order-up-to level at facility f
I_n^*	optimal global order-up-to level if n facilities are activated
q^s	quantity of scheduled transportation agents
q_t^o	quantity of on-demand transportation agents time period t
r_n	binary indicator variable which is 1 if n facilities are activated
u_{ifn}	fraction by which the safety stock of pixel i is allocated to facility f if n facilities are activated
v_n, h_0, h_1	binary indicator variables to support that r_n is 1 if n facilities are activated
\mathbf{y}	vector containing all strategic and tactical decision variables, $\mathbf{y} = \mathbf{a} \cap \mathbf{I} \cap \{q^s\}$
x_{ifvt}	fraction by which pixel i is served from facility f using a transportation agent of type v in time period t
z_n^*	optimal safety stock coefficient when demand is normally distributed with (μ, σ^2) and n facilities are activated, without accounting for the presence of B&M stores
z_n^\dagger	optimal safety stock coefficient when demand is normally distributed with (μ, σ^2) and n facilities are activated, accounting for the presence of B&M stores
\mathbf{z}	vector of safety stock decision variables

Table A.3: Notation: Decision variables

Appendix B. Proofs

To prove propositions 1 through 4, we substitute $B\chi^b$ by x . Consequently, based on Equation (8), we have

$$\gamma_n(x) = \sqrt{1 + x^2} - \frac{z_b}{z_n^*} x, \quad (\text{B.1})$$

with $x \geq 0$ since the standard deviation of B&M demand is positive.

Proposition 1 follows from $\gamma_n(0) = 1$. Furthermore, by taking the second derivative of $\gamma_n(x)$,

$$\frac{d\gamma_n(x)}{dx^2} = \frac{1}{\sqrt{1 + x^2}^{\frac{3}{2}}}. \quad (\text{B.2})$$

we see that $\frac{d\gamma_n(x)}{dx^2} \geq 0$ for $x \geq 0$, proving convexity.

A	area of service area
B	number of B&M facilities
c^a	operational cost of a scheduled transportation agent per time period
c_v^d	distance based cost for a courier of type v
c^e	per order excess inventory cost
c^l	per order late delivery cost
c_v^o	cost of summoning an on-demand transportation agent
c^s	per order global stock-out cost
c^t	per order additional transportation cost due to local stock-out
d_{if}	travel distance between pixel i and facility f
$f_{ifvt}(t')$	time that transportation agents of type v spend in time period t on orders that need to be delivered from facility f to pixel i placed in time period t' . See Equation (E.2) in Appendix Appendix E.2
I_n	network level order-up-to level, $I_n = \mu + z_n\sigma$
I_f^{\max}	maximum inventory capacity at facility f
k_{ifvt}	consolidation factor that approximates the effect of consolidating multiple orders into one trip. See Equation (E.5) in Appendix Appendix E.2
K_f^f	daily facility fixed cost for facility f
l	promised lead-time, i.e., available time to deliver order after customer request
$Q^{o\max}$	Maximum number of on-demand transportation agents that can be summoned per time period
$Q^{s\max}$	Maximum number of scheduled transportation agents that can be hired
t_{if}^d	minimum time it takes to get an order from facility f to a customer in pixel i in any time period using any transportation agent
t_{ifvt}^d	minimum time it takes to get an order from facility f to a customer in pixel i using a transportation agent of type v in time period t
t_{ifvt}^o	time a transportation agent of type v requires to serve pixel i from facility f in time period t , from the start of the trip to the end of the trip
t_{ifvt}^s	slack time when an order of a customer in pixel i is delivered from facility f using transportation agent of type v in time period t , $l - t_{ifvt}^d$
$U^s(\cdot)$	expected orders short due to global stock-out
$U^l(\cdot)$	expected orders delivered late
$U^e(\cdot)$	expected excess inventory
$U^t(\cdot)$	expected orders that incur additional transportation due to local stock-out
z^b	safety stock coefficient at B&M facilities
α	fraction of the service area within range of an arbitrary order
γ_n	safety stock coefficient scaling factor to account for B&M demand
Δ_t	length of period t
θ_c	delivery location associated to customer request c
μ	expected network-level online demand
μ_{it}	expected demand in pixel i at time period t
μ^b	expected demand at B&M facilities
ξ_v^c	carrying capacity of a courier of type v
σ	network-level standard deviation of online demand
τ_c	order arrival time of customer request c
χ^b	factor on standard deviation of the online demand such that the standard deviation of demand at B&M facilities equals $\chi^b\sigma$

Table A.4: Notation: Analytical model parameters

β	vector of correction factors
β^0	additive correction factor
β^A	correction factor on physical component of metamodel
β_f^f	correction factor on the activation of facility f
β_f^I	correction factor on the order-up-to inventory level at facility f
β^v	correction factor on the number of scheduled transportation agents
β^n	correction factor on the total number of facilities f
β_i^p	correction factor on the distribution of safety stock for pixel i
β_n^l	correction factor on the expected orders late when n facilities are activated
β_n^t	correction factor on the expected orders that incur additional transportation when n facilities are activated
β^e	correction factor on the expected excess inventory [-]
β^s	correction factor on the expected orders short [-]

Table A.5: Notation: Correction factors

$\mathcal{A}_k(\mathbf{y})$	number of additional simulations for solution \mathbf{y} in iteration k
$\hat{G}(\mathbf{y})$	the average simulation performance of solution \mathbf{y}
$g_A(\mathbf{y})$	analytical performance of solution \mathbf{y}
$\mathcal{H}(k)$	hyperbox at iteration k
\mathcal{L}	set of evaluated solutions
\mathcal{L}_n	set of evaluated solutions with n activated facilities
$\mathcal{L}(k)$	set of evaluated solutions in iteration k
$l_k^{(d)}$	lower bound of the hyperbox for coordinate d at iteration k of the algorithm
$m(\mathbf{y})$	metamodel performance of solution \mathbf{y}
$p^{\mathcal{H}}$	probability that a facility is activated when $I_f = 1$ when randomly generating solutions from the hyperbox
p_{ik}	percentage by which inventory is available in the preferred facility when serving pixel i in the current iterate of iteration k
$u_k^{(d)}$	lower bound of the hyperbox for coordinate d at iteration k of the algorithm
w_0	base weight to ensure full rank matrix
$w_k(\lambda)$	weight of solution λ in iteration k
$\tilde{\mathbf{y}}$	abbreviated decision vector including \mathbf{I} and q^s
\mathbf{y}_k^*	best solution until iteration k
$\mathbf{y}_k^{\text{meta}}$	solution to the metamodel problem in iteration k
$\mathbf{y}_k^{\text{meta-hyper}}$	solution to the hyperbox constrained metamodel problem in iteration k
Ω	feasible solution space

Table A.6: Notation: Definition of Algorithm 1

For Proposition 2, when $z_b = z_n^*$, $\lim_{x \rightarrow \infty} \gamma_n(x) = 0$. Furthermore, we know $\gamma_n(0) = 1$. Therefore, since $\gamma_n(x)$ is convex, we know that $0 \leq \gamma_n(x) \leq 1$ and therefore $0 \leq z_n^\dagger \leq z_n^*$.

For Proposition 3, when $z_b < z_n^*$, $\lim_{x \rightarrow \infty} \gamma_n(x) = \infty$. Furthermore, we find $\min_{x \geq 0, z_b < z_n^*} \gamma_n(x)$ by solving $\frac{d\gamma_n(x)}{dx} = 0$. Combining the solution to this problem, $\sqrt{1 - \frac{z_b^2}{z_n^{*2}}} > 0$, and the fact that the limit goes to ∞ , we know that there are two solutions to $\gamma_n(x) = 1$. The first solution is known, $x = 0$, and we find the second solution as $x = \frac{-2\frac{z_b}{z_n^*}}{\frac{z_b}{z_n^*}^2 - 1}$. Between those solutions $z_n^\dagger \leq z_n^*$, i.e., $\gamma_n(x) < 1$ on $[0, \frac{-2\frac{z_b}{z_n^*}}{\frac{z_b}{z_n^*}^2 - 1})$. Furthermore, the value of z_n^\dagger is bounded below by the minimum of $\gamma_n(x)$. Outside this interval, $\gamma_n(x) > 1$ and therefore $z_n^\dagger > z_n^*$.

For Proposition 5, we can follow a similar pattern as for the proof of Proposition 3. While $\lim_{x \rightarrow \infty} \gamma_n(x) = \infty$, $\min_{x \geq 0, z_b < z_n^*} \gamma_n(x) = 1$, indicating that $\gamma_n(x) \geq 1$ and $z_n^\dagger \geq z_n^*$.

For Proposition 4, we find $\lim_{x \rightarrow \infty} \gamma_n(x) = -\infty$. Since we know that $\gamma_n(0) = 1$ and $\gamma_n(x)$ is convex, we find $\gamma_n(x) \leq 1$ and $z_n^\dagger \leq z_n^*$.

Appendix C. Physical component of the metamodel

In this Appendix, we define the physical component as a MILP that explicitly models the strategic and tactical decisions, \mathbf{y} , while approximating the operational decisions introduced in Section 3.2. We start by introducing three approximations and aggregations to ensure tractability of the model.

First, we develop an expected value based deterministic analytical approximation of the stochastic decision problem, i.e., there is no uncertainty about the location or timing of demand. Second, we aggregate demand temporally and spatially. We divide the day into a set of discrete periods $t \in \mathcal{T}$ of length Δ_t and the service area in a large set of adjacent rectangular pixels $i \in \mathcal{I}$. The demand level in each time period and pixel is defined by μ_{it} , and demand is uniformly distributed within each pixel in each time period. However, this discretization allows us to capture any arbitrary temporal, across time periods, and spatial, across pixels, demand distributions. Third, we aggregate transportation capacity by computing the expected number of transportation agents required per pixel and time period, rather than modeling individual agents making individual trips. Consequently, we capture operational allocation decisions by x_{ifvt} , which captures the fraction by which pixel i is served from facility f using transportation agent type v in time period t .

Since safety stock is worthless in a deterministic world, we ensure that the model allocates safety stock to desirable facilities by distributing the optimal network-wide safety stock level determined in the first stage of the metamodel through three mechanisms: (i) we allocate safety stock to every pixel proportionally to its share of expected demand, $\mu^{-1} \sum_{t \in \mathcal{T}} \mu_{it}$; (ii) we introduce a pixel specific correction parameter β_i^p , which controls the share of safety stock allocated to a given pixel; (iii) we ensure that safety stock only gets allocated to facilities that can be reached from a particular pixel by introducing the set $\mathcal{I}_f = \{i \in \mathcal{I} | t_{if}^d \leq l\}$, where t_{if}^d is the minimum time required to deliver an order, including order processing and delivery, with any vehicle in any time period.

We refer the reader to Tables A.2 through A.6 in Appendix A for an overview of notation, and proceed by formally introducing the physical component of the metamodel, $g_A(\mathbf{y})$, as

$$g_A(\mathbf{y}) = \min_{\mathbf{a}, \mathbf{l}, \mathbf{q}^s, \mathbf{x}, \mathbf{q}^o} \sum_{f \in \mathcal{F}} K_f^f a_f + \sum_{i \in \mathcal{I}} \Delta_t c^a q^s + \sum_{i \in \mathcal{I}} \sum_{v \in \mathcal{V}} c_v^o q_i^o + \sum_{v \in \mathcal{V}} c_v^d \sum_{i \in \mathcal{I}} \sum_{f \in \mathcal{F}} d_{if} k_{ifvt} x_{ifvt} + \sum_{i \in \mathcal{I}} \sum_{t \in \mathcal{T}} c^s \mu_{it} (1 - \sum_{f \in \mathcal{F}} \sum_{v \in \mathcal{V}} x_{ifvt}) + c^e (I^* - \sum_{i \in \mathcal{I}} \sum_{t \in \mathcal{T}} \mu_{it} \sum_{f \in \mathcal{F}} \sum_{v \in \mathcal{V}} x_{ifvt}) \quad (\text{C.1})$$

$$\text{s.t.} \quad \sum_{f \in \mathcal{F}} \sum_{v \in \mathcal{V}} x_{ifvt} \leq 1, \quad i \in \mathcal{I}, t \in \mathcal{T}, \quad (\text{C.2})$$

$$I^* = \mu + \sigma \sum_{n=1}^F z_{nk}^\dagger r_n, \quad (\text{C.3})$$

$$r_n \leq n^{-1} \sum_{f \in \mathcal{F}} a_f, \quad 1 \leq n \leq F, \quad (\text{C.4})$$

$$r_n \leq v_n \quad 1 \leq n \leq F, \quad (\text{C.5})$$

$$0 \leq v_n \leq F(1 - h_0), \quad 1 \leq n \leq F, \quad (\text{C.6})$$

$$2 - n^{-1} \sum_{f \in \mathcal{F}} a_f \leq v_n, \quad 1 \leq n \leq F, \quad (\text{C.7})$$

$$2 - n^{-1} \sum_{f \in \mathcal{F}} a_f + F(1 - h_1) \geq v_n, \quad 1 \leq n \leq F, \quad (\text{C.8})$$

$$\sum_1^F r_n = 1, \quad (\text{C.9})$$

$$h_0 + h_1 = 1, \quad (\text{C.10})$$

$$I^* = \sum_{f \in \mathcal{F}} I_f, \quad (\text{C.11})$$

$$I_f \leq I_f^{\max} a_f, \quad f \in \mathcal{F}, \quad (\text{C.12})$$

$$\sum_{i \in \mathcal{I}} \sum_{i \in \mathcal{I}} \sum_{v \in \mathcal{V}} \mu_{it} x_{ifvt} + \sum_{i \in \mathcal{I}} \beta_i^p \frac{\mu_i}{\mu} \sum_{n=1}^F z_n^* \sigma u_{ifn} \leq I_f, \quad f \in \mathcal{F} / \mathcal{B}, \quad (\text{C.13})$$

$$\sum_{i \in \mathcal{I}} \sum_{i \in \mathcal{I}} \sum_{v \in \mathcal{V}} \mu_{it} x_{ifvt} + \sum_{i=1 \in \mathcal{I}} \beta_i^p \frac{\mu_i}{\mu} \sum_{n=1}^F z_n^* \sigma u_{ifn} \leq I_f + \sum_{n=1}^F r_n B^{-1} \sigma (z_n^* - z_n^\dagger)^+, \quad f \in \mathcal{B}, \quad (\text{C.14})$$

$$\sum_{f \in \mathcal{F}} u_{ifn} = r_n, \quad i \in \mathcal{I}_{fv}, 1 \leq n \leq F, \quad (\text{C.15})$$

$$\begin{aligned} & \sum_{f \in \mathcal{F}} \sum_{i \in \mathcal{I}} k_{ifvt} x_{ifvt} (l_{ifvt}^o \mu_{it} \Delta_t - \sum_{\tau=t+1}^T f_{ifv\tau}(t)) \\ & + \sum_{f \in \mathcal{F}} \sum_{i \in \mathcal{I}} \sum_{\tau=0}^{t-1} k_{ifv\tau} x_{ifv\tau} f_{ifv\tau}(\tau) \leq q^s \Delta_t, \quad v \in \{t\}, t \in \mathcal{T}, \quad (\text{C.16}) \end{aligned}$$

$$\sum_{f \in \mathcal{F}} \sum_{i \in \mathcal{I}} \mu_{it} x_{ifvt} \leq q_t^o, \quad v \in \{o\}, t \in \mathcal{T}, \quad (\text{C.17})$$

$$q_t^o \leq \Delta_t Q^{o \max}, \quad t \in \mathcal{T}, \quad (\text{C.18})$$

$$x_{ifvt} = 0, \quad i \notin \mathcal{I}_{fv}, v \in \mathcal{V}, f \in \mathcal{F}, t \in \mathcal{T}, \quad (\text{C.19})$$

$$u_{ifn} = 0, \quad i \notin \mathcal{I}_f, f \in \mathcal{F}, 1 \leq n \leq F, \quad (\text{C.20})$$

$$x_{ifvt} \geq 0, \quad i \in \mathcal{I}_{fv}, v \in \mathcal{V}, f \in \mathcal{F}, t \in \mathcal{T}, \quad (\text{C.21})$$

$$u_{ifn} \geq 0, \quad i \in \mathcal{I}_f, f \in \mathcal{F}, 1 \leq n \leq F, \quad (\text{C.22})$$

$$a_f = 1 \quad f \in \mathcal{B}, \quad (\text{C.23})$$

$$a_f, r_n \in \{0, 1\}, \quad f \in \mathcal{F}, \quad (\text{C.24})$$

$$q^s \leq Q^{s \max}, \quad (\text{C.25})$$

$$I_f, q^s, q_t^o \in \mathbb{Z}, \quad v \in \mathcal{V}, f \in \mathcal{F}. \quad (\text{C.26})$$

The objective in Equation (C.1) aims to minimize the total network cost, consisting of fixed investments in opening facilities, $K_f^f a_f$; cost of scheduled transportation capacity, $\Delta_t c^a q^s$; cost of on-demand transportation agents that are hired per delivery, $c_v^o q_t^o$; total distance based travel cost of transportation agents, $d_{if} k_{ifvt} x_{ifvt}$; cost of lost sales, $c^s \mu_{it} (1 -$

$\sum_{f \in \mathcal{F}} \sum_{v \in \mathcal{V}} x_{ifvt}$); and cost of excess inventory, $c^e(I^* - \sum_{i \in \mathcal{I}} \sum_{l \in \mathcal{T}} \mu_{il} \sum_{f \in \mathcal{F}} \sum_{v \in \mathcal{V}} x_{ifvt})$. Constraints (C.2) ensure that no more than 100% of the total demand is allocated to facility-agent combinations. Note that unallocated demand is considered as lost sales. Constraints (C.3) ensure that the network-level order-up-to level depends on the safety stock coefficient associated with the number of activated facilities through the binary indicator r_n . Constraints (C.4) through (C.10) ensure the binary indicator $r_n = 1$ if and only if n facilities are activated, i.e., if the sum of the activated facilities $\sum_{f \in \mathcal{F}} a_f = n$. Constraints (C.11) ensure the sum of the order-up-to levels at the individual facilities equals the network-level order-up-to level. Constraints (C.12) limit the order-up-to level to the inventory capacity of a facility, and ensure that the order-up-to level for facilities that are not activated equals 0. Constraints (C.13) and (C.14) ensure that the allocated demand and safety stock never exceed the order-up-to inventory levels at facilities. Note that the demand and safety stock of a pixel could be allocated to different facilities. In addition, constraints (C.14) account for the available safety stock of the B&M demand. The system-wide safety stock from the B&M network that can be leveraged by the online network if n facilities are activated, $\sigma(z_n^* - z_n^\dagger)^+$, is proportionally divided over the number of B&M stores, B . Note that safety stock can only be leveraged if $z_n^* > z_n^\dagger$. Constraints (C.16) translate the demand allocation to a number of scheduled transportation agents required to satisfy demand. Note that the demand in previous time periods influences the transportation capacity required in the current time period since there are spillover effects of orders that are being delivered or still need to be delivered. The left-hand side of these constraints computes the resulting total quantity of transportation time required in a particular time period. We define the transportation capacity overflow $f_{ifvt}(t)$ and the consolidation factor k_{ifvt} in Appendix E.2. Constraints (C.17) ensure that the number of on-demand transportation agents can handle the allocated demand, while Constraints (C.18) impose a cap on the number of on-demand agents that can be deployed per hour, based on the average deployed on-demand agents per time period. Constraints (C.19) ensure that the demand of a pixel is not allocated to a facility-transportation combination that would lead to guaranteed late delivery in time period t . To ensure this, we introduce the set $\mathcal{I}_{fvt} = \{i \in \mathcal{I} | t_{ifvt}^d \leq l\}$. Similarly, Constraints (C.20) ensure that the safety stock of a pixel is not allocated to a facility that would lead to guaranteed late delivery. Constraints (C.23) ensure that retail facilities are always activated. Lastly, Constraints (C.21) through (C.26) limit the domain of the decision variables. Particularly, Constraints (C.25) limit the maximum number of scheduled transportation agents that can be contracted.

Appendix D. Discrete SO algorithm

We propose Algorithm 1 referred to as MetaAHA+ in the following, which builds on the MetaAHA algorithm proposed by Zhou et al. (2019), which in turn extends the AHA proposed by Xu et al. (2013).

Appendix D.1. Sampling from the hyperbox.

To sample solutions from the hyperbox, we follow the asymptotically uniform sampling mechanism of AHA. Per the approach outlined in Section 4.1, we define a vector $\tilde{\mathbf{y}}$, with values for the inventory at each facility and the number of scheduled transportation agents. Next, we set the value for the activation decisions $\tilde{\mathbf{a}}$ per Algorithm 2.

Algorithm 1 MetaAHA+ Algorithm

Initialization:

- 0.1 $k = 0, \mathcal{H}(k) = \Omega$
- 0.2 $\beta_0^e = \beta_0^s = \beta_0^A = \beta_{n0}^t = \beta_{n0}^l = 1$ for $1 \leq n \leq F, \beta_{i0}^p = 1$ for $i \in \mathcal{I}, \beta_0^v = \beta_0^a = \beta_{f0}^f = \beta_{f0}^l = 0$ for $f \in \mathcal{F}$
- 0.3 Randomly generate solution $\mathbf{y}_0^* = \mathbf{y}_0 \in \mathcal{H}(k)$ and determine $\hat{G}(\mathbf{y}_0)$ through simulation
- 0.4 $\mathcal{L} = \{\mathbf{y}_0\}$

Step 1: Determine $\mathcal{L}(k)$

- 1.1 $k = k + 1$
- 1.2 Obtain r points in $\mathcal{H}(k)$ based on the sampling mechanism defined in Appendix [Appendix D.1](#)
- 1.3 Obtain $\mathbf{y}_k^{\text{meta}}$, the solution to the two-staged analytical metamodel defined in Section 3
- 1.4 Obtain $\mathbf{y}_k^{\text{meta-hyper}}$, the solution to the two-staged analytical metamodel defined in Section 3 with the additional hyperbox constraints $\mathbf{y} \in \mathcal{H}(k)$

Step 2: Simulate performance of solutions $\mathbf{y} \in \mathcal{L}(k)$

- 2.1 Determine the number of additional simulations at the current iteration $\mathcal{A}_k(\mathbf{y})$
- 2.2 Simulate and determine $\hat{G}(\mathbf{y})$ based on all current and historic simulations
- 2.3 Determine $\mathbf{y}_k^* = \text{argmin}_{\mathbf{y}} \hat{G}(\mathbf{y})$
- 2.4 Determine $\mathcal{H}(k)$ based on Equations (D.1) and (D.2)

Step 3: Check for termination criteria

- 3.1 If \mathbf{y}_k^* is a local optimum following the procedure of AHA (Xu et al. 2013): Stop.
- 3.2 If the computational budget is depleted: Stop

Step 4: Update the metamodel

- 4.1 Evaluate any solution in \mathcal{L} that has not been evaluated using the model defined by Equations (C.1) through (C.26) to determine $g_A(\mathbf{y})$
- 4.2 Update Equation (2) using Equations (13) through (15)
- 4.3 Find β^p using Equations (16) and (17)
- 4.4 Fit the metamodel parameters in the objective of the second stage of the analytical metamodel using Equation (18)

Algorithm 2 Facility activation based on Hyperbox sampling

For $f \in \mathcal{F}$:

1. If $f \in \mathcal{B}$: $a_f = 1$
 2. Else, if $u_k^f = 1$ and $I_f = 1, a_f = 1$ with probability $p^{\mathcal{H}}$, else $I_f = 0$ and $a_f = 0$
-

Step 2 of Algorithm 2 ensures that the probability of activating facilities to store just one order is controlled. If a facility has an order-up-to level of 0 in the current iterate, and the hyperbox is upperbounded by 1, according to the definition of the hyperbox there is a 50% probability that the order-up-to level is set to 1, potentially causing a dedicated online facility to be activated. Step 2 reduces this probability through p^{Tt} , which we recommend to decrease as relative cost of opening facilities increases, to reduce the probability of generating solutions that provide little additional information to learn from.

Appendix D.2. Updating the hyperbox.

To update the hyperbox we compare the current abbreviated iterate $\tilde{\mathbf{y}}_k^*$ to the set of other sampled solutions, \mathcal{L} , according to the following equations:

$$l_k^{(d)} = \begin{cases} \max_{\mathbf{y} \in \mathcal{L}, \tilde{\mathbf{y}} \neq \tilde{\mathbf{y}}_k^*} \{\tilde{y}^{(d)} : \tilde{y}^{(d)} < \tilde{y}^{*,(d)}\} & \text{if it exists,} \\ -\infty & \text{otherwise,} \end{cases} \quad (\text{D.1})$$

$$u_k^{(d)} = \begin{cases} \min_{\mathbf{y} \in \mathcal{L}, \tilde{\mathbf{y}} \neq \tilde{\mathbf{y}}_k^*} \{\tilde{y}^{(d)} : \tilde{y}^{(d)} > \tilde{y}^{*,(d)}\} & \text{if it exists,} \\ \infty & \text{otherwise.} \end{cases} \quad (\text{D.2})$$

To limit the premature convergence of AHA and spend significant computation resources on exploring a small area around a, potentially sub-optimal, local minimum, Xu et al. (2013) have combined it with the multi-start Industrial Strength COMPASS (ISC) framework of Xu et al. (2010).

Appendix D.3. Determining the Least-Squares Weight

Let $\hat{\lambda}$ be the abbreviated decision vector of solution λ . In line with Osorio and Bierlaire (2013), we define $w_k(\lambda)$ as

$$w_k(\lambda) = 1/(1 + \|\hat{\lambda} - \hat{\mathbf{y}}_k^*\|_2). \quad (\text{D.3})$$

Appendix E. Supporting parameters

In this section, we define two auxiliary variables used in the model defined by Equations (C.1) through (C.26). In particular, we define the scheduled transportation capacity overflow function $f_{ifvt}(t')$ and the consolidation parameter k_{ifvt} .

Appendix E.1. Scheduled transportation capacity overflow

The scheduled transportation capacity overflow variable ensures that agents that start a delivery in one period do not suddenly finish as soon as the period finishes, their work carries over into the next period(s). More precisely, we

consider the time an agent spends in the subsequent periods after starting a delivery in a certain period. We define $\tau_{t't}$ as the time that has passed since the start of period t' and the start of period t ,

$$\tau_{t't} = \sum_{j=t'}^{t-1} \Delta_j. \quad (\text{E.1})$$

We can define nine different cases (A to I, see Figure E.7) with potential overflow of scheduled courier capacity from period t' into period t , that we categorize into two broader categories.

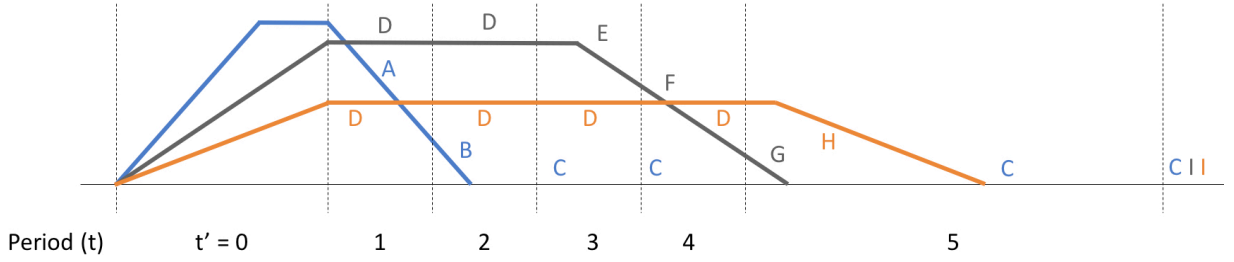


Figure E.7: Courier Overflow

First, the time a transportation agent spends on one order is smaller than the length of time period t' , therefore, agents start to free up as soon as t' is finished with rate $\mu_{it'} \Delta_t^{-1}$ (order density in orders/hr). The maximum number of transportation agents deployed simultaneously to deliver orders from time period t' is $\mu_{it'} t_{ifvt'}^o$.

- (A) At the beginning of time period t , transportation agents are finishing their job, but not all agents are finished by the end of t ($t_{ifvt'}^o \leq \Delta_{t'}$ and $\tau_{(t'+1)(t+1)} \leq t_{ifvt'}^o$). Then,

$$H_b = \mu_{it'} \Delta_t^{-1} (t_{ifvt'}^o - \tau_{(t'+1)(t)}), \quad H_e = \mu_{it'} \Delta_t^{-1} (t_{ifvt'}^o - \tau_{(t'+1)(t+1)}), \quad f_{ifvt}^A = 0.5 \Delta_t (H_b - H_e).$$

- (B) At the beginning of time period t , agents are finishing their job, and at some point in period t , all agents are done with their delivery ($t_{ifvt'}^o \leq \Delta_{t'}$ and $\tau_{(t'+1)(t)} \leq t_{ifvt'}^o \leq \tau_{(t'+1)(t+1)}$). Then,

$$H_b = \mu_{it'} \Delta_t^{-1} (t_{ifvt'}^o - \tau_{(t'+1)(t)}), \quad H_h = t_{ifvt'}^o - \tau_{(t'+1)(t)}, \quad f_{ifvt}^B = 0.5 H_b H_h.$$

- (C) All agents are done delivering items from period t' in period t ($t_{ifvt'}^o \leq \Delta_{t'}$ and $t_{ifvt'}^o \leq \tau_{(t'+1)(t)}$). Then,

$$f_{ifvt}^C = 0.$$

Second, the time agents spend on one order ($t_{ifvt'}^o$ is larger than the length of time period t' , therefore, agents start to free up after $t_{ifvt'}^o$, with rate $\mu_{it'}$. The maximum number of agents deployed simultaneously to deliver orders from time

period t' is $\mu_{it'}\Delta_{t'}$.

(D) All throughout period t , agents are busy delivering orders from period t' ($t_{ifvt'}^o \geq \Delta_{t'}$ and $\tau_{t'(t+1)} \leq t_{ifvt'}^o$). Then,

$$f_{ifvt}^D = \mu_{it'}\Delta_{t'}.$$

(E) At the beginning of time period t , all agents are still busy delivering orders from period t' , but somewhere in period t , the first agents start finishing up. However, by the end of period t , not all agents are finished yet ($t_{ifvt'}^o \leq \Delta_{t'}$ and $\tau_{t'(t)} \leq t_{ifvt'}^o \leq \tau_{t'(t+1)} \leq t_{ifvt'}^o + \Delta_{t'}$). Then,

$$S = \mu_{it'}(t_{ifvt'}^o - \tau_{t'(t)}), \quad H_b = \mu_{it'}, \quad H_e = \mu_{it'} - \mu_{it'}\Delta_{t'}^{-1}(\tau_{t'(t+1)} - t_{ifvt'}^o),$$

$$f_{ifvt}^E = 0.5(H_b + H_e)(\tau_{t'(t+1)} - t_{ifvt'}^o) + S.$$

(F) At the beginning of time period t , agents are finishing their job, but not all agents are finished by the end of t ($t_{ifvt'}^o \leq \Delta_{t'}$ and $t_{ifvt'}^o \leq \tau_{t'(t)} \leq \tau_{t'(t+1)} \leq t_{ifvt'}^o + \Delta_{t'}$). Then,

$$H_b = \mu_{it'} - \mu_{it'}\Delta_{t'}^{-1}(\tau_{t'(t)} - t_{ifvt'}^o), \quad H_e = \mu_{it'} - \mu_{it'}\Delta_{t'}^{-1}(\tau_{t'(t+1)} - t_{ifvt'}^o), \quad f_{ifvt}^F = 0.5(H_b + H_e)\Delta_{t'}.$$

(G) At the beginning of time period t , agents are finishing their job, and at some point in period t , all agents are done with their delivery ($t_{ifvt'}^o \leq \Delta_{t'}$ and $t_{ifvt'}^o \leq \tau_{t'(t)} \leq t_{ifvt'}^o + \Delta_{t'} \leq \tau_{t'(t+1)}$). Then,

$$H_b = \mu_{it'} - \mu_{it'}\Delta_{t'}^{-1}(\tau_{t'(t)} - t_{ifvt'}^o), \quad f_{ifvt}^G = 0.5H_b(t_{ifvt'}^o + \Delta_{t'} - (\tau_{t'(t)} - t_{ifvt'}^o)).$$

(H) At the beginning of time period t , all agents are still busy delivering orders from period t' , but somewhere in period t , the first agents start finishing up. By the end of period t , all agents are finished ($t_{ifvt'}^o \leq \Delta_{t'}$ and $\tau_{t'(t)} \leq t_{ifvt'}^o \leq t_{ifvt'}^o + \Delta_{t'} \leq \tau_{t'(t+1)}$). Then,

$$S = \mu_{it'}(t_{ifvt'}^o - \tau_{t'(t)}), \quad H_b = \mu_{it'}, \quad f_{ifvt}^H = 0.5H_b\Delta_{t'} + S.$$

(I) All agents are done delivering items from period t' in period t ($t_{ifvt'}^o \leq \Delta_{t'}$ and $t_{ifvt'}^o + \Delta_{t'} \leq \tau_{t'(t)}$). Then,

$$f_{ifvt}^I = 0.$$

Integrating the cases presented above, the formulation for $f_{ifvt}(t')$ leads to

$$f_{ifvt}(t') = \begin{cases} f_{ifvt}^A & \text{for } t_{ifvt}^o \leq \Delta_{t'} \text{ and } \tau_{(t'+1)(t+1)} \leq t_{ifvt}^o, \\ f_{ifvt}^B & \text{for } t_{ifvt}^o \leq \Delta_{t'} \text{ and } \tau_{(t'+1)(t)} \leq t_{ifvt}^o \leq \tau_{(t'+1)(t+1)}, \\ f_{ifvt}^C & \text{for } t_{ifvt}^o \leq \Delta_{t'} \text{ and } t_{ifvt}^o \leq \tau_{(t'+1)(t)}, \\ f_{ifvt}^D & \text{for } t_{ifvt}^o \geq \Delta_{t'} \text{ and } \tau_{t'(t+1)} \leq t_{ifvt}^o, \\ f_{ifvt}^E & \text{for } t_{ifvt}^o \leq \Delta_{t'} \text{ and } \tau_{t'(t)} \leq t_{ifvt}^o \leq \tau_{t'(t+1)} \leq t_{ifvt}^o + \Delta_{t'}, \\ f_{ifvt}^F & \text{for } t_{ifvt}^o \leq \Delta_{t'} \text{ and } t_{ifvt}^o \leq \tau_{t'(t)} \leq \tau_{t'(t+1)} \leq t_{ifvt}^o + \Delta_{t'}, \\ f_{ifvt}^G & \text{for } t_{ifvt}^o \leq \Delta_{t'} \text{ and } t_{ifvt}^o \leq \tau_{t'(t)} \leq t_{ifvt}^o + \Delta_{t'} \leq \tau_{t'(t+1)}, \\ f_{ifvt}^H & \text{for } t_{ifvt}^o \leq \Delta_{t'} \text{ and } \tau_{t'(t)} \leq t_{ifvt}^o \leq t_{ifvt}^o + \Delta_{t'} \leq \tau_{t'(t+1)}, \\ f_{ifvt}^I & \text{for } t_{ifvt}^o \leq \Delta_{t'} \text{ and } t_{ifvt}^o + \Delta_{t'} \leq \tau_{t't}. \end{cases} \quad (\text{E.2})$$

Appendix E.2. Consolidation factor

To account for the reduction in transportation capacity requirements through consolidation (i.e., assigning multiple orders to one transportation agent), we introduce a consolidation factor k_{ifvt} in the model defined by Equations (C.1) through (C.26). To approximate the effect of consolidation in pixel i , we consider the available ‘slack’ a courier has within the available time until the delivery deadline when delivering from facility f to pixel i , t_{ifvt}^s , and the potential consolidation density in orders per hour in pixel i and time period t , μ_{it}^c . We can compute t_{ifvt}^s by subtracting the time required to deliver the order (t_{ifvt}^d), including picking, courier response, traveling and loading time, from the promised delivery lead-time (l) as

$$t_{ifvt}^s = l - t_{ifvt}^d. \quad (\text{E.3})$$

Furthermore, we can compute μ_{it}^c by defining the Neighborhood of a pixel i , $\mathcal{N}(i)$ based on the maximum pixel-to-pixel consolidation distance d^c , and computing the order density in the neighborhood of pixel i as

$$\mathcal{N}(i) = \{i' | d_{ii'} \leq d^c\}, \quad \mu_{it}^c = \sum_{i' \in \mathcal{N}(i)} \sum_{s \in S} \frac{\mu_{i't}}{\Delta_t}. \quad (\text{E.4})$$

The maximum consolidation factor, i.e., the proportion of original trips required with consolidation, is the maximum of the inverse of the carrying capacity of a vehicle (ξ_v) and a function of the density of orders arriving during the ‘slack’ time. We formally define the consolidation factor as

$$k_{ifvt} = \max\left(\frac{1}{\xi_v}, \min\left(1, \frac{1}{t_{ifvt}^s * \mu_{it}^c}\right)\right). \quad (\text{E.5})$$

Appendix F. Problem Instance Definition

This appendix supports the introduction of the case study in Section 6 that is leveraged for our analysis.

Appendix F.1. Demand Instances

Systemwide Demand Density Distributions. For demand instances *uniform* (U), *concentrated* (C), and *dynamic* (D), we generate a target demand level $D_\omega > 0$ for every scenario ω based on a normal distribution with parameters $\mu = 500$ and $\sigma = 100$. Next, we generate order interarrival times through a demand level distribution that governed by the exponential distribution defined by parameter λ_ω , such that $D_\omega = \frac{T}{\lambda_\omega}$, i.e., the expected demand throughout the day is equal to D_ω . We follow a similar process for demand instance *independent* (I), except that we generate a target demand level $D_{\omega j'} > 0$ for every of the j demand areas based on normal distribution with parameters $\mu^j = \frac{\mu}{j}$ and $\sigma^j = \frac{\sigma}{\sqrt{j}}$. Consequently, we find an exponential parameter $\lambda_{\omega j'}$ that is different for every specific demand area j' .

Geographic Demand Distribution. For *uniform* (U), demand is uniformly distributed over the demand area. For both the *concentrated* (C) and *dynamic* (D) cases, we define a parameter ζ , to indicate the probability that an order belongs to a demand cluster. Consequently, with probability $1 - \zeta$, an order does not belong to the demand cluster and is uniformly distributed over the demand area. In the *concentrated* (C) case, an order assigned to the cluster is randomly located in a circle with centroid (x^c, y^c) and radius r . Similarly, in the *dynamic* (D) case, an order is assigned to a circle with radius r , but the center of the circle depends on the time (x_t^c, y_t^c) . The center of the circle moves linearly over time from (x_0^c, y_0^c) to $(x_{\text{end}}^c, y_{\text{end}}^c)$. For *independent* (I), demand in each demand area j' is independently, uniformly distributed over that particular area.

B&M Demand Distribution. We generate a target demand level $D_{b\omega} > 0$ for every B&M facility b and scenario ω based on a multivariate normal distribution with parameters $\mu_b = 100$, $\sigma_b = 20$, and $\rho_{bb'} = \frac{2}{3}$ to account for the correlation of demand between B&M stores. The data of GFR indicates that the correlation between B&M and online demand is negligible, indicating that shopping patterns for the different channels are independent processes.

Appendix F.2. Network Parameters

Inspired by Manhattan, we define the demand area as a rectangular area of 100km^2 (5km by 20km). The parameters defining the transportation agents and facilities are the same for every instance, inspired by the real-world case study. In particular, both scheduled and on-demand transportation agents are modeled as bike couriers. Scheduled agents are paid by the hour, while on-demand couriers are paid a fixed starting price per trip, in addition to a distance-based cost component. Existing B&M and dedicated online facilities differ in terms of fixed activation cost and available B&M inventory, but are otherwise identical. We generate 10 potential locations for dedicated online facilities using Algorithm 3, proposed by [Snoeck and Winkenbach \(2022\)](#). In addition, we generate 10 locations with existing B&M facilities based on a p-median problem, of which a subset is activated depending on the problem instance. We use a Euclidean distance metric throughout the analysis. While the stylized problem instances and

the network parameters are informed by a real-life case study, we believe that the results discussed in Section 6 are representative for many other real-life applications.

Facility Generation Algorithm. Since the stylized problem instances do not have actual proposed facility locations, we generate those using Algorithm 3. Note that this is just one potential mechanism to generate the facility locations, any other could be used as well. First, we generate potential facility locations by solving a p-median problem. However, in real-life cases, potential facility locations are rarely found at the optimally suggested locations, particularly in dense urban areas. To mimic this additional real-life constraint, we add a random geographical shift to the locations proposed.

Algorithm 3 Algorithm to generate facility locations for stylized problem instances ([Snoeck and Winkenbach 2022](#))

Step 1: Generate potential locations

1. Raster the demand area with dimensions X and Y into square pixels
2. Take the centroid of every pixel as demand location

Step 2: Solve p-median problem

1. Determine number (p) of potential facility locations to be included in the model
2. Solve p-median with demand locations

Step 3: Randomize locations

1. Define relative randomization as percentage z
 2. For each location i with coordinates (x_i, y_i) suggested by the p-median solution
 - Generate two uniform random numbers from $U(-1, 1)$: u_x, u_y
 - Find randomized location $(x_i + zXu_x, y_i + zYu_y)$, where X and Y are the dimensions of the area of operations
-