# Liquidity in Competitive Dealer Markets*

Peter Bank[†]     Ibrahim Ekren[‡]     Johannes Muhle-Karbe[§]

March 2, 2021

### Abstract

We study a continuous-time version of the intermediation model of Grossman and Miller [19]. To wit, we solve for the competitive equilibrium prices at which liquidity takers' demands are absorbed by dealers with quadratic inventory costs, who can in turn gradually transfer these positions to an exogenous open market with finite liquidity. This endogenously leads to transient price impact in the dealer market. Smooth, diffusive, and discrete trades all incur finite but nontrivial liquidity costs, and can arise naturally from the liquidity takers' optimization.

**Mathematics Subject Classification: (2010)** 91B26, 91B24, 91B51.

**JEL Classification:** C68, D43, G12.

**Keywords:** dealer market, segmentation, dynamic equilibrium, endogenous liquidity.

## 1 Introduction

A basic paradigm of classical financial theory is that markets are "perfectly liquid", in that arbitrary amounts can be traded immediately at the present market price. Yet, in reality liquidity is limited and generated by the interplay of strategic liquidity providers and consumers. This paper studies liquidity formation in a tractable equilibrium model, where intermediaries dynamically transfer liquidity from one market segment to another.

More specifically, we consider a risky asset that is traded in two markets. In the first market segment (the *open market*), prices and liquidity costs are given exogenously as in

the optimal execution literature following Almgren and Chriss [1]. In the second market segment (the *dealer market*), agents trade the asset competitively. Some of the agents (who we will refer to as *clients*) have exogenously given trading needs, and are therefore willing to pay a premium for the immediacy provided by the other agents (who we will call *dealers*).[1] The equilibrium price at which this second market clears is in turn driven by the agents' aggregate demand as in the model of Garleanu, Pedersen, and Pothesman [14]. The present study analyzes how agents with access to both markets dynamically intermediate between them. This extends the classical one-shot intermediation model of Grossman and Miller [19] to continuous time, where inventory management and in turn equilibrium prices reflect a tradeoff between past, present, and (expectations of) future order flow. Our paper also contributes to the literature on segmented markets [17, 33, 18], by analyzing a model where not only the positions of intermediaries in the "peripheral" dealer market are constrained, but also position adjustments in the "central" open market are costly.

Our stylized model is motivated by liquidity provision of competitive dealers in over-the-counter markets. As succinctly summarized by [7] in the context of foreign exchange markets,

> *"Dealers in over-the-counter financial markets provide liquidity to customers on a principal basis and manage the risk position that arises out of this activity in one of two ways. They may internalize a customer's trade by warehousing the risk in anticipation of future offsetting flow, or they can externalize the trade by hedging it out in the open market.[2] [..] The notion that dealers are either perfect internalizers or perfect externalizers is of course too constraining and in practice they will use to varying extent a mix of both to manage their risk."*

In our model, pure internalizers only trade in the dealer market, whereas externalization corresponds to offsetting trades in the open market. Our dynamic model in turn allows to study the optimal tradeoff between both risk management strategies. In order to permit the analysis of this complex interaction between clients, dealers, and the open market, we make a number of simplifying assumptions. First, all agents act as price takers in the dealer market as in [14]. This means that the equilibrium price is "competitive", which is reasonable for a large number of small liquidity providers, whose individual actions only have negligible effects on the overall market equilibrium.[3] Second, agents have quadratic inventory costs rather than preferences modelled by concave utility functions. Such quadratic holding costs are also used in [34, 35, 12, 29, 31], for example, because this reduced-form modelling of

---

[1]Unlike the "market makers" in [22, 25] and many more recent studies, these liquidity providers are not obliged to absorb any order flow, but trade at their discretion. Investment banks providing liquidity in foreign exchange markets are a typical example, compare [7].

[2]This refers to inter-dealer broker platforms or public electronic crossing networks, for example.

[3]Evidently, a very important but challenging direction for further research is to study extensions to a game-theoretic setting, where dealers dynamically adjust their price quotes to account for their inventories while competing with each other for the clients' order flow. For liquidation problems, a first step in this direction is undertaken in [8], where dealers strategically quote bid-ask spreads for randomly arriving clients with exogenous demand curves, and a strategic client who needs to liquidate a large position. In general, however, non-competitive equilibrium prices can be formalized in many different ways and will support many different price dynamics (cf., e.g., [35, 12]).

"inventory aversion" is considerably more tractable than risk aversion, yet still penalizes the accumulation of large and thereby risky positions. Third, the exogenous price in the open market has martingale dynamics, complemented by quadratic trading costs on the total and individual order flows. The martingale assumption, also made in [19, 16], ensures that agents do not speculate in the open market, but instead purely focus on intermediation, in line with the empirical observation that "specialists are good short-term traders but undistinguished long-term speculators" [20]. Quadratic costs on the trading rate as in [1, 16] penalize trading speed, modelling that positions can only be unwound gradually in the open market. The individual trading costs allow to distinguish how easy it is for different agents to access the open market; the benchmark example is that dealers have low or no access costs whereas retail clients have no direct access at all. The trading cost on the total order flow reflects that it becomes more difficult to, for instance, unwind a position in the open market when others are trying to do the same. For finitely many agents, this introduces a negative externality, where agents internalize the effects their trades have on their own execution prices but not on others. In order to ensure consistency with the competitive dealer market, we therefore identify the limiting cases of our results for a large number of identical small dealers. We find that the limiting model is similar to that of a representative dealer with suitably aggregated parameters.

In the setting outlined above, our main result, Theorem 3.4, identifies the unique equilibrium price in the dealer market as the solution of a linear forward-backward stochastic differential equation (FBSDE). Its explicit solution in terms of hyperbolic functions and conditional expectations of the agents' trading targets in turn reveals how liquidity is mitigated between the dealer market and the open market, and how this affects prices and optimal trading strategies. For the simplest case of an inelastic client order flow (i.e., only "noise trades") in the dealer market, we have

$$
\begin{aligned}
\text{Equilibrium Price } &= \text{Fundamental Price} \\
&\quad + \text{Holding Cost} \times \text{Expected Future Inventory.}
\end{aligned}
\tag{1.1}
$$

This adjustment is consistent with the small-risk aversion limit obtained for the model of [14] in [24, 23]. However, whereas the price impact is *permanent* in these models without access to an open market, it becomes *transient* in our model, a feature of price dynamics also documented in the empirical literature (cf., e.g., [20] and the references therein). The reason is that the client order flow can be gradually passed on to the open market here, so that the expected future demand in the formula of [23] is replaced by the optimally controlled inventory above. Put differently, the dealers in [14, 23, 24] correspond to "pure internalizers" in the terminology of [7], whereas our dealers employ both internalization and externalization in an optimal manner.

The liquidity costs implied by these equilibrium dynamics can be best understood in the case of a highly liquid open market. To wit, we show that as the trading costs in the open market tend to zero, the price in the dealer market converges to the same martingale price as in the open market. The clients' liquidity costs compared to this benchmark in turn admit simple, intuitive expressions that depend on the fluctuations of the client demand. If the latter is smooth, trading through the dealer market is approximately equivalent to trading directly in the open market at a higher cost that depends on the number of dealers and

reflects the premium that is necessary to entice the dealers to provide the necessary liquidity. This markup disappears in the limit of many small dealers, so that smooth client flow then trades at approximately the same prices as in the open market. This holds irrespective of the dealers' risk tolerances, since smooth flow can be hedged efficiently in the open market.

This changes for diffusive noise trades. Such irregular order flow is more difficult to pass on to the illiquid open market. Accordingly, the corresponding liquidity costs asymptotically scale with the square root of the trading cost in the open market, multiplied by the square root of the dealers' holding cost. The order flow at hand enters through its quadratic variation, similarly as in the reduced form model of [11]. Here, however, the trading costs of such "rough" strategies cannot be avoided as in [11, 2] by approximating it with smooth strategies. More generally, the liquidity costs implied by our model are continuous in the client demand. As a consequence, strategies of various forms are priced consistently and incur finite but nontrivial liquidity costs which reflect the strategies' regularity.

The results described so far pertain to the case of a fixed given client demand as in [14]. However, our model is tractable enough to also endogenize (some or all of) the order flow. More specifically, we can study how clients optimally track their target positions, so that their demand responds to prices. In this case, Formula (1.1) remains valid, *if* the future inventory is replaced by aggregate demand minus offsetting positions in the open market. If demand is inelastic, targets and actual positions coincide, and we recover the previous formula. In contrast, price-sensitive clients accept some deviations from their target positions, but the equilibrium price is still determined by the aggregate (expectations) of all individual target positions.

To illustrate the implications of this result and the corresponding optimal trading strategies, we consider two examples. In the first, the clients' trading targets are constant as in [12]. This leads to optimal liquidation problems similar to the ones studied in [5, 1, 32] and many more recent papers. Subsequently, we turn to an example with diffusive trading targets that correspond to the "high-frequency trading needs" considered in [26, 35].

For constant trading targets, the clients' optimal trading strategies are of a similar form as in optimal liquidation models with transient price impact [32]: isolated bulk trades combined with otherwise smooth order flow. In contrast, diffusive trading targets as in [26, 35] lead to optimal client demands with nontrivial quadratic variation. Therefore, our model consistently combines the qualitative properties of standard models for optimal liquidation, while nevertheless allowing for rapidly fluctuating inventories in line with the empirical evidence documented by [9], for example.

The effects of limited liquidity on the equilibrium price also depends on the properties of the client order flow. In the optimal liquidation example (and, more generally, for deterministic targets), illiquidity only affects price levels and expected returns, but not volatility. In contrast, random demands such as the diffusive trading targets generally also change volatility. Here, the sign of this change is determined by the correlation between trading targets and the fundamental value similarly as in [21].

This article is organized as follows. Our model for the dealer and open markets is introduced in Section 2. Section 3 in turn contains our characterization of equilibrium prices and trading strategies in this context, as well as a detailed discussion of the implications of these results. For better readability, all proofs are delegated to the appendix.

**Notation** A predictable process $X = (X_t)_{t \in [0,T]}$ belongs to $\mathcal{L}^2$ if $\mathbb{E}[\int_0^T X_t^2 dt] < \infty$ and to $\mathcal{S}^2$ if $\mathbb{E}[\sup_{t \in [0,T]} X_t^2] < \infty$. The set of square-integrable martingales is denoted by $\mathcal{M}^2$, and we write $\mathcal{H}^2$ for the semimartingales whose local martingale part belongs to $\mathcal{M}^2$ and whose finite-variation part has square-integrable total variation. For an Itô process with dynamics $dX_t = \mu_t dt + \sigma dW_t$, this holds if $\mathbb{E}[(\int_0^T |\mu_t| dt)^2 + \int_0^T \sigma_t^2 dt] < \infty$.

# 2 Model

## 2.1 Agents

We consider finitely many agents indexed by $a \in \mathcal{A} \neq \emptyset$. Agent $a$ has mass $m(a) \in (0, \infty)$, and all agents share the same beliefs and information flow described by the filtered probability space $(\Omega, \mathcal{F}, (\mathcal{F}_t), \mathbb{P})$ satisfying the usual conditions of right-continuity and completeness.

## 2.2 Financial Market

The agents invest in two assets. The first is riskless and bears no interest (for simplicity). The second asset pays an exogenous liquidating dividend $\mathscr{D}_T \in L^2(\mathcal{F}_T)$ at time $T > 0$ and can be traded in two markets.

**Dealer Market** In the first market, the agents $a \in \mathcal{A}$ competitively trade the asset at the *equilibrium price* $S \in \mathcal{H}^2$ which clears this market and matches the terminal payoff $S_T = \mathscr{D}_T$. In our main result, Theorem 3.4, we determine this price explicitly.

If some of the agents have exogenous trading needs whereas others don't, this market can be interpreted as a competitive dealer market, where the "dealers" (or "intermediaries" or "market makers") from the second group earn a premium for supplying liquidity to the first group of "clients" (or "outside customers" [19] or "end-users" [14]).

In general, if agent $a \in \mathcal{A}$ follows a trading strategy $K^a \in \mathcal{S}^2$ (specifying the number of risky assets held) in the dealer market, then this generates the expected P&L

$$\mathbb{E}\left[\int_0^T K_t^a dS_t\right] = \mathbb{E}\left[\int_0^T K_t^a dA_t\right]. \tag{2.1}$$

Here, $A$ denotes the "drift" of $S$, i.e., the predictable bounded variation part in its Doob-Meyer-decomposition.

**Open Market** The dealer market described above models liquidity provision by "pure internalizers", who absorb their customers' order flow until it is offset by future incoming orders. In real dealer markets, this is complemented by "externalization", i.e., actively "hedging out trades in the open market" [7]. In the classic model of Grossman and Miller [19] (or in the recent work of Garleanu and Pedersen [16, Section 3.2], for example) the risk inherent in this intermediation process is linked to a fixed search time needed to locate a counterparty that allows the dealers to lay off their positions. In the present study, we formulate and solve a *dynamic* version of the dealers' inventory management problem, where

the risky asset can be traded at all times on the open market, but with a trading friction that imposes a penalty for the fast liquidation of large positions.

To wit, in addition to the dealer market, we consider a second market for the risky asset, where its unaffected price process is given by its expected dividend,

$$\mathscr{D}_t := \mathbb{E}_t[\mathscr{D}_T], \quad 0 \leq t \leq T.$$

This martingale price is assumed for simplicity, because it eliminates speculation in the open market. The corresponding execution price is in turn given by $\mathscr{D}_t + \lambda \bar{u} + \frac{1}{2}\lambda^a u^a$ for the trading rate $u^a \in \mathcal{L}^2$ of agent $a \in \mathcal{A}$, that is, the speed at which her current position $U_t^a = \int_0^t u_s^a ds$ in the open market is adjusted. (The position $U^a$ then automatically belongs to $\mathcal{S}^2$.) Here, $\bar{u} := \sum_{a \in \mathcal{A}} m(a) u^a$ and the constant $\lambda \geq 0$ describes the impact that this total order flow has on the market price. In contrast, the parameter $\lambda^a$ models agent $a$'s idiosyncratic trading frictions such as search costs, that do not depend on the other agents' trading activities.[4] If the open market models interdealer trading, then the individual frictions are naturally finite for the dealers but infinite for their customers.

In summary, agent $a$'s trades in the open market generate an expected P&L of

$$\mathbb{E}\left[\mathscr{D}_T U_T^a - \int_0^T (\mathscr{D}_t + \lambda \bar{u}_t) u_t^a dt - \int_0^T \frac{\lambda^a}{2}(u_t^a)^2 dt\right]$$
$$= -\mathbb{E}\left[\int_0^T \left(\lambda u_t^{-a} u_t^a + \left(\lambda m(a) + \frac{1}{2}\lambda^a\right)(u_t^a)^2\right) dt\right]. \tag{2.2}$$

Here, $\bar{u}^{-a} := \bar{u} - m(a)u^a$ denotes the other agents' aggregate trading rate.

## 2.3 Goal Functionals

We now formulate the agents' goal functionals. Similarly as in [15, 16, 35, 6, 12], these are chosen to be of a linear-quadratic form for maximum tractability.

**Trading Targets and Inventory Costs**  The agents have an incentive to trade due to exogenous exposures to risk.[5] A particularly convenient way to model this is to fix a *target position* $\xi^a \in \mathcal{S}^2$ for each agent $a \in \mathcal{A}$, and in turn penalize the mean-squared distance to the agent's total risky position (accumulated either from dealers or in the open market):

$$\frac{1}{2}\mathbb{E}\left[\int_0^T (\xi_t^a - K_t^a - U_t^a)^2 dt\right]. \tag{2.3}$$

The importance of the inventory penalty (2.3) relative to expected trading gains and losses (2.1), (2.2) is described by the agents' *risk-tolerances* $\rho^a > 0$, $a \in \mathcal{A}$. Together, this leads to the linear-quadratic goal functional

$$J^a(K^a, u^a; \bar{u}^{-a}, S) := \mathbb{E}\left[\int_0^T K_t^a dA_t - \int_0^T \left(\lambda u_t^{-a} u_t^a + \left(\lambda m(a) + \frac{1}{2}\lambda^a\right)(u_t^a)^2\right) dt\right] \tag{2.4}$$
$$- \mathbb{E}\left[\int_0^T \frac{1}{2\rho^a}(\xi_t^a - K_t^a - U_t^a)^2 dt\right].$$

---

[4]This is similar in spirit to allowing agents to pay a quadratic cost to shorten their search times.

[5]The impact of illiquidity on trades due to heterogenous beliefs about fundamentals is studied in [28].

Given the other agents' aggregate transactions $\bar{u}^{-a}$ in the open market and a price process $S \in \mathcal{H}^2$ in the dealer market, agent $a$ maximizes this over $K^a \in \mathcal{S}^2$ and $u^a \in \mathcal{L}^2$. Our goal now is to determine the *equilibrium* price $S$ for which the dealer market clears.

# 3   Equilibrium

The agents interact with each other through the equilibrium price in the dealer market and their common price impact in the open market. With additional exogenous demands $K^N \in \mathcal{S}^2$ from noise traders in the dealer market, this leads to the following notion of equilibrium:

**Definition 3.1.** *A price process $S \in \mathcal{H}^2$ in the dealer market and trading strategies $(K^a)_{a \in \mathcal{A}}$ and $(u^a)_{a \in \mathcal{A}}$ in the dealer and the open markets, respectively, form an* equilibrium *if*

(i) *the dealer market clears, in that $K^N + \sum_{a \in \mathcal{A}} m(a) K^a = 0$;*

(ii) *the trading strategies form an (open loop) Nash equilibrium in that the strategy $(K^a, u^a)$ maximizes agent $a$'s target functional $J^a(\cdot, \cdot; \bar{u}^{-a}, S)$ over $\mathcal{S}^2 \times \mathcal{L}^2$.*

**Remark 3.2.** As in [14], we study a competitive equilibrium in the dealer market for tractability, and also because non-competitive equilibrium prices in continuous time necessarily introduce additional degrees of freedom (cf., e.g., [35, 12]).

The price-taking assumption in competitive equilibria is reasonable for a large number of small agents. In this regime, however, it is challenging to model the trading frictions in the open market in a consistent and nontrivial manner. For example, if each (small) agent simply neglects their own contribution to the execution price, this leads to a linear trading cost and in turn large trading rates that are not negligible even for small agents. To resolve this, we start from Nash competition between finitely many agents and in turn identify the competitive limit for many small agents in a second step, cf. the discussion after Theorem 3.4.

We now present our main result on existence and uniqueness result of an equilibrium price for the dealer market. To formulate this concisely, it is convenient to recall the following result from [6][Theorem A.4]:

**Lemma 3.3.** *Fix $X \in \mathcal{L}^2$ and $\Delta > 0$. Then, the unique solution $(u, U) \in \mathcal{L}^2 \times \mathcal{H}^2$ of the following linear forward-backward stochastic differential equation (FBSDE)[6]*

$$du_t = \Delta \left( U_t - X_t \right) dt + dM_t, \quad u_T = 0,$$
$$U_0 = 0, \quad dU_t = u_t dt, \tag{3.1}$$

*is given by*

$$\mathbf{u}_t^\Delta(X) := \mathbb{E}_t \left[ \int_t^T k^\Delta(t, s) X_s ds \right] - F^\Delta(t) \mathbf{U}_t^\Delta(X), \quad 0 \le t \le T, \tag{3.2}$$

$$\mathbf{U}_t^\Delta(X) := \frac{1}{\Delta} \int_0^t k^\Delta(s, t) \mathbb{E}_s \left[ \int_s^T k^\Delta(s, r) X_r dr \right] ds, \quad 0 \le t \le T,$$

---

[6]Here, the square-integrable martingale $M$ is part of the solution.

7

*for the kernel and the function*

$$k^\Delta(t,s) = \frac{\Delta \cosh(\sqrt{\Delta}(T-s))}{\cosh(\sqrt{\Delta}(T-t))}, \quad F^\Delta(t) = \sqrt{\Delta} \tanh\left(\sqrt{\Delta}(T-t)\right), \quad 0 \le s, t \le T.$$

To rule out frictionless trading in the open market, we assume that

$$\lambda + \lambda^a > 0, \quad a \in \mathcal{A},$$

which allows us to introduce the *price elasticities*

$$\eta^a = 1/(m(a)\lambda + \lambda^a), \quad \bar{\eta} = \sum_{a \in \mathcal{A}} m(a)\eta^a, \quad \eta = 1/\lambda.$$

(In the case $\lambda = 0$, it is natural to let $\eta = +\infty$ and $\Delta = \bar{\eta}/\bar{\rho}$.) We also define the agents' aggregate risk tolerance and target position:

$$\bar{\rho} = \sum_{a \in \mathcal{A}} m(a)\rho^a, \quad \bar{\xi} = \sum_{a \in \mathcal{A}} m(a)\xi^a.$$

With these ingredients, we can now formulate our main result, which characterizes the unique equilibrium price in the dealer market as well as the agents' optimal trading strategies in both the dealer market and the open market.

**Theorem 3.4.** *(i) There exists a unique equilibrium* $(S, (K^a, u^a)_{a \in \mathcal{A}}) \in \mathcal{H}^2 \times (\mathcal{S}^2 \times \mathcal{L}^2)^{\mathcal{A}}$.

*(ii) The agents' aggregate position* $\bar{U} = \sum_{a \in \mathcal{A}} m(a)U^a$ *in the open market is*

$$\bar{U}_t = \mathbf{U}_t^\Delta\left(K^N + \bar{\xi}\right), \quad 0 \le t \le T, \tag{3.3}$$

*for* $\mathbf{U}^\Delta$ *from Lemma 3.3 with*

$$\Delta = \frac{1}{\bar{\rho}}\frac{\eta\bar{\eta}}{\eta + \bar{\eta}}.$$

*Agent $a$'s share of* $\bar{U}_t$ *is* $U_t^a = \frac{\eta^a}{\bar{\eta}}\bar{U}_t$.

*(iii) The equilibrium price of the risky asset in the dealer market is*

$$S_t = \mathbb{E}_t\left[\mathscr{D}_T - \int_t^T \mu_s ds\right], \quad \text{where } \mu_t = \frac{1}{\bar{\rho}}\left(\bar{U}_t - K_t^N - \bar{\xi}_t\right), \quad 0 \le t \le T. \tag{3.4}$$

*Agent $a$'s optimal position in the dealer market is*

$$K_t^a = \xi_t^a - \frac{\eta^a}{\bar{\eta}}\bar{U}_t + \frac{\rho^a}{\bar{\rho}}\left(\bar{U}_t - K_t^N - \bar{\xi}_t\right), \quad 0 \le t \le T. \tag{3.5}$$

*Proof.* See Appendix A. ☐

The risk premium $\mu_t$ in the dealer market is thus determined by the agents' aggregate exposure $K^N + \bar{\xi} - \bar{U}$ at any one time, measured in units of their aggregate risk tolerance $\bar{\rho}$. In view of (3.5), the risk premium incentivizes agent $a$ to accept her share of this exposure in proportion to her individual risk tolerance $\rho^a$, having contributed to the aggregate risk-transfer to the open market $\bar{U}$ in proportion with her effective elasticity $\eta^a$. The asset price *fluctuations* in the dealer market (i.e., the martingale part of the asset price (3.4)) in turn depend on how uncertainty about the liquidating dividend $\mathscr{D}_T$ and future risk premia are revealed over time.

Alternatively, the first-order conditions for the optimal trading rates in the open market show that the equilibrium price from Theorem 3.4 admits the following concise representation:

$$S_t = \mathscr{D}_t + \left( \frac{1}{\eta} + \frac{1}{\bar{\eta}} \right) \bar{u}_t.$$

This means that the adjustment of the equilibrium price compared to expected dividend is determined by the agents' aggregate trading rate in the open market measured in units of the aggregate price elasticity $1/(\eta^{-1} + \bar{\eta}^{-1})$ resulting from the combination of the agents' idiosyncratic and their common impact on open market prices. To wit, suppose the agents are on aggregate buying in the open market ($\bar{u}_t > 0$) because they want to increase their net position. Then they will also be willing to purchase further risky assets in the dealer market at a premium. Despite this appealingly simple interpretation, the dependence of the equilibrium price on the model parameters is generally rather involved, since it typically depends on the past, present, and (expectations of) future demands.

## 3.1 Large-Liquidity Asymptotics for Noise Trades

In order to better understand the equilibrium prices (3.4) and their implications, we therefore first consider the simplest version of the dealer market, where $M \in \mathbb{N}$ "dealers" (i.e., agents without idiosyncratic trading targets, $\xi^d = 0$) with common masses $m(a) = 1/M$, risk tolerances $\rho^a = \rho_d > 0$, and individual trading costs $\lambda^a = 0$ absorb the demand $K^N \in \mathcal{S}^2$ of noise traders.

Our first result shows that the equilibrium price (3.4) approaches the expected dividend as the open market becomes more and more liquid for $\lambda \to 0$:

**Proposition 3.5.** *For any noise-trader demand $K^N \in \mathcal{S}^2$, the equilibrium price $S$ from Theorem 3.4 converges to the expected dividend $\mathscr{D}_t = \mathbb{E}_t[\mathscr{D}_T]$ in the particularly strong sense that*

$$\sup_{-1 \leq H \leq 1 \ predictable} \mathbb{E} \left[ \sup_{0 \leq t \leq T} \left| \int_0^t H_s d(\mathscr{D}_s - S_s) \right|^2 \right] \to 0, \quad as \ \lambda \to 0.$$

*In particular $S$ converges to $\mathscr{D}$ in the Emery topology as $\lambda \to 0$. Moreover, the corresponding wealth processes generated by the noise traders' demand satisfy*

$$\int_0^T K_t^N dS_t = \int_0^T K_t^N d\mathscr{D}_t + o(1), \quad in \ L^1 \ as \ \lambda \to 0.$$

Proposition 3.5 asserts that, as the open market becomes more and more liquid, the dealer price approaches the expected payoff of the risky asset. This is intuitive because

with vanishing trading frictions in the open market, dealers can immediately unload their inventories at this price.

Given additional structure of the noise-trader demand, we can also identify the leading-order correction term for the noise traders' wealth process, i.e., the liquidity costs implied by the dealers' nontrivial but finite risk-bearing capacity. The form of this leading-order correction term depends on the variability of the clients' demand. To illustrate this, we discuss two examples that appear frequently in applications: smooth demands $K_t^N = \int_0^t \mu_s^N ds$ and diffusive demands with Itô dynamics $K_t^N = \int_0^t \mu_s^N ds + \int_0^t \sigma_s^N dW_s$.

**Smooth Demands**  We first discuss demands that accumulate at a finite, absolutely continuous rate. In this case, the dealers could hedge their exposure perfectly by passing on their positions immediately to the open market subject to the quadratic cost $\lambda$ imposed on the corresponding trading rate. Therefore, the dealers could break even using this strategy for an execution price equal to the fundamental value plus the trading cost in the open market. However, due to the quadratic nature of the trading cost, they could achieve strictly positive profits in this case by absorbing only a fraction of the clients' demands. Accordingly, in equilibria with finitely many dealers, the dealers need to be paid an additional premium, but the latter vanishes in the limit of many small dealers.

More specifically, the subsequent result explicitly identifies the leading-order term of this liquidity cost as $(M+1)/M$ times the trading cost in the open market. This term is independent of the dealers' inventory costs, since smooth client demands can be hedged very efficiently by trading in the open market.

**Lemma 3.6.** *Suppose that $K_t^N = \int_0^t \mu_s^N ds$ for a continuous process $\mu^N \in \mathscr{S}^2$. Then, the liquidity costs generated by $K$ are*

$$\int_0^T K_t^N d\mathscr{D}_t - \int_0^T K_t^N dS_t = \lambda \frac{M+1}{M} \int_0^T \left( \mu_t^N \right)^2 dt + o(\lambda), \quad in \ L^1 \ as \ \lambda \to 0. \qquad (3.6)$$

In particular, Lemma 3.6 implies that the equilibrium trading costs in *competitive* dealer markets with many small dealers are approximately the same as in the open market.

**Diffusive demands**  Next, we turn to trading strategies with nontrivial Brownian fluctuations. These could not be implemented directly in the open market, but can be traded at a finite cost through the dealers. Since the dealers can hedge these more irregular order flows less efficiently than the smooth flows considered above, the corresponding trading costs are of a higher asymptotic order, namely $O(\sqrt{\lambda})$ instead of $O(\lambda)$ as $\lambda \to 0$. Moreover, the dealers' inventory cost now becomes visible in the leading-order term. The asymptotically crucial feature of client demand turns out to be its quadratic variation, which also appears in the reduced-form model of [11], for example:[7]

---

[7]The same scaling and the target strategy's quadratic variation also appear if such diffusive target positions are tracked optimally in markets with quadratic costs in the trading rate, cf. [27] and the references therein.

**Lemma 3.7.** *Suppose that the underlying filtration is generated by a Brownian motion $W$ and assume that the noise-trader demand has Itô dynamics,*

$$K_t^N = \int_0^t \mu_s^N ds + \int_0^t \sigma_s^N dW_s, \quad 0 \le t \le T.$$

*Here, $|K^N|^2, |\mu^N|^2, |\sigma^N|^2 \in \mathcal{H}^2$ and these processes are Malliavin differentiable in the sense of [30, p. 27]): $K_t^N, \mu_t^N, \sigma_t^N \in \mathbb{D}^{1,2}$, with continuous Malliavin derivatives*

$$s \mapsto \left( D_t\left(K_s^N\right), D_t\left(\mu_s^N\right), D_t\left(\sigma_s^N\right) \right), 0 \le t \le s \le T.$$

*Finally, suppose that $\sup_{0 \le t \le T} \mathbb{E}[\sup_{t \le s \le T}(|(D_t(K_s^N))|^2 + |(D_t(\mu_s^N))|^2)] < \infty$. Then, the liquidity costs generated by the demand $\bar{K}^N$ are*

$$\int_0^T K_t^N d\mathscr{D}_t - \int_0^T K_t^N dS_t = \sqrt{\frac{\lambda}{\rho_d} \frac{M+1}{M}} \int_0^T (\sigma_t^N)^2 dt + o(\sqrt{\lambda}), \quad in \ L^1 \ as \ \lambda \to 0.$$

**Remark 3.8.** The regularity conditions of Proposition 3.7 are satisfied, in particular, if the demand $K^N$ is the solution of a scalar stochastic differential equation whose drift and diffusion coefficients are twice continuously differentiable with bounded derivatives of orders $0, 1, 2$. In this case, the required bounds for the Malliavin derivatives follow from [30, Theorem 2.2.1].

## 3.2   Competitive Dealer Markets

Next, we apply Theorem 3.4 to study liquidity in a competitive dealer market that we model by a large number of small homogenous dealers and clients. For simplicity, suppose that there are no noise traders ($K^N = 0$), but $m_d M \in \mathbb{N}$ dealers with common risk tolerance $\rho_d$ and $m_c M$ clients with common risk tolerance $\rho_c > 0$, all with equal mass $1/((m_d + m_c)M)$. That is, the dealers and clients make up fractions

$$q_d = \frac{m_d}{m_d + m_c} \quad \text{and} \quad q_c = \frac{m_c}{m_c + m_d}$$

of the total number $(m_d + m_c)M$ of agents. This allows us to study the limiting behaviour of equilibrium prices and trading strategies as the number of agents becomes large for $M \to \infty$, while the fractions of dealers and clients remains fixed.[8] The clients have a common trading target $\xi^c \in \mathcal{S}^2$, whereas the dealers have no trading targets ($\xi^d = 0$) and therefore only trade to earn premia for providing liquidity.

As is natural for an open market that describes interdealer trading, we assume that the individual trading frictions are zero for all dealers ($\lambda^a = 0$) and infinite for the customers ($\lambda^a = \infty$). In summary, we then have

$$\bar{\rho} = q_c \rho_c + q_d \rho_d, \quad \bar{\eta} = \frac{m_d M}{\lambda}, \quad \bar{\xi} = q_c \xi^c, \quad \Delta = \frac{1}{(q_c \rho_c + q_d \rho_d)\lambda(1 + \frac{1}{m_d M})}.$$

---

[8]We are grateful to an anonymous referee for prompting us to pursue this extension with general rather than equal proportions of dealers and clients.

As the number of dealers and clients becomes large for $M \to \infty$ (with the proportions $q_d$ and $q_c$ of dealers and clients remaining fixed), $\Delta$ therefore converges to a nonzero and finite limit,

$$\Delta_\infty = \frac{1}{(q_c\rho_c + q_d\rho_d)\lambda}.$$

Let us now compare this to a market with the same proportions $q_c$ of clients and $q_d$ of dealers, but where the $m_d M$ small dealers are replaced by a single dealer with the same aggregate mass $q_d$ and risk tolerance $\rho_d$. Then the above expressions for $\bar{\rho}$, $\bar{\xi}$ remain unchanged, but we have $\bar{\eta} = 1/\lambda$ and in turn

$$\Delta_1 = \frac{1}{(q_c\rho_c + q_d\rho_d)2\lambda}.$$

As the equilibrium prices in Theorem 3.4 only depend on the aggregate trading target $\bar{\xi}$ (which is the same in both cases) and the parameter $\Delta$, we see that the market populated by many small dealers is equivalent to a market with a single representative dealer, if the liquidity cost in the competitive market is rescaled by a factor of two. Accordingly, the positive effects (for the clients) of competition (which drives down each dealer's profits) outweigh the negative effects of uncoordinated trading in the open market (which leads to excess trading because agents don't internalize the price impact costs they cause for others).

To illustrate the implications of these results on optimal trading strategies and equilibrium prices, we now specialize the discussion to two concrete examples: optimal execution as in [5, 1, 32] and diffusive trading targets as in [26, 35, 12].

**Example 3.9** (Optimal Liquidation). We first consider the simplest example where the clients' target is to sell a certain number of shares, that is, $\xi_t^c \equiv \xi^c < 0$, $0 \le t \le T$.[9] In this case, two elementary integrations show that

$$\bar{U}_t = \mathbf{U}_t^\Delta(\bar{\xi}) = \frac{1}{\Delta} \int_0^t k^\Delta(s,t) \int_s^T \left( k^\Delta(s,r) q_c \xi^c \right) dr\,ds$$

$$= \sqrt{\Delta} \cosh\left(\sqrt{\Delta}(T-t)\right) \int_0^t \frac{\sinh(\sqrt{\Delta}(T-s))}{\cosh^2(\sqrt{\Delta}(T-s))} ds\; q_c \xi^c$$

$$= \left(1 - \frac{\cosh(\sqrt{\Delta}(T-t))}{\cosh(\sqrt{\Delta}T)}\right) q_c \xi^c.$$

As a consequence, the optimal client position from (3.5) is

$$K_t^c = \xi^c + \frac{\rho_c}{q_c\rho_c + q_d\rho_d} \left(\bar{U}_t - q_c\xi^c\right)$$

$$= \xi^c \left(1 - \frac{q_c\rho_c}{q_c\rho_c + q_d\rho_d} \frac{\cosh(\sqrt{\Delta}(T-t))}{\cosh(\sqrt{\Delta}T)}\right).$$

---

[9]Complete liquidation as in [5, 1, 32] could be promoted using a quadratic liquidation penalty as in [10, 4] or enforced by a hard terminal constraint as in [3]. To ease notation, we do not pursue this here.
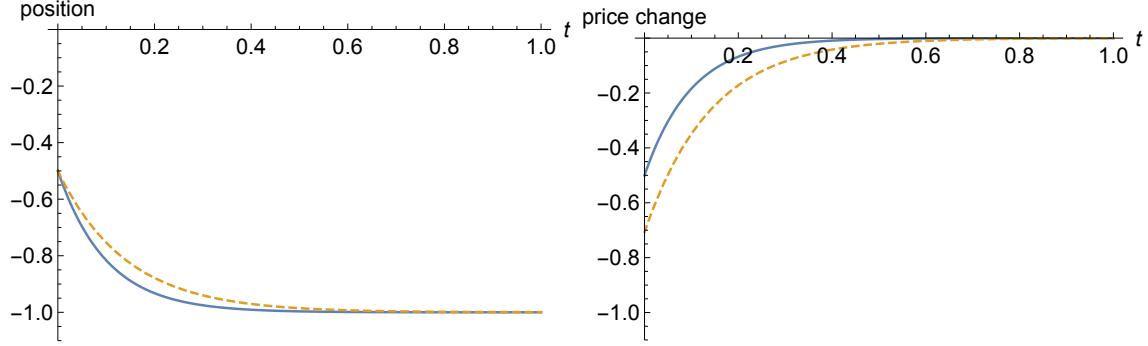
Figure 1: Optimal liquidation strategies (left panel) and price changes (right panel) for a single dealer (dotted) and infinitely many dealers (solid). Parameters are $\lambda = 0.1$, $\rho_c = \rho_d = 0.1$, $T = 1$, $\xi_c = -1$, and the dealers have half the total risk tolerance in each case.

This means that the clients use a bulk trade at time $t = 0$ to sell a fraction of their trading target equal to their share $q_c\rho_c$ of the total holding costs $q_c\rho_c + q_d\rho_d$. With an open interdealer market, they subsequently continue selling at an absolutely continuous rate, as the dealers gradually pass on their positions. Whence, the clients' optimal trading path resembles the one in the model of Obizhaeva and Wang [32] with transient price impact. The absolutely continuous trading rate is determined by the trading costs $\lambda$ in the open market, dealers' and clients' share of the total risk tolerance, and the numbers of dealers and clients through the constant $\Delta$, with low trading costs, low risk tolerances, and a large number of dealers lead to faster trading. This is illustrated in the left panel of Figure 1.

Let us now turn to the equilibrium price for the dealer market. By Theorem 3.4(iii),

$$S_t = \mathbb{E}_t[\mathscr{D}_T] - \frac{1}{\bar{\rho}}\mathbb{E}_t\left[\int_t^T (\bar{U}_s - \bar{\xi}_s)ds\right] = \mathscr{D}_t + \frac{q_c\xi^c}{q_c\rho_c + q_d\rho_d}\int_t^T \frac{\cosh(\sqrt{\Delta}(T-s))}{\cosh(\sqrt{\Delta}T)}ds$$

$$= \mathscr{D}_t + \frac{q_c\xi^c}{q_c\rho_c + q_d\rho_d}\frac{\sinh(\sqrt{\Delta}(T-t))}{\sqrt{\Delta}\cosh(\sqrt{\Delta}T)}.$$

With customers that want to liquidate a position in the dealer market ($\xi^c < 0$), the risky asset trades at a price below its expected dividend. The corresponding positive risk premium is earned by the dealers for providing liquidity to the clients as the latter liquidate their position. Like for the clients' optimal positions the transience of the price deviation from the asset's expected payoff is modulated by the constant $\Delta$. Without an open market ($\lambda = \infty$ so that $\Delta = 0$), the price impact is permanent; as the open market becomes more and more liquid, the price impact of the initial bulk trade disappears faster and faster as dealers quickly unwind their positions. Moreover, as illustrated in the right panel of Figure 1, the price impact is decreasing in the number of dealers when their proportion of the total risk tolerance in the economy is held fixed.

In this example, the volatility of the asset price in the dealer market remains unaffected because the clients' demand is deterministic. This will be different in the subsequent example with random demands.

**Example 3.10** (Diffusive Trading Targets)**.** To explore the other end of the spectrum of potential target strategies, suppose as in [26, 35, 12] that the clients have "high-frequency
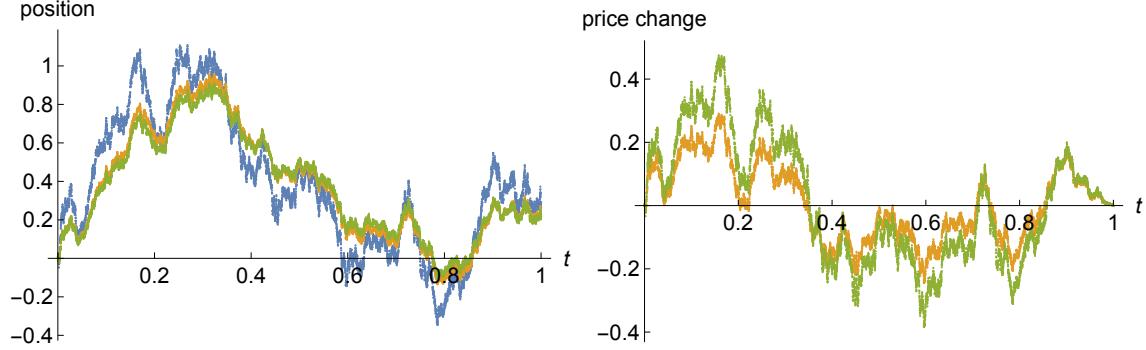
Figure 2: Left panel: simulated target position (blue) and corresponding optimal position in dealer market with infinitely many small dealers (orange) and a single dealer (green), with half the total risk tolerance in each case. Right panel: Price adjustments in equilibrium with infinitely many small dealers (orange) and a single dealer (green). Parameters are $\lambda = 0.1$, $\rho_c = \rho_d = 0.1$, $T = 1$, and $\sigma_\xi = 1$.

trading needs". This is modelled by a target position $\xi^c$ following Brownian motion with volatility $\sigma_\xi$. For such a martingale, we have

$$\bar{U}_t = \mathbf{U}_t^\Delta(\bar{\xi}) = \frac{1}{\Delta} \int_0^t k^\Delta(s,t) \int_s^T \left( k^\Delta(s,r) q_c \xi_s^c \right) dr ds$$

$$= q_c \sqrt{\Delta} \cosh\left(\sqrt{\Delta}(T-t)\right) \int_0^t \left( \frac{\sinh(\sqrt{\Delta}(T-s))}{\cosh^2(\sqrt{\Delta}(T-s))} \xi_s^c \right) ds. \qquad (3.7)$$

As a consequence, the clients' optimal position from (3.5) is a convex combination of a diffusive and a smooth component,

$$K_t^c = \frac{q_d \rho_d}{q_c \rho_c + q_d \rho_d} \xi_t^c + \frac{q_c \rho_c}{q_c \rho_c + q_d \rho_d} \sqrt{\Delta} \cosh\left(\sqrt{\Delta}(T-t)\right) \int_0^t \frac{\sinh(\sqrt{\Delta}(T-s))}{\cosh^2(\sqrt{\Delta}(T-s))} \xi_s^c ds.$$

To wit, the clients directly implement a fraction of their trading target equal to the dealers' share of the total risk tolerance. To reduce their remaining risk, they also trade in an absolutely continuous-manner. This gradually pulls the remaining deviation from their trading targets towards zero, as can be seen from the dynamics

$$dK_t^c = F^\Delta(t) \left( \xi_t^c - K_t^c \right) dt + \frac{q_d \rho_d}{q_c \rho_c + q_d \rho_d} d\xi_t^c.$$

These trading strategies are illustrated in the left panel of Figure 2, which shows that the number of dealer only has a modest impact on the clients' optimal positions relative to their target positions.

Let us now turn to the corresponding equilibrium prices in the dealer market. In view of Theorem 3.4(iii), we have

$$S_t = \mathscr{D}_t + \frac{1}{\bar{\rho}} \mathbb{E}_t \left[ \int_t^T (\bar{\xi}_s - \bar{U}_s) ds \right].$$

14

After recalling (3.7), differentiation shows that the process $\bar{\xi}_t - \bar{U}_t$ has Ornstein-Uhlenbeck-type dynamics,

$$d(\bar{\xi}_t - \bar{U}_t) = -F^\Delta(t)(\bar{\xi}_t - \bar{U}_t)dt + q_c d\xi_t^c.$$

As a consequence, the price adjustment $S_t - \mathscr{D}_t$ relative to the expected terminal dividend is

$$
\begin{aligned}
\frac{1}{\bar{\rho}}\mathbb{E}_t\left[\int_t^T (\bar{\xi}_s - \bar{U}_s)ds\right] &= \frac{(\bar{\xi}_t - \bar{U}_t)}{\bar{\rho}}\int_t^T e^{-\int_t^s F^\Delta(r)dr}ds \\
&= \frac{(\bar{\xi}_t - \bar{U}_t)}{\bar{\rho}}\int_t^T \frac{\cosh(\sqrt{\Delta}(T-s))}{\cosh(\sqrt{\Delta}(T-t))}ds = \frac{F^\Delta(t)(\bar{\xi}_t - \bar{U}_t)}{\Delta\bar{\rho}}.
\end{aligned}
$$

Differentiation in turn shows that the dynamics of the price adjustment are

$$
\begin{aligned}
d(S_t - \mathscr{D}_t) &= \frac{1}{\Delta\bar{\rho}}\left(\frac{d}{dt}F^\Delta(t) - F^\Delta(t)^2\right)(\bar{\xi}_t - \bar{U}_t)dt + \frac{F^\Delta(t)}{\Delta\bar{\rho}}q_c d\xi_t^c \\
&= -\frac{1}{\bar{\rho}}(\bar{\xi}_t - \bar{U}_t)dt + \frac{F^\Delta(t)}{\Delta\bar{\rho}}q_c d\xi_t^c \\
&= -\frac{\Delta}{F^\Delta(t)}(S_t - \mathscr{D}_t)dt + \frac{F^\Delta(t)}{\Delta\bar{\rho}}q_c d\xi_t^c.
\end{aligned}
$$

Price deviations are pulled towards zero, with fluctuations driven by the clients trading targets. Far from maturity (or in the large-liquidity limit), $F^\Delta(t) \approx \sqrt{\Delta}$ so that the price correction thus approximately has Ornstein-Uhlenbeck dynamics:

$$d(S_t - \mathscr{D}_t) \approx -\sqrt{\Delta}(S_t - \mathscr{D}_t)dt + \frac{1}{\sqrt{\Delta}}\frac{q_c}{q_c\rho_c + q_d\rho_d}d\xi_t^c$$

Accordingly, fluctuations in clients' demands have a smaller initial effect and decay faster if $\Delta$ is large for a liquid open market (small $\lambda$) as well as for many or for more risk tolerant agents. Close to maturity, price impact tends to zero and the strength $\Delta/F^\Delta(t)$ with which prices are pulled towards fundamentals explodes, so that the equilibrium price in the dealer market approaches the exogenous terminal payoff of the risky asset as illustrated in the right panel of Figure 2.

With randomly fluctuating client demand, not just the expected return but also the volatility of the equilibrium price can change relative to the expected dividend. The above calculations show that the magnitude of this effect is governed by $\approx \sigma_\xi q_c/\sqrt{\Delta}\bar{\rho}$; its sign is in turn determined by the correlation between the expected dividend and the clients' demand. If this correlation is positive, i.e., clients' demand tends to increase when expected fundamentals rise, then illiquidity (caused by an illiquid open market and few or risk-averse dealers) increases volatility. If the correlation is negative, the sign is reversed. The interpretation is that demand pressure and fundamental shocks partially offset in the second case, whereas they magnify price fluctuations in the first case. Similar comparative statics also appear in a Radner equilibrium for a market with an exogenous quadratic deadweight costs on trading [21].

## 3.3 The Effects of Segmentation

We now discuss the effects of segmentation between the dealer and open markets. To this end, we consider how the clients' optimal positions and welfare change when they gain access to the open market. Compared to the discussion in the previous section, this means that their individual trading costs in the open market now are finite; to ease notation, we focus on the case where they vanish just like for the dealers ($\lambda^a = 0$). Then, we have $\bar{\eta} = (m_c + m_d)M/\lambda$ and the corresponding parameter determining trading and equilibrium prices in the integrated market is larger than for the segmented market considered in Section 3.2,

$$\Delta^{\text{int}} = \frac{1}{(q_c\rho_c + q_d\rho_d)\lambda(1 + \frac{1}{(m_c+m_d)M})} > \frac{1}{(q_c\rho_c + q_d\rho_d)\lambda(1 + \frac{1}{m_dM})} = \Delta.$$

Let us now discuss what this implies for our concrete examples. For optimal liquidation, Theorem 3.4(i) and (iii) show that the total optimal position the clients take in the dealer and in the open market is

$$K_t^{c,\text{int}} = \frac{q_d\rho_d}{q_c\rho_c + q_d\rho_d}\xi^c + q_d\frac{\rho_c - \rho_d}{q_c\rho_c + q_d\rho_d}\bar{U}_t^{\text{int}},$$

where

$$\bar{U}_t^{\text{int}} = \left(1 - \frac{\cosh(\sqrt{\Delta^{\text{int}}}(T - t))}{\cosh(\sqrt{\Delta^{\text{int}}}T)}\right)q_c\xi^c.$$

Hence, the initial block trade remains unchanged compared to the segmented market. However, the clients subsequently build up larger positions in the integrated market with parameter $\Delta^{\text{int}}$ rather than $\Delta$. Similarly, it can be shown that also clients with diffusive trading targets share the same fraction of their Brownian shocks with the dealers, but increase the order flow to the open market if they have direct access themselves. In the corresponding equilibrium price dynamics, the initial price impact in the liquidation model disappears faster; for diffusive trading targets, price impact is both smaller and disappears faster.

Finally, let us discuss the effect of segmentation on the clients' welfare as measured by their goal functionals' (2.4). In view of Theorem 3.4(i,ii,iii), we have

$$J^{c,\text{int}} = \mathbb{E}\left[\int_0^T \left(\frac{1}{\bar{\rho}}(\bar{U}_t^{\text{int}} - \bar{\xi}_t)K_t^{c,\text{int}} - \frac{1}{2\rho_c}(\xi_t^c - \bar{U}_t^{\text{int}} - K_t^{c,\text{int}})^2 - \lambda\left(\frac{d}{dt}\bar{U}_t^{\text{int}}\right)^2\right)dt\right].$$

Similarly, we can also compute the value of the goal functional in the case the clients do *not* have access to the open market,

$$J^c = \mathbb{E}\left[\int_0^T \left(\frac{1}{\bar{\rho}}(\bar{U}_t - \bar{\xi}_t)K_t^c - \frac{1}{2\rho_c}(K_t^c - \xi_t^c)^2\right)dt\right].$$

For the optimal liquidation example, we have

$$J^{c,\text{int}} = -(\xi^c)^2\int_0^T \frac{q_cq_d}{\bar{\rho}}\frac{\cosh(\sqrt{\Delta^{\text{int}}}(T - t))}{\cosh(\sqrt{\Delta^{\text{int}}}T)}\left(1 + \frac{q_c(\rho_d - \rho_c)}{\bar{\rho}}\frac{\cosh(\sqrt{\Delta^{\text{int}}}(T - t))}{\cosh(\sqrt{\Delta^{\text{int}}}T)}\right)$$

$$+ \frac{q_c^2\rho_c}{2\bar{\rho}^2}\left(\frac{\cosh(\sqrt{\Delta^{\text{int}}}(T - t))}{\cosh(\sqrt{\Delta^{\text{int}}}T)}\right)^2 + \lambda q_c^2\Delta^{\text{int}}\left(\frac{\sinh(\sqrt{\Delta^{\text{int}}}(T - t))}{\cosh(\sqrt{\Delta^{\text{int}}}T)}\right)^2 dt.$$
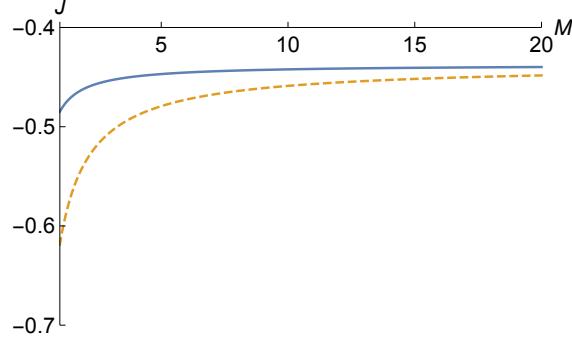
16

Figure 3: Clients' welfare with segmentation (dashed) and in the integrated market (solid) with $M$ dealers and $M$ clients ($q_c = q_d = 1/2$). Other parameters are $\lambda = 0.1$, $\rho_c = \rho_d = 0.1$, $T = 1$, and $\xi^c = -1$.

An elementary integration in turn gives

$$J^{c,\text{int}} = -\frac{(\xi^c)^2 q_c}{\bar{\rho}} \left( \frac{\tanh(\sqrt{\Delta^{\text{int}}}T)}{\sqrt{\Delta^{\text{int}}}} q_d + \left( \bar{\rho} - \frac{\rho_c}{2} \right) \frac{q_c}{4\bar{\rho}\sqrt{\Delta^{\text{int}}}} \frac{\sinh(2\sqrt{\Delta^{\text{int}}}T) + 2\sqrt{\Delta^{\text{int}}}T}{\cosh^2(\sqrt{\Delta^{\text{int}}}T)} \right)$$
$$- (\xi^c)^2 \frac{\lambda q_c^2 \sqrt{\Delta^{\text{int}}}}{4} \frac{\sinh(2\sqrt{\Delta^{\text{int}}}T) - 2\sqrt{\Delta^{\text{int}}}T}{\cosh^2(\sqrt{\Delta^{\text{int}}}T)}$$

Similarly,

$$J^c = -\frac{(\xi^c)^2 q_c}{\bar{\rho}} \int_0^T \frac{\cosh(\sqrt{\Delta}(T-t))}{\cosh(\sqrt{\Delta}T)} \left( 1 - \frac{q_c \rho_c}{2\bar{\rho}} \frac{\cosh(\sqrt{\Delta}(T-t))}{\cosh(\sqrt{\Delta}T)} \right) dt$$
$$= -\frac{(\xi^c)^2 q_c}{\bar{\rho}} \left( \frac{\tanh(\sqrt{\Delta}T)}{\sqrt{\Delta}} - \frac{q_c \rho_c}{2\bar{\rho}} \frac{\sinh(2\sqrt{\Delta}T) + 2\sqrt{\Delta}T}{4\sqrt{\Delta}\cosh^2(\sqrt{\Delta}T)} \right).$$

Now suppose the number clients and dealers is large ($M \to \infty$), whereas the proportions of each group in the total population remain fixed (fixed $m_c, m_d$ and in turn $q_c, q_d, \bar{\rho}$). Then, $\Delta, \Delta^{\text{int}} \to \Delta_\infty = 1/\bar{\rho}\lambda$ and in turn

$$\frac{J^c}{J^{c,\text{int}}} \to \frac{\tanh(\sqrt{\Delta_\infty}T) - \frac{q_c \rho_c}{2\bar{\rho}} \frac{\sinh(2\sqrt{\Delta_\infty}T) + 2\sqrt{\Delta_\infty}T}{4\cosh^2(\sqrt{\Delta_\infty}T)}}{\tanh(\sqrt{\Delta_\infty}T)q_d + \left( \bar{\rho} - \frac{\rho_c}{2} \right) \frac{q_c}{4\bar{\rho}} \frac{\sinh(2\sqrt{\Delta_\infty}T) + 2\sqrt{\Delta_\infty}T}{\cosh^2(\sqrt{\Delta_\infty}T)} + \frac{q_c}{4} \frac{\sinh(2\sqrt{\Delta_\infty}T) - 2\sqrt{\Delta_\infty}T}{\cosh^2(\sqrt{\Delta_\infty}T)}}$$
$$= \frac{\tanh(\sqrt{\Delta_\infty}T) - \frac{q_c \rho_c}{2\bar{\rho}} \frac{\sinh(2\sqrt{\Delta_\infty}T) + 2\sqrt{\Delta_\infty}T}{4\cosh^2(\sqrt{\Delta_\infty}T)}}{\tanh(\sqrt{\Delta_\infty}T)q_d - \frac{q_c \rho_c}{2\bar{\rho}} \frac{\sinh(2\sqrt{\Delta_\infty}T) + 2\sqrt{\Delta_\infty}T}{4\cosh^2(\sqrt{\Delta_\infty}T)} + \frac{q_c}{2} \frac{\sinh(2\sqrt{\Delta_\infty}T)}{\cosh^2(\sqrt{\Delta_\infty}T)}} = 1, \quad \text{as } M \to \infty.$$

(Here, we have used $\sinh(2x) = 2\sinh(x)\cosh(x)$, $\tanh(x) = \sinh(x)/\cosh(x)$, and $q_c + q_d = 1$ in the last step.) This shows that, in the competitive limit of many small clients and dealers, the welfare loss due to segmentation disappears, as illustrated in Figure 3. Maybe surprisingly, this suggests that market segmentation can be close to optimal as long as none of the agents has substantial market power.

17

# A Proof of Theorem 3.4

*Proof of Theorem 3.4.* Let $S \in \mathcal{H}^2$, $(K^a, u^a) \in \mathcal{S}^2 \times \mathcal{L}^2$, $a \in \mathcal{A}$, be *any* equilibrium. It is easy to see that each agent's target functional $J^a(K, u; \bar{u}^{-a}, S)$ is strictly concave in $(K, u) \in \mathcal{S}^2 \times \mathcal{L}^2$, and so its maximizer $(K^a, u^a)$ is uniquely determined by the first-order condition that all directional derivatives vanish. In particular, for any $K \in \mathcal{S}^2$ we must have

$$0 = \nabla_K J^a(K^a, u^a; \bar{u}^{-a}, S) = \mathbb{E}\left[\int_0^T K_t dA_t + \int_0^T \frac{1}{\rho^a}(\xi_t^a - K_t^a - U_t^a)K_t dt\right].$$

This can only hold true if $A$ is absolutely continuous with density

$$\frac{dA_t}{dt} = \mu_t = \frac{1}{\rho^a}\left(U_t^a + K_t^a - \xi_t^a\right), \quad 0 \le t \le T, \quad a \in \mathcal{A}. \tag{A.1}$$

Solving for $K^a$, we find

$$K^a = \xi^a - U^a + \rho^a \mu. \tag{A.2}$$

Hence, after aggregating over $a \in \mathcal{A}$, the market clearing condition $K^N + \sum_{a \in \mathcal{A}} m(a)K^a = 0$ implies $-K^N = \bar{\xi} - \bar{U} + \bar{\rho}\mu$ and in turn (3.4).

The first-order conditions for individual optimality also require that, for any $u \in \mathcal{L}^2$,

$$\begin{aligned}
0 =& \nabla_u J(K^a, u^a; \bar{u}^{-a}, S) \\
=& -\mathbb{E}\left[\int_0^T u_t \left(\lambda \bar{u}_t^{-a} + (2m(a)\lambda + \lambda^a)u_t^a + \mathbb{E}_t\left[\int_t^T \frac{1}{\rho^a}(K_s^a + U_s^a - \xi_s^a)\, ds\right]\right) dt\right].
\end{aligned}$$

Since this equality needs to hold for any perturbation $u \in \mathcal{L}^2$, it is equivalent to

$$\lambda \bar{u}_t^{-a} + (2m(a)\lambda + \lambda^a)u_t^a + \mathbb{E}_t\left[\int_t^T \frac{1}{\rho^a}(K_s^a + U_s^a - \xi_s^a)\, ds\right] = 0, \quad 0 \le t \le T.$$

Recalling that $\eta^a = 1/(m(a)\lambda + \lambda^a)$, the definition of $\bar{u}$ and (A.1), this implies

$$u_t^a = -\eta^a\left(\lambda \bar{u}_t + \mathbb{E}_t\left[\int_t^T \mu_s ds\right]\right), \quad 0 \le t \le T. \tag{A.3}$$

Now, we aggregate over $a \in \mathcal{A}$ to obtain

$$\bar{u}_t = -\bar{\eta}\left(\lambda \bar{u}_t + \mathbb{E}_t\left[\int_t^T \mu_s ds\right]\right), \quad 0 \le t \le T.$$

Solving for $\bar{u}$ and using the already established relation (3.4) along with $\Delta = \bar{\eta}/(\bar{\rho}(1 + \bar{\eta}\lambda))$, we find

$$\bar{u}_t = -\Delta \mathbb{E}_t\left[\int_t^T (\bar{U}_s - K_s^N - \bar{\xi}_s)\, ds\right], \quad 0 \le t \le T. \tag{A.4}$$

Therefore, the pair $(u, U) := (\bar{u}, \bar{U})$ solves the linear FBSDE (3.1) for $X = K^N + \bar{\xi}$. As this FBSDE is uniquely solved by $(\mathbf{u}^\Delta(X), \mathbf{U}^\Delta(X))$ from (3.2), this pins down $\bar{U}$ as stated in (3.3) with density $\bar{u}$ as given in (A.4). From (A.3), we infer that each individual agent's strategy $u^a$, $a \in \mathcal{A}$, is the multiple $\eta^a = 1/(m(a)\lambda + \lambda^a)$ of the same universal process. Since $\bar{u} = \sum_{a \in \mathcal{A}} m(a)u^a$, this implies $u^a = \frac{\eta^a}{\bar{\eta}}\bar{u}$ and, thus, $U^a = \frac{\eta^a}{\bar{\eta}}\bar{U}$ as claimed. Moreover, in light of (A.2) and (3.4), our observation that $U^a = \frac{\eta^a}{\bar{\eta}}\bar{U}$ allows us to write each agent's position in the dealer market in the form (3.5). As a consequence, the only candidate for an equilibrium is the one described in the present theorem.

We now show that strategies and price pinpointed above form indeed an equilibrium. For this purpose, define $\bar{U}$ by (3.3) and $\bar{u} = \frac{d}{dt}\bar{U}$ so that $(\bar{u}, \bar{U})$ solves the FBSDE (3.1). For each $a \in \mathcal{A}$ the candidate equilibrium position is then $K^a$ given by (3.5) and the candidate equilibrium trading rate in the open market is $u^a = \frac{\eta_a}{\bar{\eta}}\bar{u}$. The candidate equilibrium price is $S$ with risk premium $\mu$ given by (3.4). It is readily checked that the positions $K^a$, $a \in \mathcal{A}$, ensure market clearing. Moreover, the first order condition in $K^a$ is

$$0 = \mathbb{E}\left[\int_0^T K_t\mu_t dt + \int_0^T \frac{1}{\rho^a}(\xi_t^a - K_t^a - U_t^a)K_t dt\right], \quad \text{for all } K \in \mathcal{S}^2,$$

which is clearly satisfied by the candidate strategy $K^a$ by definition of the corresponding drift rate $\mu$. Using the first-order condition satisfied by $K^a$, the first-order condition for $u^a$ then amounts to (A.3) by the same reasoning as above. The latter condition is readily verified since $(\bar{u}, \bar{U})$ were chosen as the solution to the FBSDE (3.1). Concavity of the goal functional ensures sufficiency of the first-order conditions, and so $(K^a, u^a)$ maximizes $J^a$ as required for an equilibrium. $\square$

# B  Proofs for Section 3.1

By the definition of the equilibrium, the optimization criterion of each dealer is

$$J^a(K^a, u^a; \bar{u}, S) := \mathbb{E}\left[\int_0^T K_t^a dA_t - \int_0^T \lambda\bar{u}_t u_t^a - \frac{1}{2\rho_d}(K_t^a + U_t^a)^2 dt\right].$$

Similarly as in the Proof of Theorem 3.4, after aggregating over all the agents, the first-order optimality condition for $u^a$ is

$$\bar{u}_t = -\frac{M}{\lambda\rho_d(M+1)}\mathbb{E}_t\left[\int_t^T \left(\bar{U}_s - K_s^N\right) ds\right], \quad 0 \le t \le T.$$

Therefore $\bar{u}$ is an optimizer of the auxiliary minimization problem

$$\min_u \mathbb{E}\left[\int_0^T \lambda u_t^2 + \frac{1}{2}\frac{M}{\rho_d(M+1)}\left(K_t^N - \int_0^t u_s ds\right)^2 dt\right]. \tag{B.1}$$

*Proof of Proposition 3.5.* Set $M_t = \mathbb{E}_t\left[\int_0^T \mu_s ds\right]$ so that $d(\mathscr{D}_t - S_t) = dM_t - \mu_t dt$. By the $\varepsilon$-Young inequality, the Cauchy-Schwarz inequality and Doob's maximal inequality, for all

predictable $H$ bounded by 1, we have

$$\mathbb{E}\left[\sup_{0\le t\le T}\left|\int_0^t H_s d(\mathscr{D}_s - S_s)\right|^2\right] \le 8\mathbb{E}\left[M_T^2\right] + 2T\mathbb{E}\left[\int_0^T |\mu_s|^2 ds\right].$$

Together with the equality $M_T = \int_0^T \mu_s ds$ and a second application of the Cauchy-Schwarz inequality, it follows that

$$\mathbb{E}\left[\sup_{0\le t\le T}\left|\int_0^t H_s d(\mathscr{D}_s - S_s)\right|^2\right] \le 8\mathbb{E}\left[\left(\int_0^T \mu_s ds\right)^2\right] + 2T\mathbb{E}\left[\int_0^T |\mu_s|^2 ds\right]$$

$$\le 10T\mathbb{E}\left[\int_0^T |\mu_s|^2 ds\right] \le \frac{10T}{\bar{\rho}^2}\mathbb{E}\left[\int_0^T \left(K_s^N - \bar{U}_s\right)^2 ds\right].$$

To prove the proposition it is therefore sufficient to show that the last term converges to 0 as $\lambda \to 0$. The set $\{U = \int_0^\cdot u_s ds : u \in \mathscr{L}^2\}$ is dense in $\mathscr{L}^2$, so that there exists a sequence $u^n \in \mathscr{L}^2$ such that

$$\mathbb{E}\left[\int_0^T \left(\int_0^t u_s^n ds - K_t^N\right)^2 dt\right] \le \frac{1}{n}, \quad n = 1, 2, \dots$$

Due to the minimality condition (B.1), we have

$$\mathbb{E}\left[\int_0^T \frac{1}{2}\frac{M}{\rho_d(M+1)}(\bar{U}_t - K_t^N)^2 dt\right] \le \mathbb{E}\left[\int_0^T \left(\frac{\lambda}{2}(u_t^n)^2 + \frac{1}{2}\frac{M}{\rho_d(M+1)}\left(\int_0^t u_s^n ds - K_t^N\right)^2\right) dt\right]$$

$$\le \frac{\lambda}{2}\mathbb{E}\left[\int_0^T (u_t^n)^2 dt\right] + \frac{1}{2\rho_d n}.$$

Thus, $\mathbb{E}[\int_0^T (\bar{U}_t - K_t^N)^2 dt] \to 0$ as $\lambda \to 0$, verifying the first convergence asserted in Proposition 3.5.

To also establish the second convergence result, we apply the inequalities of Burkholder-Davis-Gundy and Hölder to obtain

$$\mathbb{E}\left[\left|\int_0^T K_t^N(dS_t - d\mathscr{D}_t)\right|\right] \le \mathbb{E}\left[\sup_{0\le t\le T}\left|\int_0^t K_s^N dM_s\right|\right] + \mathbb{E}\left[\sup_{0\le t\le T}\left|\int_0^t K_s^N \mu_s ds\right|\right]$$

$$\le C\left(\mathbb{E}\left[\left(\int_0^T (K^N)_s^2 d\langle M\rangle_s\right)^{1/2}\right] + \mathbb{E}\left[\int_0^T |K_s^N||\bar{U}_s - K_s^N| ds\right]\right)$$

$$\le C\left(\mathbb{E}\left[\sup_{s\in[0,T]} |K_s^N|\left(\langle M\rangle_T^{1/2} + \int_0^T |\bar{U}_s - K_s^N| ds\right)\right]\right)$$

$$\le C\mathbb{E}\left[\sup_{s\in[0,T]} |K_s^N|^2\right]^{1/2}\mathbb{E}\left[\int_0^T |\bar{U}_s - K_s^N|^2 ds\right]^{1/2}.$$

Here, $C > 0$ is a constant that might change from line to line but does not depend on $\lambda$. $L^1$ convergence now follows, since we have already verified above that the last term converges to 0 as $\lambda \to 0$. $\qquad\square$

*Proof of Lemma 3.6.* By (3.4), (A.4), and the integration by parts formula (using $K_0^N = \bar{u}_T$), we have

$$\int_0^T K_t^N d\mathscr{D}_t - \int_0^T K_t^N dS_t = -\lambda \frac{M+1}{M} \int_0^T K_t^N d\bar{u}_t = \lambda \frac{M+1}{M} \int_0^T \mu_t^N \bar{u}_t dt. \qquad \text{(B.2)}$$

To establish (3.6), it therefore suffices to show that

$$\int_0^T |\mu_t^N - \bar{u}_t| dt = o(1) \quad \text{in } L^1 \text{ as } \lambda \to 0. \qquad \text{(B.3)}$$

Integration by parts and (3.2) give

$$\begin{aligned}
\mu_t^N - \bar{u}_t &= \mu_t^N + F^\Delta(t)\bar{U}_t - \frac{\Delta}{\cosh(\sqrt{\Delta}(T-t))} \mathbb{E}_t\left[\int_t^T \cosh(\sqrt{\Delta}(T-s)) K_s^N ds\right] \\
&= \mu_t^N - F^\Delta(t)(K_t^N - \bar{U}_t) - \sqrt{\Delta}\mathbb{E}_t\left[\int_t^T \frac{\sinh(\sqrt{\Delta}(T-s))}{\cosh(\sqrt{\Delta}(T-t))} \mu_s^N ds\right] \\
&= \mu_t^N \left(1 - \sqrt{\Delta}\int_t^T \frac{\sinh(\sqrt{\Delta}(T-s))}{\cosh(\sqrt{\Delta}(T-t))} ds\right) \\
&\quad - F^\Delta(t)(K_t^N - \bar{U}_t) + \sqrt{\Delta}\mathbb{E}_t\left[\int_t^T \frac{\sinh(\sqrt{\Delta}(T-s))}{\cosh(\sqrt{\Delta}(T-t))} (\mu_t^N - \mu_s^N) ds\right]. \qquad \text{(B.4)}
\end{aligned}$$

Now, note that

$$\frac{\sqrt{\Delta} \int_t^T \sinh(\sqrt{\Delta}(T-s)) ds}{\cosh(\sqrt{\Delta}(T-t))} = 1 - \frac{1}{\cosh(\sqrt{\Delta}(T-t))}.$$

Together with (B.4), it follows that $K^N - \bar{U}$ satisfies the linear ODE

$$\frac{d(K_t^N - \bar{U}_t)}{dt} = \mu_t^N - \bar{u}_t = -F^\Delta(t)(K_t^N - \bar{U}_t) + w^\Delta(t) + \frac{\mu_t^N}{\cosh(\sqrt{\Delta}(T-t))}, \qquad \text{(B.5)}$$

where

$$w^\Delta(t) = \mathbb{E}_t\left[\int_t^T \frac{\sqrt{\Delta}\sinh(\sqrt{\Delta}(T-s))}{\cosh(\sqrt{\Delta}(T-t))} (\mu_t^N - \mu_s^N) ds\right].$$

Since $K_0^N = \bar{U}_0 = 0$ and by definition of the function $F^\Delta$ the explicit solution of (B.5) is

$$\begin{aligned}
K_t^N - \bar{U}_t &= \int_0^t e^{-\int_s^t F^\Delta(u)du}\left(w_s^\Delta + \frac{\mu_s^N}{\cosh(\sqrt{\Delta}(T-s))}\right) ds \\
&= \int_0^t \frac{\cosh(\sqrt{\Delta}(T-t))}{\cosh(\sqrt{\Delta}(T-s))}\left(w_s^\Delta + \frac{\mu_s^N}{\cosh(\sqrt{\Delta}(T-s))}\right) ds.
\end{aligned}$$

Together with (B.5), it follows that

$$\int_0^T |\mu_t^N - \bar{u}_t| dt \leq \sup_{u \in [0,T]} |w_u^\Delta| \int_0^T \left( 1 + \sqrt{\Delta} \int_0^t \frac{\sinh(\sqrt{\Delta}(T-t))}{\cosh(\sqrt{\Delta}(T-s))} ds \right) dt$$

$$+ \sup_{u \in [0,T]} |\mu_u^N| \int_0^T \left( \frac{1}{\cosh(\sqrt{\Delta}(T-t))} + \sqrt{\Delta} \int_0^t \frac{\sinh(\sqrt{\Delta}(T-t))}{\cosh^2(\sqrt{\Delta}(T-s))} ds \right) dt$$

$$\leq \sup_{u \in [0,T]} |w_u^\Delta| \int_0^T \left( 1 + \sqrt{\Delta} \int_0^t e^{-\sqrt{\Delta}(t-s)} ds \right) dt$$

$$+ \sup_{u \in [0,T]} |\mu_u^N| \left( \int_0^T 2 e^{-\sqrt{\Delta} t} dt + \int_0^T \sqrt{\Delta} \int_s^T \frac{\sinh(\sqrt{\Delta}(T-t))}{\cosh^2(\sqrt{\Delta}(T-s))} dt ds \right)$$

$$\leq 2T \sup_{u \in [0,T]} |w_u^\Delta| + \frac{4}{\sqrt{\Delta}} \sup_{u \in [0,T]} |\mu_u^N|. \tag{B.6}$$

Write

$$\omega(\delta) = \sup_{t,s \in [0,T], |t-s| \leq \delta} \left| \mu^N(s) - \mu^N(t) \right|$$

for the modulus of continuity of $\mu^N$. Since $t \mapsto \mu_t^N$ is continuous on the compact set $[0,T]$,

$$\omega(\delta) \to 0, \quad \text{a.s. as } \delta \to 0. \tag{B.7}$$

The definitions of $w^\Delta$ and $\omega$ and a change of variables yield the following estimate:

$$|w^\Delta(t)| \leq \frac{1}{2} \mathbb{E}_t \left[ \int_t^T \sqrt{\Delta} e^{-\sqrt{\Delta}(s-t)} \omega(s-t) ds \right] \leq \frac{1}{2} \mathbb{E}_t \left[ \int_0^\infty e^{-u} \omega \left( \frac{u}{\sqrt{\Delta}} \right) du \right] := M_t^\Delta. \tag{B.8}$$

Here, $(M_t^\Delta)_{t \in [0,T]}$ is a martingale for each $\Delta > 0$ since $|\omega(\delta)| \leq 2 \sup_{s \in [0,T]} |\mu_s^N|$ is integrable by assumption. Also note that by definition of the modulus of continuity $\omega$, the mapping $\Delta \mapsto M_t^\Delta$ is decreasing for each $t$. Define

$$M^* := \lim_{\Delta \to \infty} \sup_{t \in [0,T]} M_t^\Delta \geq 0.$$

Fix $\epsilon > 0$. Then, by the monotonicity in $\Delta$, we have

$$\mathbb{P}\left[ M^* \geq \epsilon \right] \leq \lim_{\Delta \to \infty} \mathbb{P} \left[ \sup_{t \in [0,T]} M_t^\Delta \geq \epsilon \right] \leq \lim_{\Delta \to \infty} \frac{\mathbb{E}[M_T^\Delta]}{\epsilon} = 0.$$

Here, the last equality is a consequence of (B.7), another application of the monotone convergence theorem and the integrability of right-hand side in (B.8). As a result,

$$0 \leq \limsup_{\Delta \to \infty} \sup_{t \in [0,T]} |w^\Delta(t)| \leq M^* = 0 \text{ a.s.}$$

In view of (B.6), it follows that the asserted convergence (B.3) holds in the almost-sure sense.

To show convergence in $L^1$, it therefore suffices to establish uniform integrability of (B.2). By (B.6) and (B.8),

$$\int_0^T |(\mu_t^N)^2 - \bar{u}_t \mu_t^N| dt \leq \sup_{s \in [0,T]} |\mu_s^N| \int_0^T |\mu_t^N - \bar{u}_t| dt \leq T \sup_{s \in [0,T]} |w_s^\Delta|^2 + (T + \frac{4}{\sqrt{\Delta}}) \sup_{s \in [0,T]} |\mu_s^N|^2$$

$$\leq T \sup_{s \in [0,T]} |M_s^\Delta|^2 + (T + \frac{4}{\sqrt{\Delta}}) \sup_{s \in [0,T]} |\mu_s^N|^2.$$

Observe that the right-hand side is decreasing in $\Delta$, and integrable for, e.g., $\Delta = 1$ by Doob's maximal inequality:

$$\mathbb{E}\left[T \sup_{t \in [0,T]} |M_t^1|^2 + (T + 4) \sup_{t \in [0,T]} |\mu_t^N|^2\right] \leq \mathbb{E}\left[2T |M_T^1|^2 + (T + 4) \sup_{t \in [0,T]} |\mu_t^N|^2\right] < \infty.$$

Therefore, the family, $\{\int_0^T |(\mu_t^N)^2 - u_t \mu_t^N| ds, \Delta \geq 1\}$ is uniformly integrable. Since

$$\int_0^T |(\mu_t^N)^2 - u_t \mu_t^N| ds \leq \sup_{s \in [0,T]} |\mu_s^N| \int_0^T |\mu_s^N - u_s| ds,$$

this implies that the almost sure convergence we have established for (B.3) also holds in $L^1$. □

*Proof of Lemma 3.7.* Similarly to (B.2), the liquidity costs of the clients can be written as

$$\int_0^T K_t^N d\mathscr{D}_t - \int_0^T K_t^N dS_t = -\lambda \frac{M+1}{M} \int_0^T K_t^N d\bar{u}_t$$

$$= \lambda \frac{M+1}{M} \langle \bar{u}, K^N \rangle_T + \lambda \frac{M+1}{M} \int_0^T \bar{u}_t (\mu_t^N dt + \sigma_t^N dW_t), \text{(B.9)}$$

and $\bar{u}_t = \bar{K}_t^N - F^\Delta(t)\bar{U}_t$, where $\bar{K}^N$ is defined as $\bar{K}_t^N = \mathbb{E}_t\left[\int_t^T k^\Delta(t,s) K_s^N ds\right]$. This implies that the covariation of $\bar{u}$ is the same as the one of $\bar{K}^N$. Note that by definition of $\bar{K}^N$,

$$\cosh(\sqrt{\Delta}(T-t))\bar{K}_t^N - \Delta \int_0^t \cosh(\sqrt{\Delta}(T-s)) K_s^N ds = \Delta \mathbb{E}_t\left[\int_0^T \cosh(\sqrt{\Delta}(T-s)) K_s^N ds\right]$$

$$\text{(B.10)}$$

is a square-integrable martingale. By the martingale representation theorem, it therefore can be written as a stochastic integral with respect to the Brownian motion generating the underlying filtration. The integrand in this representation can be computed using the Clark-Ocone formula. Indeed, setting

$$\Phi = \Delta \int_0^T \cosh(\sqrt{\Delta}(T-s)) K_s^N ds$$

and using the Malliavin differentiability of $\Phi$ that we prove below, the Clark-Ocone formula [30, Proposition 1.3.14] yields

$$\Phi = \mathbb{E}[\Phi] + \int_0^T \mathbb{E}_t[D_t \Phi] dW_t.$$

23

By inserting this into (B.10) and integrating by parts, we in turn obtain

$$\langle \bar{u}, K^N \rangle_T = \int_0^T \frac{\mathbb{E}_t \left[ D_t \Phi \right]}{\cosh(\sqrt{\Delta}(T-t))} \sigma_t^N dt. \tag{B.11}$$

We now show that we indeed have $\Phi \in \mathbb{D}^{1,2}$, so that the Clark-Ocone formula can be applied. Given our assumption on the square-integrability of the supremum of its Malliavin derivative, $K^N \in \mathbb{L}^{1,2,f}$, cf. [30, p. 45]. Thus, by [30, p. 45], $\Phi$ is Malliavin differentiable and it follows from the product rule that

$$D_t \Phi = \Delta \int_t^T \cosh(\sqrt{\Delta}(T-s)) D_t K_s^N ds. \tag{B.12}$$

We now expand $D_t \Phi$ and $\mathbb{E}_t[D_t \Phi]$ for $\lambda \to 0$ or, equivalently, $\Delta \to \infty$. First note that by (B.12) and the definition of the hyperbolic cosine,

$$\frac{\sqrt{\Delta}^{-1} D_t \Phi}{\cosh(\sqrt{\Delta}(T-t))} = \sqrt{\Delta} \int_t^T \frac{e^{\sqrt{\Delta}(t-s)} + e^{-\sqrt{\Delta}(2T-(s+t))}}{1 + e^{-2\sqrt{\Delta}(T-t)}} D_t K_s^N ds. \tag{B.13}$$

By continuity of $s \mapsto D_t(K_s^N)$ on $[t, T]$, some elementary integrations show that the above expression converges to $D_t K_t^N$ as $\Delta \to \infty$. In view of [30, Proposition 1.3.8], we have $D_t(\int_0^t \sigma_s^N dW_s) = \sigma_t^N$. Moreover, $D_t(\int_0^t \mu_s^N ds) = 0$, so that

$$\frac{\sqrt{\Delta}^{-1} D_t \Phi}{\cosh(\sqrt{\Delta}(T-t))} \to \sigma_t^N, \quad P\text{-a.s. as } \Delta \to \infty.$$

Next, observe that for every $t \in [0, T]$, it follows from (B.13) that

$$\sup_{\Delta > 1} \left| \frac{\sqrt{\Delta}^{-1} D_t \Phi}{\cosh(\sqrt{\Delta}(T-t))} \right| \le 2 \sup_{s \in [t,T]} |D_t K_s^N| \sup_{\Delta > 1} \left\{ \sqrt{\Delta} \int_t^T e^{\sqrt{\Delta}(t-s)} ds \right\}$$

$$\le 2 \sup_{t \le s \le T} |D_t K_s^N|. \tag{B.14}$$

Since the right-hand side is integrable by assumption, the dominated convergence theorem in turn shows

$$\mathbb{E}_t \left[ \frac{\sqrt{\Delta}^{-1} D_t \Phi}{\cosh(\sqrt{\Delta}(T-t))} \right] \to \sigma_t^N, \quad dP \times dt\text{-a.s. as } \Delta \to \infty.$$

We now show that this expansion of $D_t \Phi$ is inherited by its conditional expectation and in turn the covariation (B.11). To this end, we first use (B.14) and Young's inequality to obtain that

$$\sup_{\Delta > 1} |\sigma_t^N| \left| \mathbb{E}_t \left[ \frac{\sqrt{\Delta}^{-1} D_t \Phi}{\cosh(\sqrt{\Delta}(T-t))} \right] \right| \le \frac{1}{3} \sup_{t \in [0,T]} |\sigma_t^N|^3 + \frac{2}{3} \sup_{\Delta > 1} \left| \mathbb{E}_t \left[ \frac{\sqrt{\Delta}^{-1} D_t \Phi}{\cosh(\sqrt{\Delta}(T-t))} \right] \right|^{3/2}$$

$$\le \frac{1}{3} \sup_{t \in [0,T]} |\sigma_t^N|^3 + \frac{\sqrt{32}}{3} \mathbb{E}_t \left[ \sup_{t \le s \le T} |D_t K_s^N|^{3/2} \right].$$

24

Jensen's inequality and the integrability assumption for the supremum of the Malliavin derivative of $K^N$ yield

$$\mathbb{E}\left[\int_0^T \mathbb{E}_t\left[\sup_{t\leq s\leq T}|D_t K_s^N|^{3/2}\right]^{4/3} dt\right] \leq \mathbb{E}\left[\int_0^T \sup_{t\leq s\leq T}|D_t K_s^N|^2 dt\right] < \infty.$$

Moreover, $(\sup_{t\in[0,T]}|\sigma_t^N|^3)^{4/3} = \sup_{t\in[0,T]}|\sigma_t^N|^4$ is also integrable by assumption. Together, these two estimates show that

$$\mathbb{E}\left[\int_0^T \left(\sup_{\Delta>1}|\sigma_t^N|\left|\mathbb{E}_t\left[\frac{\sqrt{\Delta}^{-1}D_t\Phi}{\cosh(\sqrt{\Delta}(T-t))}\right]\right|\right)^{4/3} dt\right] < \infty. \qquad (B.15)$$

Since the term inside this expectation is finite, the dominated convergence theorem implies that, as $\lambda \to 0$ and in turn $\Delta \to \infty$,

$$\sqrt{\Delta}^{-1}\langle\bar{u}, K^N\rangle_T = \int_0^T \frac{\sqrt{\Delta}^{-1}\mathbb{E}_t[D_t\Phi]}{\cosh(\sqrt{\Delta}(T-t))}\sigma_t^N dt \to \int_0^T \left(\sigma_t^N\right)^2 dt, \quad P\text{-a.s.}$$

Finally, (B.15) also shows that the 4/3-th moment of $\int_0^T \frac{\sqrt{\Delta}^{-1}E_t[D_t\Phi]}{\cosh(\sqrt{\Delta}(T-t))}\sigma_t^N dt$ is bounded, uniformly for all $\Delta > 1$. Therefore, this family indexed by $\Delta > 1$ is uniformly integrable and the almost sure convergence for $\Delta \to \infty$ also holds in $L^1$.

To complete the proof, we now show that the other terms in (B.9) do not contribute at the leading order $O(\sqrt{\lambda})$, that is,

$$\lambda \int_0^T \bar{u}_t(\mu_t^N dt + \sigma_t^N dW_t) = o(\sqrt{\lambda}), \quad \text{in } L^1 \text{ as } \lambda \to 0.$$

Since $\Delta = \frac{M}{\lambda\rho_d(M+1)}$, this is implied by a bound for $\Delta^{-1/2}u^{\mathcal{K}}$. To this end, observe that the inequalities of Jensen and Burkholder-Davis-Gundy show

$$\mathbb{E}\left[(K_t^N - K_s^N)^4\right] \leq C \left(\mathbb{E}\left[\left(\int_s^t \mu_r^N dr\right)^4\right] + \mathbb{E}\left[\left(\int_s^t \sigma_r^N dW_r\right)^4\right]\right)$$

$$\leq C'\mathbb{E}\left[\sup_{0\leq u\leq T}\left\{|\mu_u^N|^4 + |\sigma_u^N|^4\right\}\right](t-s)^2,$$

for some constants $C, C' > 0$ that might only depend on $T$. For $\alpha < \frac{1}{4}$, write $R_\alpha$ for the modulus of $\alpha$-Hölder continuity of $K^N$. This quantity is well defined and satisfies $E[R_\alpha^4] < \infty$ by [13, Theorem 3.1].[10] As $K^N \in \mathcal{H}^2$ by assumption, we can define the square-integrable random variable

$$M^\alpha := \sup_{s\in[0,T]} \mathbb{E}_s[R_\alpha]\left(1 + \int_0^\infty 2e^{-u}|u|^\alpha ds\right) + \sup_{s\in[0,T]}\left|K_s^N\right| < \infty.$$

---

[10]This theorem requires an additional assumption on the iterated integral of the process $K^N$. However, a careful inspection of the proof reveals that this extra assumption is only needed to establish additional path regularity of the iterated integral and not for the path regularity of the process $K^N$ itself.

Then, we can estimate

$$\left|\Delta^{-1/2}\bar{K}_t^N - K_t^N\right| \leq \mathbb{E}_t\left[R_\alpha \int_t^T \Delta^{-1/2}k^\Delta(t,s)|t-s|^\alpha ds\right]$$

$$+ \left|1 - \int_t^T \Delta^{-1/2}k^\Delta(t,s)ds\right|\left|K_t^N\right|$$

$$\leq \mathbb{E}_t[R_\alpha]\Delta^{(1-\alpha)/2}\int_t^T \frac{2e^{\sqrt{\Delta}(T-s)}}{e^{\sqrt{\Delta}(T-t)}}|\sqrt{\Delta}(t-s)|^\alpha ds$$

$$+ \left|1 - \tanh\left(\sqrt{\Delta}(T-t)\right)\right|\left|K_t^N\right|$$

$$\leq \mathbb{E}_t[R_\alpha]\Delta^{-\alpha/2}\int_0^\infty 2e^{-u}|u|^\alpha ds + 2e^{-2\sqrt{\Delta}(T-t)}\left|K_t^N\right|$$

$$\leq \left(\Delta^{-\alpha/2} + 2e^{-2\sqrt{\Delta}(T-t)}\right)M^\alpha = C_{t,T,\alpha,\lambda}.$$

Together with the formulas for $\bar{u}, \bar{U}$ and the definition of the function $F^\Delta$ and (3.2), this estimate yields

$$\left|\Delta^{-1/2}\bar{u}_t\right| = \left|-\Delta^{-1/2}F^\Delta(t)\bar{U}_t + \Delta^{-1/2}\bar{K}_t^N\right|$$

$$= \left|-\sqrt{\Delta}\int_0^t \frac{\sinh(\sqrt{\Delta}(T-t))}{\cosh(\sqrt{\Delta}(T-s))}\Delta^{-1/2}\bar{K}_s^N ds + \Delta^{-1/2}\bar{K}_t^N\right|$$

$$\leq \left|\sqrt{\Delta}\int_0^t \frac{\sinh(\sqrt{\Delta}(T-t))}{\cosh(\sqrt{\Delta}(T-s))}K_s^N ds - K_t^N\right|$$

$$+ C_{t,T,\alpha,\lambda} + \sqrt{\Delta}\int_0^t \frac{\sinh(\sqrt{\Delta}(T-t))}{\cosh(\sqrt{\Delta}(T-s))}C_{s,T,\alpha,\lambda}ds$$

$$\leq \sqrt{\Delta}\int_0^t \frac{\sinh(\sqrt{\Delta}(T-t))}{\cosh(\sqrt{\Delta}(T-s))}\left(R_\alpha|t-s|^\alpha + C_{s,T,\alpha,\lambda}\right)ds$$

$$+ C_{t,T,\alpha,\lambda} + \left|1 - \sqrt{\Delta}\int_0^t \frac{\sinh(\sqrt{\Delta}(T-t))}{\cosh(\sqrt{\Delta}(T-s))}ds\right|\left|K_t^N\right|$$

$$\leq C_{t,T,\alpha,\lambda} + \Delta^{-\alpha/2}R_\alpha\int_0^\infty e^{-u}|u|^\alpha du + \sqrt{\Delta}\int_0^t e^{-\sqrt{\Delta}(t-s)}C_{s,T,\alpha,\lambda}ds$$

$$+ \left|1 - \sqrt{\Delta}\int_0^t \frac{\sinh(\sqrt{\Delta}(T-t))}{\cosh(\sqrt{\Delta}(T-s))}ds\right|\left|K_t^N\right|. \tag{B.16}$$

Recall the addition formula $\arctan(x) - \arctan(y) = \arctan\left(\frac{x-y}{1+xy}\right)$ for $x, y \geq 0$ and observe

26

that $|\arctan(x)| \le |x|$. As a consequence:

$$\sqrt{\Delta} \int_0^t \frac{\sinh(\sqrt{\Delta}(T-t))}{\cosh(\sqrt{\Delta}(T-s))} ds = \sinh(\sqrt{\Delta}(T-t)) \left( \arctan\left(\sinh(\sqrt{\Delta}T)\right) - \arctan\left(\sinh(\sqrt{\Delta}(T-t))\right) \right)$$

$$= \sinh(\sqrt{\Delta}(T-t)) \arctan\left( \frac{\sinh(\sqrt{\Delta}T) - \sinh(\sqrt{\Delta}(T-t))}{1 + \sinh(\sqrt{\Delta}T)\sinh(\sqrt{\Delta}(T-t))} \right)$$

$$\le \sinh(\sqrt{\Delta}(T-t)) \arctan\left( \frac{\sinh(\sqrt{\Delta}T)}{1 + \sinh(\sqrt{\Delta}T)\sinh(\sqrt{\Delta}(T-t))} \right)$$

$$\le \frac{\sinh(\sqrt{\Delta}(T-t))\sinh(\sqrt{\Delta}T)}{1 + \sinh(\sqrt{\Delta}T)\sinh(\sqrt{\Delta}(T-t))}$$

$$\le 1,$$

as well as

$$\sqrt{\Delta} \int_0^t \frac{\sinh(\sqrt{\Delta}(T-t))}{\cosh(\sqrt{\Delta}(T-s))} ds = \sqrt{\Delta} \int_0^t \frac{e^{-\sqrt{\Delta}(t-s)} - e^{-\sqrt{\Delta}(2T-s-t)}}{1 + e^{-2\sqrt{\Delta}(T-s)}} ds$$

$$\ge \sqrt{\Delta} \int_0^t \left( e^{-\sqrt{\Delta}(t-s)} - e^{-\sqrt{\Delta}(2T-s-t)} \right) \left( 1 - e^{-2\sqrt{\Delta}(T-s)} \right) ds$$

$$\ge \sqrt{\Delta} \int_0^t e^{-\sqrt{\Delta}(t-s)} - e^{-\sqrt{\Delta}(2T-2t+(t-s))} - e^{-\sqrt{\Delta}(2T-2t+3(t-s))} ds$$

$$\ge 1 - e^{-\sqrt{\Delta}t} - 2e^{-2\sqrt{\Delta}(T-t)}.$$

In view of these two estimates, (B.16) yields

$$\left| \Delta^{-1/2} \bar{u}_t \right| \le C_{t,T,\alpha,\lambda} + \Delta^{-\alpha/2} R_\alpha \int_0^\infty e^{-u} |u|^\alpha du + \sqrt{\Delta} \int_0^t e^{-\sqrt{\Delta}(t-s)} C_{s,T,\alpha,\lambda} ds$$

$$+ \left( e^{-\sqrt{\Delta}t} + 2e^{-2\sqrt{\Delta}(T-t)} \right) \sup_{s \in [0,T]} |K_s^N|$$

$$\le 4 \left( e^{-\sqrt{\Delta}t} + \Delta^{-\alpha/2} + e^{-2\sqrt{\Delta}(T-t)} \right) M^\alpha. \tag{B.17}$$

(Here, the last inequality follows from the definition of $C_{\cdot,T,\alpha,\lambda}$.) In particular, there exists a constant $C_T > 0$ only depending on $T$ such that

$$\left| \int_0^T \Delta^{-1/2} \bar{u}_t \mu_t^N dt \right| \le C_T M^\alpha \sup_{t \in [0,T]} |\mu_t^N| (\Delta^{-\alpha/2} + \Delta^{-1/2})$$

$$\le \frac{C_T}{2} \left( |M^\alpha|^2 + \sup_{t \in [0,T]} |\mu_t^N|^2 \right) (\Delta^{-\alpha/2} + \Delta^{-1/2}) \to 0,$$

as $\lambda \to 0$ and in turn $\Delta \to \infty$. By the dominated convergence theorem, this pointwise convergence also holds in $L^1$, since the upper bound in this estimate is integrable under our

assumptions. This shows that the Lebesgue integral in (B.9) is indeed of order $o(\sqrt{\lambda})$ as claimed.

The argument for the stochastic integral in (B.9) is similar. By the Burkholder-Davis-Gundy inequality, choosing $C_T > 0$ larger if necessary, we obtain

$$\mathbb{E}\left[\left|\int_0^T \Delta^{-1/2}\bar{u}_t\sigma_t^N dW_t\right|\right] \leq C_T\mathbb{E}\left[\left(\int_0^T |\Delta^{-1/2}\bar{u}_t\sigma_t^N|^2 dt\right)^{1/2}\right]$$

$$\leq 4C_T\mathbb{E}\left[|M^\alpha|\sup_{t\in[0,T]}|\sigma_t^N|\left(\int_0^T e^{-\sqrt{\Delta}t} + \Delta^{-\alpha/2} + e^{-2\sqrt{\Delta}(T-t)}dt\right)^{1/2}\right]$$

$$\leq 2C_T(\Delta^{-\alpha/2} + \Delta^{-1/2})^{1/2}\mathbb{E}\left[|M^\alpha|^2 + \sup_{t\in[0,T]}|\sigma_t^N|^2\right] \to 0,$$

as $\lambda \to 0$ and in turn $\Delta \to \infty$. Here, we have used (B.17) for the second inequality. Therefore, the stochastic integral in (B.9) also is of order $o(\sqrt{\lambda})$ in $L^1$, as $\lambda \to 0$ and the proof is complete. □

# References

[1] R. F. Almgren and N. Chriss. Optimal execution of portfolio transactions. *J. Risk*, 3:5–40, 2001.

[2] P. Bank and D. Baum. Hedging and portfolio optimization in financial markets with a large trader. *Math. Finance*, 14(1):1–18, 2004.

[3] P. Bank, M. Soner, and M. Voß. Hedging with temporary price impact. *Math. Fin. Econ.*, 11(2):215–239, 2017.

[4] P. Bank and M. Voß. Linear quadratic stochastic control problems with singular stochastic terminal constraint. *SIAM J. Control Optim.*, 56(2):672–699, 2018.

[5] D. Bertsimas and A. Lo. Optimal control of execution costs. *J. Fin. Markets*, 1(1):1–50, 1998.

[6] B. Bouchard, M. Fukasawa, M. Herdegen, and J. Muhle-Karbe. Equilibrium returns with transaction costs. *Finance Stoch.*, 22(3):569–601, 2018.

[7] M. Butz and R. Oomen. Internalisation by electronic FX spot dealers. *Quant. Finance*, 19(1):35–56, 2019.

[8] A. Capponi, A. J. Menkveld, and H. Zhang. Large orders in small markets: On optimal execution with endogenous liquidity supply. Preprint.

[9] R. Carmona and K. Webster. The self-financing equation in limit order book markets. *Finance Stoch.*, 23(3):729–759, 2013.

[10] Á. Cartea and S. Jaimungal. A closed-form execution strategy to target volume weighted average price. *SIAM J. Fin. Math.*, 7(1):760–785, 2016.

[11] U. Cetin, R. A. Jarrow, and P. Protter. Liquidity risk and arbitrage pricing theory. *Finance Stoch.*, 8(3):311–341, 2004.

[12] J.-H. Choi, K. Larsen, and D. Seppi. Equilibrium effects of TWAP and VWAP order splitting. *Math. Fin. Econ.*, to appear.

[13] P. K. Friz and M. Hairer. *A course on rough paths.* Springer, Berlin, 2014.

[14] N. Garleanu, L. Pedersen, and A. Poteshman. Demand-based option pricing. *Rev. Fin. Stud.*, 22(10):4259–4299, 2009.

[15] N. Gârleanu and L. H. Pedersen. Dynamic trading with predictable returns and transaction costs. *J. Finance*, 68(6):2309–2340, 2013.

[16] N. Gârleanu and L. H. Pedersen. Dynamic portfolio choice with frictions. *J. Econ. Theory.*, 165:487–516, 2016.

[17] D. Gromb and D. Vayanos. Equilibrium and welfare in markets with financially constrained arbitrageurs. *J. Fin. Econ.*, 66(2-3):361–407, 2002.

[18] D. Gromb and D. Vayanos. A model of financial market liquidity based on intermediary capital. *J. Eur. Econ. Assoc.*, 8(2-3):456–466, 2010.

[19] S. J. Grossman and M. H. Miller. Liquidity and market structure. *J. Finance*, 43(3):617–633, 1988.

[20] J. Hasbrouck and G. Sofianos. The trades of market makers: An empirical analysis of NYSE specialists. *J. Finance*, 48(5):1565–1593, 1993.

[21] M. Herdegen, J. Muhle-Karbe, and D. Possamaï. Equilibrium asset pricing with transaction costs. *Finance Stoch.*, to appear.

[22] T. Ho and H. R. Stoll. Optimal dealer pricing under transactions and return uncertainty. *J. Fin. Econ.*, 9(1):47–73, 1981.

[23] D. Kramkov and S. Pulido. Stability and analytic expansions of local solutions of systems of quadratic BSDEs with applications to a price impact model. *SIAM J. Fin. Math.*, 7(1):567–587, 2016.

[24] D. Kramkov and S. Pulido. A system of quadratic BSDEs arising in a price impact model. *Ann. Appl. Probab.*, 26(2):794–817, 2016.

[25] A. S. Kyle. Continuous auctions and insider trading. *Econometrica*, 53(6):1315–1335, 1985.

[26] A. W. Lo, H. Mamaysky, and J. Wang. Asset prices and trading volume under fixed transactions costs. *J. Pol. Econ.*, 112(5):1054–1090, 2004.

[27] L. Moreau, J. Muhle-Karbe, and H. M. Soner. Trading with small price impact. *Math. Finance*, 27(2):350–400, 2017.

[28] J. Muhle-Karbe, M. Nutz, and X. Tan. Asset pricing with heterogenous beliefs and illiquidity. *Math. Finance*, 30(4):1392–1421, 2020.

[29] J. Muhle-Karbe and K. Webster. Information and inventories in high-frequency trading. *Market Microstructure Liq.*, 3(02):1750010, 2017.

[30] D. Nualart. *The Malliavin calculus and related topics*. Springer, Berlin, 2006.

[31] M. Nutz and J. A. Scheinkman. Shorting in speculative markets. *J. Finance*, 75(2):995–1036, 2020.

[32] A. Obizhaeva and J. Wang. Optimal trading strategy and supply/demand dynamics. *J. Fin. Markets*, 16(1):1–32, 2013.

[33] A. Pavlova and R. Rigobon. The role of portfolio constraints in the international propagation of shocks. *Rev. Econ. Stud.*, 75(4):1215–1256, 2008.

[34] I. Rosu. Fast and slow informed trading. *J. Fin. Markets*, 43:1–30, 2019.

[35] Y. Sannikov and A. Skrzypacz. Dynamic trading: price inertia and front-running. Preprint, 2016.