

Fast and Slow Informed Trading*

Ioanid Roşu[†]

May 10, 2015

Abstract

This paper develops a model in which traders receive a stream of private signals, and differ in their information processing speed. In equilibrium, the fast traders (FTs) quickly reveal a large fraction of their information, and generate most of the volume, volatility and profits in the market. If a FT is averse to holding inventory, his optimal strategy changes considerably as his aversion crosses a threshold. He no longer takes long-term bets on the asset value, gets most of his profits in cash, and generates a “hot potato” effect: after trading on information, the FT quickly unloads part of his inventory to slower traders. The results match evidence about high frequency traders.

KEYWORDS: Trading volume, inventory, volatility, high frequency trading, price impact, mean reversion.

*Earlier versions of this paper circulated under the title “High Frequency Traders, News and Volatility.” The author thanks Kerry Back, Laurent Calvet, Thierry Foucault, Johan Hombert, Pete Kyle, Stefano Lovo, Victor Martinez, Daniel Schmidt, Dimitri Vayanos, Jiang Wang; finance seminar participants at Copenhagen Business School, HEC Paris, Univ. of Durham, Univ. of Leicester, Univ. Paris Dauphine, Univ. Madrid Carlos III, ESSEC, KU Leuven; and conference participants at the European Finance Association 2014 meetings, American Finance Association 2013 meetings, Society for Advancement of Economic Theory in Portugal, Central Bank Microstructure Conference in Norway, Market Microstructure Many Viewpoints Conference in Paris, for valuable comments. The author acknowledges financial support from the Investissements d’Avenir Labex (ANR-11-IDEX-0003/Labex Ecodec/ANR-11-LABX-0047).

[†]HEC Paris, Email: rosu@hec.fr.

1 Introduction

Today’s markets are increasingly characterized by the continuous arrival of vast amounts of information. A media article about high frequency trading reports on the hedge fund firm Citadel: *“Its market data system, for example, contains roughly 100 times the amount of information in the Library of Congress. [...] The signals, or alphas, that prove to have predictive power are then translated into computer algorithms, which are integrated into Citadel’s master source code and electronic trading program.”* (“Man vs. Machine,” CNBC.com, September 13th 2010). The sources of information from which traders obtain these signals usually include company-specific news and reports, economic indicators, stock indexes, prices of other securities, prices on various other trading platforms, limit order book changes, as well as various “machine readable news” and even “sentiment” indicators.¹

At the same time, financial markets have seen in recent years the spectacular rise of algorithmic trading, and in particular of high frequency trading.² This coincidental arrival raises the question whether or not at least some of the HFTs do process information and trade very quickly in order to take advantage of their speed and superior computing power. Recent empirical evidence suggests that this is indeed the case.³ But, despite the large role played by high frequency traders (HFTs) in the current financial landscape, there has been relatively little progress in explaining their strategies in connection with information processing.

We consider the following questions regarding HFTs: What are the optimal trading strategies of HFTs who process information? Why do HFTs account for such a large share of the trading volume? What explains the race for speed among HFTs? What are the effects of HFTs on measures of market quality, such as liquidity and price volatility?

¹ *“Math-loving traders are using powerful computers to speed-read news reports, editorials, company Web sites, blog posts and even Twitter messages—and then letting the machines decide what it all means for the markets.”* (“Computers That Trade on the News,” New York Times, December 22nd 2010).

²Hendershott, Jones, and Menkveld (2011) report that from a starting point near zero in the mid-1990s, high frequency trading rose to as much as 73% of trading volume in the United States in 2009. Chaboud, Chiquoine, Hjalmarsson, and Vega (2014) consider various foreign exchange markets and find that starting from essentially zero in 2003, algorithmic trading rose by the end of 2007 to approximately 60% of the trading volume for the euro-dollar and dollar-yen markets, and 80% for the euro-yen market.

³See Brogaard, Hendershott, and Riordan (2014), Baron, Brogaard, and Kirilenko (2014), Kirilenko, Kyle, Samadi, and Tuzun (2014), Hirschey (2013), Benos and Sagade (2013), Brogaard, Hagströmer, Nordén, and Riordan (2013).

How can HFT order flow anticipate future order flow and returns? What explains the “intermediation chains” or “hot potato” effects found among HFTs? Why do some HFTs have low inventories? Regarding the last question, some recent literature identifies HFTs as traders with both high trading volume *and* low inventories (see Kirilenko *et al.* 2014, SEC 2010). But then, a natural question arises: why would having low inventories be part of the definition of HFTs?

In this paper, we provide a theoretical model of informed trading which parsimoniously addresses these questions. Because we want to study speed differences among informed traders, we start with the standard framework of Kyle (1985), and modify it along several dimensions.⁴ First, the asset’s fundamental value is not constant but follows a random walk process, and each risk-neutral informed trader, or speculator, gradually receives signals about the asset value increments. Second, there are multiple speculators who differ in their speed, in the sense that some speculators receive their signal with a lag. Third, each speculator can trade only on lagged signals with a lag of at most m , where m is an exogenously given number.

It is the last assumption that sets our model apart from previous models of informed trading. A key effect of this assumption is to prevent the “rat race” phenomenon discovered by Holden and Subrahmanyam (1992), by which traders with identical information reveal their information so quickly, that the equilibrium breaks down at the “high frequency” limit, when the number of trading rounds approaches infinity. In our model, the speculators reveal only a fraction of their total private information, and this has a stabilizing effect on the equilibrium. Economically, we think of this assumption as equivalent to having a positive information processing cost per signal (and per trading round).⁵ Indeed, since one of our results is that the value of information decays fast, even a tiny information processing cost would make speculators optimally ignore their signals after a sufficiently large number of lags m .

⁴As in Kyle (1985), we assume that informed traders submit only market orders; this is a plausible assumption for informed HFTs (see Brogaard, Hendershott, and Riordan 2014). Also, we set the model in continuous time, which makes it easier to solve for the equilibrium.

⁵Intuitively, information processing is costly because speculators need to avoid trading on stale information, and this involves (i) constantly monitoring public information to verify that their signal has not been incorporated into the price, and (ii) extracting the predictable part of their signal from past order flow, so that speculators trade only on the unpredictable (non-stale) part.

To simplify the analysis, we restrict our attention to the particular case when $m = 1$, which we call the *benchmark model*. In this model, speculators can trade using only their current signal and its lagged value. Thus, there are two types of speculators: *fast traders* (or FTs), who observe the signal instantly; and *slow traders* (or STs), who observe the signal after one lag. The benchmark model has the advantage that the equilibrium can be described in closed form. In the Internet Appendix we verify numerically that the main results in the benchmark model carry through to the general case ($m > 1$).

Our first main result in the benchmark model is that the FTs generate most of the trading volume, volatility, and profits. To understand why, consider the decision of N fast traders about what weight to use on the last signal they have received. Because the dealer sets a price function which is linear in the aggregate order size, each FT faces a Cournot-type problem and trades such that the price impact of his order is on average $1/(N + 1)$ of his signal. That brings the expected aggregate price impact to $N/(N+1)$ of the signal, and leaves on average only $1/(N+1)$ of the signal unknown to the dealer. Thus, once the STs observe the lagged signal, they now have much less private information to exploit. Moreover, the ST profits are further diminished by competition with FTs, who also trade on the lagged signal. Empirically, Baron, Brogaard, and Kirilenko (2014) find out that the profits of HFTs are concentrated among a small number of incumbents, and the profits appear to be correlated with speed.

Our second main result is that volume, volatility and liquidity are increasing with the number of FTs. First, more competition from FTs makes the prices more informative overall, and thus increases liquidity (measured, as in Kyle 1985, by the inverse price impact coefficient). As the market is more liquid, FTs face a lower price impact, and therefore trade even more aggressively. This creates an amplification mechanism that allows the aggregate FT trading volume to be increasing roughly linearly with the number of FTs. The effect of FTs on volatility is more muted but still positive; this is because in our model price volatility is bounded above by the fundamental volatility of the asset. Empirically, in line with our theoretical results, Hendershott, Jones, and Menkveld (2011), Boehmer, Fong, and Wu (2014), and Zhang (2010) document a positive effect of HFTs on liquidity. Moreover, the last two papers find a positive effect of HFTs on volatility. We should point out, however, that our model is more likely to

apply only to the subcategory of informed, market taking HFTs, and not to all HFTs. Thus, our results should be interpreted with caution.

Our third main result in the benchmark model is the existence of anticipatory trading: the order flow of fast traders predicts the order flow of slow traders in the next period. This comes from the fact that the fast traders' signal does not fully get incorporated into the price, hence the slow traders have an incentive to use the signal in the next period, after they remove the stale (predictable) part. Anticipatory trading is therefore related to speculator order flow autocorrelation. Our model predicts that the speculator order flow autocorrelation is positive, although it is small if the number of fast traders is large. Empirically, Brogaard (2011) finds that the autocorrelation of aggregate HFT order flow is indeed small and positive. Also, using Nasdaq data on high-frequency traders, Hirschey (2013) finds that HFT order flow anticipates future order flow.

Despite being able to match several stylized facts about HFTs in our benchmark model, a few questions remain. Why do many HFTs have low inventories, both intraday and at the day close?⁶ Why do HFTs engage in “hot potato” trading (or “intermediation chains”), in which HFT pass their inventories to other traders?⁷ What is the role of speed in explaining these phenomena?

To provide some theoretical guidance on these issues, we extend our benchmark model to include one trader with inventory costs. These costs can arise from risk aversion or from capital constraints, but we take a reduced form approach and assume the costs are quadratic in inventory, with a coefficient called *inventory aversion* (see Madhavan and Smidt 1993). We call this additional trader the *Inventory-averse Fast Trader*, or IFT.⁸ We call this extension the *model with inventory management*. In addition to choosing the weight on his current signal, the IFT also chooses the rate at which he

⁶SEC (2010) characterizes HFTs by their “*very short time-frames for establishing and liquidating positions*” and argues that HFTs end “*the trading day in as close to a flat position as possible (that is, not carrying significant, unhedged positions over-night)*.” See also Kirilenko *et al.* (2014), Brogaard, Hagströmer, Nordén, and Riordan (2013), or Menkveld (2013).

⁷Weller (2014) analyzes both theoretically and empirically “intermediation chains” in which uninformed HFTs unwind inventories to slower, fundamental traders. Kirilenko *et al.* (2014) mention “hot potato effect” during the Flash Crash episode of May 6, 2010, when some HFTs would churn out their inventories very quickly to trade with other HFTs.

⁸We do not introduce more than one IFT since the model would be much more difficult to solve. The IFT is assumed fast because without slower traders it is not profitable to manage inventory.

mean reverts his inventory to zero each period. Without discussing yet optimality, suppose the IFT does inventory management, i.e., chooses a positive rate of inventory mean reversion. What are the effects of this choice?

The first effect of inventory management is that the IFT keeps essentially all his profits in cash. To see this, suppose the IFT chooses a coefficient of mean reversion of 10%. This translates into the inventory being reduced by a fraction of 10% in each trading round. Therefore, the IFT's inventory tends to become small over many rounds, and because our model is set in the high frequency limit (in continuous time), the inventory becomes in fact negligible.⁹ We call this result the *low inventory effect*.

The second effect is that the IFT no longer makes profits by betting on the fundamental value of the asset. This stands in sharp contrast to the behavior of a risk-neutral speculator, such as the fast trader in the benchmark model. Indeed, the FT accumulates inventory in the direction of his information, since he knows his signals are correlated with the asset's liquidation value. By contrast, although the IFT initially trades on his current signal, he subsequently fully reverses the bet on that signal by removing a fraction of his inventory each trading round. Thus, the IFT's direct revenue from each signal eventually decays to zero. We call this result the *information decay effect*.

The third effect of inventory management is that, in order to make a profit, the IFT must (i) anticipate the slow trading, and (ii) trade in the opposite direction to slow trading. By *slow trading* here we simply mean the part of order flow that involves the speculators' lagged signals.¹⁰ To understand this effect, consider how the IFT uses a given signal. The information decay effect means that the IFT's final revenues from betting on his signal are zero. Therefore, the IFT must benefit from inventory reversal. Since any trade has price impact, inventory reversal makes a profit only if gets pooled with order flow in the opposition direction, so that the IFT's price impact is negative. But in order to be *expected* profit, the opposite order flow must come from speculators who use lagged signals, i.e., from slow trading. We call this result the *hot potato effect*, or the *intermediation chain effect*.¹¹

⁹Formally, the inventory follows an autoregressive process, hence its variance has the same order as the variance of the signal, which at high frequencies is negligible.

¹⁰A subtle point is that slow *trading* does not need to come from actual slow *traders*. Slow trading can also arise from fast traders who use their lagged signals as part of their optimal trading strategy.

¹¹In our simplified framework, the intermediation chain only has one link, between the IFT and the

The reason behind this terminology is that the IFT’s current signal (the “potato”) produces undesirable inventory (is “hot”) and must be passed on to slower traders in order to produce a profit. Thus, speed is important to the IFT. Without slower trading, there is no hot potato effect, and the IFT makes a negative expected profit from any trading strategy that mean reverts his inventory to zero. Note also that the hot potato generates a complementarity between the IFT and slow traders: Stronger inventory mean reversion by the IFT reduces the price impact of the STs, who can trade more aggressively. But more aggressive trading by the STs allows stronger mean reversion from the IFT.

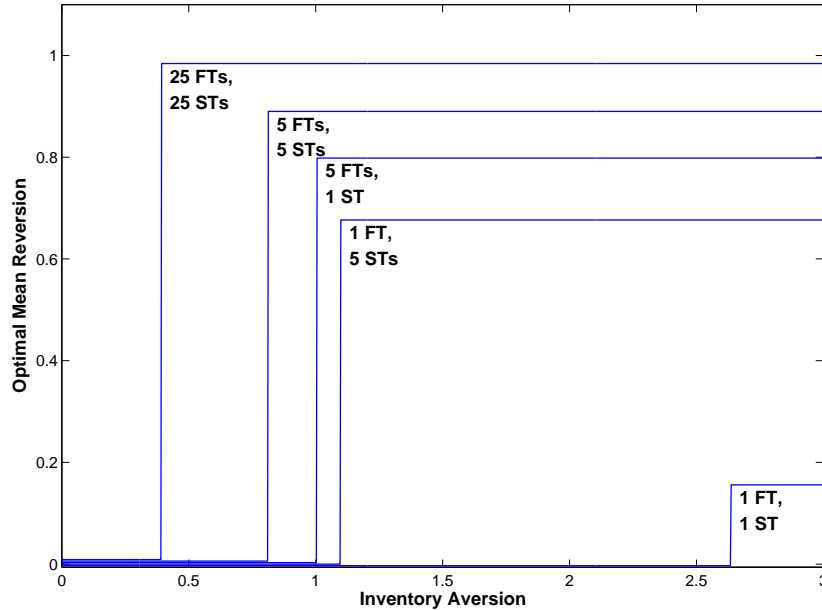
The optimal strategy of the IFT produces two contrasting types of behavior, depending on how his inventory aversion compares to a threshold. Below the threshold, the IFT behaves like a risk-neutral speculator, and lowers his inventory costs simply by reducing the weight on his signals. He does not manage inventory at all, because the information decay effect ensures that even a small but positive inventory mean reversion eventually destroys all revenues from the fundamental bets. With inventory aversion above the threshold, the IFT manages inventory and has all his profits in cash. The IFT benefits not from fundamental bets on his signals, but from the hot potato effect.

Figure 1 illustrates the optimal mean reversion for the IFT as a function of his inventory aversion coefficient. We see that, as his inventory aversion rises, the IFT changes discontinuously from the regime with no inventory mean reversion to the regime with positive inventory mean reversion. The threshold at which this discontinuity occurs depends on the number of fast traders (FTs) and slow traders (STs) in the model. This threshold is decreasing in both parameters, because the amount of slow trading is increasing in both parameters. Slow trading is clearly increasing in the number of slow traders. But it is also increasing in the number of fast traders because (i) the fast traders also use their lagged signals, and (ii) more fast traders make the market more liquid, which allows slow trading to be more aggressive.

Our results speak to the literature on high-frequency trading. One may think that in practice HFTs have very low inventories because either (i) HFTs have very high

slow traders. But we conjecture that in a model where speculators use more than one lag for their signals, the intermediation chains become longer, depending on the number of lags.

Figure 1: Optimal Inventory Mean Reversion. This figure plots the optimal mean reversion coefficient of an inventory-averse fast trader (IFT), when he competes with N_F fast traders (FTs) and N_S slow traders (STs), with $N_F, N_S \in \{1, 5, 25\}$. On the horizontal axis is the IFT's inventory aversion coefficient. The optimal mean reversion coefficient is computed using the results of Section 5, in the inventory management model with parameters N_F and $N_L = N_F + N_S$. The other parameter values are: $\sigma_w = 1, \sigma_u = 1$.



risk aversion, or (ii) HFTs do not have superior information and wish to maintain zero inventory to avoid adverse selection on their positions in the risky asset. Our results suggest that this is not necessarily the case. Indeed, Figure 1 suggests (and we rigorously prove in Proposition 6) that in the limit when the number of speculators is large, the threshold inventory aversion converges to zero, and the optimal mean reversion is close to one. In other words, even with low inventory aversion, the IFT chooses very large mean reversion. Yet, even at these high rates of mean reversion the IFT does not lose more than about 50% of his average profits from inventory management (the advantage being that he has all his profits in cash).

We predict that in practice the fast speculators are sharply divided into two categories. In both categories speculators trade with a large volume. But in one category speculators accumulate inventory by taking fundamental bets. In the other category speculators have very low inventories; they initially trade on their signals but then

quickly pass on part of their inventory to slower traders. These covariance patterns produce testable implications of our model.

The division of fast speculators in two categories appears consistent with the empirical findings of Kirilenko *et al.* (2014), who study trading activity in the E-mini S&P 500 futures during several days around the Flash Crash of May 6, 2010. The “opportunistic traders” described in their paper resembles our risk-neutral fast traders: opportunistic traders have large volume, appear to be fast, and accumulate relatively large inventories. By contrast the “high frequency traders” in their paper, while they are also fast and trade in large volume, keep very low inventories. Indeed, HFTs in their sample liquidate 0.5% of their aggregate inventories on average each second.

Related Literature

Our paper contributes to the literature on trading with asymmetric information. We show that competition among informed traders, combined with noisy trading strategies, produces a large informed trading volume and a quick information decay.¹² The market is very efficient because competition among informed traders makes them trade aggressively on their common information. This intuition is present in Holden and Subrahmanyam (1992) and Foster and Viswanathan (1996). The former paper finds that the competition among informed traders is so strong, that in the continuous time limit there is no equilibrium in smooth strategies. Our contribution to this literature is to show that there exists an equilibrium in *noisy* strategies. This rests on two key assumptions: (i) *noisy information*, i.e., speculators learn over time by observing a stream of signals, and (ii) *finite lags*, i.e., speculators only use a signal for a fixed number of lags—which is plausible if there is a positive information processing cost per signal.

Without the finite lags assumption, noisy information by itself does not generate noisy strategies, as Back and Pedersen (1998) show. Chau and Vayanos (2008), Caldentey and Stacchetti (2010), and Li (2012) find that noisy information coupled with either model stationarity or a random liquidation deadline produces strategies that are still smooth as in Kyle (1985), but towards the high frequency limit they have almost

¹²A speculator’s strategy is *smooth* if the volatility generated by that speculator’s trades is of a lower magnitude compared to the volatility from noise trading; and *noisy* if the magnitudes are the same.

infinite weight. Thus, the market in these papers is nearly strong-form efficient, which makes speculators' strategies appear noisy (there is no actual equilibrium in the limit). By contrast, in our model the market is not strong-form efficient even in the limit, yet strategies are noisy. Foucault, Hombert, and Roşu (2015) propose a model in which a single speculator receives a signal one instant before public news. The speculator's strategy is noisy, but for a different reason than in our model: the speculator optimally trades with a large weight on his forecast of the news. Yet a different mechanism occurs in Cao, Ma, and Ye (2013). In their model, informed traders must disclose their trades immediately after trading, and therefore traders optimally obfuscate their signal by adding a large noise component to their trades.

Our paper also contributes to the rapidly growing literature on High Frequency Trading.¹³ In much of this literature, it is the speed *difference* that has a large effect in equilibrium. The usual model setup has certain traders who are faster in taking advantage of an opportunity that disappears quickly. As a result, traders enter into a winner-takes-all contest, in which even the smallest difference in speed has a large effect on profits. (See for instance the model with speed differences of Biais, Foucault, and Moinas (2014), or the model of news anticipation of Foucault, Hombert, and Roşu (2015).) By contrast, our results regarding volume and volatility remain true even if all informed traders have the same speed. This is because in our model the need for speed arises endogenously, from competition among informed traders. In our model, being "slow" simply means trading on lagged signals. Since in equilibrium speculators also use lagged signals (the unanticipated part, to be precise), in some sense all traders are slow as well. Yet, it is true in our model that a genuinely slower trader makes less money, since he can only trade on older information that has already lost much of its value.

Our results regarding the optimal inventory of informed traders are, to our knowledge, new. Theoretical models of inventory usually attribute inventory mean reversion to passive market makers, who do not possess superior information, but are concerned

¹³See Biais, Foucault, and Moinas (2014), Ait-Sahalia and Saglam (2014), Budish, Cramton, and Shim (2014), Foucault, Hombert, and Roşu (2015), Du and Zhu (2014), Li (2014), Hoffmann (2014), Pagnotta and Philippon (2013), Weller (2014), Cartea and Penalva (2012), Jovanovic and Menkveld (2012), Cvitanic and Kirilenko (2010).

with absorbing order flow.¹⁴ Our paper shows that an informed investor with inventory costs (the “IFT”) can display behavior that makes him appear like a market maker, even though he only submits market orders (as in Kyle 1985). Indeed, in our model the IFT does not take fundamental bets, passes his risky inventory to slower traders (the hot potato effect), and keeps essentially all his money in cash. To obtain these results, even a small inventory aversion of the IFT suffices, but only if enough slow trading exists.

A related paper is Hirshleifer, Subrahmanyam, and Titman (1994). In their 2-period model, risk averse speculators with a speed advantage first trade to exploit their information, after which they revert their position because of risk aversion; while the slower speculators trade in the same direction as the initial trade of the faster speculators. The focus of Hirshleifer, Subrahmanyam, and Titman (1994) is different, as they are interested in information acquisition and explaining behavior such as “herding” and “profit taking.” Our goal is to analyze the inventory problem of fast informed traders in a fully dynamic context, and to study the properties of the resulting optimal strategies.

The paper is organized as follows. Section 2 describes the model setup. Section 3 solves for the equilibrium in the particular case with two categories of traders: fast and slow. Section 4 discusses the effect of fast and slow traders on various measures of market quality. Section 5 introduces an extension of the baseline model in which a new trader (the IFT) has inventory costs. Then, it analyzes the IFT’s optimal strategy and its effect on equilibrium. Section 6 concludes. All proofs are in the Appendix or the Internet Appendix. The Internet Appendix solves for the equilibrium in the general case, and analyzes several modifications and extensions of our baseline model.

2 Model

Trading for a risky asset takes place continuously over the time interval $[0, T]$, where we use the normalization:¹⁵

$$T = 1. \tag{1}$$

¹⁴See Ho and Stoll (1981), Madhavan and Smidt (1993), Hendershott and Menkveld (2014), as well as many references therein.

¹⁵To eliminate confusion with later notation, we often use T instead of 1. This way, we can denote below $t - dt$ by $t - 1$ without much confusion.

Trading occurs at intervals of length dt apart. Throughout the text, we refer to dt as representing one period, or one trading round. The liquidation value of the asset is

$$v_T = \int_0^T dv_t, \quad \text{with} \quad dv_t = \sigma_v dB_t^v, \quad (2)$$

where B_t^v is a Brownian motion, and $\sigma_v > 0$ is a constant called the *fundamental volatility*. We interpret v_T as the “long-run” value of the asset; in the high frequency world, this can be taken to be the asset value at the end of the trading day. The increments dv_t are then the short term changes in value due to the arrival of new information. The risk-free rate is assumed to be zero.

There are three types of market participants: (a) $N \geq 1$ risk neutral speculators, who observe the flow of information at different speeds, as described below; (b) noise traders; and (c) one competitive risk neutral dealer, who sets the price at which trading takes place.

Information and Speed. At $t = 0$, there is no information asymmetry between the speculators and the dealer. Subsequently, each speculator receives the following flow of signals:

$$ds_t = dv_t + d\eta_t, \quad \text{with} \quad d\eta_t = \sigma_\eta dB_t^\eta, \quad (3)$$

where $t \in (0, T]$ and B_t^η is a Brownian motion independent from all other variables. Denote by

$$w_t = \mathbb{E}(v_T \mid \{s_\tau\}_{\tau \leq t}) \quad (4)$$

the expected value conditional on the information flow until t . We call w_t the *value forecast*, or simply *forecast*. Because there is no initial information asymmetry, $w_0 = 0$. Denote by σ_w the instantaneous volatility of w_t , or the *forecast volatility*. The increment of the forecast w_t , and the forecast variance are given, respectively, by

$$dw_t = \frac{\sigma_v^2}{\sigma_v^2 + \sigma_\eta^2} ds_t, \quad \sigma_w^2 = \frac{\text{Var}(dw_t)}{dt} = \frac{\sigma_v^4}{\sigma_v^2 + \sigma_\eta^2}. \quad (5)$$

When deriving empirical implications, we call σ_w the *signal precision*, as a precise signal (small σ_η) corresponds to a large σ_w .

Speculators obtain their signal with a lag $\ell \in \{0, 1, 2, \dots\}$. A ℓ -speculator is a trader who at $t \in (0, T]$ observes the signal from ℓ periods before, $ds_{t-\ell dt}$. To simplify notation, we use the following convention:

$$\text{Notation for trading times: } t - \ell \text{ instead of } t - \ell dt. \quad (6)$$

For instance, instead of $ds_{t-\ell dt}$ we write $ds_{t-\ell}$.

Trading and Prices. At each $t \in (0, T]$, denote by dx_t^i the market order submitted by speculator $i = 1, \dots, N$ at t , and by du_t the market order submitted by the noise traders, which is of the form $du_t = \sigma_u dB_t^u$, where B_t^u is a Brownian motion independent from all other variables. Then, the aggregate order flow executed by the dealer at t is

$$dy_t = \sum_{i=1}^N dx_t^i + du_t. \quad (7)$$

The dealer is risk neutral and competitive, hence she executes the order flow at a price equal to her expectation of the liquidation value conditional on her information. Let $\mathcal{I}_t = \{y_\tau\}_{\tau < t}$ be the dealer's information set just before trading at t . The order flow at date t , dy_t , executes at

$$p_t = \mathbf{E}(v_T | \mathcal{I}_t \cup dy_t). \quad (8)$$

Together with the price, another important quantity is the dealer's expectation at t of the k -lagged signal dw_{t-k} :

$$z_{t-k,t} = \mathbf{E}(dw_{t-k} | \mathcal{I}_t). \quad (9)$$

Equilibrium Definition. In general, a trading strategy for a ℓ -speculator is a process followed by his risky asset position, x_t , which is measurable with respect to his information set $\mathcal{J}_t^{(\ell)} = \{y_\tau\}_{\tau < t} \cup \{s_\tau\}_{\tau \leq t-\ell}$. For a given trading strategy, the speculator's expected profit π_τ , from date τ onwards, is

$$\pi_\tau = \mathbf{E} \left(\int_\tau^T (v_T - p_t) dx_t | \mathcal{J}_\tau^{(\ell)} \right). \quad (10)$$

As in Kyle (1985), we focus on linear equilibria. Specifically, we consider strategies

which are linear in the unpredictable part of their signals,¹⁶

$$dw_{t-k} - z_{t-k,t}, \quad k = \ell, \ell + 1, \dots \quad (11)$$

We restrict strategies to exclude signals older than a fixed number of lags m (which is allowed to depend on the speculator's speed parameter ℓ). This assumption can be justified by costly information processing, as explained at the end of this section. Formally, the ℓ -speculator's strategy is of the form:

$$dx_t = \gamma_{\ell,t}(dw_{t-\ell} - z_{t-\ell,t}) + \gamma_{\ell+1,t}(dw_{t-\ell-1} - z_{t-\ell-1,t}) + \dots + \gamma_{m,t}(dw_{t-m} - z_{t-m,t}). \quad (12)$$

A linear equilibrium is such that: (i) at every date t , each speculator's trading strategy (12) maximizes his expected trading profit (10) given the dealer's pricing policy, and (ii) the dealer's pricing policy given by (8) and (9) is consistent with the equilibrium speculators' trading strategies.

Finally, the speculators consider the covariance structure of $z_{t-k,t}$ to be independent of their strategy. More precisely, for all $j, k \geq 0$, the speculators consider the numbers

$$Z_{j,k,t} = \text{Cov}(dw_{t-j}, z_{t-k,t}) \quad (13)$$

to depend only on j, k , and t . Thus, the covariance terms $Z_{j,k,t}$ are interpreted as being computed by the dealer, as part of her (publicly known) pricing rules.¹⁷

Model Notation. If all speculators in the model have a strategy of the form (12) with the same $m \geq 0$, we call it the *model with m lags*, and write \mathcal{M}_m . In the paper, we focus on the particular case with $m = 1$ lags. In this setup, the 0-speculators are called the *fast traders*, and the 1-speculators are called the *slow traders*; thus, we call \mathcal{M}_1 the

¹⁶Intuitively, if the strategy had a predictable component, the dealer's price would adjust and reduce the speculator's profit. We formalize this intuition in a discrete version of our model in Internet Appendix M. In the paper, however, we work in continuous time since it is easier to obtain analytical solutions. Similarly, Kyle (1985) directly assumes that the speculator's strategy in continuous time is linear in the unpredictable part of the fundamental value, $v - p_t$.

¹⁷For instance, the coefficient ρ_t in the dealer's pricing rule $z_{t-1,t} = \rho_t dy_t$ is computed using the covariance term $\text{Cov}(dw_t, dy_t)$ (see equation (A11)). Hence, even though a speculator affects dy_t by his strategy, he can consider the covariance term $\text{Cov}(dw_t, dy_t)$ to be independent of his strategy. In Internet Appendix M.3, we explore an alternative specification in which the speculator takes into account his effect on dy_t . We find, however, that the overall effect on the equilibrium coefficients is very small.

model with fast and slow traders.

If some ℓ -speculators have strategies of the form (12) with different m_ℓ , we call this the *mixed* model with m lags, where m is the maximum of all m_ℓ . We are particularly interested in the mixed model with $m = 1$ lags in which 0-speculators (fast traders) only trade on their current signal ($m_0 = 0$) and the 1-speculators (slow traders) only use their lagged signal ($m_1 = 1$). We call this the *benchmark model* with fast and slow traders, and denote it by \mathcal{B}_1 . In Section 3, we solve for the equilibrium in both \mathcal{M}_1 and \mathcal{B}_1 , and show that \mathcal{M}_1 can be regarded as a particular case of \mathcal{B}_1 .

Information Processing. The assumption that speculators cannot use lagged signals beyond a given bound can be justified by introducing an information processing cost $\delta > 0$ per individual signal and per unit of time. More precisely, we consider an alternative model in which a ℓ -speculator can use all past signals, but must pay a fixed cost $\delta_\ell dt$ each time he trades with a nonzero weight ($\gamma_{k,t}$) on his k -lagged signal (see equation 12). Then, intuitively, because the value of information decays with the lag, and the speculator does not want to accumulate too large a cost, he must stop using lagged signals beyond an upper bound. In Result 1 we show that for a particular value of δ the alternative model is equivalent to \mathcal{M}_1 .

In this paper, we do not model the exact nature of the speculators' signals and their processing costs. But, intuitively, an information processing cost per signal (and per trading round) is plausible, because in practice speculators must constantly monitor each signal in order to avoid trading on stale (predictable) information. In our model, this can be done by simply removing the predictable part ($z_{t,t-k}$) from the lagged signal (dw_{t-k}). In practice, however, speculators must monitor various sources of public information (such as news reports, economic data, or trading information in various related securities), to extract the part of the signal has not yet been incorporated into the price.

Note that an individual processing cost implicitly means that speculators cannot simply rely on free public signals, such as the price, to shortcut the learning process. This is because in reality prices may contain other relevant information about the fundamental value, along which the speculators are adversely selected. We formalize this intuition in Internet Appendix L, where we introduce an orthogonal dimension of the fundamental value, and show that trading strategies that rely on prices make an average loss.

3 Equilibrium with Fast and Slow Traders

In this section, we analyze the important case in which speculators use signals with a maximum lag of one. There are two types of speculators: (i) the *Fast Traders*, or FTs, who observe the signal with no delay (called 0-speculators in Section 2); and (ii) the *Slow Traders*, or STs, who observe the signal with a delay of one lag (called 1-speculators). The trading strategy of FTs and STs is of the form (see (12)):

$$dx_t = \gamma_t(dw_t - z_{t,t}) + \mu_t(dw_{t-1} - z_{t-1,t}), \quad t \in (0, T], \quad (14)$$

where the weight γ_t must be zero for a ST. There are two possibilities: either the FT can trade on both the current and the lagged signals, or the FT can trade only on the current signal, i.e., the FT's weight γ_t must be zero.¹⁸ The former case is the model denoted by \mathcal{M}_1 , the *model with fast and slow traders*. The latter case is the model denoted by \mathcal{B}_1 , the *benchmark model*.

In Section 3.1, we solve for the equilibrium of the model \mathcal{M}_1 in closed form. One important implication is that the FTs and STs trade identically on their lagged signal (μ_t is the same for all). Therefore, if we require the FTs to use only their current signal (as in \mathcal{B}_1) and introduce an equal number of additional STs, then the aggregate behavior remains essentially the same. Hence, the model \mathcal{M}_1 can be regarded as a particular case of \mathcal{B}_1 , and we are justified in calling \mathcal{B}_1 the *benchmark model* with fast and slow traders. This more general model can also be solved in closed form, by using essentially the same formulas as in Section 3.1. We discuss the benchmark model in Section 3.2.

3.1 The Model with Fast and Slow Traders

In this section, we solve for the equilibrium of the model \mathcal{M}_1 with fast and slow traders. From (14), the FTs have a strategy of the form $dx_t = \gamma_t(dw_t - z_{t,t}) + \mu_t(dw_{t-1} - z_{t-1,t})$, while the STs have a strategy of the same form, except that μ_t must be zero. The current signal (dw_t) is not predictable from the past order flow, hence the dealer sets $z_{t,t} = 0$. The lagged signal (dw_{t-1}) has already been used by the FTs in the previous

¹⁸Intuitively, this can occur if the FT must pay a higher processing cost per signal than the ST; see the discussion at the beginning of Section 3.2.

trading round, hence the dealer can use the past order flow to compute the predictable part $z_{t-1,t}$.¹⁹ To simplify notation, let \widetilde{dw}_{t-1} be the unanticipated part at t of the lagged signal:

$$\widetilde{dw}_{t-1} = dw_{t-1} - z_{t-1,t}. \quad (15)$$

In Theorem 1, we show that there exists a closed-form linear equilibrium of the model. The equilibrium is *symmetric*, in the sense that the FTs have identical trading strategies, and so do the STs. We also provide asymptotic results when the number N_F of fast traders is large. We say that X_∞ is the asymptotic value of a number X which depends on N_F , if the ratio X/X_∞ converges to 1 as N_F approaches infinity, and we write:

$$X \approx X_\infty \iff \lim_{N_F \rightarrow \infty} \frac{X}{X_\infty} = 1. \quad (16)$$

Let “ F ” refer to the fast traders, and “ S ” to the slow traders. Denote by N_F the number of fast traders, and by N_S the number of slow traders. We denote the total number of speculators by

$$N_L = N_F + N_S. \quad (17)$$

This is the same as the number of speculators who use their lagged signals, hence the “ L ” notation. We also call N_L the *number of lag traders*.

Theorem 1. *Consider the model \mathcal{M}_1 with $N_F > 0$ fast traders and $N_S \geq 0$ slow traders; let $N_L = N_F + N_S$. Then, there exists a symmetric linear equilibrium with constant coefficients, of the form ($t \in (0, T]$):*

$$\begin{aligned} dx_t^F &= \gamma dw_t + \mu \widetilde{dw}_{t-1}, & dx_t^S &= \mu \widetilde{dw}_{t-1}, \\ \widetilde{dw}_{t-1} &= dw_{t-1} - \rho dy_{t-1}, & dp_t &= \lambda dy_t, \end{aligned} \quad (18)$$

where the coefficients $\gamma, \mu, \rho, \lambda$ are given by:

$$\begin{aligned} \gamma &= \frac{1}{\lambda} \frac{1}{N_F + 1}, & \mu &= \frac{1}{\lambda} \frac{1}{N_L + 1} \frac{1}{1 + b}, \\ \rho &= \frac{\sigma_w}{\sigma_u} \sqrt{(1-a)(a-b^2)}, & \lambda &= \rho \frac{N_F}{N_F - b}, \end{aligned} \quad (19)$$

¹⁹In Theorem 1, we show that the dealer sets $z_{t-1,t} = \rho dy_{t-1}$ for some constant coefficient ρ .

with

$$\omega = 1 + \frac{1}{N_F} \frac{N_L}{N_L + 1}, \quad b = \frac{\sqrt{\omega^2 + 4 \frac{N_L}{N_L + 1}} - \omega}{2}, \quad a = \frac{N_F - b}{N_F + 1}. \quad (20)$$

We have the following asymptotic limits when N_F is large:

$$\omega_\infty = a_\infty = 1 \quad b_\infty = \frac{\sqrt{5} - 1}{2}, \quad \lambda_\infty = \rho_\infty = \frac{\sigma_w}{\sigma_u} \frac{1}{\sqrt{N_F}}. \quad (21)$$

The number b is increasing in both N_F and N_S . Moreover, $\omega \in [1, 2)$, $a \in (0, 1)$, $b \in [0, b_\infty)$.

One consequence of the Theorem is that FTs and STs trade with the same intensity (μ) on their lagged signals. This is true because the current signal dw_t is uncorrelated with the lagged signal \widetilde{dw}_{t-1} , which implies that the FTs and the STs get the same expression for the expected profit that comes from the lagged signal.²⁰

We now discuss some comparative statics regarding the optimal weights γ and μ (for brevity, we omit the proofs). The fast traders' optimal weight γ is decreasing in the number of fast traders, yet it is increasing in the number of slow traders. The first statement simply reflects that, when the number of fast traders is larger, these traders must divide the pie into smaller slices. The same logic applies to the coefficient on the lagged signal: μ is decreasing in both N_F and N_S , as the fast and slow traders compete in trading on their common lagged signal. This last intuition also shows that the fast traders' weight γ is increasing in the number of slow traders. Indeed, when there is more competition from slow traders, the fast traders have an incentive to trade more aggressively on their current signal, as the slow traders have not yet observed this signal.

The next Corollary helps to get more intuition for the equilibrium.

²⁰This result does not generalize to the case when there are more lags ($M > 1$). In Internet Appendix I, we see that there is a positive autocorrelation between the signals of higher lags, which reflects a more complicated covariance structure. Mathematically, this translates into the covariance matrix A having non zero entries $A_{i,j}$ when $i > j \geq 1$.

Corollary 1. *In the context of Theorem 1, we have the following formulas:*

$$\begin{aligned} \lambda \bar{\gamma} &= \frac{N_F}{N_F + 1}, & \lambda \bar{\mu} &= \frac{1}{1 + b} \frac{N_L}{N_L + 1}, \\ \frac{\text{Var}(\widetilde{dw}_t)}{dt} &= (1 - a) \sigma_w^2 = \frac{1 + b}{N_F + 1} \sigma_w^2, & \frac{\text{Cov}(\widetilde{dw}_t, w_t)}{dt} &= \frac{1 - a}{1 + b} \sigma_w^2 = \frac{\sigma_w^2}{N_F + 1}. \end{aligned} \tag{22}$$

The first equation in (22) implies that $\lambda \bar{\gamma} dw_t = \frac{N_F}{N_F + 1} dw_t$, which shows that most of the current signal (dw_t) is incorporated into the price by the fast traders. The intuition comes from the Cournot nature of the equilibrium. Indeed, when trading on the current signal, the benefit of each of each FT increases linearly with the intensity of trading γ on his signal; while the price at which he eventually trades increases linearly with the aggregate quantity demanded. Given that the price impact of the other $N_F - 1$ fast traders aggregates to $\frac{N_F - 1}{N_F + 1} dw_t$, the FT is a monopsonist against the residual supply curve, and trades such that his price impact is half of $\frac{2}{N_F + 1} dw_t$, i.e., his price impact equals $\frac{1}{N_F + 1} dw_t$.

After incorporating $\frac{N_F}{N_F + 1} dw_t$ in trading round t , the fast traders must compete with the slow traders for the remaining $\frac{1}{N_F + 1} dw_t$ in the next trading round. As explained before, the speculators must trade a multiple of the unanticipated part of the lagged signal, $\widetilde{dw}_t = dw_t - \rho dy_t$. Thus, when trading on the lagged signal, the benefit of each speculator—fast or slow—increases linearly with the intensity of trading μ , and is proportional to the covariance $\text{Cov}(\widetilde{dw}_t, w_t)$. At the same time, each speculator faces a price that increases linearly with the aggregate quantity demanded, and which is proportional to the lagged signal variance $\text{Var}(\widetilde{dw}_t)$. The argument is now similar to the Cournot one above, except that everything gets multiplied by the ratio $\text{Cov}(\widetilde{dw}_t, w_t) / \text{Var}(\widetilde{dw}_t)$, which according to (22) is equal to $1/(1 + b)$. This justifies the second equation in (22). It also implies that in the case of the lagged signal only a fraction $1/(1 + b)$ of it is incorporated by the speculators into the price.

We use the results in Theorem 1 to compute the expected profits of the fast traders and the slow traders.

Proposition 1. *In the context of Theorem 1, the expected profit of the FTs and STs at*

$t = 0$ from their equilibrium strategies are given, respectively, by:

$$\begin{aligned}\frac{\pi^F}{\sigma_w^2} &= \frac{\gamma}{N_F + 1} + \frac{1}{N_F + 1} \frac{\mu}{N_L + 1}, \\ \frac{\pi^S}{\sigma_w^2} &= \frac{1}{N_F + 1} \frac{\mu}{N_L + 1}.\end{aligned}\tag{23}$$

The ratio of slow profits to fast profits is therefore

$$\frac{\pi^S}{\pi^F} = \frac{1}{1 + \frac{(N_L+1)^2(1+b)}{N_F+1}} \implies \frac{\pi^S}{\pi^F} \approx \frac{N_F}{(N_F + N_S)^2} \frac{1}{1 + b_\infty}.\tag{24}$$

Thus, even if there is only one ST ($N_S = 1$), the ST profits are small compared to the FT profits. The reason is that FTs trade also on their lagged signals, and thus compete with the STs. Indeed, FTs compete for trading on dw_t only among themselves, while they also compete with the STs for trading on the lagged signal \widetilde{dw}_{t-1} .

Finally, Proposition 1 gives an estimate for the information processing cost δ that would be sufficient to discourage speculators from trading on lagged signals beyond one, if that were not imposed by the model. We state the following numerical result.

Result 1. *Consider the alternative model setup with N_F fast speculators and N_S slow speculators, in which each speculator can use past signals at any lag, but must pay for each signal (used with nonzero weight) an information processing cost of*

$$\delta = \frac{1}{N_F + 1} \frac{\mu}{N_L + 1} \sigma_w^2.\tag{25}$$

Then, the alternative model is equivalent to the model with fast and slow traders (\mathcal{M}_1).

3.2 The Benchmark Model

We now consider the *benchmark model* \mathcal{B}_1 , in which the fast traders use only the current signal, while the slow traders use only the lagged signal.²¹ The strategies of the FTs

²¹As in Result 1, $\mathcal{M}_{0,1}$ is equivalent to an alternative setup with information processing costs, in which (i) the STs pay the cost δ from (25), while (ii) the FTs pay a cost slightly higher than δ . Indeed, if a FT paid δ , he would be indifferent between using his lagged signal and not using it; while with a slightly higher cost, he would be strictly worse off and would ignore his lagged signal.

and STs are, respectively, of the form

$$dx_t^F = \gamma_t dw_t, \quad dx_t^S = \mu_t \widetilde{dw}_{t-1}, \quad (26)$$

where $\widetilde{dw}_{t-1} = dw_{t-1} - \rho_t dy_{t-1}$. The dealer sets the price using the rule $dp_t = \lambda_t dy_t$. Let $N_F \geq 1$ be the number of FTs and $N_L \geq 0$ the number of STs.

The next result shows that the model \mathcal{M}_1 with N_F fast traders and N_S slow traders produces essentially the same outcome as the benchmark model \mathcal{B}_1 with N_F fast traders and $N_L = N_F + N_S$ slow traders.

Corollary 2. *Consider (a) the model \mathcal{M}_1 with $N_F \geq 1$ fast traders and $N_S \geq 0$ slow traders; and (b) the benchmark model \mathcal{B}_1 with N_F fast traders and $N_L = N_F + N_S$ slow traders. Then, the equilibrium coefficients γ , μ , λ , ρ in the two models are identical.*

This Corollary is obtained by simply following the proof of Theorem 1 to solve for the equilibrium in the \mathcal{B}_1 model. The key step is to observe that in Theorem 1 the fast trader’s choice of μ is the same as the slow trader’s choice of μ , and therefore it does not matter who does the optimization, as long as the total number of speculators using their lagged signal is the same.

We finally note that the benchmark model \mathcal{B}_1 with $N_F > 0$ fast traders and N_L slow traders has two important particular cases:

- If $N_L \geq N_F$, \mathcal{B}_1 is equivalent to the model \mathcal{M}_1 with N_F fast traders and $N_S = N_L - N_F$ slow traders;
- If $N_L = 0$, \mathcal{B}_1 is the model \mathcal{M}_0 (with 0 lags).

4 Market Quality with Fast and Slow Traders

In this section, we study the effect of fast and slow trading on various measures of market quality. The setup is the benchmark model \mathcal{B}_1 with $N_F \geq 1$ fast traders and $N_L \geq 0$ slow traders. In this context, “fast trading” is the speculators’ aggregate trading on their current signal, and “slow trading” is the speculators’ aggregate trading on their lagged signal. The measures of market quality analyzed are illiquidity (measured by the price

impact coefficient), trading volume, price volatility, price informativeness, the speculator participation rate, and the speculator's order flow autocorrelation. The main conclusion of this section is that fast trading has the strongest effect on most of our measures of market quality, while slow trading has a relatively smaller effect. The only measure that depends crucially on slow trading is the speculators' order flow autocorrelation, which becomes positive only in the presence of slow trading. This is shown to be related to anticipatory trading: the order flow coming from fast traders anticipates the order flow coming from the slow traders in the next period.

4.1 Measures of Market Quality

We first decompose the aggregate speculator order flow into fast trading and slow trading. Denote by $d\bar{x}_t$ be the aggregate speculator order flow. Let $\bar{\gamma}$ be the aggregate weight on the current signal (dw_t), and $\bar{\mu}$ the aggregate weight on the lagged signal (\widetilde{dw}_{t-1}). We decompose the aggregate speculator order flow $d\bar{x}_t$ into two components: *fast trading*, which represents the aggregate trading on the current signal; and *slow trading*, which represents the aggregate trading on the lagged signal:

$$d\bar{x}_t = \underbrace{\bar{\gamma} dw_t}_{\text{Fast Trading}} + \underbrace{\bar{\mu} \widetilde{dw}_{t-1}}_{\text{Slow Trading}}, \quad \text{with } \bar{\gamma} = N_F \gamma, \quad \bar{\mu} = N_L \mu. \quad (27)$$

As in Theorem 1, we define $b = \rho \bar{\mu}$. We call b the *slow trading coefficient*. Then, slow trading exists (is nonzero) only if the number of traders who use their lagged signal is positive, or equivalently if $b > 0$:

$$\text{Slow Trading exists} \quad \iff \quad N_L > 0 \quad \iff \quad b > 0. \quad (28)$$

Note that the case when there is no slow trading coincides with the model \mathcal{M}_0 with 0 lags from Section 2. In that model, N_F fast traders use only their current signal.

We now define the measures of market quality. Recall that the dealer sets a price that changes in proportion to the total order flow $dy = d\bar{x}_t + du_t$:

$$dp_t = \lambda dy_t = \lambda \left(\bar{\gamma} dw_t + \bar{\mu} \widetilde{dw}_{t-1} + du_t \right), \quad (29)$$

First, as it is standard in the literature, we define *illiquidity* to be the price impact coefficient λ . Thus, the market is considered illiquid if the price impact of a unit of trade is large, i.e., if the coefficient λ is large.

Second, we define *trading volume* as the infinitesimal variance of the aggregate order flow dy_t :

$$TV = \sigma_y^2 = \frac{\text{Var}(dy_t)}{dt}. \quad (30)$$

We argue that this is a measure of trading volume. Indeed, in each trading round the actual aggregate order flow is given by dy_t . Thus, one can interpret trading volume as the absolute value of the order flow: $|dy_t|$. From the theory of normal variables, the average trading volume is given by $\mathbf{E}(|dy_t|) = \sqrt{\frac{2}{\pi}} \sigma_y$. With our definition $TV = \sigma_y^2$, we observe that TV is monotonic in $\mathbf{E}(|dy_t|)$, and thus TV can be used a measure of trading volume. Using (29), we compute the trading volume in our model by the formula

$$TV = \bar{\gamma}^2 \sigma_w^2 + \bar{\mu}^2 \sigma_{\tilde{w}}^2 + \sigma_u^2, \quad \text{with} \quad \sigma_{\tilde{w}}^2 = \frac{\text{Var}(\tilde{d}w_t)}{dt}. \quad (31)$$

The trading volume measure TV can be decomposed into the speculator trading volume and the noise trading volume:

$$TV = TV^s + TV^n, \quad \text{with} \quad TV^s = \bar{\gamma}^2 \sigma_w^2 + \bar{\mu}^2 \sigma_{\tilde{w}}^2, \quad TV^n = \sigma_u^2. \quad (32)$$

Third, we define *price volatility* σ_p to be the square root of the instantaneous price variance:

$$\sigma_p = \left(\frac{\text{Var}(dp_t)}{dt} \right)^{1/2}. \quad (33)$$

From (29), it follows that the instantaneous price variance can be computed simply as the product of the illiquidity measure λ and the trading volume $TV = \sigma_y^2$. Thus,

$$\sigma_p^2 = \lambda^2 TV = \lambda^2 \left(\bar{\gamma}^2 \sigma_w^2 + \bar{\mu}^2 \sigma_{\tilde{w}}^2 + \sigma_u^2 \right). \quad (34)$$

Fourth, we define *price informativeness* as a measure inversely related to the forecast error variance $\Sigma_t = \text{Var}((w_t - p_{t-1})^2)$. Thus, if prices are informative, they stay close to the forecast w_t , i.e., the variance Σ_t is small. In Internet Appendix I, in the general

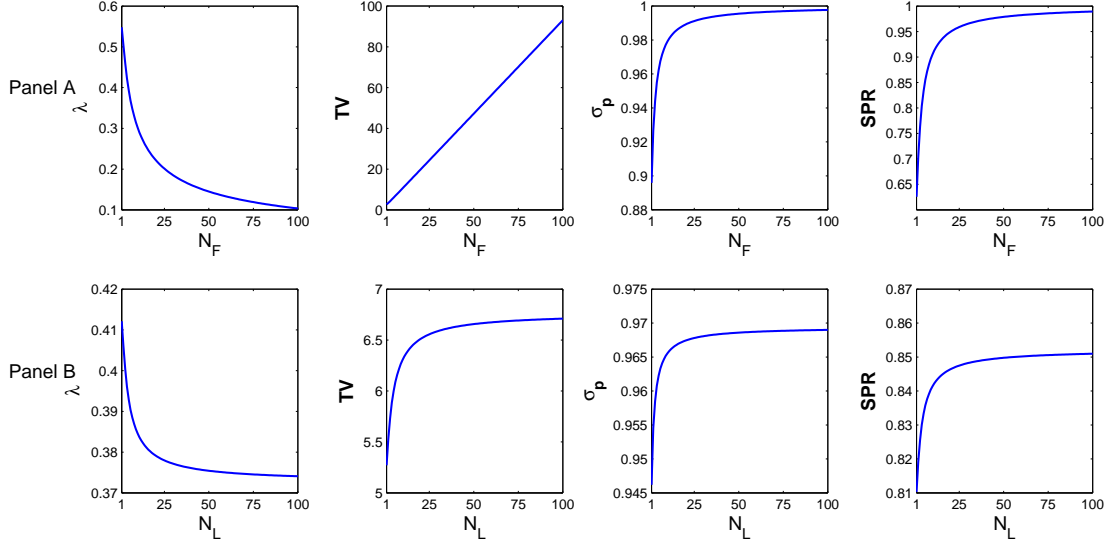
model with at most m lagged signals (\mathcal{M}_m) we show that Σ_t evolves according to $\Sigma'_t = \sigma_w^2 - \sigma_p^2$, where σ_p^2 is the price variance (Proposition I.1). Therefore, since Σ'_t is inversely monotonic in the price variance, we do not use it as a separate measure of market quality.

Fifth, the *speculator participation rate* is defined as the ratio of speculator trading volume over total trading volume:

$$SPR = \frac{TV^s}{TV} = \frac{\bar{\gamma}^2 \sigma_w^2 + \bar{\mu}^2 \sigma_w^2}{\bar{\gamma}^2 \sigma_w^2 + \bar{\mu}^2 \sigma_w^2 + \sigma_u^2}. \quad (35)$$

SPR can also be interpreted as the fraction of price variance due to the speculators.

Figure 2: Market Quality with Fast and Slow Traders. This figure plots the following measures of market quality: (i) illiquidity λ ; (ii) trading volume TV ; (iii) price volatility σ_p ; and (iv) speculator participation rate SPR . Panel A plots the dependence of the four market quality measures on the number of fast traders N_F , while taking the number of slow traders $N_L = 5$. Panel B plots the dependence of the four market quality measures on N_L , while taking $N_F = 5$. The other parameters are $\sigma_w = 1$, $\sigma_u = 1$.



4.2 Comparative Statics on Market Quality

We now give explicit formulas for our measures of market quality. As before, we use asymptotic notation when N_F is large: $X \approx Y$ stands for $\lim_{N_F \rightarrow \infty} \frac{X}{Y} = 1$.

Proposition 2. *Consider the benchmark model with $N_F \geq 1$ fast traders and $N_L \geq 0$ slow traders. Then, the price impact coefficient, trading volume, price volatility, and speculator participation rate satisfy:*

$$\begin{aligned} \lambda &= \frac{\sigma_w}{\sigma_u} \frac{\sqrt{(1+b)(a-b^2)}}{\sqrt{N_F+1}} \frac{N_F}{N_F-b}, & TV &= \sigma_u^2 (N_F+1) \frac{a}{(1+b)(a-b^2)}, \\ \sigma_p^2 &= \sigma_w^2 \frac{N_F^2}{(N_F+1)(N_F-b)}, & SPR &= a + \frac{b^2(1+b)}{N_F-b}, \end{aligned} \quad (36)$$

where $b^2 + b(1 + \frac{1}{N_F} \frac{N_L}{N_L+1}) = \frac{N_L}{N_L+1}$, and $a = \frac{N_F-b}{N_F+1}$.

Panel A of Figure 2 shows how the four measures of market quality vary with the number of fast traders N_F , while holding the number of slow traders N_L constant. Panel B of Figure 2 shows how the four measures of market quality vary with N_L , while holding N_F constant. We find that all four market quality measures vary in the same direction with respect to N_F and N_L . Nevertheless, the number of fast traders has a much stronger effect on these measures than the number of slow traders.

To get more intuition about the effect of fast trading on market quality, we consider the simplest case, when $N_L = 0$. Since all speculators trade only on their current signal, this case coincides with the model \mathcal{M}_0 as defined in Section 2. In this model there is no slow trading ($\bar{\mu} = 0$), hence the slow trading coefficient b is zero. Moreover, $a = \frac{N_F-b}{N_F+1} = \frac{N_F}{N_F+1}$. Thus, we can solve the model \mathcal{M}_0 by simply using Proposition 2. Nevertheless, it is instructive to solve for the equilibrium of \mathcal{M}_0 independently.

Proposition 3. *Consider the model \mathcal{M}_0 , with N_F fast traders whose trading strategy is of the form $dx_t = \gamma_t dw_t$. Then, the optimal coefficient γ is constant and equal to $\gamma = \frac{1}{\lambda} \frac{1}{N_F+1} = \frac{\sigma_u}{\sigma_w} \frac{1}{\sqrt{N_F}}$. The price impact coefficient, trading volume, price volatility, and speculator participation rate satisfy, respectively,*

$$\lambda = \frac{\sigma_w}{\sigma_u} \frac{\sqrt{N_F}}{N_F+1}, \quad TV = \sigma_u^2 (N_F+1), \quad \sigma_p^2 = \sigma_w^2 \frac{N_F}{N_F+1}, \quad SPR = \frac{N_F}{N_F+1}. \quad (37)$$

Using Proposition 3, we now discuss in more detail the effect of the number N_F of fast traders on the measures of market quality. First, we note by quickly inspecting the formulas in Proposition 3, that we obtain the same qualitative results as those displayed

in Figure 2. Namely, illiquidity is decreasing in N_F , while the other three measures are increasing in N_F .

An important consequence of Proposition 3 is that in our model the speculator participation rate can be made arbitrarily close to 1 if the number of fast traders is large. Thus, noise trading volatility is only a small part of the total volatility. This stands in sharp contrast for instance with the models of Kyle (1985) or Back, Cao, and Willard (2000), in which virtually *all* instantaneous price volatility is generated by the noise traders at the high frequency limit (in continuous time).

The market is more efficient when the number of fast traders is large. Indeed, in the proof of Proposition 3 we show that the rate of change of the forecast error variance Σ' is constant and equal to $\frac{\sigma_w^2}{N_F+1}$. Since by assumption there is no initial informational asymmetry ($\Sigma_0 = 0$), it follows that $\Sigma_t \leq \frac{\sigma_w^2}{N_F+1}$ for all t . In other words, the price stays close to the fundamental value at all times. Thus, a larger number N_F of fast traders, rather than destabilizing the market, makes the market more efficient.

The trading volume TV strongly increases with the number of fast traders. This occurs because of the competition among FTs make them trade more aggressively. By trading more aggressively, FTs reveal more information, which as we see later lowers the traders' price impact. This has an amplifier effect on the trading aggressiveness, such that the trading volume grows essentially linearly in the number of speculators (see equation (37)). Moreover, the speculator participation rate SPR also increases in N_F , since SPR is the fraction of trading volume caused by the speculators.

Surprisingly, a larger number of fast traders make the market more liquid, as more information is revealed when there are more competing speculators. This appears to be in contradiction with the fact that more informed trading should increase the amount of adverse selection. To understand the source of this apparent contradiction, note that illiquidity is measured by the price impact λ of one unit of volume. But, while the trading volume TV strongly increases in N_F in an unbounded way, its price impact is bounded by magnitude of the signal dw_t .²² Thus, the price impact *per unit of volume* actually decreases, indicating that prices are more informative. This makes the market

²²In Internet Appendix I, we make this intuition rigorous in the general case; see the discussion surrounding Proposition I.4.

overall more liquid. This result is consistent with the empirical studies of Zhang (2010), Hendershott, Jones, and Menkveld (2011), and Boehmer, Fong, and Wu (2014).

To understand the effect of fast traders on the price volatility σ_p , consider the pricing formula $dy_t = \lambda dy_t$, which implies $\sigma_p^2 = \lambda^2 TV$. There are two effects of N_F on the price volatility σ_p . First, the trading volume TV strongly increases in N_F , which has a positive effect on σ_p . Second, price impact λ decreases in N_F , which has a negative effect on σ_p . The first effect is slightly stronger than the second, hence the net effect is that price volatility σ_p increases in N_F . This result is consistent with the empirical studies of Boehmer, Fong, and Wu (2014) and Zhang (2010).

A few caveats are in order. First, all these studies analyze the effects of HFT activity, where activity is proxied either by turnover or by intensity of order-related message traffic, and *not* by the number of HFTs present in the market. An answer to this concern is that, as we have seen, trading volume does increase in N_F . Second, in our paper we do not model *passive* HFTs, that is, HFTs that offer liquidity via limit orders. Therefore, it is possible that an increase in the number of passive HFTs decreases price volatility, which would cancel the opposite effect of the number of *active* HFTs. For instance, Hasbrouck and Saar (2012) document a negative effect of HFTs on volatility, possibly because they also consider passive HFTs, which by providing liquidity have a stabilizing effect on price volatility. Moreover, Chaboud, Chiquoine, Hjalmarsson, and Vega (2014) find essentially no relation. In our model, the dependence of price volatility on N_F is weak, which may explain the mixed results in the empirical literature.

Next, we discuss how the various measures of market quality depend on the speculators' signal precision σ_w . Note that, according to equation (5), the signal precision is related to the fundamental volatility σ_v by a monotonic relation: $\sigma_w = \frac{\sigma_v}{(1 + \sigma_\eta^2 / \sigma_v^2)^{1/2}}$. Therefore, we only analyze the dependence of market quality on signal precision, while keeping in mind that these results apply equally to the fundamental volatility.

The price volatility σ_p increases in signal precision, indicating that speculators trade more aggressively when they have a more precise signal. Indeed, σ_p is the volatility of dp_t , which is the price impact of the aggregate order flow. In particular, the order flow coming from the FTs has an aggregate price impact which is proportional to dw_t .²³

²³From Proposition 3, the FTs' order flow equals $\lambda N_F \gamma dw_t = \frac{N_F}{N_F + 1} dw_t$.

Thus, price volatility increases in the signal precision.

A larger signal precision σ_w generates more adverse selection between fast traders and the dealer, hence the illiquidity λ is increasing in the signal precision. However, the trading volume TV is independent of σ_w . To get some intuition for this result, note that $TV = \frac{\sigma_p^2}{\lambda^2}$. Since both the numerator and denominator increase with signal precision, the net effect is ambiguous. Proposition 3 shows that the two effects exactly offset each other.

4.3 Order Flow Autocorrelation and Anticipatory Trading

We start by analyzing the autocorrelation of the components of the order flow. Since the dealer is competitive and risk neutral, the total order flow dy_t has zero autocorrelation. But because the dealer cannot identify the part of the order flow that comes from speculators, the *speculator* order flow can in principle be autocorrelated.

As in Section 4.1, in the benchmark model with fast and slow traders, the aggregate speculator order flow decomposes into its fast trading and slow trading components:

$$d\bar{x}_t = \underbrace{d\bar{x}_t^F}_{\text{Fast Trading}} + \underbrace{d\bar{x}_t^S}_{\text{Slow Trading}}, \quad \text{with} \quad d\bar{x}_t^F = \bar{\gamma} dw_t, \quad d\bar{x}_t^S = \bar{\mu} \widetilde{dw}_{t-1}, \quad (38)$$

with $\bar{\gamma} = N_F \gamma$ and $\bar{\mu} = N_L \mu$. As before, we say that slow trading exists if $b = \rho \bar{\mu} > 0$, or equivalently $N_L > 0$.

We define *speculator order flow autocorrelation* by $\text{Corr}(d\bar{x}_t, d\bar{x}_{t+1})$. Because $d\bar{x}_{t+1}^F$ is orthogonal to both components of $d\bar{x}_t^F$, we obtain the decomposition:

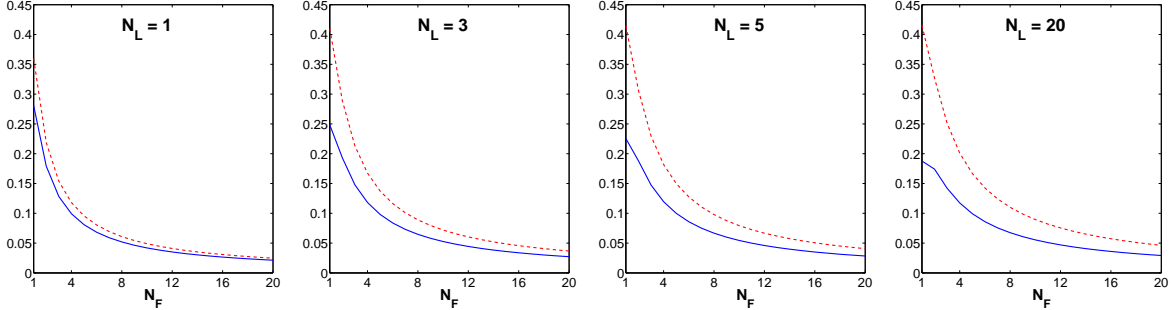
$$\rho_{\bar{x}} = \text{Corr}(d\bar{x}_t, d\bar{x}_{t+1}) = \underbrace{\frac{\text{Cov}(d\bar{x}_t^F, d\bar{x}_{t+1}^S)}{\text{Var}(d\bar{x}_t)}}_{\text{Anticipatory Trading}} + \underbrace{\frac{\text{Cov}(d\bar{x}_t^S, d\bar{x}_{t+1}^S)}{\text{Var}(d\bar{x}_t)}}_{\text{Expectation Adjustment}}. \quad (39)$$

We denote the *anticipatory trading* part by ρ_{AT} and the *expectation adjustment* part by ρ_{EA} . The first component arises because fast trading at t anticipates slow trading at $t + 1$. Indeed, there is a positive correlation between fast trading at t and slow trading at $t + 1$ ($\bar{\mu} \widetilde{dw}_t$). The second component arises because slow trading at $t + 1$ is based on lagged signals, adjusted by subtracting the dealer's expectation which incorporates past

lagged signals. Because of this expectation adjustment, we see below that the slow order flow is negatively autocorrelated. Formally, slow trading at $t + 1$ ($\bar{\mu}\widetilde{dw}_t$) is proportional to the lagged signal minus dealer's expectation, $\widetilde{dw}_t = dw_t - \rho dy_t$. But the dealer's expectation is proportional on the total order flow at t , which includes the previous slow trading ($dy_t = \bar{\gamma} dw_t + \bar{\mu} \widetilde{dw}_{t-1} + du_t$). We compute:

$$\rho_{\bar{x}} = \rho_{AT} + \rho_{EA}, \quad \text{with} \quad \rho_{AT} = \bar{\mu}\bar{\gamma} \frac{\text{Var}(dw_t)}{\text{Var}(d\bar{x}_t)}, \quad \rho_{EA} = -\rho\bar{\mu}^3 \frac{\text{Var}(\widetilde{dw}_{t-1})}{\text{Var}(d\bar{x}_t)}. \quad (40)$$

Figure 3: Speculator Order Flow Autocorrelation. This figure plots the speculator order flow autocorrelation $\rho_{\bar{x}}$ (solid line) and the anticipatory trading component ρ_{AT} (dashed line) as a function of the number N_F of fast traders. The four graphs correspond to four values of the number N_L of speculators using their lagged signal: $N_L = 1, 3, 5, 20$.



Proposition 4. Consider the benchmark model with $N_F \geq 1$ fast traders and $N_L \geq 0$ slow traders. Then, the speculator order flow autocorrelation and its components satisfy

$$\rho_{\bar{x}} = \frac{b(b+1)(a-b^2)}{a^2 + b^2(1-a)} \frac{1}{N_F + 1}, \quad \frac{\rho_{AT}}{\rho_{\bar{x}}} = \frac{a}{a-b^2}, \quad \frac{\rho_{EA}}{\rho_{\bar{x}}} = -\frac{b^2}{a-b^2}, \quad (41)$$

where a and b are as in Proposition 2. Moreover, $\rho_{\bar{x}}$ is strictly positive if and only if there exists slow trading, i.e., $N_L > 0$.

One implication of Proposition 4 is that, as long as there exists slow trading, the speculator order flow autocorrelation $\rho_{\bar{x}}$ is nonzero. To understand why, note that both the anticipatory trading component and the expectation adjustment component depend on the existence of slow trading. Formally, if there is no slow trading, $\bar{\mu} = 0$ implies that both components of the speculator order flow autocorrelation are zero.

Figure 3 shows how the speculator order flow autocorrelation ($\rho_{\bar{x}}$) and its anticipatory trading component (ρ_{AT}) depend on the number of fast traders (N_F) for four different values of the number of slow traders ($N_L = 1, 3, 5, 20$). We see that both $\rho_{\bar{x}}$ and ρ_{AT} are decreasing in N_F . Indeed, when the number of fast traders is large, there is only $\frac{1}{N_F+1}$ of the signal left in the next period for the slow traders. Hence, one should expect the autocorrelation to decrease by the order of $\frac{1}{N_F+1}$, which is indeed the case. For instance, when $N_L = 5$, we see that the speculator order flow autocorrelation is 22.56% when there is one FT, but decreases to 2.84% when there are 20 FTs. Our results are consistent with the empirical literature on HFTs. For instance, Brogaard (2011) finds that the autocorrelation of aggregate HFT order flow is small but positive.

The anticipatory trading component ρ_{AT} is increasing in the number of slow traders N_L (to see this, fix for instance $N_F = 10$ in each of the four graphs in Figure 3). The intuition is simple: when the number of slow traders is larger, fast trading in each period can better predict the slow trading the next period, hence the correlation ρ_{AT} is larger. Using Nasdaq data on high-frequency traders, Hirschey (2013) finds that HFT order flow anticipates non-HFT order flow. But Nasdaq defines HFTs along several criteria including the use of large trading volume and low inventories. In our model, these are indeed the characteristics of fast traders, but not those of slow traders (see the next section for a discussion about traders' inventories). Thus, if in our model we classified fast traders as HFTs and slow traders as non-HFTs, our previous results would imply that HFT order flow anticipates non-HFT order flow.

5 Inventory Management

In this section, we analyze the inventory problem of fast traders. Because the benchmark model cannot address this problem (when speculators are risk-neutral, their inventory follows a random walk), we modify the model by introducing an additional trader with inventory costs.²⁴ We call this new trader the *Inventory-averse Fast Trader*, or IFT, and the resulting setup the *model with inventory management*, or the *model with an IFT*.

²⁴Introducing more than one inventory-averse trader makes the problem considerably more complicated, as the number of state variables increases.

To get intuition for the model with inventory management, we first solve for the optimal strategy of the IFT in a partial equilibrium framework, taking as fixed the behavior of the other speculators and the dealer. The solution of this problem is provided in closed form. Then, we continue with a general equilibrium analysis. We show that the equilibrium reduces to a non-linear equation in one variable, which can be solved numerically. We then study the properties of the general equilibrium, and the effect of the inventory management on market quality.

5.1 Model

To define the model with inventory management, we consider a setup similar to the benchmark model, but we replace one risk-neutral fast trader with an inventory-averse fast trader (IFT). Specifically, the IFT maximizes an expected utility U of the form (recall that $T = 1$):

$$U = \mathbb{E} \left(\int_0^T (v_1 - p_t) dx_t \right) - C_I \mathbb{E} \left(\int_0^T x_t^2 dt \right), \quad (42)$$

where x_t is his inventory in the risky asset, and $C_I > 0$ is a constant. We call C_I the *inventory aversion coefficient*. We do not identify the exact source of inventory costs for the IFT, but these can be thought to arise either from capital constraints or from risk aversion.

In this model, there are N_F fast traders, N_L slow traders, and one IFT. The equilibrium concept is similar to the linear equilibrium from Section 2. But, because the inventory problem is very difficult in a more general formulation, we assume directly that the speculators' strategies have constant coefficients, and that the dealer has pricing rules as in the benchmark model. Thus, the fast trader $i = 1, \dots, N_F$ and the slow trader $j = 1, \dots, N_L$ have strategies, respectively, of the form:

$$dx_{i,t}^F = \gamma_i dw_t, \quad dx_{j,t}^S = \mu_j \widetilde{dw}_{t-1}. \quad (43)$$

The dealer has pricing rules of the form:

$$dp_t = \lambda dy_t, \quad z_{t-1,t} = \rho dy_{t-1}, \quad (44)$$

where dy_t is the aggregate order flow at t , and $z_{t-1,t} = E_t(dw_{t-1})$ is the dealer's expectation of the current signal given the past order flow. The coefficient λ is chosen so that the dealer breaks even, meaning that her expected profit is zero.²⁵

Since the IFT has quadratic inventory costs, it is plausible to expect that his optimal trading strategy is linear in the inventory.²⁶ Therefore, we assume that the IFT's strategy is of the following type:

$$dx_t = -\Theta x_{t-1} + G dw_t, \quad (45)$$

with constant coefficients $\Theta \in [0, 2)$ and $G \in \mathbb{R}$. Equivalently, the IFT's inventory x_t follows an $AR(1)$ process

$$x_t = \phi x_{t-1} + G dw_t, \quad \phi = 1 - \Theta, \quad (46)$$

with autoregressive coefficient $\phi \in (-1, 1]$.²⁷

If $\Theta > 0$, in each trading round the IFT removes a fraction Θ of his current inventory, with the goal of bringing his inventory eventually to zero. One measure of how quickly the inventory mean reverts to zero is the *inventory half life*. This is defined as the average number of periods (of length dt) that the process needs to halve the distance from its mean, i.e.,

$$\text{Inventory Half Life} = \frac{\ln(1/2)}{\ln(\phi)} dt = \frac{\ln(1/2)}{\ln(1 - \Theta)} dt. \quad (47)$$

²⁵Note that because of inventory management, the aggregate order flow is no longer completely unpredictable by the dealer. Nevertheless, the only source of predictability is the IFT's inventory, and, as we prove later, this inventory in equilibrium is very small because of fast mean reversion. Moreover, not being able to properly compute the expectation of IFT's inventory does not mean that the dealer loses money. Indeed, we have assumed that the dealer chooses λ so that her expected profit is zero.

²⁶This is standard in the literature. See for instance Madhavan and Smidt (1993), but also Hendershott and Menkveld (2014), or Ho and Stoll (1981).

²⁷A standard result is that the $AR(1)$ process becomes explosive (with infinite mean and variance) if ϕ is outside $[-1, 1]$, or equivalently if Θ is outside $[0, 2]$.

Hence, the inventory half life is of the order of dt . This in practice can be short (minutes, seconds, milliseconds), which means that when $\Theta > 0$ the IFT does very quick, “real-time” inventory management.

We end this section with a brief discussion of the different types of inventory management. In Section 5.2 we will see that there is a discontinuity between the cases $\Theta = 0$ and $\Theta > 0$. To explain this discontinuity, we introduce a new case, $\Theta = 0_+$, in which the IFT mean reverts his inventory, but much more smoothly (formal details are below). It turns out that this intermediate inventory management regime indeed connects continuously the other two. Thus, there are three different cases (regimes):

- $\Theta = 0$, the *neutral regime*: the IFT’s strategy is of the form $dx_t = Gdw_t$, similar to the strategy of a (risk-neutral) fast trader.
- $\Theta > 0$, the *fast regime*: the IFT’s strategy is of the form $dx_t = -\Theta x_{t-1} + Gdw_t$. In this regime, the inventory half life is of the order of dt .
- $\Theta = \theta dt$, the *smooth regime*: the IFT’s strategy is of the form $dx_t = -\theta x_{t-1} dt + Gdw_t$, with $\theta \in (0, \infty)$.²⁸ In this regime, the inventory half life $\frac{\ln(1/2)}{\ln(1-\theta dt)} dt = \frac{\ln(2)}{\theta}$, which is much larger than the inventory half life in the fast regime.

The smooth regime is discussed in detail in Internet Appendix K. We find that indeed the smooth regime connects continuously the cases $\Theta = 0$ (neutral regime) with the case $\Theta > 0$ (fast regime).²⁹ However, we show that the smooth regime is not optimal for the IFT when there is enough slow trading (this is true for instance if the $N_L \geq 2$ and $N_F \geq 1$). Therefore, in the rest of the paper we assume that there is enough slow trading, and ignore the smooth regime.

5.2 Optimal Inventory Management

In this section, we do a partial equilibrium analysis, and solve for the optimal strategy of the IFT while fixing the behavior of the other players. This allows us to get insight

²⁸This is called an Ornstein-Uhlenbeck process.

²⁹More precisely, $\theta = 0$ in the smooth regime coincides with $\Theta = 0$; while the limit when $\theta \nearrow \infty$ in the smooth regime coincides with the limit when $\Theta \searrow 0$ in the fast regime.

about the IFT's behavior, without having to do a full equilibrium analysis. We leave this more general analysis to Section 5.3.

Consider the inventory management model with one IFT, N_F fast traders and N_L slow traders. Let γ, μ be the coefficients arising from the strategies of the FTs and STs (not necessarily optimal), and λ, ρ the coefficients from the dealer's pricing rules. Define additional *model coefficients* by:

$$R = \frac{\lambda}{\rho}, \quad \gamma^- = N_F \gamma, \quad \bar{\mu} = N_L \mu, \quad a^- = \rho \gamma^-, \quad b = \rho \bar{\mu}. \quad (48)$$

Proposition 5 analyzes the inventory management regime, where by definition the IFT has a trading strategy with positive mean reversion ($\Theta > 0$). For this result, the strategy need not be optimal.

Proposition 5. *In the inventory management model, let $dx_t = -\Theta x_{t-1} + Gdw_t$ be the IFT's strategy (not necessarily optimal), with $\Theta > 0$. Suppose $b \in (-1, 1)$. Then, the IFT has zero inventory costs, and all his expected profits are in cash. His expected profit π satisfies:*

$$\pi = \lambda \left(\bar{\mu} G \frac{1 - a^-}{1 + \phi b} - G^2 \frac{b + \frac{1}{1+\phi}}{1 + \phi b} \right) \sigma_w^2. \quad (49)$$

Because the IFT reduces his inventory by a fraction $\Theta > 0$ in each trading round, his inventory decays exponentially on average. As our model is set at the high frequency limit (in continuous time), the decrease in IFT inventory is very quick, and the inventory remains infinitesimal at all times.³⁰

In general, the expected profit π of any speculator satisfies:

$$\pi = \mathbb{E} \int_0^T (v_T - p_t) dx_t = \underbrace{\mathbb{E}(v_T x_T)}_{\text{Risky Component}} + \underbrace{\mathbb{E} \int_0^T (-p_t) dx_t}_{\text{Cash Component}}. \quad (50)$$

The *risky component* is the expected profit due to the accumulation of inventory in the risky asset. This does not translate into cash profits until the liquidation date, $T = 1$. The *cash component* is the expected profit that comes from changes in the cash account due to trading. Because the IFT has an infinitesimal inventory, his risky component of

³⁰Mathematically, the average squared inventory $\mathbb{E}(x_t^2)$ is of the order of dt ; see equation (A32).

profits is negligible. Hence, all IFT profits come from the cash component, as stated in Proposition 5.

As a result of keeping all his profits in cash, the behavior of the IFT is very different than the behavior of a risk-neutral speculator. Indeed, while the risk-neutral speculator trades directly on his private information, the IFT benefits only indirectly, from timing his trades and unloading his inventory to slower traders.

To understand why the IFT behaves differently, suppose he observes a new signal dw_t . Initially, the IFT trades on his signal (Gdw_t), but subsequently he fully reverses his trade by unloading a positive fraction of his inventory each period. Therefore, the only way for the IFT to make money is to ensure that the inventory reversal is done at a profit. This can occur for instance if the IFT expects that when he sells, other traders buy even more, and as a result his overall price impact is negative. But this is only possible if there exist slow traders, whose lagged signals can be predicted by the IFT.

In general, the expected profit of a speculator who manages inventory satisfies the following formula:³¹

$$\pi = \mathbb{E} \int_0^T x_{t-1} dp_t. \quad (51)$$

Thus, inventory management is profitable only when the speculator can use his past inventory (x_{t-1}) to forecast the current price change (dp_t). In particular, this formula explains why the IFT trades at $t - 1$ an amount Gdw_{t-1} even though he knows that subsequently he will fully reverse his trade. He trades like this because his signal dw_{t-1} anticipates the slow trading at t , which in turn affects dp_t . To make this intuition more precise, the next result specializes the formula (51) to our model.

Corollary 3. *In the context of Proposition 5, the IFT's expected profit satisfies:*

$$\pi = \mathbb{E} \int_0^T x_{t-1} \left(\lambda \bar{\mu} \widetilde{dw}_{t-1} \right) - \lambda \Theta \mathbb{E} \int_0^T x_{t-1}^2. \quad (52)$$

If there is no slow trading ($\bar{\mu} = 0$), the IFT's makes negative expected profits.

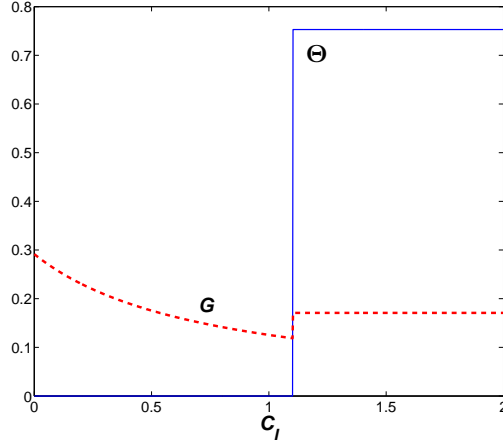
Using Corollary 3, we see that the IFT's speculative trade Gdw_{t-1} is part of the

³¹For the IFT, see equation (A35). The result is true in general when the speculator has infinitesimal inventory. Indeed, if we integrate $d(x_t p_t) = p_t dx_t + x_{t-1} dp_t$, we get $x_T p_T$, which is zero in expectation since x_T is infinitesimal. Hence, $\pi = \mathbb{E} \int_0^T (-p_t) dx_t = \mathbb{E} \int_0^T x_{t-1} dp_t$.

IFT's inventory x_{t-1} , and is also correlated with the price change dp_t via the slow trading component of the order flow $\bar{\mu} \widetilde{dw}_{t-1}$. Without slow trading ($\bar{\mu} = 0$), there is no correlation, hence no revenue source for the IFT. Therefore, the IFT makes negative profits on average, as he loses from the price impact of his trades.

The main result of this section describes the IFT's optimal strategy when there is enough slow trading, i.e., the slow trading coefficient b is above a threshold.

Figure 4: Optimal IFT Inventory Management. This figure plots the coefficients of the IFT's optimal trading strategy ($dx_t = -\Theta x_{t-1} + G dw_t$) in the inventory management model with $N_F = 5$ fast traders and $N_L = 5$ slow traders. On the horizontal axis is the IFT's inventory aversion, C_I . The parameter values are $\sigma_w = 1$, $\sigma_u = 1$. For the model coefficients, we use the equilibrium values from Section 5.3: $a^- = 0.7088$, $b = 0.5424$, $\lambda = 0.3782$, $\rho = 0.3439$. The formulas for G , Θ , and \bar{C}_I are from Theorem 2.



Theorem 2. *In the inventory management model, suppose the model coefficients satisfy $0 \leq a^-, b < 1$ and $\lambda, \rho > 0$. In addition, suppose $b > \frac{\sqrt{17}-1}{8} = 0.3904$.³² Let $\bar{C}_I = 2\lambda \left(\frac{(1-Ra^-)^2(1+\sqrt{1-b})^2}{R^2b(1-a^-)^2} - 1 \right)$. Then, if $C_I < \bar{C}_I$, the optimal strategy of the IFT is to set*

$$\Theta = 0, \quad G = \frac{1 - Ra^-}{2\lambda + C_I}. \quad (53)$$

If $C_I > \bar{C}_I$, the optimal strategy of the IFT is to set

$$\Theta = 2 - \frac{\sqrt{1-b}}{b} \in (0, 2), \quad G = \frac{1 - a^-}{2\rho \left(1 + \frac{1}{\sqrt{1-b}} \right)}. \quad (54)$$

³²In equilibrium (section 5.3) we have the following numerical results: the condition $b < 1$ is always satisfied, and the condition $b > \frac{\sqrt{17}-1}{8}$ is equivalent to having (i) $N_L \geq 2$ and (ii) $N_L \geq 6$ if $N_F = 0$.

Thus, there are two different types of behavior (regimes) for the IFT, depending on how his inventory aversion compares to a threshold value (\bar{C}_I).

- (*Neutral regime*) If the inventory aversion coefficient is small (below \bar{C}_I), the IFT sets $\Theta = 0$ and controls his inventory by choosing his weight G . As his inventory aversion gets larger, the IFT reduces his inventory costs by decreasing G . The tradeoff is that a smaller G also reduces expected profits. The behavior of the IFT when $\Theta = 0$ is essentially the same as the behavior of a FT.
- (*Fast regime*) If the inventory aversion is large (above \bar{C}_I), the IFT manages his inventory by choosing a positive mean reversion coefficient ($\Theta > 0$). There is no longer a tradeoff between expected profit and inventory costs, as the IFT has zero inventory costs. Hence, the IFT chooses the weight G and the mean reversion Θ to maximize expected profit (more details below).

Theorem 2 implies that a small change in IFT's inventory aversion can have a large effect on the IFT's behavior. Figure 4 plots the coefficients of the optimal strategy when $N_F = 5$, $N_S = 5$. We see that when the IFT's inventory aversion rises above the threshold $\bar{C}_I = 0.1021$, his optimal mean reversion coefficient jumps from $\Theta = 0$ to $\Theta = 0.7530$. Also, his optimal weight jumps from $G = 0.1186$ (the left limit of G at the threshold) to $G = 0.1708$ (the constant value of G above the threshold).

The sharp discontinuity between the two regimes arises because the IFT has zero inventory costs in the fast regime ($\Theta > 0$). Let $U_{\Theta > 0}^*$ and $U_{\Theta = 0}^*$ be the maximum expected utility of the IFT respectively when he manages inventory versus when he does not. Because the IFT has zero inventory costs in the fast regime, $U_{\Theta > 0}^*$ does not depend on C_I ; while in the neutral regime $U_{\Theta = 0}^*$ is decreasing in C_I and is precisely equal to $U_{\Theta > 0}^*$ at the threshold value $C_I = \bar{C}_I$.³³ This implies that the neutral regime is optimal when C_I is below the threshold, while the fast regime is optimal when C_I is above the threshold.

Note that a necessary condition for Theorem 2 is the existence of enough slow trading. Formally, the slow trading coefficient must be larger than the threshold $b = 0.3904$, which numerically is true for instance if there are $N_F \geq 1$ fast traders and $N_L \geq 2$ slow traders. If slow trading is below the threshold, we show that a similar analysis holds,

³³Formally, these statements follow from equations (A43) and (A39) in the Appendix.

but with the fast regime replaced by the *smooth regime*, in which the IFT still manages inventory but with a strategy of the form: $dx_t = -\theta x_t dt + Gdw_t$. (See Section J.1 in the Internet Appendix.) To simplify presentation, we assume that there is enough slow trading, and ignore the smooth regime in the rest of the paper.

Using our results, we predict that in practice fast speculators are sharply divided into two categories. In both categories speculators generate large trading volume. But in one category the speculators make fundamental bets and accumulate inventories, while in the other category speculators mean revert their inventories very quickly, and keep their profits in cash. Our results appear consistent with the “opportunistic traders” and the “high frequency traders” described in Kirilenko *et al.* (2014). Both opportunistic traders and HFTs have large volume and appear to be fast. But while opportunistic traders have relatively large inventories, the HFTs in their sample (during several days around the Flash Crash of May 6, 2010) liquidate 0.5% of their aggregate inventories on average each second. This implies that HFT inventories have an $AR(1)$ half life of a little over 2 minutes.

We finish this section with a brief discussion of how the IFT’s optimal strategy is correlated with slow trading. Corollary 3 shows that if there is no slow trading, the IFT cannot make positive profits. Theorem 2 shows that with enough slow trading, the IFT can manage inventory and make positive profits (see equation (A43) in the Appendix). In the previous discussion, we have argued that this is possible only if the IFT trades in the opposite direction to the slow trading. We now prove this is indeed the case.

Corollary 4. *In the context of Theorem 2, suppose the IFT is sufficiently averse ($C_I > \bar{C}_I$). Denote by $d\bar{x}_t^S = \bar{\mu}\widetilde{dw}_{t-1}$ the slow trading component of the speculator order flow. Then, the IFT’s optimal strategy is negatively correlated with slow trading:*

$$\text{Cov}(dx_t, d\bar{x}_t^S) = -\Theta \text{Cov}(x_{t-1}, d\bar{x}_t^S) < 0. \quad (55)$$

We call this phenomenon the *hot potato effect*, or the *intermediation chain effect*. The intuition is that the IFT’s current signal generates undesirable inventory and must be passed on to slower traders in order to produce a profit. The passing of inventory can be thought as the beginning of an intermediation chain. Kirilenko *et al.* (2014) and

Weller (2014) document such hot potato effects among high frequency traders.

5.3 Equilibrium Results

In this section, we solve for the full equilibrium of the inventory management model. For simplicity, we assume that the IFT is *sufficiently averse*, meaning that his inventory aversion is above a certain threshold (formally, above the threshold value \bar{C}_I from Theorem 2). Then, the solution can be expressed almost in closed form, except for the slow trading coefficient b , which satisfies a non-linear equation in one variable.

Theorem 3. *Consider the inventory management model with one sufficiently averse IFT, N_F fast traders, and N_L slow traders. Suppose there is an equilibrium in which the speculators's strategies are: $dx_t = -\Theta x_{t-1} + Gdw_t$ (the IFT), $dx_t^F = \gamma dw_t$ (the FTs), $dx_t^S = \mu \widetilde{dw}_{t-1}$ (the STs); and the dealer's pricing rules are: $dp_t = \lambda dy_t$, $\widetilde{dw}_t = dw_t - \rho dy_t$. Denote the model coefficients R , a^- , b as in (48). Suppose $\frac{\sqrt{17}-1}{8} < b < 1$. Then, the equilibrium coefficients satisfy equations (A44)–(A46) from the Appendix.*

Conversely, suppose the equations (A44)–(A46) have a real solution such that $\frac{\sqrt{17}-1}{8} < b < 1$, $a < 1$, $\lambda > 0$. Then, the speculators' strategies and the dealer's pricing rules with these coefficients provide an equilibrium of the model.

Rather than relying on numerical results to study the equilibrium, we start by providing asymptotical results when the number of FTs and STs is large. The advantage is that the asymptotic results can be expressed in closed form, and thus help provide a clearer intuition for the equilibrium. Let \bar{C}_I be the threshold aversion from Theorem 2. Let π be the expected profit of a sufficiently averse IFT ($C_I \geq \bar{C}_I$), and $\pi^{C_I=0}$ the maximum expected profit of a risk-neutral IFT ($C_I = 0$), where the behavior of the other speculators and the dealers is taken to be the same. Let γ_0 the benchmark FT weight, and $\pi_0^F = \frac{\gamma_0}{N_F+2} \sigma_w^2$ the benchmark profit of a FT, as in Proposition (1). We use the asymptotic notation: $X \approx X_\infty$ stands for $\lim_{N_F, N_L \rightarrow \infty} \frac{X}{X_\infty} = 1$.

Proposition 6. *Consider (i) the inventory management model with one sufficiently averse IFT, N_F fast traders, and N_L slow traders, and (ii) the benchmark model with $N_F + 1$ fast traders and N_L slow traders. Then, the equilibrium coefficients γ , μ , λ , ρ*

are asymptotically equal across the two models when N_F and N_L are large. Also, $a \approx 1$, $b \approx b_\infty = 0.6180$, and we have the following asymptotic formulas:

$$\begin{aligned} \Theta &\approx 1, & \frac{G}{\gamma_0} &\approx 1 - b_\infty = 0.3820, & \frac{\pi}{\pi_0^F} &\approx 2b_\infty - 1 = 0.2361, \\ \frac{\pi}{\pi^{C_I=0}} &\approx \frac{4}{5} b_\infty = 49.44\%, & \bar{C}_I &\approx \frac{1+5b_\infty}{2} \lambda_\infty \approx 2.0451 \frac{\sigma_w}{\sigma_u} \frac{1}{\sqrt{N_F+1}}. \end{aligned} \quad (56)$$

The first implication of Proposition 6 is that model with inventory management is asymptotically the same as the benchmark model when both N_F and N_L are large. This is not surprising, since when there are many other speculators, the IFT has a relatively smaller and smaller role in the limit.

The behavior of the IFT is more surprising. First, when there are many other speculators, the IFT's inventory mean reversion becomes extreme (Θ approaches 1). This means that the IFT's inventory half life becomes essentially zero, as the IFT removes most of his inventory each period. This extreme mean reversion is possible because the existence of a sufficient amount of slow trading allows the hot potato effect to generate positive profits for the IFT. Furthermore, the equation $\pi \approx 49.44\% \times \pi^{C_I=0}$ implies that even under extreme inventory mean reversion ($\Theta = 1$) the IFT can trade so that he only loses on average only about 50% of his maximum expected profits when he has zero inventory aversion.³⁴

The equation $\bar{C}_I \approx 2.0451 \frac{\sigma_w}{\sigma_u} \frac{1}{\sqrt{N_F+1}}$ implies that the threshold inventory aversion above which the IFT chooses to mean revert his inventory becomes very small when the number of competing fast traders is large. This is perhaps counterintuitive, since one may think that the IFT chooses fast inventory mean reversion because he has very high inventory aversion. This is not the case, however. Indeed, even when the IFT has small inventory aversion, a sufficient amount of slow trading is enough to convince the IFT to engage in very fast inventory mean reversion. This is because inventory management is a zero/one proposition. Once the IFT engages in inventory management ($\Theta > 0$), any profits from fundamental bets become zero, and the hot potato effect is the sole source of profits.

³⁴This recalls the saying attributed to Joseph Kennedy (the founder of the Kennedy dynasty) that “I would gladly give up half my fortune if I could be sure the other half would be safe.”

We now compare the IFT with the other speculators. For the IFT, we consider the following variables: (i) IFT's trading volume, measured by the his order flow variance $TV_x = \text{Var}(dx_t)/dt$, as in Section 4.1, (ii) IFT's order flow autocorrelation, $\rho_x = \text{Corr}(dx_t, dx_{t+1})$; and (iii) $\beta_{x,\bar{x}^S} = \text{Cov}(dx_t, d\bar{x}_t^S) / \text{Var}(d\bar{x}_t^S)$, which is the regression coefficient of the IFT's strategy (dx_t) on the slow trading component ($d\bar{x}_t^S$). We are also interested in the individual FT volume, TV_{x^F} ; the aggregate FT volume, $TV_{\bar{x}^F}$; and the aggregate ST volume, $TV_{\bar{x}^S}$; the aggregate FT order flow autocorrelation, $\rho_{\bar{x}^F}$; and the aggregate ST order flow autocorrelation, $\rho_{\bar{x}^S}$.

The next result computes all these quantities, and provides asymptotic results when both N_F and N_L are large. Some of these results provide new testable implications, regarding the relationship between trading volume, order flow covariance, and inventory.

Proposition 7. *In the context of Theorem 3, consider a sufficiently averse IFT. Then, the variables defined above satisfy the following formulas:*

$$\begin{aligned} \frac{TV_x}{TV_{x^F}} &= \frac{2G^2}{(1+\phi)\gamma^2} \approx 4 - 6b_\infty = 0.2918, & \frac{TV_{\bar{x}^S}}{TV_{\bar{x}^F}} &= \frac{b^2(1-a)}{(a^-)^2} \approx \frac{b_\infty}{N_F + 1}, \\ \rho_x &= -\frac{\Theta}{2} \approx -\frac{1}{2}, & \rho_{\bar{x}^F} &= 0, & \rho_{\bar{x}^S} &\approx -b_\infty = -0.6180, \\ \beta_{x,\bar{x}^S} &= -\frac{\Theta(1-a^-)}{2b(1+2\sqrt{1-b})} \approx -\frac{3+b_\infty}{5(N_F+1)} = -\frac{0.7236}{N_F+1}. \end{aligned} \tag{57}$$

The last result illustrates the hot potato effect. The IFT's order flow has a negative beta on the STs' aggregate order flow, which means that the IFT and the STs trade in opposite directions. As the number of FTs becomes larger, there is more information released to the public by the trades of the fast traders, hence there is less room for slow trading. As a result, the hot potato effect is less intense when there is a large number of FTs.

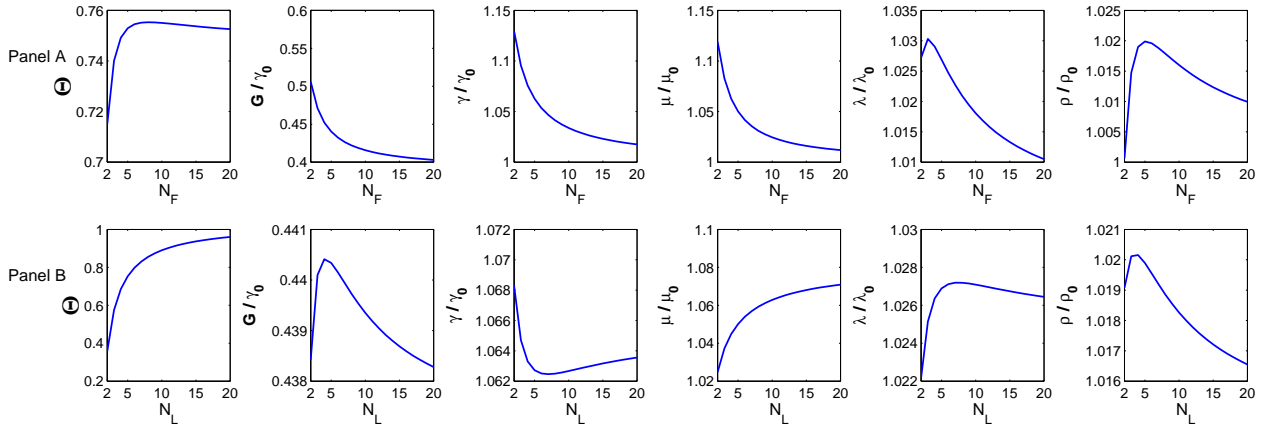
Proposition 7 implies that in the limit when N_F and N_L are large, the IFT's trading volume is about 30% of the individual FT trading volume. This implies that the IFT's trading volume is comparable to that of a regular FT. By contrast, just as in the benchmark model, the volume coming from STs is much smaller than the volume coming from FTs. This confirms our intuition that in an empirical analysis which selects traders based on volume, the IFT and the FTs are in the category with large trading volume,

while the STs are in the category with small trading volume.

If we compare order flow autocorrelations, we see that the IFT is similar to the STs, but not to the FTs. Indeed, the IFT and the STs have negative and large order flow autocorrelation. By contrast, the FTs have zero order flow autocorrelation.³⁵ Finally, if we compare inventories, the IFT has infinitesimal inventory, while the variance of the other speculators' inventory increases over time.³⁶ Nevertheless, the STs' inventories are smaller relative to the FTs' inventories, since the STs have smaller volume.

We now present some numerical results for the equilibrium coefficients. Figure 5 plots the equilibrium coefficients (Θ , G , γ , μ , λ , ρ). We normalize some variables X in the inventory management model by the corresponding variable X_0 in the benchmark model. Panel A of Figure 5 plots the variables against N_F , while holding N_L constant. Panel B plots the same variables against N_L , while holding N_F constant.³⁷

Figure 5: Equilibrium Coefficients with Inventory Management. This figure plots the equilibrium coefficients that arise in the inventory management model. Some variables X are normalized by the corresponding variable X_0 in the benchmark model. The coefficients are Θ , G , γ , μ , λ , ρ . Panel A plots the dependence of the six variables on the number of fast traders N_F , while taking the number of slow traders $N_S = 5$. Panel B plots the dependence of the six variables on N_S , while taking $N_F = 5$. The other parameters are $\sigma_w = 1$, $\sigma_u = 1$.



As expected, we find that the mean reversion coefficient Θ is increasing in the number

³⁵Even if we allowed FTs to trade on lagged signals, one can see that the FTs would still have very small order flow autocorrelation (of the order of $\frac{1}{N_F+1}$) because of their large trading volume.

³⁶For the IFT, $\text{Var}(x_t) = \frac{G^2}{1-\phi^2} \sigma_w^2 dt$ (see equation (A25) in the Appendix); while for the FT, $\text{Var}(x_t^F) = t\sigma_w^2$, as the FT's inventory follows a random walk.

³⁷We consider $N_F, N_L \geq 2$. The reason is that in order to apply Theorem 2, we need to have $b > \frac{\sqrt{17}-1}{8}$. This is true in equilibrium if $N_F, N_L \geq 2$.

of slow traders N_L . This is because the IFT needs slow traders in order to make profits. The IFT's weight G is less than half the benchmark weight γ_0 , indicating that the IFT shifts towards inventory management in order to make profits. This leaves more room for fundamental profits, which explains why both the FTs and the STs are better off with inventory management than in the benchmark model (γ/γ_0 and μ/μ_0 are both above one), despite the price impact λ being larger than in the benchmark (we see that $\lambda/\lambda_0 > 1$). The reason why the market is more illiquid in the inventory management model is that the IFT trades much less intensely on his signal (G is less than half of γ_0), and therefore the informational efficiency is lower. To see directly that the market is less informationally efficient in the inventory management model, we use the fact that in our model price volatility is a proxy for informational efficiency (see the discussion in Section 4). Then, we verify numerically that indeed $\sigma_p/\sigma_{p,0} < 1$, which implies that with inventory management the market is less informationally efficient.

6 Conclusion

We have presented a theoretical model in which traders continuously receive signals over time about the value of an asset, but only use each signal for a finite number of lags (which can be justified by an information processing cost per signal). We have found that competition among speculators reveals much private information to the public, and the value of information decays fast. Therefore, a trader who is just one instant slower than the other traders loses the majority of the profits by being slow. Another consequence is that the market is very efficient and liquid. As a feedback effect, because of the small price impact (high market liquidity), the informed traders are capable of trading even more aggressively. In equilibrium, the fast speculator trading volume is very large and dominates the overall trading volume. We have also considered an extension of the model in which a fast speculator, called the inventory-averse fast trader (IFT), has quadratic inventory costs. We find that a sufficiently averse IFT has a very different behavior compared to a risk-neutral fast trader. The IFT keeps his profits in cash, makes no fundamental bets on the value of the risky asset, and quickly passes his inventory to “slow traders,” who use their lagged signals. This hot potato effect is

possible because the existence of slower traders more than reverses the price impact of the IFT.

Appendix A. Proofs

Notation Preliminaries

Recall that $t + 1$ is notation for $t + dt$, and

$$T = 1. \tag{A1}$$

In general, a tilde above a symbol denotes normalization by σ_w . For instance, if σ_u is the instantaneous volatility of the noise trader order flow, and σ_y the instantaneous volatility of the total order flow, we denote by

$$\tilde{\sigma}_u = \frac{\sigma_u}{\sigma_w}, \quad \tilde{\sigma}_y = \frac{\sigma_y}{\sigma_w}, \quad \text{with} \quad \sigma_u^2 = \frac{\text{Var}(du_t)}{dt}, \quad \sigma_y^2 = \frac{\text{Var}(dy_t)}{dt}. \tag{A2}$$

Consider a trading strategy dx_t for $t \in (0, T]$, where $T = 1$. Denote by $\tilde{\pi}$ the normalized expected profit at $t = 0$ from the strategy dx_t :

$$\tilde{\pi} = \frac{1}{\sigma_w^2} \mathbb{E} \left(\int_0^T (w_t - p_t) dx_t \right). \tag{A3}$$

For variances and covariances, a tilde the symbol means normalization by both σ_w^2 and dt . For instance, denote by

$$\widetilde{\text{Var}}(\widetilde{dw}_t) = \frac{\text{Var}(\widetilde{dw}_t)}{\sigma_w^2 dt} = A_t, \quad \widetilde{\text{Cov}}(w_t, \widetilde{dw}_t) = \frac{\text{Cov}(w_t, \widetilde{dw}_t)}{\sigma_w^2 dt} = B_t. \tag{A4}$$

Proof of Theorem 1. We look for an equilibrium with the following properties: (i) the equilibrium is symmetric, in the sense that the FTs have identical trading strategies, and the same for the STs; (ii) the equilibrium coefficients are constant with respect to time.

To solve for the equilibrium, in the first step we take the dealer's pricing functions as given, and solve for the optimal trading strategies for the FTs and STs. In the second

step, we take the speculators' trading strategies as given, and we compute the dealer's pricing functions. In Section 2, we have assumed that the speculators take the signal covariance structure as given (see equation (13)). In the current context, this means that the speculators take the following covariances as given and constant:

$$A_t = \widetilde{\text{Var}}(\widetilde{dw}_t) = \frac{\text{Var}(\widetilde{dw}_t)}{\sigma_w^2 dt}, \quad B_t = \widetilde{\text{Cov}}(w_t, \widetilde{dw}_t) = \frac{\text{Cov}(w_t, \widetilde{dw}_t)}{\sigma_w^2 dt} \quad (\text{A5})$$

Thus, in the rest of the Appendix we consider that the dealer also sets A and B , in addition to setting λ and ρ .

Speculators' Optimal Strategy (γ, μ)

Since we search for an equilibrium with constant coefficients, we assume that the speculators take as given the dealer's pricing rules $dp_t = \lambda dy_t$ and $z_{t-1,t} = \rho dy_{t-1}$, and also the covariances $A = \widetilde{\text{Var}}(\widetilde{dw}_t)$ and $B = \widetilde{\text{Cov}}(w_t, \widetilde{dw}_t)$.

Consider a FT, indexed by $i = 1, \dots, N_F$. He chooses $dx_t^i = \gamma_t^i dw_t + \mu_t^i \widetilde{dw}_{t-1}$, and assumes that at each $t \in (0, T]$, the price satisfies:

$$dp_t = \lambda dy_t, \quad \text{with} \quad dy_t = (\gamma_t^i + \gamma_t^{-i}) dw_t + (\mu_t^i + \mu_t^{-i}) \widetilde{dw}_{t-1} + du_t, \quad (\text{A6})$$

where the superscript “ $-i$ ” indicates the aggregate quantity from the other speculators. Since dw_t and \widetilde{dw}_{t-1} are both orthogonal on the public information set \mathcal{I}_t , and $p_{t-1} \in \mathcal{I}_t$, it follows that dx_t^i is orthogonal to p_{t-1} as well. The normalized expected profit of FT i at $t = 0$ satisfies:

$$\begin{aligned} \tilde{\pi}^F &= \frac{1}{\sigma_w^2} \mathbf{E} \int_0^T \left(w_t - p_{t-1} - \lambda \left((\gamma_t^i + \gamma_t^{-i}) dw_t + (\mu_t^i + \mu_t^{-i}) \widetilde{dw}_{t-1} + du_t \right) \right) dx_t^i \\ &= \gamma_t^i - \lambda \gamma_t^i (\gamma_t^i + \gamma_t^{-i}) + \mu_t^i B - \lambda \mu_t^i (\mu_t^i + \mu_t^{-i}) A. \end{aligned} \quad (\text{A7})$$

This is a pointwise optimization problem, hence it is enough to consider the profit at $t = 0$, and maximize the expression over γ_t^i and μ_t^i . The solution of this problem is $\lambda \gamma_t^i = \frac{1 - \lambda \gamma_t^{-i}}{2}$, and $\lambda \mu_t^i = \frac{B/A - \lambda \mu_t^{-i}}{2}$. The ST $j = 1, \dots, N_S$ solves the same problem, only that his coefficient on dw_t is $\gamma_t^j = 0$. Thus, all γ 's are equal for the FTs, and all μ 's are equal for the FTs and STs. We also find that they are constant, and since

$N_L = N_F + N_S$, we have

$$\gamma = \frac{1}{\lambda} \frac{1}{1 + N_F}, \quad \mu = \frac{B/A}{\lambda} \frac{1}{1 + N_L}. \quad (\text{A8})$$

Dealer's Pricing Rules (λ, ρ, A, B)

The dealer takes the speculators' strategies as given, and assumes that the aggregate order flow is of the form:

$$dy_t = du_t + \bar{\gamma} dw_t + \bar{\mu} \widetilde{dw}_{t-1}, \quad \text{with} \quad \bar{\gamma} = N_F \gamma, \quad \bar{\mu} = N_L \mu. \quad (\text{A9})$$

Moreover, the dealer assumes that, in their trading strategy, the speculators set:

$$\widetilde{dw}_{t-1} = dw_{t-1} - \rho^* dy_{t-1}. \quad (\text{A10})$$

Naturally, later we require that in equilibrium the dealer's pricing coefficient ρ coincides with the coefficient ρ^* used by the speculators.

Since the order flow dy_t is orthogonal to the dealer's information set \mathcal{I}_t , the dealer sets $\lambda_t, \rho_t, A_t, B_t$ such that the following equations are satisfied:

$$\begin{aligned} \lambda_t &= \frac{\widetilde{\text{Cov}}(w_t, dy_t)}{\widetilde{\text{Var}}(dy_t)} = \frac{\bar{\gamma} + \bar{\mu} B_{t-1}}{\sigma_{y,t}^2}, & dp_t &= \lambda_t dy_t, \\ \rho_t &= \frac{\widetilde{\text{Cov}}(dw_t, dy_t)}{\widetilde{\text{Var}}(dy_t)} = \frac{\bar{\gamma}}{\sigma_{y,t}^2}, & \widetilde{dw}_t &= dw_t - \rho_t dy_t, \\ \sigma_{y,t}^2 &= \widetilde{\text{Var}}(dy_t^2) = \tilde{\sigma}_u^2 + \bar{\gamma}^2 + \bar{\mu}^2 A_{t-1}, & & \\ B_t &= \widetilde{\text{Cov}}(w_t, dw_t - \rho^* dy_t) = (1 - \rho^* \bar{\gamma}) - \rho^* \bar{\mu} B_{t-1}, & & \\ A_t &= \widetilde{\text{Var}}(dw_t - \rho^* dy_t) = 1 - 2\rho^* \bar{\gamma} + (\rho^*)^2 \sigma_{y,t}^2 \\ &= 1 - 2\rho^* \bar{\gamma} + (\rho^*)^2 (\tilde{\sigma}_u^2 + \bar{\gamma}^2) + (\rho^*)^2 \bar{\mu}^2 A_{t-1}. \end{aligned} \quad (\text{A11})$$

Consider the last equation in (A11), $A_t = 1 - 2\rho^* \bar{\gamma} + (\rho^*)^2 (\tilde{\sigma}_u^2 + \bar{\gamma}^2) + (\rho^*)^2 \bar{\mu}^2 A_{t-1}$, which is a recursive equation in A_t . Then, Lemma A.1 (below) implies that A does not depend on t , as long as $|\rho^* \bar{\mu}| < 1$. But, since the dealer takes the speculators' strategies as given, we can use the equilibrium condition $\rho^* \bar{\mu} = b \in (0, 1)$. The same method shows that

B does not depend on t . Moreover, Lemma A.1 can be used to compute the constant values of A and B :

$$A = \frac{(1 - \rho^* \bar{\gamma})^2 + (\rho^*)^2 \tilde{\sigma}_u^2}{1 - (\rho^* \bar{\mu})^2}, \quad B = \frac{1 - \rho^* \bar{\gamma}}{1 + \rho^* \bar{\mu}}. \quad (\text{A12})$$

Then, equation (A11) shows that λ , ρ , $\tilde{\sigma}_y$ are independent on t as well.

Equilibrium Conditions

We now use the equations derived above to solve for the equilibrium values of γ , μ , λ , $\rho = \rho^*$, A , B , $\tilde{\sigma}_y$. Denote by

$$a = \rho \bar{\gamma}, \quad b = \rho \bar{\mu}, \quad R = \frac{\lambda}{\rho}. \quad (\text{A13})$$

From (A12) we have $A = \frac{(1-a)^2 + \rho^2 \tilde{\sigma}_u^2}{1-b^2}$. Then, substitute A in $\tilde{\sigma}_y^2 = \tilde{\sigma}_u^2 + \bar{\gamma}^2 + \bar{\mu}^2 A$ from (A11), to obtain $\rho^2 \tilde{\sigma}_y^2 = \frac{\rho^2 \tilde{\sigma}_u^2 + (a^2 + b^2 - 2ab^2)}{1-b^2}$. To summarize,

$$B = \frac{1-a}{1+b}, \quad A = \frac{(1-a)^2 + \rho^2 \tilde{\sigma}_u^2}{1-b^2}, \quad \rho^2 \tilde{\sigma}_y^2 = \frac{\rho^2 \tilde{\sigma}_u^2 + (a^2 + b^2 - 2ab^2)}{1-b^2}. \quad (\text{A14})$$

Using (A11), we get $R = \frac{\lambda}{\rho} = \frac{\bar{\gamma} + \bar{\mu} B}{\bar{\gamma}} = \frac{a+b \frac{1-a}{1+b}}{a} = \frac{a+b}{a(1+b)}$. Also, the equation for ρ implies $\rho = \frac{\bar{\gamma}}{\tilde{\sigma}_y^2} = \frac{\rho a}{\rho^2 \tilde{\sigma}_y^2}$. Using the formula for $\rho^2 \tilde{\sigma}_y^2$ in (A14), we compute $\rho^2 \tilde{\sigma}_u^2 = (1-a)(a-b^2)$. Using this formula, we obtain $\rho^2 \tilde{\sigma}_y^2 = a$ and $A = 1-a$. To summarize,

$$R = \frac{\lambda}{\rho} = \frac{a+b}{a(1+b)}, \quad \rho^2 \tilde{\sigma}_u^2 = (1-a)(a-b^2), \quad \rho^2 \tilde{\sigma}_y^2 = a, \quad A = 1-a. \quad (\text{A15})$$

From (A8), we have $\frac{N_F}{N_F+1} = \lambda \bar{\gamma} = \frac{\lambda}{\rho} a = \frac{a+b}{1+b}$. From this, $a = \frac{N_F - b}{N_F + 1}$, and $B = \frac{1-a}{1+b} = \frac{\frac{1+b}{N_F+1}}{1+b} = \frac{1}{N_F+1}$. Also, $\frac{B}{A} \frac{N_L}{N_L+1} = \lambda \bar{\mu} = \frac{\lambda}{\rho} b = \frac{b(a+b)}{a(1+b)}$. Since $\frac{B}{A} = \frac{1}{1+b}$, we have $\frac{N_L}{N_L+1} = \frac{b(a+b)}{a}$, or $\frac{a}{b(1+b)} \frac{N_L}{N_L+1} = \frac{a+b}{1+b}$. The two formulas for $\frac{a+b}{1+b}$ imply $b(1+b) \frac{N_F}{N_F+1} = a \frac{N_L}{N_L+1}$. To summarize,

$$a = \frac{N_F - b}{N_F + 1}, \quad B = \frac{1}{N_F + 1}, \quad b(1+b) \frac{N_F}{N_F + 1} = \frac{N_F - b}{N_F + 1} \frac{N_L}{N_L + 1}. \quad (\text{A16})$$

From $\frac{\lambda}{\rho} a = \frac{N_F}{N_F+1}$ and $a = \frac{N_F - b}{N_F + 1}$, we get $\frac{\lambda}{\rho} = \frac{N_F}{N_F - b}$, as stated.

From (A16), we obtain the quadratic equation $b^2 + b\omega = \frac{N_L}{N_L+1}$, with $\omega = 1 + \frac{1}{N_F} \frac{N_L}{N_L+1}$. One solution of this quadratic equation is $b = \frac{\omega + (\omega + 4\frac{N_L}{N_L+1})^{1/2}}{2} \geq 1$, which leads to a negative $\tilde{\sigma}_y^2$ (see (A14)). Thus, we must choose the other solution, $b = \frac{-\omega + (\omega + 4\frac{N_L}{N_L+1})^{1/2}}{2} \geq 0$. Let $b_\infty = \frac{\sqrt{5}-1}{2}$. Since $b_\infty^2 + b_\infty = 1$ and $\omega \geq 1$, we have $b_\infty^2 + b_\infty\omega \geq 1$. Moreover, since $b^2 + b\omega = \frac{N_L}{N_L+1} < 1$, we get $b^2 + b\omega < b_\infty^2 + b_\infty\omega$. But the function $b^2 + b\omega$ is strictly increasing in b when $b \geq 0$, hence we obtain $b < b_\infty$. Thus, $b \in [0, b_\infty)$, as stated in the Theorem. We also obtain $a = \frac{N_F - b}{N_F + 1} \in (0, 1)$. The proof of the exact formulas in (19) is now complete.

We now derive the asymptotic formulas in (19). When N_F is large, note that $a = \frac{N_F}{N_F - b} \approx a_\infty = 1$, $\omega = 1 + \frac{1}{N_F} \frac{N_L}{N_L+1} \approx \omega_\infty = 1$. Therefore, we also get $b \approx b_\infty = \frac{\sqrt{5}-1}{2}$. One can now verify that the formulas for γ_∞ , μ_∞ , λ_∞ , and ρ_∞ are as stated in (19).

We now show how b depends on N_F and N_L (the dependence on N_S is the same as the dependence on $N_L = N_F + N_S$). Consider the function $F(\beta, \omega) = \sqrt{\omega^2 + 4\beta} - \omega$, and note that $b = F(\beta, \omega)/2$, with $\beta = \frac{N_L}{N_L+1}$ and $\omega = 1 + \frac{\beta}{N_F}$. We compute $\frac{\partial b}{\partial N_F} = \frac{\partial \beta}{\partial N_L} = \frac{1}{(N_L+1)^2}$, $\frac{\partial \omega}{\partial N_F} = -\frac{N_L(N_L+1) - N_F}{N_F^2(N_L+1)^2} < 0$, $\frac{\partial \omega}{\partial N_L} = \frac{1}{N_F(N_L+1)^2} > 0$. Also, $\frac{\partial F}{\partial \beta} = \frac{2}{\sqrt{\omega^2 + 4\beta}} > 0$, and $\frac{\partial F}{\partial \omega} = \frac{\beta}{\sqrt{\omega^2 + 4\beta}} - 1 = -\frac{b}{\sqrt{\omega^2 + 4\beta}} < 0$. Then, $\frac{\partial(2b)}{\partial N_F} = \frac{\partial F}{\partial \beta} \cdot \frac{\partial \beta}{\partial N_F} + \frac{\partial F}{\partial \omega} \cdot \frac{\partial \omega}{\partial N_F} > 0$, and $\frac{\partial(2b)}{\partial N_L} = \frac{\partial F}{\partial \beta} \cdot \frac{\partial \beta}{\partial N_L} + \frac{\partial F}{\partial \omega} \cdot \frac{\partial \omega}{\partial N_L} = \frac{1}{(N_L+1)^2 \sqrt{\omega^2 + 4\beta}} (2 - \frac{b}{N_F}) > 0$, where the last inequality follows from $b \in (0, 1)$.

We end the analysis of the equilibrium conditions, by proving several more useful inequalities for a and b . Denote by $\beta_F = \frac{N_F}{N_F+1}$ and recall that $\beta = \frac{N_L}{N_L+1}$. Then, b satisfies the quadratic equation $b^2 + b\omega = \beta$, with $\omega = 1 + \frac{\beta}{N_F}$. Now start with the straightforward inequality $\beta < \beta_F + 1$, and multiply it by β_F . We get $\beta\beta_F < \beta_F^2 + \beta_F$. Since $\beta_F = 1 - \frac{\beta_F}{N_F}$, we get $\beta(1 - \frac{\beta_F}{N_F}) < \beta_F^2 + \beta_F$, or equivalently $\beta < \beta_F^2 + \beta_F(1 + \frac{\beta_F}{N_F})$. Since $b^2 + b\omega = \beta$ and $\omega = 1 + \frac{\beta_F}{N_F}$, we get $b^2 + b\omega < \beta_F^2 + \beta_F\omega$. Because the function $f(x) = x^2 + x\omega$ is increasing in $x \in (0, 1)$, we have $b < \beta_F = \frac{N_F}{N_F+1}$. This inequality is equivalent to $N_F - b > N_F b$. Dividing by $N_F + 1$, we get $a = \frac{N_F - b}{N_F + 1} > \frac{N_F b}{N_F + 1} = b\beta_F$. But we have already seen that $\beta_F > b$, hence $a > b\beta_F > b^2$. To summarize,

$$b < \frac{N_F}{N_F + 1}, \quad a > b^2. \quad (\text{A17})$$

Lemma A.1 can now be used to show that the coefficients A and B are constant. Indeed, in the proof of the Theorem, we have seen that both A_t and B_t satisfy recursive equations of the form $X_t = \alpha + \beta X_{t-1}$, with $\beta \in (-1, 1)$. Then, Lemma A.1 implies that X_t converges to a fixed number $\frac{\alpha}{1-\beta}$, regardless of the starting point. But, since we work in continuous time, and $t + 1$ actually stands for $t + dt$, the convergence occurs in an infinitesimal amount of time. Thus, X_t is constant for all t , and that constant is equal to $\frac{\alpha}{1-\beta}$. \square

We now state the Lemma that is used in the proof of Theorem 1.

Lemma A.1. *Let $X_1 \in \mathbb{R}$, and consider a sequence $X_t \in \mathbb{R}$ which satisfies the following recursive equation:*

$$X_t - \beta X_{t-1} = \alpha, \quad t \geq 2. \quad (\text{A18})$$

Then the sequence X_t converges to $\bar{X} = \frac{\alpha}{1-\beta}$, regardless of the initial value of X_1 , if and only if $\beta \in (-1, 1)$.

Proof. First, note that \bar{X} is well defined as long as $\beta \neq 1$. If we denote by $Y_t = X_t - \bar{X}$, the new sequence Y_t satisfies the recursive equation $Y_t - \beta Y_{t-1} = 0$. We now show that Y_t converges to 0 (and \bar{X} is well defined) if and only if $\beta \in (-1, 1)$. Then, the difference equation $Y_t - \beta Y_{t-1} = 0$ has the following general solution:

$$Y_t = C\beta^t, \quad t \geq 1, \quad \text{with } C \in \mathbb{R}. \quad (\text{A19})$$

Then, Y_t is convergent for any values of C if and only if all $\beta \in (-1, 1]$. But in the latter case, $1 - \beta = 0$, which makes \bar{X} nondefined. \square

Proof of Corollary 1. In the proof of Theorem 1, equation (A8) implies $\lambda\bar{\gamma} = \frac{N_F}{N_F+1}$, $\lambda\bar{\mu} = \frac{B}{A} \frac{N_L}{N_L+1}$. But from (A14) and (A15), we have $\frac{B}{A} = \frac{1}{1+b}$, which proves the first row in (22). The second row in (22) just rewrites the formulas for A and B from equations (A14) and (A15). \square

Proof of Proposition 1. From Corollary 1, $\lambda\bar{\gamma} = \frac{N_F}{N_F+1}$ and $\lambda\bar{\mu} = \frac{B}{A} \frac{N_L}{N_L+1}$. From (A7),

the equilibrium normalized expected profit of the FT is

$$\tilde{\pi}^F = \gamma - \lambda\gamma\bar{\gamma} + \mu B - \lambda\mu\bar{\mu}A = \gamma\left(1 - \frac{N_F}{N_F + 1}\right) + B\mu\left(1 - \frac{N_L}{N_L + 1}\right) \quad (\text{A20})$$

From (A16), $B = \frac{1}{N_F + 1}$, which proves the desired formula for π^F . The profit of the ST is the same as for the FT, but with $\gamma = 0$. The last statement now follows from the asymptotic results in Theorem 1. \square

Justification of Result 1. According to Proposition 1, δdt is the expected profit that speculators get per unit of time dt from trading on their lagged signal (\widetilde{dw}_{t-1}) . Given that all speculators break even on this lag, they would not trade on any signal with a larger lag, as this would cost them the same (δ), but would bring a lower profit. For this last statement we use the results of Internet Appendix I (Proposition I.3), where we show numerically and asymptotically that the profit generated by lagged signals is decreasing in the number of lags. \square

Proof of Proposition 2. Since $1 - a = \frac{1+b}{N_F + 1}$, equation (19) implies that $\lambda = \rho \frac{N_F}{N_F - b} = \frac{\sigma_w}{\sigma_u} \sqrt{(1-a)(a-b^2)} \frac{N_F}{N_F - b} = \frac{\sigma_w}{\sigma_u} \frac{\sqrt{(1+b)(a-b^2)}}{\sqrt{N_F + 1}} \frac{N_F}{N_F - b}$, which proves the first equation in (36).

By definition, the trading volume is $TV = \sigma_y^2$. From (A15), $TV = \sigma_y^2 = \tilde{\sigma}_y^2 \sigma_w^2 = \frac{a\sigma_w^2}{\rho^2}$. From (19), $\rho^2 = \frac{\sigma_w^2}{\sigma_u^2} (1-a)(a-b^2)$, hence $TV = \sigma_u^2 \frac{a}{(1-a)(a-b^2)}$. Substituting $1 - a = \frac{1+b}{N_F + 1}$, we get $TV = \sigma_u^2 (N_F + 1) \frac{a}{(1+b)(a-b^2)}$, which proves the second equation in (36).

The price volatility is $\sigma_p^2 = \lambda^2 TV = \left(\frac{\lambda}{\rho}\right)^2 \rho^2 TV = \left(\frac{\lambda}{\rho}\right)^2 a\sigma_w^2$. From (19), $\frac{\lambda}{\rho} = \frac{N_F}{N_F - b}$, hence $\sigma_p^2 = \left(\frac{N_F}{N_F - b}\right)^2 \frac{N_F - b}{N_F + 1} \sigma_w^2 = \frac{N_F^2}{(N_F + 1)(N_F - b)} \sigma_w^2$, which proves the third equation in (36).

The speculator participation rate is $SPR = \frac{\bar{\gamma}^2 \sigma_w^2 + \bar{\mu}^2 \sigma_w^2}{TV} = \frac{\rho^2 (\bar{\gamma}^2 \sigma_w^2 + \bar{\mu}^2 \sigma_w^2)}{a\sigma_w^2}$. Since $\rho\bar{\gamma} = a$, $\rho\bar{\mu} = b$, and $\sigma_w^2 = (1-a)\sigma_w^2$, we get $SPR = \frac{a^2 + b^2(1-a)}{a}$. This proves the last equation in (36), since $\frac{1-a}{a} = \frac{1+b}{N_F - b}$. \square

Proof of Proposition 3. As in Theorem 1, we start with the FT's choice of optimal trading strategy. Each FT $i = 1, \dots, N_F$ observes dw_t , and chooses $dx_t^i = \gamma_t^i dw_t$ to maximize the expected profit:

$$\pi_0 = \mathbb{E} \left(\int_0^T \left(w_t - p_{t-1} - \lambda_t (dx_t^i + dx_t^{-i} + du_t) \right) dx_t^i \right) = \int_0^T \gamma_t^i \sigma_w^2 dt - \lambda_t \gamma_t^i (\gamma_t^i + \gamma_t^{-i}) \sigma_w^2 dt, \quad (\text{A21})$$

where the superscript “ $-i$ ” indicates the aggregate quantity from the other FTs. This is a pointwise quadratic optimization problem, with solution $\lambda_t \gamma_t^i = \frac{1 - \lambda_t \gamma_t^{-i}}{2}$. Since this is true for all $i = 1, \dots, N_F$, the equilibrium is symmetric and we compute $\gamma_t = \frac{1}{\lambda_t} \frac{1}{1 + N_F}$.

The dealer takes the FTs’ strategies as given, thus assumes that the aggregate order flow is of the form $dy_t = du_t + N_F \gamma_t dw_t$. To set λ_t , the dealer sets p_t such that $dp_t = \lambda_t dy_t$, with $\lambda_t = \frac{\text{Cov}(w_t, dy_t)}{\text{Var}(dy_t)} = \frac{N_F \gamma_t \sigma_w^2}{\sigma_u^2 + N_F^2 \gamma_t^2 \sigma_w^2}$. This implies $\lambda_t^2 \sigma_u^2 + (N_F \gamma_t \lambda_t)^2 \sigma_w^2 = N_F \gamma_t \lambda_t \sigma_w^2$. But $N_F \lambda_t \gamma_t = \frac{N_F}{N_F + 1}$. Hence, $\lambda_t^2 \sigma_u^2 + \left(\frac{N_F}{N_F + 1}\right)^2 \sigma_w^2 = \frac{N_F}{N_F + 1} \sigma_w^2$, or $\lambda_t^2 \sigma_u^2 = \frac{N_F}{(N_F + 1)^2} \sigma_w^2$, which implies the formula $\lambda = \frac{\sigma_w}{\sigma_u} \frac{\sqrt{N_F}}{N_F + 1}$. We then compute $\gamma_t = \frac{1}{\lambda_t} \frac{N_F}{1 + N_F} = \frac{\sigma_u}{\sigma_w} \frac{1}{\sqrt{N_F}}$.

We have $TV = \sigma_y^2 = N_F^2 \gamma^2 \sigma_w^2 + \sigma_u^2$. But $N_F \gamma = \frac{\sigma_u}{\sigma_w} \sqrt{N_F}$, hence $TV = \sigma_u^2 (1 + N_F)$. Next, $\sigma_p^2 = \lambda^2 TV = \frac{\sigma_w^2}{\sigma_u^2} \frac{N_F}{(N_F + 1)^2} \sigma_u^2 (N_F + 1) = \sigma_w^2 \frac{N_F}{N_F + 1}$. Also, $SPR = \frac{TV - \sigma_u^2}{TV} = \frac{\sigma_u^2 (N_F + 1) - \sigma_u^2}{\sigma_u^2 (N_F + 1)} = \frac{N_F}{N_F + 1}$.

Finally, we compute Σ' . From the formula above for λ , we get $\text{Var}(dp_t) = \lambda^2 \text{Var}(dy_t) = \lambda \text{Cov}(w_t, dy_t) = \text{Cov}(w_t, dp_t)$. Since $\Sigma_t = \text{Var}(w_t - p_{t-1}) = \mathbf{E}((w_t - p_{t-1})^2)$, we compute $\Sigma'_t = \frac{1}{dt} \mathbf{E}(2(dw_{t+1} - dp_t)(w_t - p_{t-1}) + (dw_{t+1} - dp_t)^2) = -2 \frac{\text{Cov}(w_t, dp_t)}{dt} + \sigma_w^2 + \frac{\text{Var}(dp_t)}{dt} = \sigma_w^2 - \sigma_p^2 = \frac{\sigma_w^2}{N_F + 1}$. \square

Proof of Proposition 4. We use the formulas from the Proof of Theorem 1. Since \widetilde{dw}_t is orthogonal on dy_t , we have $\widetilde{\text{Cov}}(\widetilde{dw}_t, dw_t) = \widetilde{\text{Cov}}(\widetilde{dw}_t, \widetilde{dw}_t) = A = 1 - a = \frac{1+b}{N_F+1}$. Then, $\widetilde{\text{Cov}}(\widetilde{dw}_t, \widetilde{dw}_{t-1}) = \widetilde{\text{Cov}}(dw_t - \rho \bar{\gamma} dw_t - \rho \bar{\mu} \widetilde{dw}_{t-1}, \widetilde{dw}_{t-1}) = -\rho \bar{\mu} A$. Therefore,

$$\begin{aligned} \widetilde{\text{Cov}}(d\bar{x}_{t+1}, d\bar{x}_t) &= \widetilde{\text{Cov}}(\bar{\gamma} dw_{t+1} + \bar{\mu} \widetilde{dw}_t, \bar{\gamma} dw_t + \bar{\mu} \widetilde{dw}_{t-1}) = \bar{\mu} \bar{\gamma} A + \bar{\mu}^2 (-bA) \\ \widetilde{\text{Var}}(d\bar{x}_t) &= \widetilde{\text{Var}}(\bar{\gamma} dw_t + \bar{\mu} \widetilde{dw}_{t-1}) = \bar{\gamma}^2 + \bar{\mu}^2 A. \end{aligned} \quad (\text{A22})$$

By multiplying both the numerator and denominator by ρ^2 , we compute

$$\rho_{\bar{x}} = \frac{\bar{\mu} \bar{\gamma} A}{\bar{\gamma}^2 + \bar{\mu}^2 A} - \frac{b \bar{\mu}^2 A}{\bar{\gamma}^2 + \bar{\mu}^2 A} = \frac{ab(1-a)}{a^2 + b^2(1-a)} - \frac{b^3(1-a)}{a^2 + b^2(1-a)} = \rho_{AT} + \rho_{EA}. \quad (\text{A23})$$

Then, $\rho_{\bar{x}} = \frac{ab-b^3}{a^2+b^2(1-a)} (1-a) = \frac{(a-b^2)b}{a^2+b^2(1-a)} \frac{1+b}{N_F+1}$, which implies the desired formulas.

We now prove that $\rho_{\bar{x}} > 0$ if and only if there exists slow trading. When there is no slow trading, $b = \rho \bar{\mu} = 0$, hence $\rho_{\bar{x}} = 0$. When there is slow trading, we show that $\rho_{\bar{x}} = \frac{b(b+1)(a-b^2)}{a^2+b^2(1-a)} \frac{1}{N_F+1} > 0$. Indeed, we have $b > 0$, $a < 1$, and from equation (A17), $a - b^2 > 0$. \square

Proof of Proposition 5. If x_t is the IFT's inventory in the risky asset, denote by

$$\begin{aligned}\Omega_t^{xx} &= \frac{\mathbb{E}(x_t^2)}{\sigma_w^2 dt}, & \Omega_t^{xe} &= \frac{\mathbb{E}(x_t(w_t - p_t))}{\sigma_w^2 dt}, & X_t &= \frac{\mathbb{E}(x_t \widetilde{dw}_t)}{\sigma_w^2 dt} \\ \Omega_t^{xw} &= \frac{\mathbb{E}(x_t w_t)}{\sigma_w^2 dt}, & \Omega_t^{xp} &= \frac{\mathbb{E}(x_t p_t)}{\sigma_w^2 dt}, & Z_t &= \frac{\mathbb{E}(x_{t-1} dy_t)}{\sigma_w^2 dt}.\end{aligned}\tag{A24}$$

Since $\Theta > 0$, we have $\Theta \in (0, 2)$, or $\phi = 1 - \Theta \in (-1, 1)$. From (46), x_t satisfies the recursive equation $x_t = \phi x_{t-1} + G dw_t$. We compute $\Omega_t^{xx} = \frac{\mathbb{E}((x_t)^2)}{\sigma_w^2 dt} = \frac{\mathbb{E}((\phi x_{t-1} + G dw_t)^2)}{\sigma_w^2 dt} = \phi^2 \Omega_{t-1}^{xx} + G^2$. Since $\phi^2 \in (-1, 1)$, we apply Lemma A.1 to the recursive formula $\Omega_t^{xx} = \phi^2 \Omega_{t-1}^{xx} + G^2$. Then, Ω_t^{xx} is constant and equal to:

$$\Omega^{xx} = \frac{G^2}{1 - \phi^2} = \frac{G^2}{\Theta(1 + \phi)},\tag{A25}$$

which is the usual variance formula for the $AR(1)$ process. The order flow at t is $dy_t = -\Theta x_{t-1} + \bar{\gamma} dw_t + \bar{\mu} \widetilde{dw}_{t-1} + du_t$, with $\bar{\gamma} = \gamma^- + G$. Then, Z_t is a function of X_{t-1} :

$$Z_t = \frac{\mathbb{E}(x_{t-1} dy_t)}{\sigma_w^2 dt} = -\Theta \Omega_{t-1}^{xx} + \bar{\mu} X_{t-1} = -\frac{G^2}{1 + \phi} + \bar{\mu} X_{t-1}.\tag{A26}$$

The recursive formula for X_t is $X_t = \frac{\mathbb{E}(x_t \widetilde{dw}_t)}{\sigma_w^2 dt} = \frac{\mathbb{E}((\phi x_{t-1} + G dw_t)(dw_t - \rho dy_t))}{\sigma_w^2 dt} = -\phi \rho Z_t + G - G \rho \bar{\gamma} = -\phi \rho \bar{\mu} X_{t-1} + \phi \frac{\rho G^2}{1 + \phi} + G - G \rho \bar{\gamma} = -\phi b X_{t-1} + G(1 - a^-) - \frac{\rho G^2}{1 + \phi}$. By assumption, $0 \leq b < 1$, hence $\phi b \in (-1, 1)$. Lemma A.1 implies that X_t is constant and equal to

$$X = \frac{G(1 - a^-) - \frac{\rho G^2}{1 + \phi}}{1 + \phi b}.\tag{A27}$$

From (A26), Z_t is also constant and satisfies:

$$Z = \bar{\mu} X - \frac{G^2}{1 + \phi} = \bar{\mu} G \frac{1 - a^-}{1 + \phi b} - G^2 \frac{b + \frac{1}{1 + \phi}}{1 + \phi b}.\tag{A28}$$

We are interested in $\Omega_t^{xe} = \Omega_t^{xw} - \Omega_t^{xp}$. The recursive equation for Ω_t^{xw} is $\Omega_t^{xw} = \frac{\mathbb{E}(x_t w_t)}{\sigma_w^2 dt} = \frac{\mathbb{E}((\phi x_{t-1} + G dw_t)(w_{t-1} + dw_t))}{\sigma_w^2 dt} = \phi \Omega_{t-1}^{xw} + G$. Since $\phi \in (-1, 1)$, Lemma A.1 implies that Ω_t^{xw} is constant and equal to

$$\Omega^{xw} = \frac{G}{\Theta}.\tag{A29}$$

The recursive formula for Ω_t^{xp} is $\Omega_t^{xp} = \frac{\mathbb{E}(x_t p_t)}{\sigma_w^2 dt} = \frac{\mathbb{E}((\phi x_{t-1} + G dw_t)(p_{t-1} + \lambda dy_t))}{\sigma_w^2 dt} = \phi \Omega_{t-1}^{xp} + \lambda \phi Z + \lambda G \bar{\gamma}$. Lemma A.1 implies that Ω_t^{xp} is constant and equal to $\Omega^{xp} = \frac{\lambda \phi Z + \lambda G \bar{\gamma}}{\Theta}$. It follows that $\Omega_t^{xe} = \Omega_t^{xw} - \Omega_t^{xp}$ is constant and satisfies:

$$-\Theta \Omega^{xe} = -(\Theta \Omega^{xw} - \Theta \Omega^{xp}) = -G + \lambda \phi Z + \lambda G \bar{\gamma}. \quad (\text{A30})$$

The IFT's expected profit satisfies $\pi_{\Theta > 0} = \mathbb{E} \int_0^T (w_t - p_t) dx_t = \mathbb{E} \int_0^T (w_{t-1} - p_{t-1} + dw_t - \lambda dy_t)(G dw_t - \Theta x_{t-1})$. Hence, the IFT's normalized expected profit is $\tilde{\pi}_{\Theta > 0} = \int_0^T (-\Theta \Omega^{xe} + G - \lambda \bar{\gamma} G + \lambda \Theta Z) dt$. If we use the formula (A30) for Ω^{xe} , we obtain:

$$\tilde{\pi}_{\Theta > 0} = \lambda Z = \lambda \left(\bar{\mu} G \frac{1 - a^-}{1 + \phi b} - G^2 \frac{b + \frac{1}{1 + \phi}}{1 + \phi b} \right), \quad (\text{A31})$$

where the second equality comes from (A28). This proves (49).

We now show that the inventory costs are zero, which implies that the IFT's expected utility is the same as his expected profit. According to equation (A25), $\Omega_t^{xx} = \frac{G^2}{\Theta(1 + \phi)}$ is constant. Then, by the definition (A24) of Ω_t^{xx} , we have

$$\mathbb{E}(x_t^2) = \frac{G^2}{\Theta(1 + \phi)} \sigma_w^2 dt, \quad (\text{A32})$$

which implies that the expected squared inventory of the IFT is infinitesimal, and therefore becomes zero when integrated up over $[0, 1]$ ($dt^2 = 0$). Hence, from the definition (42), the inventory costs of the IFT are $C_I \int_0^T \mathbb{E}(x_t^2) dt = 0$.

To show that all IFT's expected profits are in cash, consider the decomposition

$$\pi_{\Theta > 0} = \mathbb{E} \int_0^T (v_T - p_t) dx_t = \mathbb{E} \int_0^T w_t dx_t - \mathbb{E} \int_0^T p_t dx_t, \quad (\text{A33})$$

which is the same as (50). We need to show that the first term (the risky component) is zero. From (A29), $\Omega_t^{xw} = \frac{G}{\Theta}$. Thus, we compute

$$\frac{\mathbb{E} \int_0^T w_t dx_t}{\sigma_w^2} = \frac{\mathbb{E} \int_0^T (w_{t-1} + dw_t)(-\Theta x_{t-1} + G dw_t)}{\sigma_w^2} = -\Theta \Omega^{xw} + G = 0. \quad (\text{A34})$$

which implies that the risky component is indeed zero. This finishes the proof. \square

Proof of Corollary 3. From (A24), $Z_t = \frac{\mathbb{E}(x_{t-1}dy_t)}{\sigma_w^2 dt}$, which implies $\mathbb{E}(x_{t-1}dy_t) = Z_t \sigma_w^2 dt$. From this, $\mathbb{E} \int_0^T x_{t-1} dp_t = \lambda \int_0^T Z_t \sigma_w^2 dt = \lambda Z \sigma_w^2$, since Z_t is constant. But from (A31), the IFT's expected profit is $\pi_{\Theta > 0} = \lambda Z \sigma_w^2$. Therefore,

$$\pi_{\Theta > 0} = \mathbb{E} \int_0^T x_{t-1} dp_t. \quad (\text{A35})$$

Now, write $dp_t = \lambda dy_t = \lambda(-\Theta x_{t-1} + \bar{\gamma} dw_t + \bar{\mu} \widetilde{dw}_{t-1} + du_t)$. Since x_{t-1} is orthogonal to dw_t and du_t , we get $dp_t = \lambda(\bar{\mu} \widetilde{dw}_{t-1} - \Theta x_{t-1})$. If we substitute this formula in (A35), we obtain (52). \square

Proof of Theorem 2. Let $\Theta = 0$. Then, the IFT's strategy is of the form $dx_t = Gdw_t$. We compute the IFT's expected profit $\pi_{\Theta=0} = \mathbb{E} \int_0^T (w_t - p_t) dx_t = \mathbb{E} \int_0^1 (w_{t-1} - p_{t-1} + dw_t - \lambda dy_t) (Gdw_t) = \mathbb{E} \int_0^1 (dw_t - \lambda dy_t) (Gdw_t) = \mathbb{E} \int_0^1 (dw_t - \lambda \bar{\gamma} dw_t) (Gdw_t) = G(1 - \lambda \bar{\gamma}) \sigma_w^2$. But $\lambda \bar{\gamma} = \lambda G + \lambda \gamma^- = \lambda G + Ra^-$. The normalized IFT's expected profit is:

$$\tilde{\pi}_{\Theta=0} = G(1 - \lambda \bar{\gamma}) = G(1 - Ra^-) - \lambda G^2. \quad (\text{A36})$$

To compute the IFT's inventory costs, denote by $\Omega_t^{xx} = \frac{\mathbb{E}(x_t^2)}{\sigma_w^2 dt}$. We compute $\frac{d\Omega_t^{xx}}{dt} = \frac{1}{\sigma_w^2 dt} \mathbb{E}(2x_{t-1} dx_t + (dx_t)^2) = \frac{1}{\sigma_w^2 dt} \mathbb{E}(2Gx_{t-1} dw_t + G^2 (dw_t)^2) = G^2$. Since $\Omega_0^{xx} = 0$, the solution of this first order ODE is $\Omega_t^{xx} = tG^2$, for all $t \in [0, 1]$. Hence, the inventory costs are equal to

$$C_I \mathbb{E} \int_0^1 x_t^2 dt = C_I G^2 \int_0^1 t dt = \frac{C_I}{2} G^2, \quad (\text{A37})$$

From (A36) and (A37), the IFT's normalized expected utility when $\Theta = 0$ is:

$$\tilde{U}_{\Theta=0} = G(1 - Ra^-) - G^2 \left(\lambda + \frac{C_I}{2} \right). \quad (\text{A38})$$

The function $\tilde{U}_{\Theta=0}$ attains its maximum at $G = \frac{1 - Ra^-}{2\lambda + C_I} = \frac{1 - Ra^-}{2\lambda(1 + \frac{C_I}{2\lambda})}$, as stated in the Theorem. The maximum value is:

$$\tilde{U}_{\Theta=0}^{\max} = \frac{(1 - Ra^-)^2}{2(2\lambda + C_I)}. \quad (\text{A39})$$

Let $\Theta > 0$, which is equivalent to $\phi = 1 - \Theta \in (-1, 1)$. In the proof of Proposition 5, we have already computed the IFT's expected profit (see (49)) and showed that the IFT's inventory costs are zero. Hence, the IFT's expected utility is the same as his expected profit, and satisfies $\tilde{U}_{\Theta > 0} = \tilde{\pi}_{\Theta > 0} = \frac{\lambda}{\rho} \left(bG \frac{1-a^-}{1+\phi b} - \rho G^2 \frac{b+\frac{1}{1+\phi}}{1+\phi b} \right)$. The first order condition with respect to G implies that at the optimum $G = \frac{b(1-a^-)}{2\rho \left(b+\frac{1}{1+\phi} \right)}$, as stated in the Theorem. The second order condition for a maximum is $\lambda \frac{b+\frac{1}{1+\phi}}{1+\phi b} > 0$, which follows from $\lambda > 0$, $b \in [0, 1)$, and $\phi \in (-1, 1)$. For the optimum G , the normalized expected utility (profit) of the IFT is:

$$\tilde{U}_{\Theta > 0} = \frac{(Rb(1-a^-))^2}{4\lambda(1+\phi b) \left(b+\frac{1}{1+\phi} \right)}. \quad (\text{A40})$$

We now analyze the function

$$f(\phi) = (1+\phi b) \left(b+\frac{1}{1+\phi} \right) \implies f'(\phi) = \frac{b^2(1+\phi)^2 + b - 1}{(1+\phi)^2}. \quad (\text{A41})$$

The polynomial in the numerator has two roots:

$$\phi_1 = -1 + \frac{\sqrt{1-b}}{b} \quad \phi_2 = -1 - \frac{\sqrt{1-b}}{b}. \quad (\text{A42})$$

By assumption $b < 1$, hence both roots are real. Clearly, we have $\phi_2 < -1$. We show that $\phi_1 \in (-1, 1)$. First, note that ϕ_1 is decreasing in b . For $b = 1$ we have $\phi_1 = -1$; while for $b = \frac{\sqrt{17}-1}{8}$ (which satisfies $4b^2 + b = 1$) we have $\phi_1 = 1$. Since by assumption $\frac{\sqrt{17}-1}{8} < b < 1$, it follows that indeed $\phi_1 \in (-1, 1)$. Thus, $f'(\phi)$ is negative on $(-1, \phi_1)$ and positive on $(\phi_1, 1)$. Hence, $f(\phi)$ attains its minimum at $\phi = \phi_1$, which implies that the normalized expected utility $\tilde{U}_{\Theta > 0}$ from (A40) attains its maximum at $\phi = \phi_1$, or $\Theta = 2 - \frac{\sqrt{1-b}}{b}$, as stated in the Theorem. The maximum value (over both G and Θ) is:

$$\tilde{U}_{\Theta > 0}^{\max} = \frac{(Rb(1-a^-))^2}{4\lambda b(1+\sqrt{1-b})^2}. \quad (\text{A43})$$

To determine the cutoff value for the inventory aversion coefficient C_I , we set $\tilde{U}_{\Theta=0}^{\max} = \tilde{U}_{\Theta > 0}^{\max}$. From (A39) and (A43), algebraic manipulation shows that the cutoff value is

$\bar{C}_I = 2\lambda \left(\frac{(1-Ra^-)^2(1+\sqrt{1-b})^2}{R^2b(1-a^-)^2} - 1 \right)$, as stated in the Theorem. \square

Proof of Corollary 4. Let $\Theta > 0$. We are in the context of Theorem 2, where $b > \frac{\sqrt{17}-1}{8} > 0$ and $\rho > 0$, hence $\bar{\mu} = \frac{b}{\rho} > 0$. The IFT's strategy is $dx_t = -\Theta x_{t-1} + Gdw_t$, while the slow trading component is $d\bar{x}_t^S = \bar{\mu}d\bar{w}_{t-1}$. Since dw_t is orthogonal to $\bar{d}w_{t-1}$, $\text{Cov}(dx_t, d\bar{x}_t^S) = -\Theta \text{Cov}(x_{t-1}, d\bar{x}_t^S) = -\Theta \bar{\mu} \text{Cov}(x_{t-1}, \bar{d}w_{t-1})$. This proves the equality in (55). Since $\Theta > 0$ and $\bar{\mu} > 0$, it remains to prove the inequality $\text{Cov}(x_{t-1}, \bar{d}w_{t-1}) > 0$. But $\text{Cov}(x_{t-1}, \bar{d}w_{t-1}) = X\sigma_w^2 dt$ (see (A24)). From (A27), $X = \frac{G(1-a^-) - \frac{\rho G^2}{1+\phi}}{1+\phi}$. Substituting the optimal G and $\phi = 1 - \Theta$ from Theorem 2, we obtain $X = \frac{(1-a^-)^2}{4(1+\sqrt{1-b})}$. As in Theorem 2, $a^-, b \in [0, 1)$, hence $X > 0$ and the proof is complete. \square

Proof of Theorem 3. Consider the following implicit equation in b

$$\frac{2b(1+b)(2B+1)}{n_L} = \frac{Q}{B^2(a^-+b)} + \frac{3bB+2b^2B-1-b}{b}(1-a^-) - 2, \quad (\text{A44})$$

where the following substitutions are made:³⁸

$$\begin{aligned} B &= \frac{1}{\sqrt{1-b}}, & q &= (B+1) \left(2(B^2-1) - n_F(3B^2-2) \right), \\ a^- &= \frac{-q \pm \sqrt{q^2 + n_F B^5 ((4-n_F)B + 2(2-n_F))}}{B^2((4-n_F)B + 2(2-n_F))}, & (\text{A45}) \\ Q &= B^3(a^-)^2 + 2(3B^3 + 3B^2 - 2B - 1)a^- + (B^3 + 2B^2 - 2). \end{aligned}$$

We write the equations for the other coefficients:

$$\begin{aligned} R &= \frac{4(B+1)B^2(a^-+b)}{Q}, & a &= \frac{(2B+1)a^-+1}{2(B+1)} \\ \rho^2 &= \left((a-b^2) + \frac{2bB-1}{2B+1}(1-a) \right) (1-a) \frac{\sigma_w^2}{\sigma_u^2}, & \lambda &= R\rho \\ \Theta &= 2 - \frac{\sqrt{1-b}}{b}, & G &= \frac{1-a}{\rho(2B+1)}, & \gamma &= \frac{a^-}{\rho N_F}, & \mu &= \frac{b}{\rho N_L}. \end{aligned} \quad (\text{A46})$$

The proof is now left to Internet Appendix J (see Sections J.4 and J.5). \square

Proof of Proposition 6. See Internet Appendix J (Section J.5). \square

³⁸To be rigorous, we have included the case when a^- is negative. However, numerically this case never occurs in equilibrium, because it leads to $\lambda < 0$, which contradicts the FT's second order condition (J56) in Internet Appendix J.

Proof of Proposition 7. See Internet Appendix J (Section J.5). □

REFERENCES

- [1] AÏT-SAHALIA, YACINE, AND MEHMENT SAGLAM (2014): “High Frequency Traders: Taking Advantage of Speed,” Working Paper.
- [2] BACK, KERRY, HENRY CAO, AND GREGORY WILLARD (2000): “Imperfect Competition among Informed Traders,” *Journal of Finance*, 55, 2117–2155.
- [3] BACK, KERRY, AND HAL PEDERSEN (1998): “Long-Lived Information and Intraday Patterns,” *Journal of Financial Markets*, 1, 385–402.
- [4] BARON, MATTHEW, JONATHAN BROGAARD, AND ANDREI KIRILENKO (2014): “Risk and Return in High Frequency Traders,” Working Paper.
- [5] BENOS, EVANGELOS, AND SATCHIT SAGADE (2013): “High-Frequency Trading Behaviour and Its Impact on Market Quality: Evidence from the UK Equity Market,” Working Paper.
- [6] BIAIS, BRUNO, THIERRY FOUCAULT, AND SOPHIE MOINAS (2014): “Equilibrium Fast Trading,” *Journal of Financial Economics*, forthcoming.
- [7] BOEHMER, EKKEHART, KINGSLEY FONG, AND JULIE WU (2014): “International Evidence on Algorithmic Trading,” Working Paper.
- [8] BROGAARD, JONATHAN (2011): “High Frequency Trading and Its Impact on Market Quality,” Working Paper.
- [9] BROGAARD, JONATHAN, BJÖRN HAGSTRÖMER, LARS NORDÉN, AND RYAN RIORDAN (2013): “Trading Fast and Slow: Colocation and Market Quality,” Working Paper.
- [10] BROGAARD, JONATHAN, TERRENCE HENDERSHOTT, AND RYAN RIORDAN (2014): “High-Frequency Trading and Price Discovery,” *Review of Financial Studies*, 27, 2267–2306.
- [11] BUDISH, ERIC, PETER CRAMTON, AND JOHN SHIM (2014): “The High-Frequency Trading Arms Race: Frequent Batch Auctions as a Market Design Response,” Working Paper.
- [12] CAO, HUINING HENRY, YUAN MA, AND DONGYAN YE (2013): “Disclosure, Learning, and Coordination,” Working paper.
- [13] CALDENTEY, RENÉ, AND ENNIO STACCHETTI (2010): “Insider Trading with a Random Deadline,” *Econometrica*, 78, 245–283.
- [14] CARTEA, ÁLVARO, AND JOSÉ PENALVA (2012): “Where is the Value in High Frequency Trading?,” *Quarterly Journal of Finance*, 2, 1–46.
- [15] CHABOUD, ALAIN, BENJAMIN CHIQUOINE, ERIK HJALMARSSON, AND CLARA VEGA (2014): “Rise of the Machines: Algorithmic Trading in the Foreign Exchange Market,” *Journal of Finance*, 69, 2045–2084.
- [16] CHAU, MINH, AND DIMITRI VAYANOS (2008): “Strong-Form Efficiency with Monopolistic Insiders,” *Review of Financial Studies*, 18, 2275–2306.
- [17] CVITANIĆ, JAKŠA, AND ANDREI KIRILENKO (2010): “High Frequency Traders and Asset Prices,” Working Paper.
- [18] DU, SONGZI, AND HAOXIANG ZHU (2014): “Welfare and Optimal Trading Frequency in Dynamic Double Auctions,” Working Paper.
- [19] FOSTER, DOUGLAS, AND S. VISWANATHAN (1996): “Strategic Trading When Agents Forecast the Forecast of Others,” *Journal of Finance*, 51, 1437–1478.

- [20] FOUCAULT, THIERRY, JOHAN HOMBERT, AND IOANID ROȘU (2015): “News Trading and Speed,” *Journal of Finance*, forthcoming.
- [21] HASBROUCK, JOEL, AND GIDEON SAAR (2013): “Low-Latency Trading,” *Journal of Financial Markets*, 16, 646–679.
- [22] HENDERSHOTT, TERRENCE, CHARLES JONES, AND ALBERT MENKVELD (2011): “Does Algorithmic Trading Improve Liquidity?,” *Journal of Finance*, 66, 1–33.
- [23] HENDERSHOTT, TERRENCE, AND ALBERT MENKVELD (2014): “Price Pressures,” *Journal of Financial Economics*, 114, 405–423.
- [24] HIRSCHHEY, NICHOLAS (2013): “Do High-Frequency Traders Anticipate Buying and Selling Pressure?,” Working Paper.
- [25] HIRSHLEIFER, DAVID, AVANIDHAR SUBRAHMANYAM, AND SHERIDAN TITMAN (1994): “Security Analysis and Trading Patterns When Some Investors Receive Information Before Others,” *Journal of Finance*, 49, 1665–1698.
- [26] HO, THOMAS, AND HANS STOLL (1981): “Optimal Dealer Pricing Under Transactions and Return Uncertainty,” *Journal of Financial Economics*, 9, 47–73.
- [27] HOFFMANN, PETER (2014): “A Dynamic Limit Order Market with Fast and Slow Traders,” *Journal of Financial Economics*, 113, 156–169.
- [28] HOLDEN, CRAIG, AND AVANIDHAR SUBRAHMANYAM (1992): “Long-Lived Private Information and Imperfect Competition,” *Journal of Finance*, 47, 247–270.
- [29] JOVANOVIĆ, BOYAN, AND ALBERT MENKVELD (2012): “Middlemen in Limit-Order Markets,” Working Paper.
- [30] KIRILENKO, ANDREI, ALBERT KYLE, MEHRDAD SAMADI, AND TUGKAN TUZUN (2014): “The Flash Crash: The Impact of High Frequency Trading on an Electronic Market,” Working Paper.
- [31] KYLE, ALBERT (1985): “Continuous Auctions and Insider Trading,” *Econometrica*, 53, 1315–1335.
- [32] LI, SU (2012): “Speculative Dynamics I: Imperfect Competition, and the Implications for High Frequency Trading,” Working Paper.
- [33] LI, WEI (2014): “High Frequency Trading with Speed Hierarchies,” Working Paper.
- [34] LYONS, RICHARD (1997): “A Simultaneous Trade Model of the Foreign Exchange Hot Potato,” *Journal of International Economics*, 42, 275–298.
- [35] MADHAVAN, ANANTH, AND SEYMOUR SMIDT (1993): “An Analysis of Changes in Specialist Inventories and Quotations,” *Journal of Finance*, 48, 1595–1628.
- [36] PAGNOTTA, EMILIANO, AND THOMAS PHILIPPON (2013): “Competing on Speed,” Working Paper.
- [37] SEC (2010): “Concept Release on Equity Market Structure,” Release No. 34-61358; File No. S7-02-10.
- [38] SUBRAHMANYAM, AVANIDHAR, AND SHERIDAN TITMAN (1999): “The Going-Public Decision and the Development of Financial Markets,” *Journal of Finance*, 54, 1045–1082.
- [39] WELLER, BRIAN (2014): “Intermediation Chains,” Working Paper.
- [40] ZHANG, X. FRANK (2010): “High Frequency Trading, Stock Volatility, and Price Discovery,” Working Paper.