

Inventory, Market Making, and Liquidity in OTC Markets*

Assa Cohen[†] Mahyar Kargar[‡] Benjamin Lester[§] Pierre-Olivier Weill[¶]

November 30, 2023

Abstract

We develop a search-theoretic model of a dealer-intermediated over-the-counter market. Our key departure from the literature is to assume that, when a customer meets a dealer, the dealers can only sell assets that they already own. Hence, in equilibrium, dealers choose to hold *inventory*. We derive the equilibrium relationship between dealers' cost of holding assets on their balance sheets, their optimal inventory holdings, and various measures of liquidity, including bid-ask spreads, trade size, volume, and turnover. Using transaction-level data from the corporate bond market, we calibrate the model to quantitatively assess the impact of post-crisis regulations on dealers' inventory costs, liquidity, and welfare.

Keywords: Over-the-counter markets, intermediation, liquidity, dealer inventory, financial regulation

JEL Classification: G11, G12, G21.

*We thank Jack Bao, Briana Chang, Tim Johnson, Charlie Kahn, Lucie Lebeau, Mariano Palleja, Dejanir Silva; conference participants at AEA 2023, SED 2021, SFS Cavalcade 2022, Search and Matching in Macro and Finance Virtual Seminar Series, the Virtual Finance Theory Seminar; and seminar participants at CMU, the Federal Reserve Bank of Atlanta, the Federal Reserve Bank of Philadelphia, FGV, UCI, UCLA, UCSB, UIUC, Wisconsin, and Wharton Macro Lunch. Lian Chen provided expert research assistance. The views expressed here are those of the authors and do not necessarily reflect the views of the Federal Reserve Bank of Philadelphia or the Federal Reserve System. All errors are our own responsibility.

[†]Yeshiva University Sy Syms School of Business. Email: assac@sas.upenn.edu

[‡]University of Illinois at Urbana-Champaign. Email: kargar@illinois.edu

[§]Federal Reserve Bank of Philadelphia. Email: benjamin.lester@phil.frb.org

[¶]UCLA, NBER and CEPR. Email: poweill@econ.ucla.edu

1 Introduction

In many over-the-counter (OTC) markets, dealer banks provide liquidity through their willingness to hold *inventory*: they absorb assets onto their balance sheets when investors need to sell quickly, and they use these assets to fulfill investors’ buy orders without delay. After the Global Financial Crisis (GFC) of 2007-2008, several regulations were introduced that increased the cost to dealers of holding inventory.¹ Not surprisingly, dealers responded by reducing their inventory holdings. For instance, according to data from the Flow of Funds, the share of outstanding corporate bonds and non-agency mortgage-backed securities held by broker-dealers fell from 2-3% in 2006 to less than 1% in 2018.² At the same time, both market participants and academics alike argued that post-GFC regulations posed a threat to market liquidity.³ Since maintaining liquid financial markets is crucial for a well-functioning economy, understanding and quantifying the effects of post-GFC regulations on market liquidity and welfare has emerged as a central challenge.

In this paper, we develop a structural model of dealer-intermediated OTC markets in order to meet this challenge. Our starting point is the benchmark search-theoretic framework developed by [Duffie, Gârleanu, and Pedersen \(2005\)](#) and extended to allow for arbitrary preferences and asset holdings by [Lagos and Rocheteau \(2009\)](#) and [Gârleanu \(2009\)](#). However, a key abstraction in these papers is that inventories do not play any economic role for market making. Indeed, in these models, dealers *never* hold inventory—they merely enable customers to access a frictionless market, for which they charge a fee. Our innovation is to bring inventories back into market making with a simple and, arguably, quite natural constraint: we assume that a dealer can only sell to customers the assets that she currently holds in inventory. As a result of this “inventory-in-advance” constraint, a central feature of our model is that dealers choose an optimal amount of inventory in order to provide liquidity to their customers. Technically, this model is more difficult to analyze than its

¹These regulations include the 2010 Basel III framework, which introduced enhanced capital and liquidity requirements, along with the so-called “Volcker rule,” which reduced implicit government guarantees (thus increasing banks’ funding costs) and began monitoring banks’ inventory holdings in concert with the regulation’s ban on proprietary trading.

²Source: Table L.213 of the Federal Reserve’s Flow of Funds, shown in Figure 6.

³See the extensive discussions by [Thakor \(2012\)](#), [Duffie \(2017\)](#), [Bessembinder, Jacobsen, Maxwell, and Venkataraman \(2018\)](#), and the many references therein.

predecessors because it no longer admits closed form solutions. However, using standard recursive methods, we can characterize the equilibrium and study how dealers' optimal inventory holdings depend on various features of the economic environment, such as the frequency and variation in investors' preferences—which captures their trading needs—and the flow (dis)utility dealers receive from holding assets themselves—which captures the effects of regulatory costs imposed by policymakers. Hence, our model provides a structural framework to evaluate the effects of various regulatory or technological changes in OTC markets on asset prices, transaction costs, trading volume, dealers' profits, and investors' surplus.

Then, we apply the framework to the secondary market for U.S. corporate bonds, to quantitatively assess the impact of rising inventory costs associated with post-GFC regulations. The quantitative analysis proceeds in two steps.

First, using transaction-level data—and focusing on transactions above \$1 million made by institutional investors—we calibrate the model to match several target moments constructed from the corporate bond market data before the GFC. Given the parameter values implied by our calibration, we find that the inventory-in-advance constraint had a relatively modest effect on equilibrium outcomes, relative to an environment where dealers did not face such a constraint. For example, we find that welfare loss in the equilibrium with inventory constraints is approximately 44% larger than an environment without inventory constraints (but with search and bargaining frictions). Intuitively, the implied cost of holding inventory before the GFC was relatively small, and hence dealers held sufficient inventory to fulfill most customer-buy orders in full.

Second, we quantify the implicit cost of regulations to dealers. More specifically, holding all other structural parameters fixed, we increase dealers' cost of holding assets to levels consistent with the aggregate decline in dealers' inventory observed in the data. We find that dealers' inventory costs must increase tenfold, from about 4% to almost 40% of the asset coupon—or, dividing by $1/0.05 = 20$, from 0.2% to 2% of the asset's face value. As a result of this increase in dealers' balance sheet costs, our model predicts that trading costs (i.e., bid-ask spreads) rise by approximately 40%

relative to the pre-GFC benchmark, from 10 to 14 basis points (bps). This represents nearly 80% of the increase that we observe in our data across the same time periods.

An important advantage of our structural approach is that it allows us to go beyond the analysis of trading costs and generate additional predictions that would be difficult to make in reduced-form models. For one, it allows us to measure welfare, which is crucial for distinguishing between distributional effects—such as shifts in the share of surplus that accrues to dealers vs. customers—and distortions to the efficient allocation. Indeed, we find that the welfare cost of frictions, stemming from the combination of search and inventory constraints, increased substantially after the introduction of post-GFC regulations, from 1.25% to 2.4% of the total gains from trade. Moving beyond the welfare of investors trading in OTC market, we turn to the cost of capital and calculate the model-implied change in liquidity yield spread. We find that this spread increased by a factor of 2.5, going from 2 to 5 bps.

Related literature

Our theory contributes to the the literature that uses search-theoretic models of trade to study OTC markets. Many of these papers build off of the basic framework developed in [Duffie, Gârleanu, and Pedersen \(2005\)](#), including the important contributions by [Lagos and Rocheteau \(2009\)](#) and [Gârleanu \(2009\)](#), who extend the basic framework to accommodate arbitrary preferences and asset holdings. Importantly, in most papers within this literature, dealers are assumed to have unfettered access to a frictionless, inter-dealer market, which obviates the need for any dealer to hold inventory. Such papers include, but are not limited to, [Feldhütter \(2012\)](#), [Lester, Rocheteau, and Weill \(2015\)](#), [Milbradt \(2017\)](#), [Pagnotta and Philippon \(2018\)](#), and [Lagos and Zhang \(2020\)](#), [Kargar, Passadore, and Silva \(2020\)](#), [Pinter and Üslü \(2021\)](#), [Palleja \(2022\)](#), and [Li \(2023\)](#). See [Weill \(2020\)](#) for a thorough review of the literature. Of course, the result that dealers hold no inventories makes these models more tractable, highlighting the important role of search and bargaining frictions in the determination of prices and allocations. However, it also makes them ill-suited to study dealers' incentives to hold inventory and provide liquidity in response to various changes in the economic

environment, and the consequences for asset prices, transaction costs, trade size, volume, and welfare.

In the literature on search-based OTC markets, several papers have proposed models of dealers' inventory management. In [Weill \(2007\)](#) and [Lagos, Rocheteau, and Weill \(2008\)](#), for example, dealers find it optimal to hold inventories in anticipation of aggregate fluctuations in customers' demand. However, in both environments, optimal inventory holdings are always zero in the long run, that is, in the non-stochastic steady state. In our model, inventories play a non-trivial economic role even in the non-stochastic steady state, which we believe is an important feature for studying the long-run decline in inventories between 2008 and 2018. Our work is also related to [An \(2018\)](#), who shows that, despite the presence of holding costs, imperfectly competitive dealers have incentive to hold inventories in order to gain market power with their customers. [Tse and Xu \(2021\)](#) also develop a model where dealers (with different trading capacity) carry inventory in order to rationalize empirical observations about inter-dealer trades in OTC markets. Subsequent attempts to incorporate inventory into OTC models of trade include [Diao, Dudley, and Sun \(2023\)](#) and [Dyskant, Silva, and Sultanum \(2023\)](#).

There are also a number of papers in which all agents, including those who play the role of dealers, trade in decentralized markets. See, e.g., [Hugonnier, Lester, and Weill \(2020, 2022\)](#), [Shen, Wei, and Yan \(2021\)](#), [Üslü \(2019\)](#), [Farboodi, Jarosch, and Shimer \(2022\)](#), [Farboodi, Jarosch, and Menzio \(2017\)](#), [Bethune, Sultanum, and Trachter \(2022\)](#), [Yang and Zeng \(2019\)](#), and [Nosal, Wong, and Wright \(2019\)](#). In these models, since agents face a short-selling constraint, those who play the role of dealers must hold inventories. Though these models have proven useful in studying the determinants of *inter-dealer* market structure and trading patterns, we assume instead that the inter-dealer market is centralized. This simplification allows us to focus our analysis more squarely on the issue at hand; to derive new, testable implications regarding, e.g, the relationship between dealers' inventory costs and the distribution of trade size.

In addition to our substantive contribution, we also make several methodological contributions. Indeed, the inventory-in-advance constraint implies that our model no longer admits closed-form

solutions for the value functions and distributions. Hence, to characterize equilibria, we adapt the standard recursive methods of [Stokey and Lucas \(1989\)](#) to our environment to formally establish key properties of the equilibrium. See also [Rocheteau, Weill, and Wong \(2018\)](#) and [Choi and Rocheteau \(2021\)](#) for related methodological contributions in the New Monetarist literature.

Outside of search-based models, there is also, of course, a celebrated literature on inventory management by dealers, starting with [Amihud and Mendelson \(1980\)](#), [Ho and Stoll \(1981, 1983\)](#), and [Mildenstein and Schleef \(1983\)](#). Relative to this literature, our main contribution is to consider a model in which customers' supply and demand are derived from explicit, dynamic optimization problems, subject to search frictions. This enables us to quantify the gains from trade created by the inter-dealer market, and offer a welfare analysis of post-GFC regulations.

Because of its quantitative focus, our work is also related to papers who structurally estimate models of OTC markets, either search-based as in [Feldhütter \(2012\)](#), [Gavazza \(2016\)](#), [Brancaccio, Li, and Schurhoff \(2017\)](#), [Hendershott, Li, Livdan, and Schürhoff \(2020\)](#), [Liu \(2020\)](#), [Pinter and Üslü \(2021\)](#), [Brancaccio and Kang \(2022\)](#), or network-based as in [Gofman \(2014, 2017\)](#), and [Eisfeldt, Herskovic, Rajan, and Siriwardane \(2023\)](#). We contribute to this literature by developing a new model and focusing on a different market phenomena.

Finally, given the focus of our application, our paper is related to several recent empirical studies that have attempted to identify the effect of post-crisis regulations on market liquidity, including [Tebbi and Xiao \(2019\)](#), [Bao, O'Hara, and Zhou \(2018\)](#), [Bessembinder, Jacobsen, Maxwell, and Venkataraman \(2018\)](#), [Dick-Nielsen and Rossi \(2019\)](#), and [Choi, Huh, and Shin \(2023\)](#). By studying this issue within the context of a structural equilibrium model, our analysis complements these existing empirical exercises in several important ways. First, by calibrating our model to match moments before and after the introduction of new regulations, we are able to infer the implicit cost of these regulations on dealers; this cost is difficult to measure directly and, to the best of our knowledge, such an estimate is new to the literature. Second, while existing empirical studies based on difference-in-difference regressions identify "local" effects of new regulations on a particular measure of liquidity, such as price impact, our model allows us to explore the broader

implications of policy for the behavior of customers and dealers, and the subsequent implications for a variety of outcomes, both observable (such as bid-ask spreads, trade size, or volume) and unobservable (such as the time customers wait to complete their trade). Third, and perhaps most important, our structural equilibrium model provides natural measures of welfare, along with the opportunity to perform counterfactuals, which is crucial for evaluating the quantitative impact of policy.

The remainder of the paper has two parts. In Section 2, we describe the model, show that an equilibrium exists, and study analytically a number of its properties. In Section 3, we calibrate the model to the U.S. Corporate Bond market and study the welfare impact of post-GFC regulation.

2 The Model

We consider a continuous time, infinite horizon model of an over-the-counter asset market in the spirit of Gârleanu (2009) and Lagos and Rocheteau (2009). There are two types of infinitely-lived agents: a measure of customers normalized to one and a measure $\mu > 0$ of dealers. There is one asset that is durable, perfectly divisible, and in fixed supply, $s > 0$.

We assume that customers have stochastically varying preferences defined over the quantity of asset they hold, and a numéraire consumption good. In particular, let $u(q, \delta) + c$ denote a customer's flow utility, where $q \geq 0$ denotes the units of asset the customer holds, δ denotes her current preferences for assets, and c denotes her net consumption (or production if negative) of the numéraire good. We assume that $u(q, \delta)$ is strictly increasing and strictly concave in $q > 0$, continuously differentiable, and satisfies the Inada conditions $\lim_{q \rightarrow 0} u_q(q, \delta) = +\infty$ and $\lim_{q \rightarrow \infty} u_q(q, \delta) = 0$. We also assume that $u_q(q, \delta)$ is strictly increasing in δ , where u_q denotes the partial derivative with respect to q . Hence a larger preference shock δ creates a stronger demand for the asset.

Preference shocks arrive at rate γ , at which time a new δ' is drawn according to the cumulative distribution function (CDF) $F(\delta')$.⁴ We assume that the CDF has support included in some compact

⁴Micro-foundations for such a specification have been provided earlier in the literature. For example, under

interval $[\underline{\delta}, \bar{\delta}]$ but otherwise make no other restriction; in particular, the CDF can be discrete (as in, e.g., [Lagos and Rocheteau, 2009](#)), continuous, or a mixture of the two. For simplicity, we assume that dealers have linear preferences that do not change over time: a dealer receives flow utility $vq + c$ from holding q units of the asset in inventory and consuming c units of the numéraire good. All agents discount the future at rate $r > 0$.

Dealers have continuous access to a frictionless, competitive market where they can buy or sell any amount of the asset at price $P > 0$. Customers do not meet each other and trade directly. Instead, customers meet a randomly chosen dealer at independent Poisson arrival times with intensity λ . If there are gains from trade, the two bargain over the terms of trade. We denote by $\theta \in [0, 1]$ the dealers' bargaining power.

Our key departure from the existing literature is an inventory-in-advance constraint: when a dealer meets a customer, she can buy any quantity of assets from the customer, but she can only sell assets that she currently holds in inventories. After completing a transaction, a dealer can then access the inter-dealer market and rebalance her portfolio, either selling the assets she just accumulated or buying assets to restore an optimal level of inventory.

2.1 Customers

Let $V(q, \delta)$ denote the maximum attainable expected discounted utility of a customer with current asset holdings q and preferences δ . The Hamilton-Jacobi-Bellman (HJB) equation for $V(q, \delta)$ can be written as:

$$rV(q, \delta) = u(q, \delta) + \gamma \mathbb{E}^F [V(q, \delta') - V(q, \delta)] + \lambda [V(q', \delta) - V(q, \delta) - P(q' - q) - \phi]. \quad (1)$$

where $\mathbb{E}^F [\cdot]$ denotes the expectation with respect to the CDF $F(\delta')$. The interpretation of the HJB equation is standard: the customer enjoys the flow utility $u(q, \delta)$ until one of two events occurs. First, at rate γ , a preference shock arrives, at which time a new δ' is drawn from $F(\delta')$. Second,

appropriate specification, $u(q, \delta)$ represents the flow certainty equivalent of holding q units of the asset. See [Weill \(2020\)](#) for a survey.

at rate λ , the customer has the opportunity to trade with a dealer. At this time, the dealer transfers $q' - q$ units of the asset in exchange for the payment $P(q' - q) + \phi$. This payment is comprised of the cost (or revenue) of purchasing (selling) the asset at the inter-dealer price, $P(q' - q)$, plus an intermediation fee, ϕ .

A customer in state (q, δ) and a dealer holding $i \geq 0$ units of the asset choose a pair $(\hat{q}, \hat{\phi})$ to maximize the Nash product

$$[V(\hat{q}, \delta) - V(q, \delta) - P(\hat{q} - q) - \hat{\phi}]^{1-\theta} \hat{\phi}^\theta,$$

subject to the inventory-in-advance constraint

$$0 \leq \hat{q} \leq q + i. \tag{2}$$

Maximizing with respect to \hat{q} reveals that the optimal post-trade asset holding, q' , maximizes the trade surplus,

$$q' \in \arg \max V(\hat{q}, \delta) - V(q, \delta) - P(\hat{q} - q) \tag{3}$$

subject to (2). Given the value q' that solves this program, the transfer ϕ is set so that the dealer appropriates a fraction θ of the maximized joint surplus:

$$\phi = \theta [V(q', \delta) - V(q, \delta) - P(q' - q)]. \tag{4}$$

In what follows, we will adopt the usual convention of using a lower case i to denote an individual dealer's inventory and an upper case I to denote the choice of other dealers. Therefore, in a symmetric, steady-state equilibrium in which $i = I$, substituting (3) and (4) into the HJB

equation yields

$$rV(q, \delta) = u(q, \delta) + \gamma \mathbb{E}^F [V(q, \delta') - V(q, \delta)] \\ + \lambda(1 - \theta) \max_{0 \leq q' \leq q+I} \{V(q', \delta) - V(q, \delta) - P(q' - q)\}.$$

Informally differentiating with respect to q and applying the envelope condition yields

$$rV_q(q, \delta) = u_q(q, \delta) + \gamma \mathbb{E}^F [V_q(q, \delta') - V_q(q, \delta)] \\ + \lambda(1 - \theta) [\max \{V_q(q + I, \delta), P\} - V_q(\delta, q)].$$

Let $\Sigma(q, \delta) \equiv V_q(q, \delta) - P$ denote the marginal trade surplus, i.e., the marginal value to a customer of an additional unit of asset, net of the inter-dealer price. We can rewrite the expression above as

$$[r + \gamma + \lambda(1 - \theta)] \Sigma(q, \delta) = u_q(\delta, q) - rP + \gamma \mathbb{E}^F [\Sigma(q, \delta')] \\ + \lambda(1 - \theta) \max \{\Sigma(q + I, \delta), 0\}. \quad (5)$$

Since equation (5) characterizes $\Sigma(q, \delta)$ for any given (P, I) , it will sometimes be helpful to make this dependence explicit by writing $\Sigma(q, \delta | P, I)$; otherwise, we will suppress this dependence to simplify notations. Notice that the environment of [Lagos and Rocheteau \(2009\)](#), where there is no inventory-in-advance constraint, corresponds to the case where $I \rightarrow \infty$ and the final term in (5) disappears. Our first Proposition studies the fixed point equation defined by (5).

Proposition 1. *Equation (5) admits a unique, continuous solution $\Sigma(\cdot)$ that has the following properties:*

1. *it is the basis of a solution to the HJB equation (1).*
2. *it is strictly increasing in δ , strictly decreasing in q and P , and weakly decreasing in I ;*
3. *for all $\delta \in [\underline{\delta}, \bar{\delta}]$, $\lim_{q \rightarrow 0} \Sigma(q, \delta) = \infty$;*
4. *there exists \hat{q} such that, for all $\delta \in [\underline{\delta}, \bar{\delta}]$ and $q > \hat{q}$, $\Sigma(\delta, q) < 0$.*

The first point states that a value function $V(q, \delta)$ solving the HJB equation (1) can be constructed based on $\Sigma(q, \delta)$; the details are in the appendix. Importantly, the construction confirms that the envelope condition that we used informally earlier indeed holds. The properties in the last three points are inherited from the flow marginal value, $u_q(q, \delta) - rP$, except for one: the marginal surplus is decreasing in aggregate inventories, I . Indeed, if I is smaller, customers anticipate that the inventory-in-advance constraint is more likely to bind in the future. This makes it harder for them to accumulate assets, reduces their asset holding and raises their marginal value for the asset.

The function $\Sigma(\cdot)$ entirely characterizes a customer's optimal trading behavior. To see this, note that if the customer and the dealer were unconstrained by inventories, then they would trade to the "target holding" $q^*(\delta | P, I)$ such that the marginal trade surplus is equal to zero, i.e.,

$$\Sigma(q^*(\delta | P, I), \delta | P, I) = 0.$$

Proposition 1 ensures that this equation has a unique solution. It also implies some intuitive relationships between an individual customer's current state, the aggregate state, and the customer's target asset holdings. In particular, as one might expect, the customer's target asset position, $q^*(\delta | P, I)$, is increasing in his idiosyncratic preference shock δ and decreasing in the price P . A new feature of our model is induced by the inventory constraints: customers now have incentive to acquire additional assets out of precautionary motives. This incentive grows stronger as I declines. That is, since an additional unit of the asset is more valuable when dealers hold less inventory, ceteris paribus, $q^*(\delta | P, I)$ is decreasing in I .

2.2 Dealers

Let $\Phi(q, \delta)$ denote the joint distribution of asset holdings and preference shocks across customers. We characterize this distribution below and note for now that optimal trading behavior implies that its support is included in $[0, \bar{q}] \times [\underline{\delta}, \bar{\delta}]$, for some $\bar{q} > q^*(\bar{\delta})$. Using the Nash bargaining solution,

we can write the dealer's (flow) profit function as

$$r\Pi(i) = (v - rP)i + \frac{\lambda}{\mu} \theta \int_{(q', \delta')} \max_{0 \leq q'' \leq q' + i} \{V(q'', \delta') - V(q', \delta') - P(q'' - q')\} d\Phi(q', \delta'),$$

where we use $\int_{(q', \delta')}$ to denote the integral over $(q', \delta') \in [0, \bar{q}] \times [\underline{\delta}, \bar{\delta}]$. Hence, the dealer's objective function has two components: the flow payoff from owning i units of the asset, $(v - rP)i$; and the expected capital gains from trading with a randomly selected customer.

Lemma 1. *For any $\Phi(q, \delta)$, P , and I , the profit function $\Pi(i)$ is concave and continuously differentiable in i , with derivative:*

$$\frac{d\Pi}{di}(i) = v - rP + \frac{\lambda}{\mu} \theta \int_{(q', \delta')} \max\{\Sigma(q' + i, \delta'), 0\} d\Phi(q', \delta').$$

The expression for the derivative of the profit function is intuitive. The first term is the direct flow utility that a dealer enjoys by holding a marginal unit of the asset. The second term is the user cost: what the dealer has to pay per unit of time to hold a marginal unit of the asset. The third term is the marginal impact of increasing inventory on intermediation profits. Indeed, the dealer meets customers with intensity λ/μ and appropriates a fraction θ of the marginal trading surplus created by increasing inventories, which is equal to $\max\{\Sigma(q' + i, \delta'), 0\}$ with a customer of type (q', δ') . Notice in particular that this marginal surplus is strictly positive if and only if $q' + i < q^*(\delta)$, that is, if and only if it relaxes a binding inventory-in-advance constraint and helps the customer to trade closer to the target.

The first-order condition for the dealer's optimal inventory holdings is simply

$$\Pi'(i) \leq 0, \text{ with equality if } i > 0. \tag{6}$$

Note that a solution to (6) requires $rP \geq v$; if $v > rP$, then dealers would have incentive to acquire infinite inventory. Hence, in equilibrium, the price will adjust to incorporate the dealers' flow value from holding the asset *and* the marginal benefit of increasing inventory on intermediation profits.

As in our analysis of the customer's optimal asset position, the properties of $\Sigma(\cdot)$ also allow for some natural, partial equilibrium comparative statics with respect to an individual dealer's optimal inventory holdings. In particular, given the behavior of all other agents (and, hence, aggregate variables), one can easily show that an individual dealer's optimal i is increasing in the rate at which he meets customers, λ/μ , and the fraction of the trading surplus he extracts through bargaining, θ .

2.3 The steady-state distribution

We now characterize the steady state distribution $\Phi(q, \delta)$. In a steady state, the gross outflow from any Borel set B of $[0, \bar{q}] \times [\underline{\delta}, \bar{\delta}]$ must be equal to the gross inflow:

$$(\gamma + \lambda)\Phi(B) = \int_{(q, \delta)} \left(\gamma \int_{\delta'} \mathbb{I}_{\{(q, \delta') \in B\}} dF(\delta') + \lambda \mathbb{I}_{\{(\min\{q^*(\delta), q+I\}, \delta) \in B\}} \right) d\Phi(q, \delta).$$

The left-hand side is the gross outflow: customers leave the set B when they change utility, with intensity γ , or when they trade, with intensity λ . The right-hand side is the gross inflow. It states that a customer of type (q, δ) may transition into the set B in two ways. First, with intensity γ and probability $dF(\delta')$, she draws the new utility shock δ' and her new type (q, δ') belongs to B . Second, with intensity λ , she receives a trading opportunity, and her new type $(\min\{q^*(\delta), q+I\}, \delta)$ belongs to B .

Dividing both sides by $(\gamma + \lambda)$ we see that the steady state distribution solves the following fixed point problem

$$\Phi = T^*[\Phi], \text{ where } T^*[\Phi](B) = \int_{(q, \delta)} \mathbb{P}(q, \delta, B) d\Phi(q, \delta) \quad (7)$$

and

$$\mathbb{P}(q, \delta, B) = \frac{\gamma}{\lambda + \gamma} \int \mathbb{I}_{\{(q, \delta') \in B\}} dF(\delta') + \frac{\lambda}{\lambda + \gamma} \mathbb{I}_{\{(\min\{q^*(\delta), q+I\}, \delta) \in B\}}. \quad (8)$$

The function $\mathbb{P}(q, \delta, B)$ is the transition probability function for the state of a customer when she draws a new preference shock or receives a trading opportunity. After checking appropriate regularity conditions, one can apply Theorems 11.12 and 12.3 in [Stokey and Lucas \(1989\)](#) to establish the following results.

Proposition 2. *Assume that $P > 0$ and $I > 0$. Then, there exists a unique steady-state distribution $\Phi(q, \delta | P, I)$. This distribution has the following properties:*

1. *it is decreasing in P in that, for any bounded function $h(q, \delta)$ that is increasing in q , the function $P \mapsto \int h(q, \delta) d\Phi(q, \delta | P, I)$ is decreasing in P ;*
2. *it is weakly continuous in (P, I) ;*
3. *given any initial condition, Φ_0 , the sequence $T^{*n}[\Phi_0] \rightarrow \Phi$ strongly.*

The monotonicity property implies that the law of demand holds in steady state: higher prices are associated with lower aggregate asset holdings by customers. Together with the continuity property, it is crucial to our equilibrium existence proof. The strong convergence result is useful to compute moments since it allows us to calculate any stationary moment by successive iteration.

2.4 Equilibrium

An equilibrium is made up of the following objects: a marginal trade surplus function $\Sigma(q, \delta | P, I)$, an optimal inventory holdings of each dealer I , a joint distribution of asset holdings and preferences $\Phi(q, \delta | P, I)$, and an inter-dealer price P . These objects must satisfy the following conditions:

- (i) $\Sigma(q, \delta | P, I)$ solves the Bellman equation (5) given I and P ;
- (ii) $i = I$ solves the dealer's optimality condition (6) given Φ and P ;
- (iii) $\Phi(q, \delta | P, I)$ is the stationary distribution solving (7);
- (iv) The asset market clears

$$\int_{(q', \delta')} q' d\Phi(q', \delta' | P, I) + \mu I = 0. \tag{9}$$

It is easy to construct an equilibrium with $I = 0$ by choosing a v sufficiently small. In such an equilibrium, there is no trade: when $I = 0$, customers can never buy and, thus, for the market to clear, the equilibrium inter-dealer price must be small enough to ensure that no customer finds it optimal to sell. Hence, in the steady-state the asset is randomly allocated across customers, so that $\Phi(q, \delta) = \Psi(q)F(\delta)$ for some distribution Ψ of asset holdings. Assuming that the support of the distribution of asset holdings has an upper bound \bar{q} , this implies that P should be chosen so that the customer with the strongest incentive to sell chooses to hold on to her asset, i.e., $\Sigma(\bar{q}, \underline{\delta} | P, 0) \geq 0$. Finally, v has to be small enough so that, given the distribution Ψ and the marginal trade surplus function, Σ , the dealer's first-order condition (6) holds with inequality.

Next, we characterize the set of v such that there exists an equilibrium with active intermediation and trade: $I > 0$. Clearly, some v belongs to this set if and only if there exists some $I > 0$ and some $P > 0$ such that the market clearing condition (9) is satisfied and the dealer's first-order condition (6) holds with equality. Because the stationary distribution is decreasing and weakly continuous in P , it can be shown that the market-clearing condition implies a unique market clearing price given I , denoted by $P(I)$. Plugging this price into the dealer's first-order condition at equality, one obtains in turns that $v = V(I)$, where

$$V(I) \equiv rP(I) - \frac{\lambda}{\mu} \theta \int_{(q', \delta')} \max \{ \Sigma(q' + I, \delta' | P(I), I), 0 \} d\Phi(q', \delta' | P(I), I). \quad (10)$$

Hence, the set of v such that there exists an equilibrium with $I > 0$ is simply the range of the function $V(I)$ above, for all values of I that may arise in an equilibrium with active intermediation, that is, for all $I \in (0, s/\mu)$.

The range of I is unbounded above because $V(I) \rightarrow \infty$ as $I \rightarrow s/\mu$. This is true for two reasons. First, when $I \rightarrow s/\mu$ customers hold almost no assets, which implies that their marginal utility and the inter-dealer price $P(I)$ go to infinity. Second, since dealers' inventory are bounded away from zero but customers' asset holdings go to zero, the inventory-in-advance constraint never binds. Hence, $V(I) = rP(I) \rightarrow \infty$. In the Appendix, we show that the marginal trade surplus—and,

hence, $V(I)$ —is bounded below as $I \rightarrow 0$. Together with continuity, this implies that $V(I)$ remains bounded below over $(0, s/\mu)$. The next theorem summarizes.

Theorem 1. *There exists a \underline{v} such that an equilibrium with active intermediation exists if $v > \underline{v}$ and does not exist if $v < \underline{v}$.*

While we do not know whether there are multiple equilibria, it is easy to study this question numerically. In particular, one sees that multiple equilibria arise whenever there is a region of I where the function $V(I)$ is decreasing; since $\lim_{I \rightarrow s/\mu} V(I) = +\infty$, the Intermediate Value Theorem implies that there are several I mapping to the same v , which is just another way to say that the same v can be associated with several I .

2.5 A key property of the inventory constraint: Asymmetry

Before proceeding to our quantitative analysis, we highlight an important qualitative implication of introducing an inventory constraint into an otherwise-standard model of OTC trade. Namely, since the inventory constraint only has a direct impact on purchases, but not on sales, it creates asymmetries that are unique to our model relative to [Lagos and Rocheteau \(2009\)](#).

Consider first a customer who seeks to purchase from a dealer who holds i assets in inventories. While she ideally wants to trade up to her target $q^*(\delta)$, the inventory constraint implies that it is not always feasible. Instead, she will trade so as to be as close as possible to the target given the constraint; that is, $q' = \min\{q^*(\delta), q + i\}$. Clearly, sales are never constrained by inventories in this way: a customer who holds $q > q^*(\delta)$ does not face a constraint on the size of her trade and is able to reach her target in one transaction. Figure 1 illustrates. Therefore, relative to an otherwise identical trading opportunity without a constraint ($i = \infty$), a customer purchase is smaller but a customer sale is not.

A similar asymmetry holds for proportional transaction costs – the proportional markup/markdown charged by dealers over the inter-dealer price. Consider first a customer purchase $q^*(\delta) > q$ and recall equation 4, which states that, with Nash bargaining, the trading fee ϕ is a constant share of

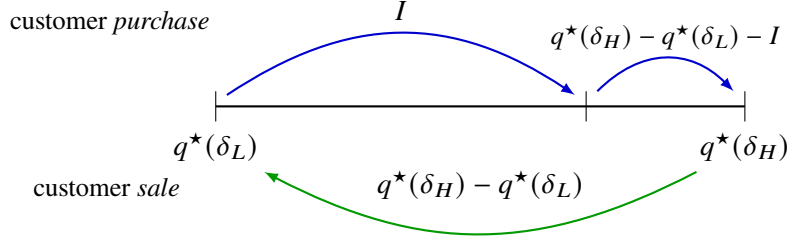


Figure 1. An illustration of the asymmetry between purchases and sales in the special case in which the distribution of preference shocks has a discrete support and can take only two values $\delta_L < \delta_H$, with parameters such that $I < q^*(\delta_H) - q^*(\delta_L) < 2I$.

the surplus. Hence, as a function of *marginal* surplus, the proportional transaction cost is equal to:

$$\frac{\theta}{P(\min\{q^*(\delta), q + i\} - q)} \int_q^{\min\{q^*(\delta), q+i\}} \Sigma(x, \delta) dx.$$

Clearly, since the marginal surplus is decreasing, the transaction cost above decreases in i . In other words, relative to an otherwise identical trading opportunity without a constraint, a customer-buyer pays a larger transaction cost. A customer-seller, on the other hand, is not constrained by inventories and pays the same transaction cost as she would in the absence of an inventory constraint.

Corollary 1. *Relative to otherwise identical trading opportunities without an inventory constraint, customer buyers trade smaller quantities and pay larger proportional transaction cost. Customer sellers trade identical quantities and pay identical proportional transaction costs.*

One may wonder about the empirical implications of this Corollary: does it imply that customer purchases are smaller and more expensive than customer sales? The answer to this question is not obvious since, in principle, the size and the cost of purchases and sales may differ *even without* inventory constraints.

However, in Lemma 4 of Appendix B, we establish that, in an otherwise identical model *without* an inventory constraint, purchases and sales have the same average size. Consequently, in the model *with* a constraint, we find numerically that customer purchases are smaller, on average, than customer sales. Since, in the aggregate, the total quantity of assets purchased and sold are equal, this means that our model predicts that the number of customer purchases will be larger than

that of sales. These asymmetries in trade size and the number of trades are well-known empirical features of several major OTC markets (see, for example, [Green, Hollifield, and Schürhoff, 2007](#)). Hence, introducing inventory constraints is not just a conceptual contribution, but it also brings the model closer to the data along this important dimension.

Asymmetries in transaction costs, on the other hand, can go either way. The reason is that, in the absence of an inventory constraint, proportional transaction costs can be different for purchases than sales. For example, under their preferred parameterization, the model of [Duffie, Gârleanu, and Pedersen \(2005\)](#) implies that transaction costs are zero for purchases and strictly positive for sales. In [Appendix B](#), we show that, in the model of [Lagos and Rocheteau \(2009\)](#) with an isoelastic utility function, $u(q, \delta) = \frac{q^{1-1/\eta}}{1-1/\eta} \delta$, value-weighted transaction costs for purchases are strictly smaller than those for sales if and only if $\eta > 2$. In these cases, we find numerically that the inventory in advance constraint generally raises the transaction cost for purchases, but not always by a sufficient amount to make them larger than the transaction costs for sales.

3 Quantitative Exercise

In this section, we use our model to quantitatively evaluate the market impact of post-GFC regulations that increased dealers' balance sheet costs. Relative to a purely empirical approach that focuses on traditional liquidity measures, such as transaction costs, our structural approach enables to evaluate welfare. This is crucial: for example, we can imagine a scenario where post-GFC regulations increased transaction costs without impacting asset allocation, in which case it would have only had redistributive effects. Our analysis proceeds in two steps. First, we calibrate our model to match moments from the corporate bond market before the GFC. Then, we infer the change in ν that is consistent with the observed decline in dealers' inventory holdings post-GFC. Using this implicit change in dealers' inventory costs, we evaluate the effects of post-GFC regulations on liquidity, prices, allocations, and welfare. Our calibration generates a rise in trading costs consistent with observed data. Furthermore, it reveals that regulations have had a substantial negative impact

on asset allocation: pre-GFC, the welfare loss created by the OTC market, relative to the first best, was about 1.25%, post-GFC this loss has increased to 2.4%.

3.1 Data

We use the academic version of the Trade Reporting and Compliance Engine (TRACE) database of US corporate bond transactions, made available by the Finance Industry Regulatory Authority (FINRA). The raw TRACE data provides detailed information on all secondary market transactions self-reported by FINRA member dealers. These include bond's CUSIP, trade execution time and date, transaction price (\$100 = par), the volume traded (in multiple of par), a buy/sell indicator, and flags for dealer-to-customer and inter-dealer trades. Unlike the public version, the academic TRACE does not censor trade volume at \$5 million (for investment grade bonds) or \$1 million (for high-yield bonds). The academic version also contains masked dealer identities as well as transactions in privately traded Rule 144A bonds that are not disseminated to the public.

Dealers are required to correct errors in previously reported trades with flags corresponding to trade cancellations, modifications, or reversals. We use the standard cleansing algorithm described in [Dick-Nielsen \(2009, 2014\)](#) and [Dick-Nielsen and Poulsen \(2019\)](#) to remove these self-reported errors. Our TRACE sample starts in July 2002 and covers transactions until June 2020. We exclude the COVID-19 crisis period in March and April 2020 from our sample, since our maintained assumption of a non-stochastic steady state is not appropriate for such a turbulent period (see [Kargar, Lester, Lindsay, Liu, Weill, and Zúñiga, 2021](#) for an empirical study).

We collect issue credit ratings and bond characteristics from Mergent Fixed Income Securities Database (FISD). We drop all bonds not contained in the FISD and only consider CUSIPs in TRACE identified by FISD as fixed-coupon US corporate debentures and US corporate bank notes with non-missing maturity dates and amounts outstanding. We also exclude bonds with equity-like and special features.⁵ Furthermore, we exclude trades associated with new issuances and also

⁵Following earlier work, we exclude all bonds that are convertible, puttable, exchangeable, preferred, asset-backed, secured lease obligations, unit deals, or Yankee bonds. Additionally, we do not consider bonds with variable coupons or sinking funds, or those issued in a foreign currency or as part of unit deals.

remove transactions that happen within 90 days of the traded bond issuance. This ensures that trading activity in the sample closely aligns with the non-stochastic steady state envisioned by the model.

Our model features a representative asset and a representative customer. As is well known, in reality, there is substantial heterogeneity in both these dimensions in the corporate bond market. Hence, we apply the following two filters to control for heterogeneity while keeping the sample economically relevant. First, we concentrate on trades for investment-grade (IG) bonds, which exhibit more homogeneous liquidity properties and represent the vast majority of daily trading volume in TRACE.⁶ Second, we argue that it would be unreasonable to require that our representative customer model rationalizes the vast difference in trade size between institutional and retail investors. For this reason, we focus on trades larger than a threshold of \$1 million, which are more likely to originate from institutional rather than retail investors.

Finally, in practice, dealers provide liquidity via “agency” and “risky-principal” trades. In an agency trade, the dealer acts as a match maker between buyers and sellers, and never actually owns the asset being traded. In contrast, in a risky-principal trade, dealers buy and sell on their own account, absorbing sell orders onto their balance sheet and fulfilling buy orders by reducing their inventory holdings. Since trading costs for these two types of trades have been shown to be quite different (see e.g., [Choi, Huh, and Shin, 2023](#))—and, in our theory, inventories are a key input into the provision of liquidity services—it is natural for us to focus exclusively on risky-principal transactions. We identify these transactions in the data using the procedure described in [Kargar, Lester, Lindsay, Liu, Weill, and Zúñiga \(2021\)](#).

Table 1 reports summary statistics for the daily number and volume of inter-dealer, customer-bought and customer-sold trades for our final, filtered sample. Note that, while the total volume of customer buys and sells are very similar (approximately \$3.9 billion), we observe, on average, more customer buys than customer sells, which is a key qualitative prediction of our model with a (binding) inventory-in-advance constraint.

⁶In the third quarter of 2023, IG bonds represent approximately 85% of the average total daily trading volume of publicly traded US corporate bonds. Source: U.S. corporate bond statistics from SIFMA.

[Table 1 about here.]

3.2 Calibration to Pre-GFC corporate bond market: Targets

We set the discount factor, r , equal to 5%. We assume that customers have an iso-elastic utility function of the form

$$u(q, \delta) = \frac{q^{1-1/\eta}}{1-1/\eta} \delta.$$

In addition, we assume that the preference shock, δ , is an iid draw from a discretized log-normal distribution, $F(\delta)$. Given our choice of an iso-elastic utility function, the model is homogeneous in s , the per-capita supply of the asset.⁷ Hence, we normalize s such that the asset supply held by customers is one, and the aggregate dealer inventories represent 2% of the total asset supply, similar to the dealer sector's pre-GFC corporate bond holding share. Finally, we normalize the mean of $F(\delta)$ such that the price of the asset equals $1/r$ in a Walrasian equilibrium in which the supply held by customers is one, i.e., the same aggregate quantity of asset they hold in our calibration.

The intensity of contact between customers and dealers is typically difficult to calibrate, because the TRACE data does not offer direct evidence about customers' search process. We rely on the work of [Kargar, Lester, Plante, and Weill \(2023\)](#), who leverage proprietary data from an electronic trading platform to measure customers' time to trade in the U.S. corporate bond market. Following their estimate, we set λ so that a customer contacts a dealer every 3 days.

Given the assumptions above, there are five parameters that we need to calibrate: the variance of the preference shocks, σ_δ^2 ; the arrival rate of preference shocks, γ ; the elasticity parameter for the customers' utility function, η ; the dealers' bargaining power, θ ; the dealers' utility parameter, ν , and the measure of active dealers, μ . In what follows, we describe the five target moments that

⁷Specifically, suppose that we scale the supply by the constant κ , i.e., $\tilde{s} = \kappa s$. Then, scaling preference shocks by the same factor, $\tilde{\delta} = \kappa^{1/\sigma} \delta$, renders the marginal utilities the same for all investors if they scale their holdings by κ . As a result, the equilibrium price remains unchanged, as do the holdings of customers and dealers relative to the aggregate supply.

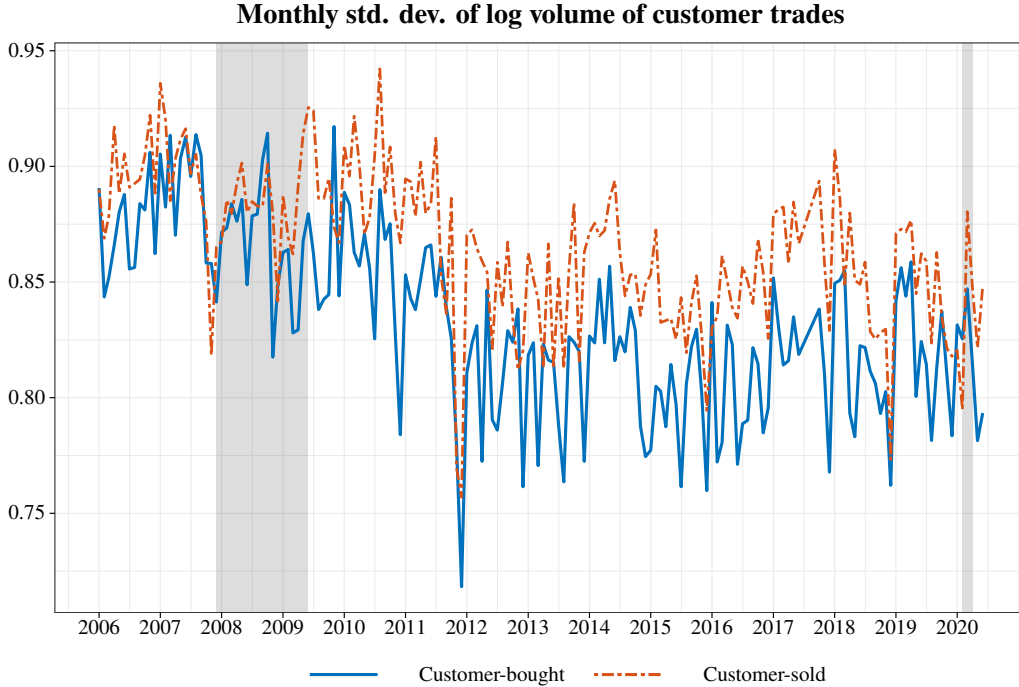


Figure 2. Monthly standard deviation of log trading volume for customer-to-dealer trades in percentage points. We restrict the sample to subset of trades involving risky-principal trades of investment-grade bonds with a volume exceeding \$1 million. The vertical shaded bars indicate NBER recessions. Sources: Academic TRACE and FISD.

we use to determine these parameters. Though in general the target moments and parameters are determined simultaneously, we try to connect each moment to the parameter it affects most directly.

The variance of preference shocks. The dispersion in preference shocks determines customers' equilibrium asset holdings and the size of trades they execute when their asset holdings differ from their target portfolios. Hence, to help identify σ_θ^2 , we target the standard deviation of log trade size. In our sample of IG bond trades with par value exceeding \$1 million, we find that the monthly standard deviation of log size is about 0.9. Figure 2 plots the monthly standard deviation of log trade size for customer-bought and customer-sold trades.

The frequency of preference shock. This parameter is a key determinant of how frequently customers want to buy or sell. Therefore, to help determine γ , we target the turnover of assets that customers purchase, generated by trades above of size greater than \$1 million. This is calculated by dividing the total quarterly trading volume by the quarterly average of the amount outstanding of

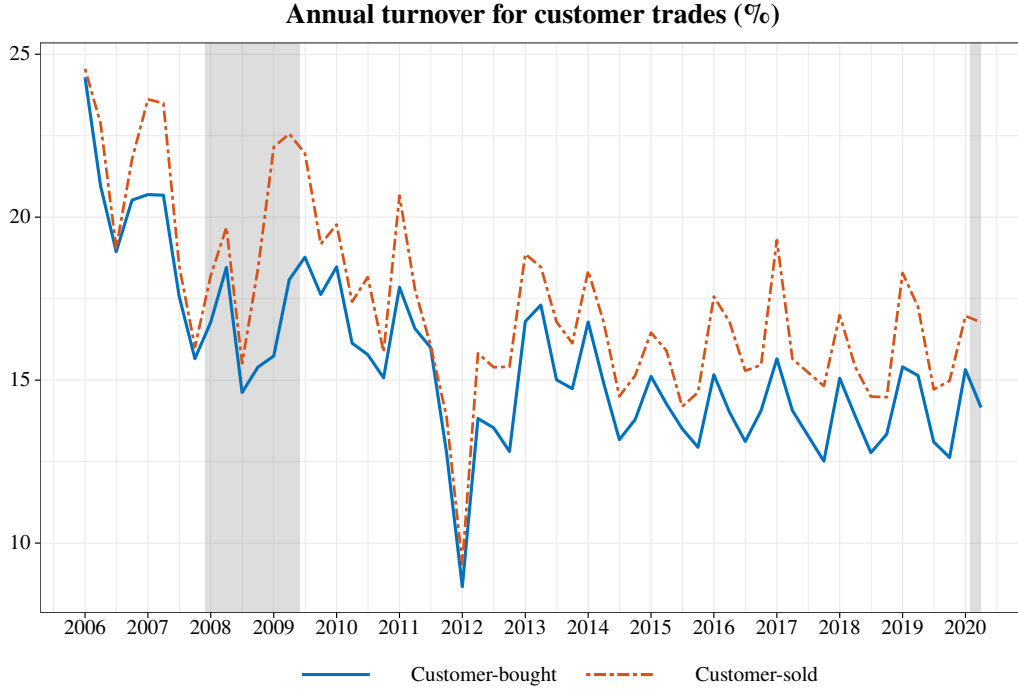


Figure 3. Annual turnover for customer-to-dealer trades in percentage points. We restrict the sample to subset of trades involving risky-principal trades of investment-grade bonds with a volume exceeding \$1 million. The vertical shaded bars indicate NBER recessions. Sources: Academic TRACE and FISD.

bonds for dealer-to-customer transactions in our sample. We find that the turnover is approximately 20% annually . Figure 3 plots the annual turnover for customer-bought and customer-sold trades.

Dealers’ bargaining power and the elasticity of the customer’s utility function. To help determine these parameters we focus on proportional transaction costs. On the one hand, dealers’ bargaining power determines the overall level of transaction costs paid by customers. On the other hand, the elasticity of the customer’s utility function is one determinant of the asymmetry between transaction costs for customer purchases and customers’ sales. This is demonstrated formally in Lemma 5 of the appendix: it shows that, without an inventory-in-advance constraint, a lower value of η tends to make the transaction costs for customers purchases larger than for sales.

Now turning to measurement, our empirical target is based on the value-weighted two-way trading cost proposed by Choi, Huh, and Shin (2023):

$$2Q \times \frac{\text{traded price} - \text{reference price}}{\text{reference price}}, \quad (11)$$

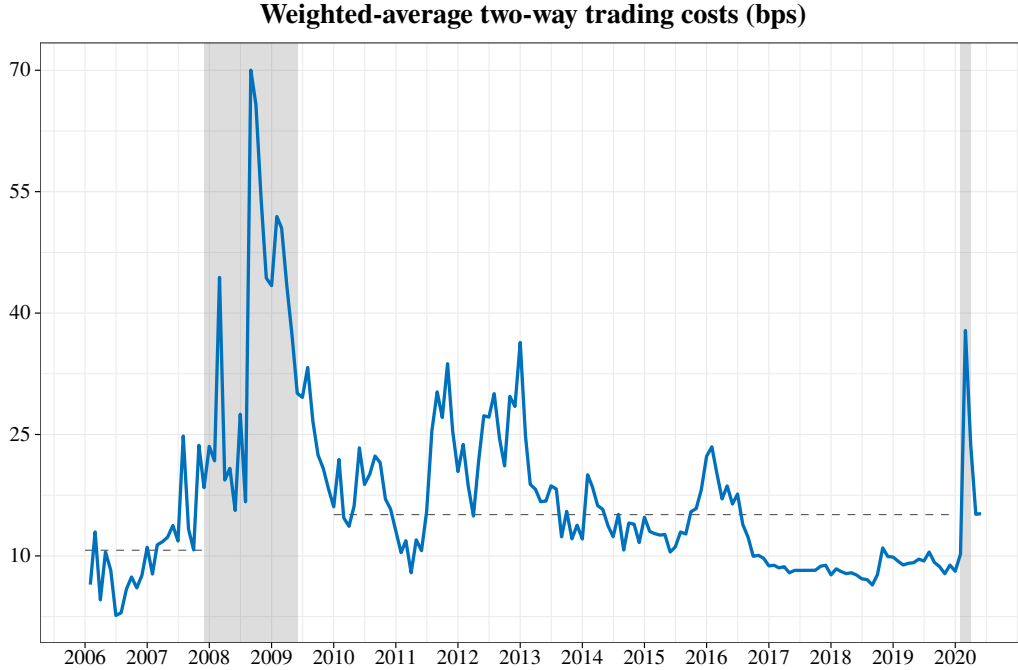


Figure 4. Monthly size-weighted average of customer-bought and customer-sold two-way (round-trip) trading costs proposed by [Choi, Huh, and Shin \(2023\)](#) from equation (11). We restrict the sample to subset of trades involving risky-principal trades of investment-grade bonds with a volume exceeding \$1 million. The horizontal dashed lines represent subsample averages for the pre-GFC (2006–2007) and post-GFC (2010–2019) periods. The vertical shaded bars indicate NBER recessions. Sources: Academic TRACE and FISD.

where Q is equal to +1 for a customer *buy* and -1 for a customer *sell*. For each customer trade, a “reference price” is calculated as the volume-weighted average price of inter-dealer trades larger than \$100,000 in the same bond-day, excluding inter-dealer trades executed within 15 minutes. The measure is calculated at the trade level for all customer trades classified as risky principal, and then calculated at the bond-day level by taking the volume-weighted average of trade level spreads.

We find that, for risky-principal customer trades involving IG bonds with a volume exceeding \$1 million, two-way trading costs from [Choi, Huh, and Shin \(2023\)](#) for buy and sell transactions in the pre-GFC periods are 10.8 bps and 9.6 bps, respectively. Figure 4 plots the volume-weighted average round-trip (two-way) customer-bought and customer-sold transactions costs from [Choi, Huh, and Shin \(2023\)](#).

Measure and utility flow of dealers. The measure of dealers μ and their utility flow v jointly determine how much inventory they hold on aggregate, $\mu \times I$. Moreover, dealers’ individual choice

of inventory, I , determines how many purchases customers need to make in order to reach their target holding. Hence, we propose two targets that help determine μ and ν : the total asset holdings of the dealer sector before the GFC, as a share of outstanding assets; and the ratio of the number of customer-sell transactions to the number of customer-buy transactions.

To set a target for the asset holdings of the dealer sector prior to GFC, we rely on data from the Federal Reserve's Flow of Funds. Figure 6 plots the share of corporate and foreign bonds held by security brokers and dealers from the Federal Reserve's Flow of Funds. Up until 2002, this share was slightly above 2%, then it dramatically increased during the years leading up to the GFC, only to drop sharply to levels below 1% after the GFC. The substantial increase leading up to the GFC may be partly attributed to non-agency mortgage-backed securities (MBS), which are included in the Flow of Funds accounting but not relevant to our quantitative exercise. For this reason, for the pre-GFC calibration, we take the dealer sector asset holding to be 2% of the aggregate asset supply.

In Figure 5, we plot the ratio of the number of customer-sold trades to the number of customer-bought trades in our sample, after adjusting for order flow imbalance (see Appendix C). We calculate that, prior to the GFC, this ratio averaged about 0.76.

3.3 Calibration to pre-GFC corporate bond market: Outcomes

Table 2 presents the target moments and Table 3 reports our calibrated parameters. All targets are matched nearly exactly.

[Table 2 about here.]

Since the Walrasian price of the asset would be $1/r$ and dealers are risk neutral, it is natural to interpret the flow utility that dealers receive from holding a unit of this bond as $\nu = 1 - \tau$, where τ denotes the (flow) inventory cost to dealers of holding the asset on their balance sheet. Hence, if we think of the asset as a consol bond, our calibration implies that inventory costs before the GFC were approximately 3.99% of the bond's coupon.

[Table 3 about here.]

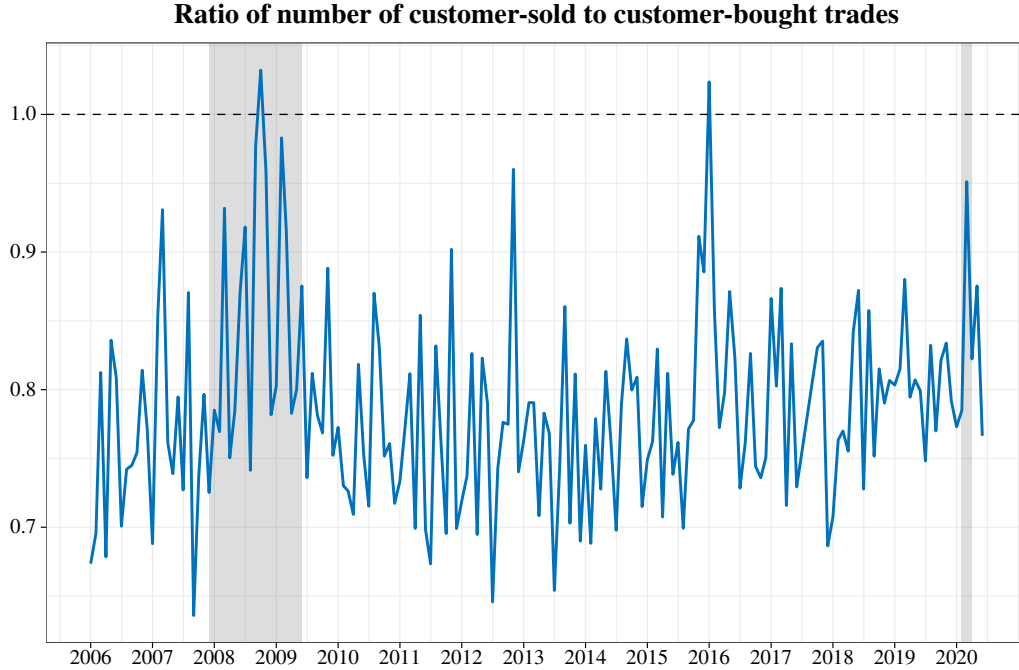


Figure 5. Ratio of the monthly number of customer-sold to customer-bought trades. We adjust this ratio for order imbalance, as described in Appendix C. We restrict the sample to subset of trades involving risky-principal trades of investment-grade bonds with a volume exceeding \$1 million. The horizontal dashed line at 1 represents the case in which the number of customer-bought trades equals the number of customer-sold trades, as is the case in a model without the inventory constraint. The vertical shaded bars indicate NBER recessions. Sources: Academic TRACE and FISD.

Figure 7 plots the (log of) customers’ target asset holdings, $q^*(\delta)$, along with each dealers’ equilibrium inventory holdings, I . To highlight the effects of the inventory constraint in our framework, we also plot the target asset holdings in an environment without the inventory constraint (i.e., the environment of Lagos and Rocheteau, 2009, with the same parameters, assuming that the asset supply is equal to one). The dashed horizontal line is the inventory holdings of dealers: hence, the inventory constraint only binds for those customers who receive the large enough preference shock. Also note that, at the scale of the figure, the target asset holdings are indistinguishable from the targets in the equilibrium without inventory constraint (they are in fact slightly larger).

Introducing an inventory constraint creates a small increase in the bid-ask spread charged by dealers: given the parameter values that emerge from our calibration, the trading costs in the no-constraint environment would be approximately 0.5 bps smaller. As trading costs rise, the customers’ valuation for the asset declines, which puts downward pressure on the inter-dealer price.

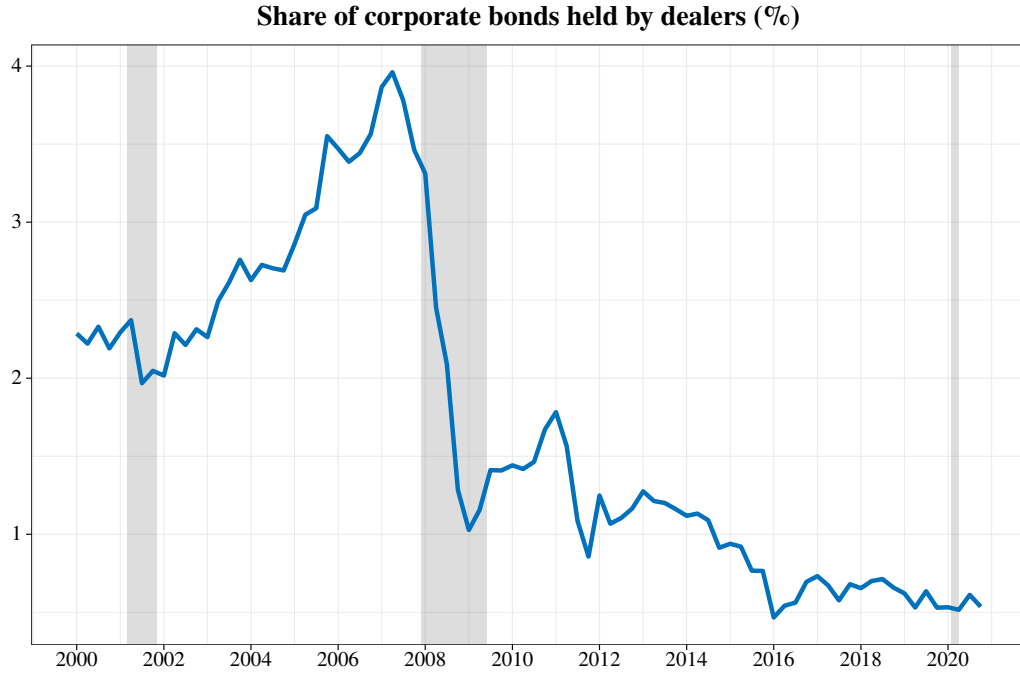


Figure 6. Share of corporate and foreign bond holdings for security broker-dealers. The vertical shaded bars indicate NBER recessions. Source: Table L.213 of the Federal Reserve’s Z.1: Financial Accounts of the United States (the Flow of Funds).

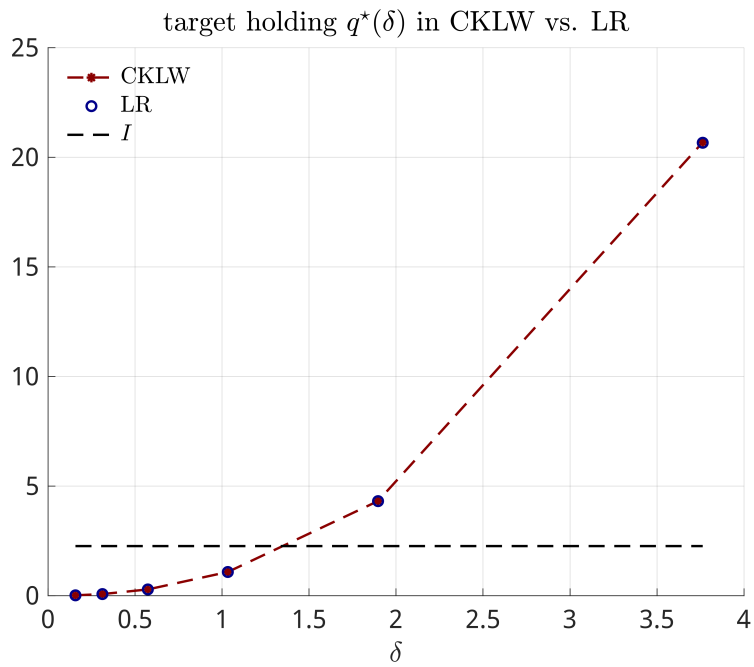


Figure 7. Target asset holdings, $q^*(\delta)$, with inventory constraints (CKLW) and without (LR). The horizontal line is I in CKLW.

However, the presence of an inventory constraint also puts upward pressure on the price, because precautionary incentives increase customers' demand for the asset. In equilibrium, we find that the forces putting downward pressure on the price dominate, as the inter-dealer price in our benchmark model is slightly lower than in the model without an inventory constraint.

Overall, however, we find that the presence of an inventory constraint in the pre-GFC economy had quite mild effects on equilibrium prices and target holdings, relative to an environment where dealers are not required to hold inventory in order to intermediate trade. To get a sense of the welfare cost of the inventory constraint, we calculate the gains from trade that are realized in equilibrium relative to the gains from trade in a frictionless environment. More precisely, we consider the following measure of welfare *loss*:

$$L = \frac{W_{fb} - W_{eqm}}{W_{fb} - W_{aut}}, \quad (12)$$

where W_{fb} , W_{eqm} , and W_{aut} , denote total welfare in the (first best) frictionless environment, in equilibrium, and in autarky, respectively. This measures the fraction of gains from trade lost by the market in equilibrium relative to the first best. As a point of reference, we find that the environment without an inventory constraint—but with search and bargaining frictions—creates a welfare loss of 0.87%. In our environment, where dealers must hold inventory in order to sell, the welfare loss is 1.25%.

3.4 The effects of rising inventory costs

We now study the impact of increasing dealers' cost of holding assets (by decreasing ν) so as to create a threefold decline in their inventories, from 2% to $2\%/3 \simeq 0.66\%$ of aggregate bond supply. This decline is in line with Figure 9, where we observe that dealers' inventories were around 2% of aggregate bond supply in early 2000, and dropped to around 0.6% by 2020.

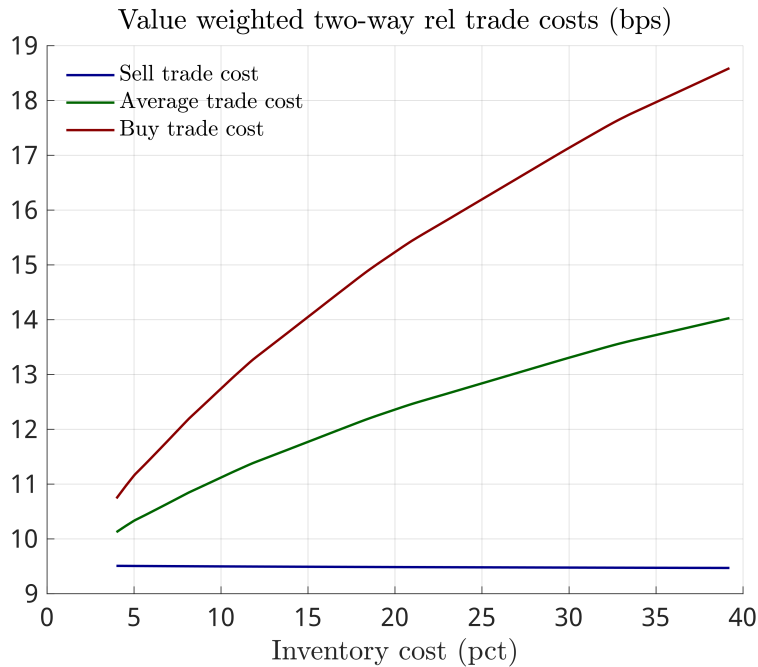


Figure 8. Value-weighted trading costs

Trading costs. Figure 8 shows that the resulting increase in the value-weighted trading cost was about 4 bps, nearly 80% of the increase we observed in the data. The Figure also reveals a striking asymmetry: trading costs increase for customer purchases, but not for customer sales. This quantitative finding is consistent with the intuition given in Section 2.5. In the data, we observe a similar asymmetry, although it is much less pronounced than in the model: two-way transaction costs increased by 5.4 bps for customer purchases and by 4.8 bps for customer sales.

Three structural measurement exercises. One advantage of a structural model is to facilitate measurements that would be difficult to make based on a purely reduced-form econometric approach.

First, our model provides a measure of the implicit cost of regulation on dealers. Figure 9 plots dealers’ aggregate inventories as a fraction of total supply, $\mu I/s$ against the implicit inventory cost τ . It shows that, in order to engineer a threefold decline in dealers’ inventory holdings—from 2% to $2\%/3 \approx 0.66\%$ —inventory costs need to increase almost tenfold, from approximately 4% in the

pre-crisis calibration to approximately 40%, as a fraction of the asset coupon, holding all other parameters fixed.

Second, in Figure 10, we illustrate the effect of post-GFC regulation on welfare. Namely, we plot the welfare loss, defined in (12), as the inventory cost rises. As shown by the plain blue curve in the figure, an increase in inventory cost τ from approximately 4% to 17% increases the welfare loss from 1.25% to about 2.4%, an increase of nearly 100%. One caveat of this calculation, of course, is that it abstracts from potential benefits that derive from greater financial stability.

One may argue that our measure of welfare is narrow because we wrote a partial equilibrium model: it focuses on the welfare of those investors who participate (directly and indirectly) in OTC markets. One could also argue that regulation has a welfare impact through other channels, too, such as the effects on firms' cost of capital. One way to measure the impact on the cost of capital is to calculate the liquidity yield spread of the asset implied by our model. Recall that, given the normalization of preference shocks explained above, the frictionless price is equal to $1/r$, the present value of a riskless consol bond with a coupon equal to 1. Hence, it is natural to define the liquidity yield spread based on the following pricing condition: the present value of this consol bond, at rate equal to r plus the liquidity yield spread, should be equal to P , the price of the consol bond in our theoretical OTC market. This gives:

$$\text{liquidity yield spread} = \frac{1}{P} - r.$$

Figure 11 shows that as inventory costs rise, the yield spread also increases. Prior to the GFC, the liquidity yield spread is about 2 bps and, after the GFC, rises to more than 5 bps. Hence, according to our calibration, the liquidity component of firms' cost of capital has increased dramatically after the crisis.

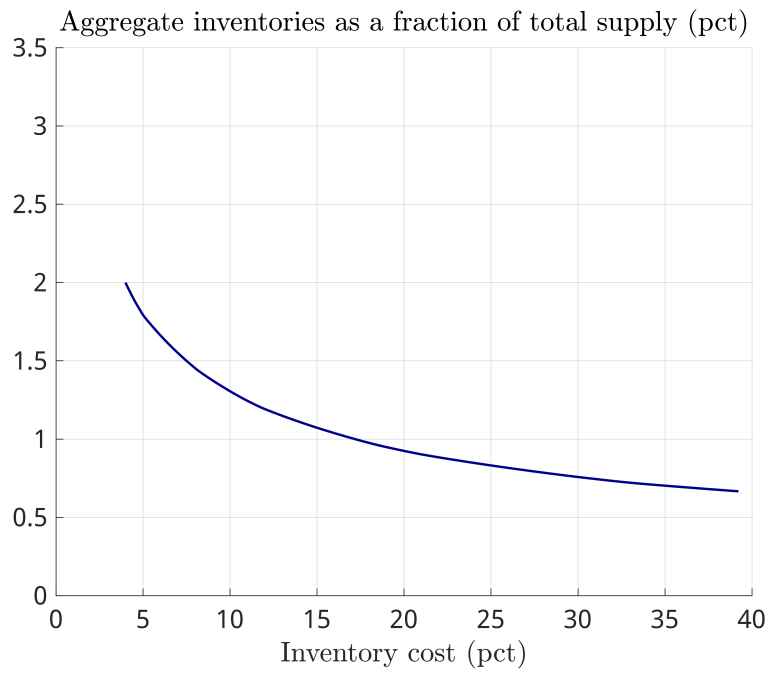


Figure 9. Dealers' aggregate inventory as a percentage of total supply

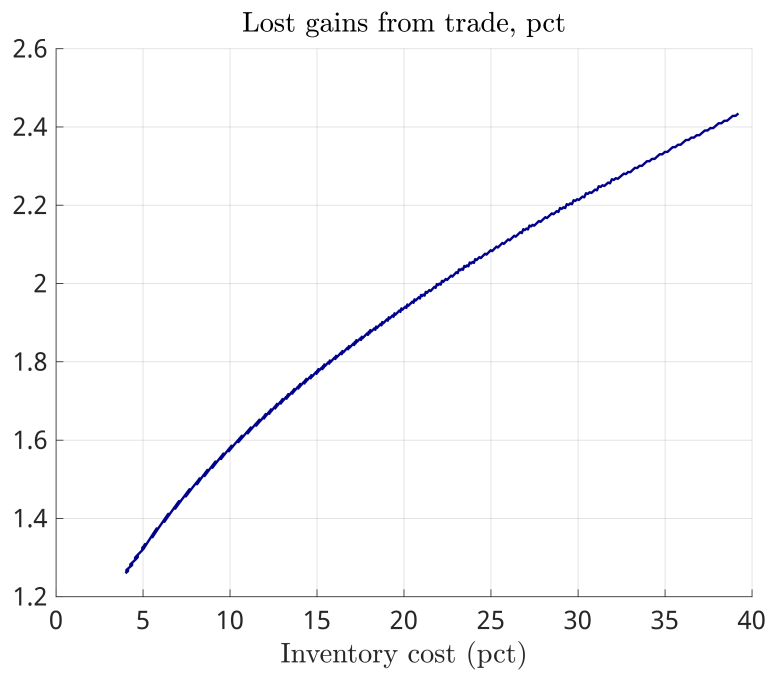


Figure 10. Equilibrium welfare loss

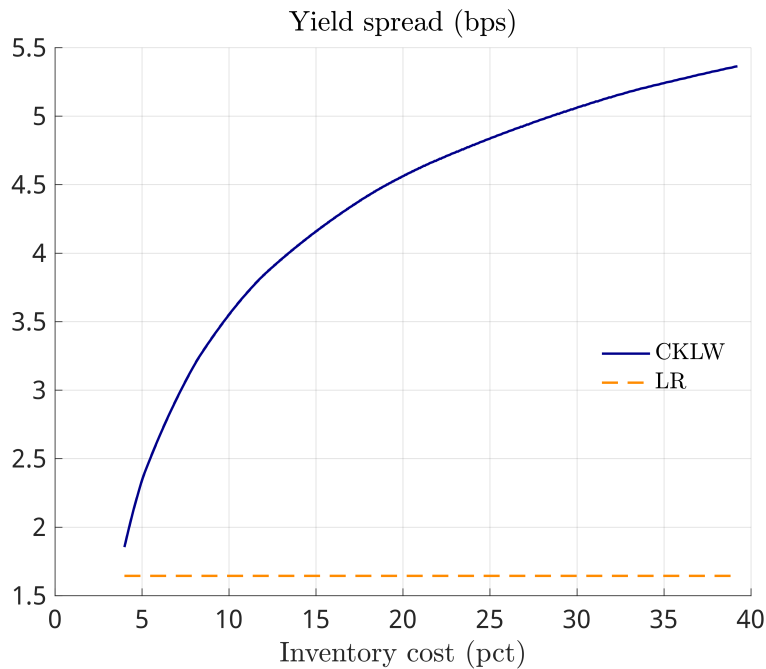


Figure 11. Yield spread with inventory in advance constraint (CKLW) and without (LR)

4 Conclusion

We extend the standard search-theoretic model of dealer-intermediated OTC markets, in which dealers never hold inventory, by introducing a simple and natural “inventory-in-advance” constraint, which makes inventory a necessary input to intermediation. We characterize the equilibrium, and study how dealers’ optimal inventory choice depends on inventory costs. We calibrate the model to transaction-level data from the corporate bond market and analyzed the welfare impact of rising inventory costs associated with post-crisis regulations. We measure the welfare loss as the fraction of total gains from trade that the OTC market fails to generate. We find that the rise in inventory costs substantially increases the welfare loss, from about 1.25% to approximately 2.4% of the total gains from trade.

References

- Amihud, Yakov, and Haim Mendelson, 1980, Dealership market: Marketmaking with inventory, *Journal of Financial Economics* 8, 31–53.
- An, Yu, 2018, Competing with inventory in dealership markets, Working paper, Johns Hopkins University.
- Bao, Jack, Maureen O’Hara, and Alex Zhou, 2018, The volcker rule and market making in times of stress, *Journal of Financial Economics* 130, 95–113.
- Bessembinder, Hendrik, Stacey Jacobsen, William Maxwell, and Kumar Venkataraman, 2018, Capital commitment and illiquidity in corporate bonds, *Journal of Finance* 73, 1615–1661.
- Bethune, Zachary, Bruno Sultanum, and Nicholas Trachter, 2022, An information-based theory of financial intermediation, *Review of Economic Studies* 89, 2381–2444.
- Bogachev, Vladimir I., 2007, *Measure Theory* (Springer-Verlag).
- Brancaccio, Giulia, and Karam Kang, 2022, Search frictions and product design in the municipal bond market.
- Brancaccio, Giulia, Dan Li, and Norman Schurhoff, 2017, Learning by trading: The case of the U.S. market for municipal bonds, Working paper, Cornell University.
- Carter, Michael, and Bruce Van Brunt, 2000, *The Lebesgue-Stieltjes Integral: a Practical Introduction* (Springer-Verlag, New York).
- Choi, Jaewon, Yesol Huh, and Sean Seunghun Shin, 2023, Customer liquidity provision: Implications for corporate bond transaction costs, *Management Science*, forthcoming .
- Choi, Michael, and Guillaume Rocheteau, 2021, New monetarism in continuous time: Methods and applications, *Economic Journal* 658–696.
- Diao, Chengjie, Evan Dudley, and Hongfei Amy Sun, 2023, Search and inventory in the over-the-counter market, Working paper, Queen’s University.

- Dick-Nielsen, Jens, 2009, Liquidity biases in TRACE, *Journal of Fixed Income* 19, 43–55.
- Dick-Nielsen, Jens, 2014, How to clean enhanced TRACE data, Working paper, CBS.
- Dick-Nielsen, Jens, and Thomas Kjær Poulsen, 2019, How to clean academic TRACE data, Working paper, CBS and BI Norwegian.
- Dick-Nielsen, Jens, and Marco Rossi, 2019, The cost of immediacy for corporate bonds, *Review of Financial Studies* 32, 1–41.
- Duffie, Darrell, 2017, Post-crisis bank regulations and financial market liquidity, Banca d'Italia, Thirteenth Paolo Baffi Lecture on Money and Finance.
- Duffie, Darrell, Nicolae Gârleanu, and Lasse H. Pedersen, 2005, Over-the-counter markets, *Econometrica* 73, 1815–1847.
- Dyskant, Lucas B., André C. Silva, and Bruno Sultanum, 2023, Dealer costs and customer choice, Working paper, Richmond Fed.
- Eisfeldt, Andrea L., Bernard Herskovic, Sriram Rajan, and Emil Siriwardane, 2023, OTC intermediaries, *Review of Financial Studies* 36, 615–677.
- Farboodi, Maryam, Gregor Jarosch, and Guido Menzio, 2017, Intermediation as rent extraction, Working Paper 24171, National Bureau of Economic Research.
- Farboodi, Maryam, Gregor Jarosch, and Robert Shimer, 2022, The emergence of market structure, *The Review of Economic Studies* 90, 261–292.
- Feldhütter, Peter, 2012, The same bond at different prices: Identifying search frictions and selling pressures, *Review of Financial Studies* 25, 1155–1206.
- Folland, Gerald B., 1999, *Real Analysis: Modern Techniques and Their Applications* (John Wiley & Sons Inc.).
- Gârleanu, Nicolae, 2009, Portfolio choice and pricing in illiquid markets, *Journal of Economic Theory* 144, 532–564.

- Gavazza, Alessandro, 2016, An empirical equilibrium model of a decentralized asset market, *Econometrica* 84, 1755–1798.
- Gofman, Michael, 2014, A network-based analysis of over-the-counter markets, Working paper, University of Rochester.
- Gofman, Michael, 2017, Efficiency and stability of a financial architecture with too-interconnected-to-fail institutions, *Journal of Financial Economics* 124, 113–14.
- Green, Richard, Burton Hollifield, and Norman Schürhoff, 2007, Dealer intermediation and price behavior in the aftermarket for new bond issues, *Journal of Financial Economics* 86, 643–682.
- Hendershott, Terrence, Dan Li, Dmitry Livdan, and Norman Schürhoff, 2020, Relationship trading in over-the-counter markets, *Journal of Finance* 75, 683–734.
- Ho, Thomas, and Hans R. Stoll, 1981, Optimal dealer pricing under trading transactions and return uncertainty, *Journal of Financial Economics* 9, 47–73.
- Ho, Thomas, and Hans R. Stoll, 1983, The dynamics of dealer markets under competition, *Journal of Finance* 38, 1053–1074.
- Hugonnier, Julien, Benjamin Lester, and Pierre-Olivier Weill, 2020, Frictional intermediation in over-the-counter markets, *Review of Economic Studies* 87, 1432–1469.
- Hugonnier, Julien, Benjamin Lester, and Pierre-Olivier Weill, 2022, Heterogeneity in decentralized asset markets, *Theoretical Economics* 17, 1313–1356.
- Kargar, Mahyar, Benjamin Lester, David Lindsay, Shuo Liu, Pierre-Olivier Weill, and Diego Zúñiga, 2021, Corporate bond liquidity during the covid-19 crisis, *Review of Financial Studies* 34, 5352–5401.
- Kargar, Mahyar, Benjamin Lester, Sébastien Plante, and Pierre-Olivier Weill, 2023, Sequential search for corporate bonds, Working Paper 31904, National Bureau of Economic Research.
- Kargar, Mahyar, Juan Passadore, and Dejanir Silva, 2020, Liquidity and risk in OTC markets: A theory of asset pricing and portfolio flows, Working paper, UIUC.

- Lagos, Ricardo, and Guillaume Rocheteau, 2009, Liquidity in asset markets with search frictions, *Econometrica* 77, 403–426.
- Lagos, Ricardo, Guillaume Rocheteau, and Pierre-Olivier Weill, 2008, Crashes and recoveries in illiquid markets, Working paper, NYU, UCI, and UCLA.
- Lagos, Ricardo, and Shengxing Zhang, 2020, Turnover liquidity and the transmission of monetary policy, *American Economic Review* 110, 1635–1672.
- Lester, Benjamin, Guillaume Rocheteau, and Pierre-Olivier Weill, 2015, Competing for order flow in OTC markets, *Journal of Money, Credit and Banking* 47, 77–126.
- Li, Jessica S., 2023, Strategic bargaining and portfolio choice in intermediated markets, Working paper, University of Chicago.
- Liu, Shuo, 2020, Dealer’s search intensity in U.S. corporate bond markets, Working paper, Tsinghua University.
- Milbradt, Konstantin, 2017, Asset heterogeneity in over-the-counter markets, Working paper, Kellogg School of Management.
- Mildenstein, Eckart, and Harold Schlee, 1983, The optimal pricing policy of a monopolistic marketmaker in equity market, *Journal of Finance* 38, 218–231.
- Nosal, Ed, Yuet-Yee Wong, and Randall Wright, 2019, Intermediation in markets for goods and markets for assets, *Journal of Economic Theory* 183, 876–906.
- Pagnotta, Emiliano S., and Thomas Philippon, 2018, Competing on speed, *Econometrica* 86, 1067–1115.
- Palleja, Mariano J., 2022, Filters for the academic trace database, Working paper, UCLA.
- Pinter, Gabor, and Semih Üslü, 2021, Comparing search and intermediation frictions across markets, Working paper, Bank of England and Johns Hopkins University.
- Rocheteau, Guillaume, Pierre-Olivier Weill, and Tsz-Nga Wong, 2018, A tractable model of monetary exchange with ex post heterogeneity, *Theoretical Economics* 13, 1369–1423.

- Shen, Ji, Bin Wei, and Hongjun Yan, 2021, Financial intermediation chains in an over-the-counter market, *Management Science* 67, 4623–4642.
- Stokey, Nancy L., and Robert E. Lucas, 1989, *Recursive Methods in Economic Dynamics* (Harvard University Press, Cambridge).
- Thakor, Anjan V., 2012, The economic consequences of the Volcker rule, Report by the US Chamber of Commerce Center for Capital Market Competitiveness.
- Trebbi, Francesco, and Kairong Xiao, 2019, Regulation and market liquidity, *Management Science* 65, 1949–1968.
- Tse, Chung-Yi, and Yujing Xu, 2021, Inter-dealer trades in OTC markets—who buys and who sells?, *Review of Economic Dynamics* 39, 220–257.
- Üslü, Semih, 2019, Pricing and liquidity in decentralized asset markets, *Econometrica* 87, 2079–2140.
- Weill, Pierre-Olivier, 2007, Leaning against the wind, *Review of Economic Studies* 74, 1329–1354.
- Weill, Pierre-Olivier, 2020, The search theory of over-the-counter markets, *Annual Review of Economics* 12.
- Yang, Ming, and Yao Zeng, 2019, Coordination and fragility in liquidity provision, Working paper, UCL and Wharton.

Tables

Table 1. Summary statistics. This table provides mean, standard deviation, median, 25th, 75th, and 95th percentiles of the average daily number of trades and volume by counterparty type, all years. The “daily num.” variables refer to the daily number of trades and the “daily vol.” variables refer to the average total daily volume, in millions USD. “customer” trades refer to trades between a dealer and a customer which represent the sum of “customer-bought” and “customer-sold” trades. The sample is from the academic version of TRACE and runs from July 2002 to the end of June 2020. Our sample only includes trades for investment-grade bonds with size exceeding \$1 million and excludes the COVID-19 crisis period in March and April 2020. All agency transactions, where dealers act as match makers have been removed. Rule 144A bonds for which trades not disseminated to the public are excluded. We filter the sample as described in the main text.

	Mean	Std. dev.	Q25	Q50	Q75	Q95
Daily num. inter-dealer	844.40	469.91	584	817.50	1,048	1,663.30
Daily num. customer	988.94	522.80	681	965	1,252.75	1,866.65
Daily num. customer-bought	511.40	272.60	353	497	646	985.55
Daily num. customer-sold	477.54	257.75	322.25	462	614	907.20
Daily vol. interdealer (\$m)	2,603.91	1,568.39	1,699.70	2,455.07	3,224.50	5,442.98
Daily vol. customer (\$m)	4,516.78	2,389.89	3,154.19	4,384.17	5,697.99	8,484.16
Daily vol. customer-bought (\$m)	2,200.03	1,173.49	1,539.42	2,145.60	2,737.46	4,165.64
Daily vol. customer-sold (\$m)	2,316.75	1,255.36	1,591.49	2,234.40	2,979.39	4,429.07

Table 2. Calibration targets.

This table presents moments from the TRACE data, which serve as calibration targets for the model. The pre-GFC period spans from 2006 to 2007, while the post-GFC period covers 2010 to 2019. We use the values from the pre-GFC period as our calibration targets.

Moment	Pre-GFC	Post-GFC
Num. customer-sold/num. customer-bought trades	0.7621	0.7794
Two-way spread, customer-sold (bps)	9.5160	13.6105
Two-way spread, customer-bought (bps)	10.8928	16.9578
Two-way spread for customer trades (bps)	10.1066	15.1040
Monthly std. dev. of log trade size, customer-bought	0.8802	0.8180
Monthly std. dev. of log trade size, customer-sold	0.8943	0.8562
Annual turnover, customer-bought (%)	19.9163	14.5082
Annual turnover, customer-sold (%)	21.2388	16.2388

Table 3. Values of calibrated parameters.

This table reports parameter values used in calibrating the model with associated empirical targets in the TRACE data. The target values are reported in Table 2.

Parameter	Value	Target (target value)
σ_δ^2 Dispersion in preference shocks	0.2277	Std. dev. of log trade size (0.89)
θ Dealers' bargaining power	0.6415	Avg. two-way customer trading cost, buy (10.8 bps)
η Elasticity of customers' utility	2.296	Avg. two-way customer trading cost, sell (9.5 bps)
γ Preference shock intensity	0.3902	Annual turnover for customer trades (20%)
ν Flow utility of dealers	0.9601	Dealer sector's pre-GFC bond holding share (2%)
μ Measure of dealers	0.0088	No. customer-sold / No. customer-bought (0.76)

Appendix

A Omitted Proofs

A.1 Proof of Proposition 1

Existence, uniqueness, and continuity. Let $x = (\delta, q, P, I)$ and $X = [\underline{\delta}, \bar{\delta}] \times (0, \infty) \times (0, \infty) \times [0, \infty)$. For any strictly positive q' and P' , let $X' = [\underline{\delta}, \bar{\delta}] \times [q', \infty) \times (0, P'] \times [0, \infty)$. Now consider the set $C_b(X')$ bounded continuous function of $x \in X'$, equipped with the sup norm. For any $h \in C_b(X')$, define the operator:

$$\begin{aligned} T[h](q, \delta \mid P, I) &= \frac{u_q(q, \delta) - rP}{r + \gamma + \lambda(1 - \theta)} + \frac{\gamma}{r + \gamma + \lambda(1 - \theta)} \mathbb{E}^F [h(q, \delta' \mid P, I)] \\ &\quad + \frac{\lambda(1 - \theta)}{r + \gamma + \lambda(1 - \theta)} \max\{h(q + I, \delta \mid P, I), 0\}. \end{aligned}$$

Since the first term $u_q(q, \delta) - rP$ is bounded on the domain X' , so is $T[h]$. Moreover, since $T[h]$ is the sum of continuous functions, it is also continuous. Hence, the operator T maps $C_b(X')$ into itself. Next, one easily verifies that it satisfies the Blackwell sufficient condition for a contraction (see Theorem 3.3 in [Stokey and Lucas, 1989](#), henceforth SLP), with modulus of contraction $(\gamma + \lambda(1 - \theta))/(r + \gamma + \lambda(1 - \theta))$. An application of the Contraction Mapping Theorem 3.2 in SLP establishes uniqueness of a bounded and continuous solution over any X' . Given uniqueness, this solution can be extended uniquely over the entire set X by letting $q' \rightarrow 0$ and $P' \rightarrow \infty$.

Conversely, if we consider any solution of the HJB defined over the domain X , then its restriction over the domain X' satisfies the HJB as well and so must coincide with the solution we constructed above.

Monotonicity. The operator T preserves the the following weak monotonicity properties: if h is increasing in δ and decreasing in (q, P, I) , then so is $T[h]$. Since weak monotonicity properties are preserved by passing to the limit, they are inherited by the fixed point. Now note that the first term of $T[h]$, $u_q(\delta, q) - rP$ is in fact strictly increasing in δ , and strictly decreasing in (q, P) . Hence, the fixed point, $\Sigma = T[\Sigma]$, also has these strict monotonicity properties.

$\Sigma(q, \delta)$ goes to infinity as $q \rightarrow 0$. Given that the third term of the Bellman equation is positive it follows that, for any q :

$$\Sigma(q, \underline{\delta}) \geq \frac{u_q(q, \underline{\delta}) - rP}{r + \gamma + \lambda(1 - \theta)} + \frac{\gamma}{r + \gamma + \lambda(1 - \theta)} \Sigma(q, \underline{\delta}) \Rightarrow \Sigma(q, \underline{\delta}) \geq \frac{u_q(q, \underline{\delta}) - rP}{r + \lambda(1 - \theta)}, \quad (13)$$

where we omitted the dependence of Σ on (P, I) for notational convenience. Since the utility function satisfies Inada condition, $\lim_{q \rightarrow 0} \Sigma(q, \underline{\delta}) = +\infty$. Since $\Sigma(q, \delta)$ is increasing in δ , it follows that $\lim_{q \rightarrow 0} \Sigma(q, \delta) = +\infty$ for all δ as well.

$\Sigma(q, \delta) < 0$ for q large enough. Let \hat{q} denote the solution of $u_q(\bar{\delta}, \hat{q}) = rP$. Evaluating $T[\Sigma]$ at $(\bar{\delta}, \hat{q})$ and keeping in mind that Σ is a fixed point, we obtain:

$$\Sigma(\hat{q}, \bar{\delta}) = \frac{\gamma}{r + \gamma + \lambda(1 - \theta)} \mathbb{E}^F [\Sigma(\hat{q}, \bar{\delta})] + \frac{\lambda(1 - \theta)}{r + \gamma + \lambda(1 - \theta)} \max\{\Sigma(\hat{q} + I, \bar{\delta}), 0\},$$

Now using that Σ is increasing in δ and that $\Sigma(\bar{\delta}, \hat{q} + I) \leq \Sigma(\bar{\delta}, \hat{q})$, we obtain

$$\Sigma(\hat{q}, \bar{\delta}) \leq \frac{\gamma}{r + \gamma + \lambda(1 - \theta)} \Sigma(\hat{q}, \bar{\delta}) + \frac{\lambda(1 - \theta)}{r + \gamma + \lambda(1 - \theta)} \max\{\Sigma(\hat{q}, \bar{\delta}), 0\},$$

implying that $\Sigma(\hat{q}, \bar{\delta}) \leq 0$. Given that Σ is strictly decreasing in q and strictly increasing in δ , the result follows.

Using $\Sigma(q, \delta)$ to construct the value function $V(q, \delta)$. The properties established in the previous paragraph implies that the equation $\Sigma(q, \delta) = 0$ has a unique solution, which we denote by $q^*(\delta)$. Fix some $q_0 \in (0, \infty)$ and let $\delta \mapsto V(q_0, \delta)$ denote the solution of:

$$rV(q_0, \delta) = u(q_0, \delta) + \gamma \mathbb{E}^F [V(q_0, \delta') - V(q_0, \delta)] + \lambda(1 - \theta) \int_{q_0}^{\min\{q^*(\delta), q_0 + I\}} \Sigma(x, \delta) dx.$$

The existence and uniqueness of such a function is guaranteed by standard contraction-mapping arguments. Our guess for the value function at any (δ, q) is:

$$V(q, \delta) = V(q_0, \delta) + \int_{q_0}^q \Sigma(x, \delta) dx + P(q - q_0). \quad (14)$$

Note that, since $\Sigma(q, \delta)$ is strictly decreasing, it follows that $V(q, \delta)$ is strictly concave. Next, we verify that $V(q, \delta)$ constructed above solves the HJB equation (1). Namely, multiplying the above equation by r and substituting in the HJB equation for $V(q_0, \delta)$ and $\Sigma(x, \delta)$, we have:

$$\begin{aligned} rV(q, \delta) &= u(q_0, \delta) + \gamma \mathbb{E}^F [V(q_0, \delta') - V(q_0, \delta)] + \lambda(1 - \theta) \int_{q_0}^{\min\{q^*(\delta), q_0 + I\}} \Sigma(x, \delta) dx \\ &+ \int_{q_0}^q (u_q(x, \delta) - rP) dx + \gamma \int_{q_0}^q \mathbb{E}^F [\Sigma(x, \delta') - \Sigma(x, \delta)] dx \\ &+ \lambda(1 - \theta) \int_{q_0}^q [\Sigma(\min\{q^*(\delta), x + I\}, \delta) - \Sigma(x, \delta)] dx + rP(q - q_0). \end{aligned}$$

Adding the first term on the first line, the first term on the second line, and the third term on the third line, we obtain $u(q_0, \delta)$, that is, the first term on the right-hand side of the HJB equation (1). Adding the second term on the first line together with the second term on the second line, we obtain $\gamma \mathbb{E}^F [V(q_0, \delta') - V(q_0, \delta)]$, that is, the second term on the right-hand side of the HJB equation (1). Grouping the last two other terms

together, we obtain:

$$\begin{aligned}
& \lambda(1 - \theta) \left[\int_{q_0}^{\min\{q^*(\delta), q_0+I\}} \Sigma(x, \delta) dx + \int_{q_0}^q [\Sigma(\min\{q^*(\delta), x+I\}, \delta) - \Sigma(x, \delta)] dx \right] \\
&= \lambda(1 - \theta) \left[\int_{q_0}^{\min\{q^*(\delta), q_0+I\}} \Sigma(x, \delta) dx + \int_{q_0+I}^{q+I} \Sigma(\min\{q^*(\delta), x\}, \delta) dx + \int_q^{q_0} \Sigma(x, \delta) dx \right] \\
&= \lambda(1 - \theta) \left[\int_{q_0}^{\min\{q^*(\delta), q_0+I\}} \Sigma(x, \delta) dx + \int_{\min\{q^*(\delta), q_0+I\}}^{\min\{q^*(\delta), q+I\}} \Sigma(x, \delta) dx + \int_q^{q_0} \Sigma(x, \delta) dx \right] \\
&= \lambda(1 - \theta) \int_q^{\min\{q^*(\delta), q+I\}} \Sigma(x, \delta) dx = \lambda(1 - \theta) (V(q', \delta) - V(q) - P(q' - q))
\end{aligned}$$

where $q' \equiv \min\{q^*(\delta), q + I\}$. In the above, the second line obtains by change of variable, and the third line because, by definition of $q^*(\delta)$, $\Sigma(\min\{q^*(\delta), x\}) = 0$ for all $x \geq q^*(\delta)$. The fourth line follows by piecing the three integral together and using our definition of $V(q, \delta)$ in equation (14). The last step is to verify that q' maximizes surplus subject to $0 \leq q' \leq q + I$, which follows immediately since $V(q', \delta)$ is strictly concave.

A.2 Proof of Lemma 1

The integrand in the dealer's profit function is the maximized trade surplus:

$$\begin{aligned}
g(i, q', \delta') &= \max_{0 \leq q'' \leq q'+i} \{V(q'', \delta') - V(q', \delta') - P(q'' - q')\} \\
&= V(\min\{q^*(\delta'), q' + i\}, \delta') - V(q', \delta') - P(\min\{q^*(\delta'), q' + i\} - q'),
\end{aligned}$$

where the second equality follows from Section 2.1, where we established that the optimum is attained for $q'' = \min\{q^*(\delta), q' + i\}$. Since, by its construction in Section 2.1, the value function is continuously differentiable in q' , a direct calculation reveals that g is continuously differentiable in i with derivative:

$$0 \leq g_i(i, q', \delta') = \Sigma(\min\{q^*(\delta), q' + i\}, \delta') = \max\{\Sigma(q' + i, \delta'), 0\} \leq \max\{\Sigma(q', \delta'), 0\}.$$

The inequality on the right-hand side shows that $|g_i|$ is bounded by an integrable function of (q', δ') . Hence, an application of Theorem 2.27 in Folland (1999) shows that

$$\int_{(q', \delta')} g(i, q', \delta') d\Phi(q', \delta')$$

is differentiable with respect to i and that its derivative is obtained by differentiating under the integral sign. The result follows.

A.3 Proof of Proposition 2

For this proof, pick some $0 < \underline{P} < \bar{P}$ and $0 < \underline{I} < \bar{I}$. Now, for all $(P, I) \in [\underline{P}, \bar{P}] \times [\underline{I}, \bar{I}]$, define the transition probability function (8) over $[0, \bar{q}] \times [\underline{\delta}, \bar{\delta}]$ where $\bar{q} > q^*(\delta | \underline{I}, \underline{P})$. One sees that \bar{q} has been chosen sufficiently large so that, for all $(P, I) \in [\underline{P}, \bar{P}] \times [\underline{I}, \bar{I}]$, the stationary distribution will belong to $[0, \bar{q}] \times [\underline{\delta}, \bar{\delta}]$.

Existence, uniqueness, and strong convergence. We rely on Lemma 11.11 and Theorem 11.12 in SLP, which provide a sufficient condition for existence of the operator T^{*N} to be a contraction mapping for some N and guarantee that the desired properties hold. The sufficient condition, labeled “condition M” by SLP, is that there exists some $\varepsilon > 0$ and some integer N such that, for any Borel set $B \subseteq [0, \bar{q}] \times [\underline{\delta}, \bar{\delta}]$, $\mathbb{P}^N(q, \delta, B) \geq \varepsilon$ for all (q, δ) or $\mathbb{P}^N(q, \delta, B^c) \geq \varepsilon$ for all (q, δ) . In our setting condition M follows because through their trades, customers are always able to reach their target holdings in an uniformly bounded number of trades, so that they eventually transition to the “diagonal” set:

$$D \equiv \{(q^*(\delta), \delta) : \delta \in [\underline{\delta}, \bar{\delta}]\} \quad (15)$$

Specifically, consider any Borel set B of $[0, \bar{q}] \times [\underline{\delta}, \bar{\delta}]$ and pick N such that $(N-1)I \geq q^*(\bar{\delta})$, so a customer reaches her target holding in at most $N-1$ successive trades with dealers. Then we have that, for all (q, δ) :

$$\mathbb{P}^N(q, \delta, B) \geq \mathbb{P}^N(q, \delta, B \cap D) \geq \frac{\gamma}{\gamma + \lambda} F((B \cap D)_\delta) \left(\frac{\lambda}{\gamma + \lambda} \right)^{N-1},$$

and where we use the short-hand A_δ to denote the set of δ such that $(q, \delta) \in A$ for some q (the “ δ -section” of the set A). In words, the above inequality states that the probability of reaching B in N transitions is greater than the probability of reaching the intersection of B with the diagonal set, which is itself greater than the probability of first drawing a preference shock in the δ -section of $B \cap D$, and receiving $N-1$ successive trading opportunities, which is sufficient to reach a point in $B \cap D$. Likewise, for all (q, δ) :

$$\mathbb{P}^N(q, \delta, B^c) \geq \mathbb{P}^N(q, \delta, B^c \cap D) \geq \frac{\gamma}{\gamma + \lambda} F([\underline{\delta}, \bar{\delta}] \setminus (B \cap D)_\delta) \left(\frac{\lambda}{\gamma + \lambda} \right)^{N-1},$$

since $B^c \cap D = D \setminus (B \cap D)$. Since either $F((B \cap D)_\delta) \geq 1/2$ or $F([\underline{\delta}, \bar{\delta}] \setminus (B \cap D)_\delta) \geq 1/2$, condition M holds for

$$\varepsilon = \frac{1}{2} \frac{\gamma}{\gamma + \lambda} \left(\frac{\lambda}{\gamma + \lambda} \right)^{N-1}.$$

Weak continuity with respect to (P, I) . One obtains weak continuity with respect to $(P, I) \in [\underline{P}, \bar{P}] \times [\underline{I}, \bar{I}]$ by an application of Theorem 12.13 in SLP. The first condition of the theorem is that the state space is compact, which is true here by assumption. The second condition is that the transition probability function is weakly continuous in (q, δ, P, I) , which follows by an application of Theorem 12.3, keeping in mind that

$q^*(\delta | P, I)$ is continuous in (δ, P, I) since it uniquely solve $\Sigma(q, \delta | P, I) = 0$, where $\Sigma(q, \delta | P, I)$ is continuous in (q, δ, P, I) . The third condition is the that the operator T^* has a unique fixed point for all (P, I) , which we established in the previous paragraph.

Monotonicity in P . Consider, for any bounded and measurable function, the conditional expectation operator:

$$\begin{aligned} T[g](q, \delta | P, I) &= \int_{(q', \delta')} g(q', \delta' | P, I) \mathbb{P}(q, \delta, dq', d\delta' | P, I) \\ &= \frac{\gamma}{\gamma + \lambda} \int_{\delta'} g(q, \delta' | P, I) dF(\delta') + \frac{\lambda}{\gamma + \lambda} g(\min\{q^*(\delta | P, I), q + I\}, \delta | P, I). \end{aligned}$$

Since $q^*(\delta | P, I)$ is decreasing in P , one sees that the operator preserves the following joint monotonicity property: for any bounded measurable function $g(q, \delta | P, I)$ that is increasing in q and decreasing in P , then $T[g](q, \delta | P, I)$ is also increasing in q and decreasing in P . By induction, it follows that this is also true for the n -transitions ahead conditional expectation: $T^n[g](q, \delta | P, I)$ is increasing in q and decreasing in P as well. In particular, if $P' \geq P$:

$$T^n[g](q, \delta | P', I) \leq T^n[g](q, \delta | P, I).$$

Given the strong convergence result established before, we can pass to the limit and obtain

$$\int g(q', \delta' | P, I) d\Phi(q', \delta' | P, I) \leq \int g(q', \delta' | P, I) d\Phi(q', \delta' | P, I),$$

as claimed.

A.4 Proof of Theorem 1

Lower and upper bounds on target holdings. Equation (13) evaluated at $q^*(\underline{\delta})$ implies that $u_q(q^*(\underline{\delta}), \underline{\delta}) - rP \leq 0$, from which it follows that, for all $\delta \in [\underline{\delta}, \bar{\delta}]$:

$$q^*(\delta) \geq u_q^{-1}(rP, \underline{\delta}). \tag{16}$$

Now consider the Bellman equation for $\Sigma(q, \bar{\delta})$ evaluated at $q \geq q^*(\delta)$ so that $\max\{\Sigma(q + I, \bar{\delta}), 0\} = 0$:

$$\Sigma(q, \bar{\delta}) \leq \frac{u_q(q, \delta) - rP}{r + \gamma + \lambda(1 - \theta)} + \frac{\gamma}{r + \gamma + \lambda(1 - \theta)} \Sigma(q, \bar{\delta}) \implies \Sigma(q, \bar{\delta}) \leq \frac{u_q(q, \delta) - rP}{r + \lambda(1 - \theta)}.$$

Letting $q \downarrow q^*(\bar{\delta})$ we obtain that $u_q(q^*(\bar{\delta}), \bar{\delta}) - rP \geq 0$, implying that for all $\delta \in [\underline{\delta}, \bar{\delta}]$

$$q^*(\delta) \leq u_q^{-1}(rP, \bar{\delta}). \quad (17)$$

Market-clearing given inventory. We now establish that, given some $I \in (0, s/\mu)$, there is a unique price $P(I)$ such that the market-clearing condition (9) holds. First, since the stationary distribution of asset holdings $\Phi(q', \delta' | P, I)$ is weakly continuous and decreasing in P , it follows that the left-hand side of the market-clearing condition is continuous and decreasing in P as well. Second, the lower and the upper bound of equations (16) and (17) imply that

$$u_q^{-1}(rP, \underline{\delta}) \leq \int_{(q', \delta')} q' d\Phi(q', \delta' | P, I) \leq u_q^{-1}(rP, \bar{\delta}).$$

Together with the Inada condition, this means that the amount of asset held by customers goes to infinity as $P \rightarrow 0$, and to zero as $P \rightarrow \infty$. Hence, an application of the Intermediate Value Theorem implies that the market clearing equation (9) has at least one solution.

To establish uniqueness we show that the left-hand side of (9) is a strictly decreasing function of P . To do so, first note that, by stationarity, trading between customers and dealers keeps the amount of asset held by customers sector constant:

$$\int_{(q', \delta')} q' d\Phi(q', \delta' | P, I) = \int_{(q', \delta')} \min\{q^*(\delta' | P, I), q' + I\} d\Phi(q', \delta' | P, I) \quad (18)$$

Hence for any two $0 < P_1 < P_2$:

$$\begin{aligned} & \int_{(q', \delta')} q' d\Phi(q', \delta' | P_1, I) - \int_{(q', \delta')} q' d\Phi(q', \delta' | P_2, I) \\ &= \int_{(q', \delta')} \min\{q^*(\delta' | P_1, I), q' + I\} d\Phi(q', \delta' | P_1, I) - \int_{(q', \delta')} \min\{q^*(\delta' | P_2, I), q' + I\} d\Phi(q', \delta' | P_2, I) \\ &\geq \int_{(q', \delta')} \left(\min\{q^*(\delta' | P_1, I), q' + I\} - \min\{q^*(\delta' | P_2, I), q' + I\} \right) d\Phi(q', \delta' | P_1, I) \\ &\geq \int_{(q', \delta') \in D} \left(\min\{q^*(\delta' | P_1, I), q' + I\} - \min\{q^*(\delta' | P_2, I), q' + I\} \right) d\Phi(q', \delta' | P_1, I) \\ &= \int_{(q', \delta') \in D} \left(q^*(\delta' | P_1, I) - q^*(\delta' | P_2, I) \right) d\Phi(q', \delta' | P_1, I) > 0, \end{aligned}$$

where: the equality on the second line follows from (18); the inequality on the third line follows from the fact that the stationary distribution is decreasing in P ; the inequality on the fourth line follows because target holdings are decreasing in P so that the integrand is positive; the equality on the last line follows because, on the diagonal set D which we defined earlier in (15), $(q', \delta') = (q^*(\delta'), \delta')$. Finally, the strict inequality on the last line follows because target holdings are strictly decreasing in P and the diagonal set D has strictly positive measure. Indeed, for any N such that $NI \geq q^*(\bar{\delta})$, it takes at most N consecutive trading

opportunities to reach the diagonal set from any (q', δ') in the support of the stationary distribution. Hence $\mathbb{P}(q', \delta', D) \geq (\lambda/(\gamma + \lambda))^N$. Now using stationarity we have that

$$\Phi(D) = T^* [\Phi] (D) = \int_{(q', \delta')} \mathbb{P}(q', \delta', D) d\Phi(q', \delta') \geq \left(\frac{\lambda}{\lambda + \gamma} \right)^N > 0.$$

Taken together, we obtain that the market-clearing condition has a unique solution $P(I)$. Given uniqueness of a solution and given the continuity of the market-clearing condition, it follows that $P(I)$ is continuous.

The set of v consistent with active intermediation. Let $V(I)$ denote the function defined in equation (10). Note that, in any equilibrium, $P > 0$ implies by (16) that $q^*(\underline{\delta}) > 0$ and so that $\mu I < s$. Thus, the set of v consistent with active intermediation is equal to the range of $V(I)$ over the open interval $(0, s/\mu)$. As $I \rightarrow s/\mu$, customers' average asset holding must go to zero, implying that the same is true for the smallest of customer's asset holdings, i.e., $q^*(\underline{\delta} | P(I), I) \rightarrow 0$. Therefore, (16) implies that $u_q^{-1}(rP(I), \underline{\delta}) \rightarrow 0$ and, given Inada conditions, that $P(I) \rightarrow \infty$. It thus follows from (17) that $q^*(\bar{\delta} | P(I), I) \rightarrow 0$ and that, for all (q', δ') in the support of the distribution $q' + I \geq q^*(\underline{\delta} | P(I), I) + I > q^*(\bar{\delta} | P(I), I) \geq q^*(\bar{\delta} | P(I), I)$. Hence, the inventory in advance never binds and $V(I) = rP(I)$. We conclude that $\lim_{I \rightarrow s/\mu} V(I) = +\infty$.

Next, recall that since the marginal trade surplus is decreasing in asset holdings and decreasing in preference shocks, we have that

$$\Sigma(q', \delta') \leq \Sigma(q^*(\underline{\delta} | P(I), I), \bar{\delta})$$

for all (q', δ') in the support of the stationary distribution. Keeping in mind that $\Sigma(q^*(\underline{\delta} | P(I), I), \bar{\delta}) \geq 0$, we obtain from the Bellman equation of the marginal trade surplus that:

$$\Sigma(q^*(\underline{\delta} | P(I), I), \bar{\delta}) \leq \frac{u_q(q^*(\underline{\delta} | P(I), I), \bar{\delta}) - rP(I)}{r + \gamma + \lambda(1 - \theta)} + \frac{\gamma + \lambda(1 - \theta)}{r + \gamma + \lambda(1 - \theta)} \Sigma(q^*(\underline{\delta} | P(I), I), \bar{\delta}).$$

Simplifying and using the lower bound for asset holdings in equation (16), we obtain

$$\Sigma(q', \delta') \leq \frac{u_q \left(u_q^{-1}(rP(I), \underline{\delta}), \bar{\delta} \right) - rP(I)}{r}.$$

Now the upper bound on asset holdings (17), together with market clearing, implies $rP(I) \leq u_q(s - \mu I, \bar{\delta})$ otherwise all asset holdings would be less than $s - \mu I$. Likewise, $rP(I) \geq u_q(s - \mu I, \underline{\delta})$. Taken together, we obtain an upper bound on the marginal trade surplus:

$$\Sigma(q', \delta') \leq \frac{u_q \left(u_q^{-1}(u_q(s - \mu I, \bar{\delta}), \underline{\delta}), \bar{\delta} \right) - u_q(s - \mu I, \underline{\delta})}{r}.$$

This implies that the marginal trade surplus remains bounded above as $I \rightarrow 0$ and, in turns, that $V(I)$ remains

bounded below as $I \rightarrow 0$. Altogether this means that:

$$\underline{y} \equiv \inf_{I \in (0, s/\mu)} V(I) > -\infty,$$

and we are done.

B The Model Without Inventory Constraints

To better understand the unique implication of our model, we study the model without inventory-in-advance constraints. Lemma 2 and 3 are versions of the results of Lagos and Rocheteau (2009) and Pinter and Üslü (2021). But Lemma 4 and 5, on the symmetry properties of this model, are new.

The results of this section are useful for several reasons. First, they provide natural initial conditions for our computations. Second, they help motivate the use of certain moments for parameter identification – in particular, they highlight the role played by the elasticity of the utility function in generating asymmetries in value-weighted transaction costs for purchases and sales. Third, they allow us to highlight unique properties of the model *with* inventory-in-advance constraint: in particular, in Lemma 4, we show that, without inventory constraints, the number of purchases and the number of sales are equal. Hence, in our model, any asymmetry in the number of sales and purchases is due to the presence of an inventory constraint.

B.1 Marginal surplus

In the model without inventory constraint, the Bellman equation for the marginal surplus reduces to:

$$(r + \gamma + \lambda(1 - \theta))\Sigma(q, \delta) = u_q(q, \delta) - rP + \gamma\mathbb{E}^F [\Sigma(q, \delta')].$$

Taking expectations on both sides we obtain that:

$$\mathbb{E}^F [\Sigma(q, \delta')] = \frac{\mathbb{E}^F [u_q(q, \delta')] - rP}{r + \lambda(1 - \theta)}$$

Plugging back into the Bellman equation, we obtain after a few lines of algebra, we obtain the following results:

Lemma 2. *In the model without inventory-in-advance constraint, the marginal surplus is*

$$(r + \lambda(1 - \theta))\Sigma(q, \delta) = \frac{r + \lambda(1 - \theta)}{r + \gamma + \lambda(1 - \theta)}u_q(q, \delta) + \frac{\gamma}{r + \gamma + \lambda(1 - \theta)}\mathbb{E}^F [u_q(q, \delta')] - rP.$$

In particular, for all $\eta > 0$ let $u(q, \delta) = \frac{q^{1-1/\eta}\delta}{1-1/\eta}$ if $\eta \neq 1$ and $u(q, \delta) = \log(q)\delta$ if $\eta = 1$. Then, the marginal surplus solves

$$(r + \lambda(1 - \theta))\Sigma(q, \delta) = D(\delta)q^{-1/\eta} - rP,$$

where

$$D(\delta) \equiv \frac{r + \lambda(1 - \theta)}{r + \gamma + \lambda(1 - \theta)}\delta + \frac{\gamma}{r + \gamma + \lambda(1 - \theta)}\mathbb{E}^F [\delta'].$$

B.2 Stationary distribution

In this Section, we derive a closed form solution for the steady-state distribution of asset holdings and preference shocks. [Lagos and Rocheteau \(2009\)](#) first proposed a solution for the case of a discrete distribution of preference shocks, and [Pinter and Üslü \(2021\)](#) extended this result to the case of an arbitrary distribution.

Lemma 3. *In the model without inventory-in-advance constraint, the cumulative measure of customers with holding less than q and utility type less than δ is:*

$$\Phi(q', \delta') = \frac{\gamma}{\gamma + \lambda} F(\delta^*(q')) F(\delta') + \frac{\lambda}{\gamma + \lambda} F(\min\{\delta^*(q'), \delta'\}). \quad (19)$$

In particular, for any integrable function $h(q', \delta')$, it holds that:

$$\begin{aligned} & \int_{(q', \delta')} h(q', \delta') d\Phi(q', \delta') \\ &= \frac{\gamma}{\gamma + \lambda} \int_{(\delta, \delta')} h(q^*(\delta), \delta') dF(\delta) dF(\delta') + \frac{\lambda}{\gamma + \lambda} \int_{\delta'} h(q^*(\delta'), \delta') dF(\delta'). \end{aligned} \quad (20)$$

Perhaps the formula that is easiest to interpret is (20), which allows to calculate any moment of the distribution. It shows that the joint distribution of asset holding can be viewed as a mixture of two distributions. First, with probability $\gamma/(\gamma + \lambda)$, the distribution is random: a customer of type δ' is endowed with a random asset holding, $q^*(\delta)$, where δ is drawn according to the distribution F . Second, with probability $\lambda/(\gamma + \lambda)$, the distribution is perfect: a customer of type δ' is endowed with her target asset holding, $q^*(\delta')$.

Proof of equation (19). We characterize $\Phi(q, \delta)$ in two steps. First, we derive the measure of customers with holdings less than q ,

$$\int_{(q', \delta')} \mathbb{I}_{\{q' \leq q\}} d\Phi(q', \delta').$$

In steady state, the gross outflow from the set of customers with holding less than q must equal the inflow:

$$\lambda \int_{(q', \delta')} \mathbb{I}_{\{q' \leq q\}} d\Phi(q', \delta') = \lambda \int_{\delta'} \mathbb{I}_{\{q^*(\delta') \leq q\}} dF(\delta')$$

The left-hand side is the gross outflow, created by customers who contact the market with current holding less than q . The right-hand side is the gross inflow, generated by all customers with *optimal holding* less than q who contact dealers. Clearly, $q^*(\delta) \leq q$ if and only if $\delta \leq \delta^*(q)$ where $\delta^*(q) \equiv (q^*)^{-1}(q)$. Hence, the above steady-state equation writes:

$$\int_{(q', \delta')} \mathbb{I}_{\{q' \leq q\}} d\Phi(q', \delta') = F(\delta^*(q)).$$

This preliminary step facilitates the derivation of the entire distribution. Indeed, the outflow-inflow equation for $\Phi(q, \delta)$ can now be written:

$$(\gamma + \lambda)\Phi(q, \delta) = \gamma F(\delta^*(q))F(\delta) + \lambda \int \mathbb{I}_{\{\delta' \leq \delta \text{ and } q^*(\delta) \leq q\}} dF(\delta').$$

The left-hand side is the gross outflow, created by all customers with type less than δ and holding less than q who either change type or contact dealers. The first term on the right-hand side is the gross inflow created by customers with holding less than q who draw a new type less than δ . The second term is the gross inflow created by trade with dealers: customers with utility type less than δ and optimal holding $q^*(\delta)$ less than q . Recalling the definition of $\delta^*(q)$, equation (19) follows.

Proof of equation (20). By Theorem 3.6.1 in Bogachev (2007), it follows that:

$$\int_{q', \delta'} h(q', \delta') d\Phi(q', \delta') = \int_{\delta, \delta'} h(q^*(\delta), \delta') d\Psi(\delta, \delta'),$$

where

$$\begin{aligned} \Psi(\delta, \delta') &= \Phi(q^*(\delta), \delta') \\ &= \frac{\gamma}{\gamma + \lambda} F(\delta)F(\delta') + \frac{\lambda}{\gamma + \lambda} F(\min\{\delta, \delta'\}) \\ &= \frac{\gamma}{\gamma + \lambda} \int_{\underline{\delta}}^{\delta} \int_{\underline{\delta}}^{\delta'} dF(x) dF(y) + \frac{\lambda}{\gamma + \lambda} \int_0^{\delta} dF(x) \mathbb{I}_{\{x \leq \delta'\}} \\ &= \frac{\gamma}{\gamma + \lambda} \int_{\underline{\delta}}^{\delta} \int_{\underline{\delta}}^{\delta'} dF(x) dF(y) + \frac{\lambda}{\gamma + \lambda} \int_0^{\delta} dF(x) \int_0^{\delta'} d\mathbb{I}_{\{x \leq y\}}, \end{aligned}$$

where the last line follows because $\mathbb{I}_{\{x \leq \delta'\}} = \int_0^{\delta'} d\mathbb{I}_{\{x \leq y\}}$, where $y \mapsto d\mathbb{I}_{\{x \leq y\}}$ is the Dirac measure centered at point x . We thus obtain that:

$$d\Psi(\delta, \delta') = \frac{\gamma}{\gamma + \lambda} dF(\delta) dF(\delta') + \frac{\lambda}{\gamma + \lambda} dF(\delta) d\mathbb{I}_{\{\delta \leq \delta'\}},$$

and equation (20) follows.

B.3 Buy-sell symmetry

In this section, we show that, in the model with no inventory constraints, the number of purchases and sales are equal. With a continuum of customers, the natural measure of “number of purchases” is the flow of

purchases per unit of time:

$$\lambda \int_{(q', \delta')} \mathbb{I}_{\{q' < q^*(\delta')\}} d\Phi(q', \delta').$$

Using equation (20) with $h(q', \delta') = \mathbb{I}_{\{q^*(\delta') > q'\}}$, we obtain that the flow of purchase writes:

$$\frac{\lambda\gamma}{\gamma + \lambda} \int_{(\delta, \delta')} \mathbb{I}_{\{q^*(\delta) < q^*(\delta')\}} dF(\delta) dF(\delta') + \frac{\lambda^2}{\gamma + \lambda} \int_{\delta'} \mathbb{I}_{\{q^*(\delta') < q^*(\delta')\}} dF(\delta').$$

Since the target holding function is strictly increasing, the indicator in the first integral simplifies to $\mathbb{I}_{\{\delta < \delta'\}}$. Moreover, it is clear that the indicator in the second integral is zero. Hence, the flow of purchase is:

$$\int_{(\delta, \delta')} \mathbb{I}_{\{\delta < \delta'\}} dF(\delta) dF(\delta') = \frac{\lambda\gamma}{\lambda + \gamma} \int_{\delta} dF(\delta) (1 - F(\delta)) = \frac{\lambda\gamma}{2(\lambda + \gamma)} \left(1 - \sum_{\delta \in [\underline{\delta}, \bar{\delta}]} \Delta F(\delta)^2 \right).$$

where $\Delta F(\delta) = F(\delta) - F(\delta-)$ and the last equality follows from the integration by part formula for functions of bounded variations, which can be found in Theorem 6.2.2 of [Carter and Van Brunt \(2000\)](#). Following the same steps we obtain the flow of sales:

$$\lambda \int_{(q', \delta')} \mathbb{I}_{\{q' > q^*(\delta')\}} d\Phi(q', \delta') = \frac{\lambda\gamma}{\lambda + \gamma} \int_{(\delta, \delta')} \mathbb{I}_{\{\delta > \delta'\}} dF(\delta) dF(\delta'),$$

which is clearly equal to the flow of purchase. Taking stock:

Lemma 4. *In the model without inventory-in-advance constraint, the flow of purchases and the flow of sales are both equal to*

$$\frac{\lambda\gamma}{2(\lambda + \gamma)} \left(1 - \sum_{\delta \in [\underline{\delta}, \bar{\delta}]} \Delta F(\delta)^2 \right).$$

Correspondingly, the average trade size of a sale and of a purchase are also equal.

B.4 Transaction cost asymmetry

We now investigate a different source of asymmetry: the average proportional transaction cost incurred by customers who purchase the asset vs. those who want to sell it. We show that, in the model without an inventory-in-advance constraint, the transaction costs are in general asymmetric. In addition the direction of the asymmetry depends on the elasticity of the utility function $u(q, \delta)$.

Let W denote the total value of purchases which, by market clearing, must be equal to the total value of

sales. The value weighted proportional transaction costs for purchase writes:

$$\begin{aligned} \text{TC}_p &= \int_{(q', \delta')} \mathbb{I}_{\{q^*(\delta') > q'\}} \frac{P(q^*(\delta') - q')}{W} \theta \frac{V(q^*(\delta'), \delta') - V(q', \delta') - P(q^*(\delta') - q')}{q^*(\delta') - q'} d\Phi(q', \delta') \\ &= \frac{\gamma}{\gamma + \lambda} \frac{\theta P}{W} \int_{\delta < \delta'} [V(q^*(\delta'), \delta') - V(q^*(\delta), \delta') - P(q^*(\delta') - q^*(\delta))] dF(\delta) dF(\delta'), \end{aligned}$$

where the second equality follows from equation (20). This formula shows that the value weighted transaction cost is proportional to the total surplus generated by purchases. Likewise, we obtain that the value weighted proportional transaction costs for sales writes:

$$\begin{aligned} \text{TC}_s &= \int_{(q', \delta')} \mathbb{I}_{\{q^*(\delta') < q'\}} \frac{P(q' - q^*(\delta'))}{W} \theta \frac{V(q^*(\delta'), \delta') - V(q', \delta') - P(q^*(\delta') - q')}{q' - q^*(\delta')} d\Phi(q', \delta') \\ &= \frac{\gamma}{\lambda + \gamma} \frac{\theta P}{W} \int_{\delta' < \delta} [V(q^*(\delta'), \delta') - V(q^*(\delta), \delta') - P(q^*(\delta') - q^*(\delta))] dF(\delta) dF(\delta') \\ &= \frac{\gamma}{\lambda + \gamma} \frac{\theta P}{W} \int_{\delta < \delta'} [V(q^*(\delta), \delta) - V(q^*(\delta'), \delta) - P(q^*(\delta) - q^*(\delta'))] dF(\delta) dF(\delta') \\ &= - \frac{\gamma}{\lambda + \gamma} \frac{\theta P}{W} \int_{\delta < \delta'} [V(q^*(\delta'), \delta) - V(q^*(\delta), \delta) - P(q^*(\delta') - q^*(\delta))] dF(\delta) dF(\delta'). \end{aligned}$$

where the second to last line renames the variables of integration, replacing δ by δ' and vice versa.

Subtracting the above expressions for TC_p and TC_s , and writing the surplus as an integral of marginal surplus, we see that transaction costs are larger for purchases than for sales, $\text{TC}_p > \text{TC}_s$, if

$$\int_{q^*(\delta)}^{q^*(\delta')} [\Sigma(q, \delta) + \Sigma(q, \delta')] dq > 0, \quad (21)$$

for all $\delta < \delta'$ and, vice versa, $\text{TC}_p < \text{TC}_s$ if the above expression is strictly negative for all $\delta < \delta'$.

For intuition consider any pair of target holdings $q^*(\delta)$ and $q^*(\delta')$, $\delta < \delta'$. A customer will purchase the quantity $q^*(\delta') - q^*(\delta)$ if her current asset holding is $q^*(\delta)$ and her current utility type is δ' . With Nash-bargaining, transaction costs are proportional to the surplus, which can be calculated by integrating below the marginal surplus curve, $\Sigma(q, \delta')$. Likewise, a customer will sell the same quantity $q^*(\delta') - q^*(\delta)$ when her current asset holding is $q^*(\delta')$ and her current utility type is δ . In that case, the transaction cost is obtained by integrating below the marginal surplus curve $-\Sigma(q, \delta)$. Condition (21) ensures that the integral is larger for the purchase than for the sale. Given buy-sell symmetry, there is an equal number of purchases and sales of this particular quantity, and the result follows.

Our main result in this section is:

Lemma 5. Suppose that the utility function is isoelastic, $u(q, \delta) = \frac{q^{1-1/\eta}}{1-1/\eta} \delta$. Then:

$$\text{TC}_p \begin{cases} > \text{TC}_s & \text{if } \eta < 2 \\ = \text{TC}_s & \text{if } \eta = 2 \\ < \text{TC}_s & \text{if } \eta > 2. \end{cases}$$

With an iso-elastic utility function, the marginal surplus has a simple closed-form solution shown in Lemma 2. Then, condition (21) writes:

$$\frac{1}{2} (D(\delta) + D(\delta')) \left(\frac{q^*(\delta')^{1-1/\eta}}{1-1/\eta} - \frac{q^*(\delta)^{1-1/\eta}}{1-1/\eta} \right) - rP (q^*(\delta') - q^*(\delta)) > 0.$$

Using that target holdings have zero marginal surplus, we have that $D(\delta)q^*(\delta)^{-1/\eta} = D(\delta')q^*(\delta')^{-1/\eta} = rP$, implying that target holdings are

$$q^*(\delta) = \left(\frac{D(\delta)}{rP} \right)^{-\eta} \quad \text{and} \quad q^*(\delta') = \left(\frac{D(\delta')}{rP} \right)^{-\eta}.$$

Plugging back, and assuming for now that $\eta \neq 1$, we can factor out the price rP and, after letting $x \equiv D(\delta')/D(\delta)$, we find that condition (21) holds for all $\delta' > \delta$ if and only if

$$f(x) > 0 \text{ for all } x > 1, \text{ where } f(x) \equiv \frac{x+1}{2} \frac{x^{\eta-1} - 1}{1-1/\eta} - (x^\eta - 1).$$

Taking derivatives twice, we obtain that:

$$\begin{aligned} \frac{df}{dx} &= \frac{\eta}{2} \left(\frac{x^{\eta-1} - 1}{\eta - 1} + x^{\eta-2} - x^{\eta-1} \right) \\ \frac{d^2f}{dx^2} &= \frac{\eta}{2} x^{\eta-3} (2 - \eta) (x - 1). \end{aligned}$$

Hence, $df/dx(x) = 0$ when $x = 1$, is strictly increasing in x if $\eta < 2$, is identically equal to zero if $\eta = 2$, and is strictly decreasing if $\eta > 2$. Since $f(1) = 0$ as well, it thus follows that, for $x > 1$, $f(x) > 0$ if $\eta < 2$, $f(x) = 0$ if $\eta = 2$ and $f(x) < 0$ if $\eta > 2$. The result follows. The log case $\eta = 1$ can be addressed separately.

C Sell-to-Buy Ratio: Adjustment for Order Imbalance

Suppose that, over some period of time, we index customer purchases by $b \in \{1, 2, \dots, N_B\}$ and customer sales by $s \in \{1, 2, \dots, N_S\}$. The total quantity of customer purchases and sales are:

$$Q_B = \sum_{b=1}^{N_B} q_b, \quad \text{and} \quad Q_S = \sum_{s=1}^{N_S} q_s,$$

where q is the trade quantity. The average trade size for purchases and sales are:

$$\bar{Q}_B = \frac{1}{N_B} \sum_{b=1}^{N_B} q_b, \quad \text{and} \quad \bar{Q}_S = \frac{1}{N_S} \sum_{s=1}^{N_S} q_s.$$

If the market is “balanced” during that time period, i.e., $Q_B = Q_S$ and dealers do not accumulate net inventories, we have:

$$N_B \bar{Q}_B = N_S \bar{Q}_S \implies \frac{N_S}{N_B} = \frac{\bar{Q}_B}{\bar{Q}_S}.$$

However, in the data, the market is not always balanced. Suppose, for example, that $Q_S > Q_B$. Then, the quantity $Q_S - Q_B > 0$ is absorbed on dealers’ balance sheets to clear the market, and a quantity Q_B is sold to customers. We further assume that the average of trade size is the same for both portions of the sale volume, $Q_S - Q_B$ and Q_B . Let \hat{N}_S and \hat{N}_B be the number of customer sales and purchases that do *not* stay on dealers’ balance sheets. Then, we have:

$$\hat{N}_S = N_S \frac{Q_B}{Q_S}, \quad \text{and} \quad \hat{N}_B = N_B.$$

In general, to adjust for order imbalance, we redefine the number of sales and purchases, N_S and N_B , as:

$$\hat{N}_S = N_S \times \min \left\{ \frac{Q_B}{Q_S}, 1 \right\}, \quad \text{and} \quad \hat{N}_B = N_B \times \min \left\{ \frac{Q_S}{Q_B}, 1 \right\}.$$

We can see that, by construction,

$$\frac{\hat{N}_S}{\hat{N}_B} = \frac{\bar{Q}_B}{\bar{Q}_S}.$$

In Figure 5, we plot the adjusted ratio $\frac{\hat{N}_S}{\hat{N}_B}$.