

How Rigged Are Stock Markets? Evidence from Microsecond Timestamps

Robert P. Bartlett, III^{*}
University of California, Berkeley

Justin McCrary^{**}
University of California, Berkeley, NBER

Abstract:

We examine the incidence of “SIP latency arbitrage” strategies using new timestamp data from the two Securities Information Processors (SIPs). On average, the SIPs report quote updates from stock exchanges 1.13 milliseconds after they occur. However, liquidity-taking orders *gain* on average \$0.0002 per share when priced at the SIP-reported national best bid or offer (NBBO) rather than the NBBO calculated using exchanges’ direct data feeds. Trading surrounding SIP-priced trades shows little evidence that fast traders initiate these liquidity-taking orders to pick-off stale quotes. These findings contradict widespread claims that fast traders systematically exploit traders who transact at the SIP NBBO.

Draft Date: January 29, 2017

JEL codes: G10, G15, G18, G23, G28, K22

Keywords: latency arbitrage, high-frequency trading; SIP; market structure

Statement of Financial Disclosure and Conflict of Interest: Neither author has any financial interest or affiliation (including research funding) with any commercial organization that has a financial interest in the findings of this paper. The authors are grateful to the University of California, Berkeley School of Law, for providing general faculty research support.

^{*} rbartlett@berkeley.edu, 890 Simon Hall, UC Berkeley, Berkeley CA 94720. Tel: 510-542-6646.

^{**} jmccrary@berkeley.edu, 586 Simon Hall, UC Berkeley, Berkeley CA 94720. Tel: 510-643-6252.

1. Introduction

Concerns over the different speeds at which market participants access information and the resulting potential for adverse selection in financial markets have occupied center stage in recent years. In particular, the emergence of low-latency trading strategies that can exploit sub-second information asymmetries has led not just to economic research, but also to extensive regulatory scrutiny, litigation, and the formation of a new stock exchange.¹ Describing high frequency trading (HFT) as “one of the greatest threats to public confidence in the markets,” New York attorney general Eric Schneiderman in 2014 launched a series of high profile lawsuits against dark pools, exchanges, and HFT firms.² Regulators from the Federal Bureau of Investigation,³ to the Commodity Futures Trading Commission,⁴ to the Securities and Exchange Commission (SEC) have all brought pressure to bear on HFT.⁵

An important recent paper proposes a fundamental change to market design—frequent batch auctions (FBAs)—that would eliminate the trading gains from obtaining information microseconds (or nanoseconds) before others and presumably end the “high-frequency trading arms race” (Budish, Cramton, and Shim (2015)). The central advantage of FBAs is to determine a single market price of a security at discrete time periods rather than in continuous time, thus eliminating the race to exploit inter-market information arbitrage possibilities. However, eliminating the current system of competing stock exchanges in favor of a single market using FBAs would seem to require a complete overhaul of Regulation National Market System (Reg NMS), which is likely not realistic.⁶ Proposals to implement FBAs have accordingly been

¹ The Investor’s Exchange (IEX) obtained regulatory approval to operate as a public exchange in June 2016. IEX and its CEO Brad Katsuyama featured prominently in Michael Lewis’ *Flash Boys*.

² Speech by Eric Schneiderman on March 18, 2014 at New York Law School, “Insider Trading 2.0 – A New Initiative to Crack Down on Predatory Practices.”

³ Scott Patterson and Michael Rothfeld, “FBI Investigates High-Speed Trading,” *Wall Street Journal*, March 31, 2014. Available at <http://www.wsj.com/articles/SB10001424052702304886904579473874181722310>, last accessed December 17, 2016.

⁴ Douwe Miedema, “U.S. Futures Regulator CFTC Probing Speed Traders,” *Reuters Business News*, April 3, 2014. Available at <http://www.reuters.com/article/us-hedgefunds-speed-trading-cftc-idUSBREA321QU20140403>, last accessed December 17, 2016.

⁵ John McCrank, “Exclusive: SEC Targets 10 Firms in High Frequency Trading Probe—SEC Document,” *Reuters Business News*, July 17, 2014. Available at <http://www.reuters.com/article/us-sec-investigation-highfrequencytrading-idUSKBN0FM2TW20140717>, last accessed December 17, 2016. No doubt a contributor to the pressure on HFTs was its vivid description in *Flash Boys*.

⁶ As stated in its 2005 adopting release, Reg NMS has as a central objective a goal of promoting competition among exchanges and avoiding “a totally centralized system that loses the benefits of vigorous competition and innovation among individual markets.”

limited to single-venue solutions that allow inter-market arbitrage to persist.⁷ This highlights how the market design flaws identified by Budish, Cramton, and Shim (2015) are likely to persist in the future given political constraints.

As emphasized by Budish, Cramton, and Shim (2015), FBAs would simultaneously solve two very different problems: the waste of the speed-based arms race *and* the asymmetric access to publicly-available market information, whereby some participants receive public data fractions of a second before others. Understanding the current HFT controversy in terms of the *arms race*, on the one hand, and *asymmetric access*, on the other, is useful for unpacking the political economy of the debate. For while the arms race may raise concerns over the utility of allowing asset prices to move within ever shorter time frames, it is in many ways as old as financial markets themselves.⁸ In contrast, there is far greater controversy over asymmetric access to market information, which for many, including attorney general Schneiderman, has the hallmarks of insider trading.⁹ And instances of this asymmetric access, such as Thomson Reuters giving HFT firms early access to consumer survey data, have in fact been the subject of government enforcement actions.¹⁰

In this paper, we examine “SIP latency arbitrage,” a seemingly arcane aspect of the modern financial market that has emerged as arguably the most controversial use of low-latency, asymmetric access to information in financial markets. Significantly, the source of this asymmetric access is structurally embedded in the regulation of trade and quote data arising from U.S. trading venues, making it legally permissible even though it creates a heightened risk of adverse selection in equity trading.

To describe more precisely what “SIP latency arbitrage” is, we must first describe briefly some institutional aspects of modern equities markets. By design, price discovery in U.S. equity

⁷ Budish, Cramton, and Shim (2014) propose to implement FBAs within an alternative trading system (ATS). Single-venue on-demand auctions now exist in two forms. Centralized On-Demand Auctions (CODA) Markets, Inc. runs CODA, an ATS focused on on-demand auctions for small and block trades, and the Chicago Stock Exchange (a public exchange) also now runs on-demand “SNAP auctions.”

⁸ The Rothschilds are said to have used carrier pigeons to communicate the news of Waterloo in time for London trading. See, for example, http://www.moaf.org/publications-collections/financial-history-magazine/111/_res/id=sa_File1/Plundered_by_Harpies.pdf

⁹ See, for example, Rachel Abrams, “Attorney General Vows to Crack Down on ‘Insider Trading 2.0,’” *New York Times*, January 9, 2014.

¹⁰ For background and description of the eventual settlement agreement, see “A.G. Schneiderman Secures Agreement By Thomson Reuters To Stop Offering Early Access To Market-Moving Information,” Press Release, July 8, 2013. Available at <https://ag.ny.gov/press-release/ag-schneiderman-secures-agreement-thomson-reuters-stop-offering-early-access-market>, accessed on January 3, 2017.

trading relies on a system of interlinked, competing market centers through which liquidity providers display the prices at which they are willing to buy and sell securities in hopes of drawing trading interest. Consistent with this design, U.S. regulations mandate that all trading centers disclose quote updates and trades occurring on a venue to two centralized Securities Information Processors (“SIPs”) where members of the public can obtain consolidated pricing information. Yet exchanges also sell the exact same—or better—information to those willing to pay for it, allowing clients to avoid waiting for the SIPs to process and disseminate the data. Accessing this so-called “direct feed” data accordingly allows some market participants to obtain right now the data the public will obtain in a bit. The delay is miniscule and becomes smaller every year, but any discrepancy creates the potential for risk-free arbitrage. The financial markets are thus currently organized as a two-tiered system where a trading firm can access public information before it is in fact public.

The discrepancy between direct feed data and the SIP data is particularly relevant because trading rules benchmark trades to the national best bid or offer (the “NBBO”) available across exchanges. For instance, brokers’ best execution duties encourage brokers to fill retail market orders at (or better than) the NBBO. Likewise, dark pools allow clients to post orders to buy or sell that are “pegged” to the prevailing NBBO. If brokers or dark pools determine the NBBO from the SIPs, which was the historical practice, the latency with which the SIPs process the NBBO relative to a trader using direct data feeds thus generates an information asymmetry, allowing fast traders to choose whether to trade or not at NBBO prices they know to be stale.

At the same time, any such arbitrage play requires that there be liquidity required to execute the play—and it is possible that few such arbitrage plays occur. Consequently, the scope of SIP latency arbitrage is fundamentally an empirical question. Until recently, however, understanding the extent to which these opportunities actually arise has been hampered by the absence of detailed information concerning the speed advantage of traders who use exchanges’ proprietary data rather than data from the SIPs.

We use new timestamp data provided by the two SIPs to conduct the first market-wide analysis of the latency with which the SIPs process quote and trade data, and we present new results regarding the economic significance of SIP latency.¹¹ For ease of computation, we focus

¹¹ Our paper builds on important prior work in this area focusing on BATS trading (Ding, Hanna, and Hendershott (2014)).

on all trades involving the Dow Jones 30 during the first eleven months of these new reporting requirements. We show how to reconstruct for each trade in our sample the NBBO that prevailed on the SIP (the “SIP NBBO”) at the microsecond in which the trade occurred, along with the NBBO that was theoretically possible were there no latency at all in transmitting quote updates (the “Direct NBBO”). Reconstruction of this “direct feed” NBBO is made possible by the fact that for each quote update from an exchange, the new timestamp data includes the time at which a quote update was released by the exchange matching engine and therefore available for distribution over an exchange’s direct proprietary data feed.¹²

To preview our specific findings, we first document descriptively that the mean time gap between the time a quote update is recorded by an exchange matching-engine and the time it is processed by a SIP is now just 1.13 milliseconds. Mean latency for processing trades, however, is approximately 20 times higher, clocking in at 22.84 milliseconds.¹³ While latencies are small on average, we document long right-hand tails for both quote and trade reports. For instance, more than 2% of all quote updates in our sample transactions from the Nasdaq BX and the Chicago Stock Exchange have latencies exceeding 10 milliseconds.¹⁴

As noted, we also use the new timestamp data to explore empirically the economic significance of these latencies. We first consider the extent to which liquidity takers and liquidity providers are differentially affected by SIP reporting latencies. Somewhat surprisingly, both classes of traders are commonly alleged to be injured by SIP reporting latencies, often at the hand of the other. For instance, a widely-followed Department of Justice investigation into retail market-making firms Citadel and KCG is reportedly premised on the allegation that market makers filling marketable orders at (or within) the SIP-generated NBBO may do so at stale

¹² As emphasized below, the Direct NBBO is a hypothetical construct that approximates what traders actually observe if they subscribe to exchanges’ direct data feeds. See Section 3 for discussion.

¹³ The slower processing time for trades largely reflects the fact that nearly one-third of trades occur in non-exchange venues whereas quote updates are disseminated by exchange matching engines. Excluding trades executed in non-exchange venues, mean reporting latency for trades is less than 1 millisecond.

¹⁴ As we discuss in greater detail below, we also document that the variation exhibited by quote and trade latencies reflects the institutional structure of SIP reporting obligations, with quote and trade reports in Tape A securities released by the NYSE matching engine arriving at the NYSE-SIP in Mahwah, New Jersey, almost instantaneously, while quote and trade reports in Tape A securities occurring on the Nasdaq matching engine in Cartaret, New Jersey over 188 microseconds more slowly (due to the approximate time it takes light to travel 35 miles). The data confirm this basic institutional prediction, strengthening our confidence in the quality of measurement.

prices to the disadvantage of retail investors using marketable orders.¹⁵ At the same time, a central thesis in the widely-followed book *Flash Boys* is rooted in the strategic use of marketable orders by HFT firms to “pick off” resting limit orders that have been pegged to stale NBBO prices.¹⁶

Overall, our analysis suggests SIP reporting latencies generate remarkably little scope for exploiting the informational asymmetries available to subscribers to exchanges’ direct data feeds, regardless of whether trading is targeted at liquidity takers or at liquidity providers. Indeed, with respect to liquidity takers, on a size-weighted basis, liquidity-taking trades in our sample that were priced at either the SIP NBB or the SIP NBO *gained* on average \$0.0002 per share by having their trades priced at the SIP NBBO rather than the Direct NBBO. Moreover, approximately 97% of trades within our sample occur at a time when the SIP NBBO and Direct NBBO match. This simple fact highlights the low probability that the choice of NBBO benchmark matters at all for liquidity-taking trades at the best ask or best offer. And even among the 3% of trades in our sample where SIP-pricing affected a trade’s profitability, less than 10% left a liquidity-taking trader in a worse position. Although surprising in light of contemporary debates about equity market structure, this finding makes sense. The NBBO will often change in response to serial buy (sell) orders so that late-arriving buy (sell) orders benefit from the stale SIP quotes that have yet to reflect the new trading interest.

To be sure, these findings suggest some *liquidity providers* might suffer an avoidable economic loss by trading at the SIP NBBO rather than at the Direct NBBO. However, we find little evidence that these trades are the result of fast traders using market orders to “pick off” stale limit orders priced at the SIP NBBO to earn risk-free profits. Specifically, our analysis exploits the fact that such an arbitrage play would require a pair of trades, and we find that at

¹⁵ For instance, suppose a direct feed showed the NBBO changing from \$10.00 x \$10.01 to \$9.99 x \$10.00, while the SIP’s NBBO remained at \$10.00 x \$10.01. A broker might fill buy orders by selling to them at \$10.01 (the stale NBO reflected in the SIP NBBO) rather than at \$10.00 (the NBO shown in its direct feed).

¹⁶ As an illustration of this behavior, consider the following example given in Fox, Glosten & Rauterberg (2015). In it, an institutional investor posts to a dark venue a midpoint buy order for a security when the NBBO is \$161.11 x \$161.15 so that an incoming market order to sell would result in this order being filled at \$161.13. However, if the exchange holding the best ask subsequently decreases its displayed quote from \$161.15 to \$161.12 while the midpoint order rests in the dark pool, a fast trader can detect the new NBBO before the dark venue, providing it a momentary opportunity to send an immediate-or-cancel sell order to the dark venue that will execute at the stale midpoint of \$161.13. Upon receiving confirmation, the fast trader can cover the resulting short position by sending a marketable buy order to an exchange to execute at the new national best bid of \$161.12, producing a penny of risk-free profit. In the meantime, the institutional investor—rather than buying at \$161.115, the actual midpoint—buys at \$161.13.

most 0.8% of these liquidity-taking trades could be part of such a strategy. Equally important, while our sample includes over \$4 trillion of trades, we estimate that liquidity providers trading at the SIP NBBO could have saved just \$11 million in lost profits had they transacted at the Direct NBBO instead. This latter finding substantially undercuts the likelihood that stale quote arbitrage generates sufficient economic rents to explain the high speed arms race.

Finally, we also assess the extent to which SIP reporting latencies can affect a venue's trade execution statistics which are routinely used by brokers to route orders. As we show below, any divergence between the SIP and Direct NBBOs creates the possibility for conflicting trade execution measures depending on which NBBO a venue chooses to use as its pricing benchmark. In this regard, aside from the economic costs of SIP reporting latencies on trader welfare, SIP reporting latencies can independently undermine the reliability of a venue's published trade execution statistics.¹⁷ Our results show, however, the low likelihood that a trading center's choice of NBBO benchmark can meaningfully affect their trade execution performance metrics. Specifically, calculating effective spreads using the Direct NBBO rather than the SIP NBBO changes effective spreads by less than 1.9 percentage points for exchange trades and less than a half percentage point for all non-exchange trades.

In summary, our results show that, whatever the economic significance of SIP latency in the past, SIP latency does not currently play a meaningful role in creating profitable arbitrage opportunities. More generally, absent other documented evidence of legally permissible asymmetric access affecting trading outcomes, we posit that the controversy surrounding HFT is best resolved as a question about the utility of allowing the speed-based arms race to continue.

This paper is most closely related to two recent studies of latency arbitrage. Wah and Wellman (2013) estimate the prevalence of latency arbitrage opportunities created by market fragmentation when two or more exchanges create a crossed market (i.e., when the best bid on one exchange creates a NBB that is greater than the NBO). However, their analysis is based on simulated data from an agent-based model, while our approach is empirical. More relevant to our empirical analysis of stale quote arbitrage is Ding, Hanna & Hendershott (2014). Using proprietary data feeds from select exchanges, they study the latency between NBBO updates provided by the publicly-available SIP and NBBO updates calculated using direct data feeds for

¹⁷ Returning to the example in n. 15, by using the SIP NBO of \$10.01, the broker would report an effective spread of just \$0.01 (twice the difference between the trade price of \$10.01 and the midpoint of the SIP NBBO) rather than the actual effective spread of \$0.03.

a trader based at BATS exchange in Secaucus, New Jersey. For such a trader, they find that price dislocations between the two observed NBBOs average 3.4 cents and last on average 1.5 milliseconds. Using a single trading day for Apple, Inc., they use these estimates to conclude that a fast trader could theoretically earn up to \$32,000 over the course of the trading day by trading against stale orders in dark pools based on the volume of off-exchange trades. This estimate, however, assumes each off-exchange trade is made during a period of price dislocation. Our data, in contrast, permits analysis of how many trades are actually made during a period of price dislocation across both exchange and non-exchange venues, enabling a precise estimate of the probability that a trade is adversely affected by latency arbitrage. Our data also permits an estimate of the trading gains and losses traders experience by having their trades priced at the SIP NBBO. Consequently, our results establish that such fast traders are not likely to be as highly compensated as the analysis in Ding, Hanna, and Hendershott (2014) suggests.

Finally, while our results establish that there is little scope in equity markets currently for latency arbitrage arising from stale SIP quotes, we caution that these results should not be over-interpreted. In particular, our results do not rule out other types of latency arbitrage that might be prevalent in the current environment. Nor do our results rule out the possibility that latency arbitrage arising from stale SIP quotes might have been prevalent in the quite recent past (e.g., 2014), for the simple reason that our data are not available until mid-2015. Nonetheless, our results do clarify that a popular narrative regarding stale-quote arbitrage would appear to be scarcely relevant to markets in 2015-2016, and they provide the first broad-based evidence on the extent of quote, trade, and NBBO latency using the SIPs' new microsecond timestamps.

The remainder of this paper is organized as follows. Section 2 provides institutional details regarding the rules governing the dissemination of trade and quote data and the theoretical advantage they provide to fast traders. Section 3 summarizes the new microsecond timestamps and sample selection choices. Section 4 presents our empirical estimates of trade and quote reporting latencies. Section 5 examines the economic consequence to liquidity takers and liquidity providers of having trades priced at the SIP NBBO rather than the Direct NBBO. Section 5 also analyzes how differences between these two NBBOs can affect a trading center's trade performance statistics. Section 6 concludes.

2. Institutional Background

At present, there are three national market plans governing the dissemination of quote and trade data for National Market System (NMS) equity securities. These three plans are required by Rule 603 of Regulation National Market System (Reg. NMS) and reflect the historical structure of U.S. equity markets.¹⁸ For trades in NYSE-listed securities (“Tape A” securities) and securities listed on regional exchanges and their successors (“Tape B” securities), the Consolidated Trade Association (“CTA”) Plan requires all exchanges and FINRA to report last sale information to the Securities Industry Automation Corporation (“SIAC”), a subsidiary of the NYSE which acts as the central SIP for any transaction in Tape A and Tape B securities. The Consolidated Quotation (“CQ”) Plan similarly obligates exchanges and FINRA to report to the SIAC any change in the best bid or best offer (including changes to the number of shares) currently available on each trading venue for Tape A and Tape B securities, which the SIAC uses to calculate the NBBO for these securities.¹⁹ For transactions in Nasdaq-listed securities (“Tape C” securities), the Unlisted Trading Privileges (“UTP”) Plan governs reporting obligations for both trades and quotations. Under this plan, exchanges and FINRA must provide trade and quote updates in any Tape C securities to Nasdaq, which operates as the SIP for transactions in these securities. We refer to the SIP managed by the SAIC as the “NYSE SIP” and the SIP managed by Nasdaq as the “Nasdaq SIP.”

While the trade reporting plans initially focused on exchange-based trades, the SEC has required since March 2007 that all off-exchange transactions be reported to a formal FINRA-managed Trade Reporting Facility (a “FINRA TRF”) (O’Hara & Ye, 2011). At present, FINRA manages two facilities operated separately by the NYSE and Nasdaq, with the Nasdaq facility receiving the vast majority of trade reports from non-exchange venues.²⁰ In combination with FINRA’s trade reporting obligations under the CTA and UTP Plans, this SEC reporting

¹⁸ Rule 603 requires that all exchanges and FINRA “act jointly pursuant to one or more effective national market system plans to disseminate consolidated information, including a national best bid and national best offer, on quotations for and transactions in [National Market System] stocks.” Adopted pursuant to Section 11A of the Securities Exchange Act of 1934 (the “Exchange Act”), Rule 603 reflects Section 11A’s mandate that the SEC develop rules that ensure trading data historically published by exchanges and broker-dealers for their customers is made available to all investors “on terms which are not unreasonably discriminatory.”

¹⁹ FINRA operates an Alternative Display Facility (the “FINRA ADF”) through which non-exchange venues (such as an electronic communications network, or “ECN”) might choose to disseminate quotations from their subscribers. At present, no venue disseminates any quotations through the FINRA ADF.

²⁰ For instance, in unreported results, we find that 87.27% of non-exchange trades within our sample were reported to the TRF operated by the Nasdaq, and 12.73% were reported to the NYSE TRF.

requirement for FINRA members means that off-exchange trades made through a broker-dealer internalizer or in a dark pool are now effectively segregated and reported to the appropriate SIP as having been executed at a FINRA TRF.

In addition to sending market data to the SIPs for consolidation, exchanges and FINRA TRFs are also permitted to sell access to the same transaction data directly to customers through proprietary data feeds.²¹ Importantly, the SEC has interpreted Rule 603 to require only that exchanges *transmit* data to the SIPs no later than they transmit data through their proprietary data feeds.²² This implies that traders subscribing to a direct feed avoid the inevitable latency arising from the SIPs' obligation to consolidate and process transaction information before disseminating it.

To establish the magnitude of this delay, Table 1 provides processing times for trade and quote information disclosed by both SIPs from 2014 through the second quarter of 2016.²³ For Tape A and B securities, the time between receipt of a transaction report by the NYSE SIP and its subsequent dissemination of that report averaged 410 microseconds for trades and 450 microseconds for quote updates. Processing times for Tape C securities were slightly higher at 700 microseconds and 750 microseconds, respectively. A trader subscribing to an exchange's direct feed can accordingly avoid these processing-related latencies when receiving the exchange's transaction data.²⁴

In addition to allowing exchanges to sell their direct feed data, the SEC also allows exchanges to sell co-location services. These services allow customers to place their computer

²¹ Exchanges submit to the SEC for review specific proposals to offer proprietary feeds. Fees for accessing these feeds must also be reviewed by the SEC.

²² See *In re NYSE LLC*, Exchange Act Release No. 34-67857, at 2 (Sept. 14, 2012). In adopting Reg NMS, the SEC similarly noted that while Rule 603 requires exchanges that offer proprietary feeds to do so on terms that are fair and reasonable and not unreasonably discriminatory, "Rule 603(a) will not require a market center to synchronize the delivery of its data to end-users with delivery of data by a Network processor to end-users." This SEC guidance accordingly permits subscribers of exchange data to *receive* this data before a SIP so long as the exchange *releases* the data to the subscriber no sooner than it does for the SIP.

²³ Data in Table 1 for Tape A and Tape B securities can be found at [https://www.nyse.com/publicdocs/ctaplan/notifications/trader-update/CTA%20SIP%20Q16%20Consolidated%20Data%20Operating%20Metrics%20Report%20\(7-13-16%20Update\).pdf](https://www.nyse.com/publicdocs/ctaplan/notifications/trader-update/CTA%20SIP%20Q16%20Consolidated%20Data%20Operating%20Metrics%20Report%20(7-13-16%20Update).pdf). Data for Tape C securities can be found at <http://www.utpplan.com/DOC/UTP%202015-Q4%20Stats%20with%20Processor%20Stats.pdf>.

²⁴ The secular decline in processing-related latencies shown in Table 1 reflect several initiatives of both SIPs. As Tabb (2016) summarizes, "The quality of SIP data over the past few years has improved and is scheduled to dramatically improve again in the near future. Latencies for the SIP currently are approximately 500 microseconds, but they are scheduled to decrease to 50 microseconds by year end and to less than 25 microseconds within a year. By any account, the SIPs (UTP and CTA) have done a yeoman's job improving SIP latency and robustness over the past few years."

servers in close physical proximity to the exchanges' matching engines to minimize the transit time of the exchanges' market data. For Tape A and B securities, co-location accordingly allows a trader to avoid the additional latency a transaction report experiences when traveling from a market center to the NYSE SIP in the NYSE's Mahwah, New Jersey datacenter (the same datacenter housing the NYSE's matching engine); for Tape C securities, it avoids the latency a report experiences when traveling to the Nasdaq SIP's processing platform in Carteret, New Jersey (the same datacenter housing Nasdaq's matching engine).²⁵

In light of widespread concerns about the advantages these direct feeds provide fast traders, SEC Chair Mary Jo White stated in Congressional testimony (White, 2015) that she had "asked the exchanges and the SIPs to incorporate a time stamp in their data feeds to facilitate greater transparency on the issue of data latency." We use these new timestamps in the analyses below.

3. Data and Sample Selection

We obtain all trade and quote reports published by the two SIPs for the common stock of firms listed within the Dow Jones 30 as of August 1, 2015. We focus on the Dow Jones 30 in light of popular claims that high frequency trading firms are "overwhelmingly interested in heavily traded" securities (Lewis, 2004: p. 115). Our sample period commences with the full implementation of the new microsecond timestamps on August 6, 2015 (the first full day on which exchanges complied with the new reporting requirements) and ends on June 30, 2016.²⁶ To ensure that all quotes and trades occur during the trading day after the opening cross and before the closing auction, we subset the data to exclude quotes and trades occurring before 9:45:00 and after 15:44:59.999999.²⁷ Because we use the quote data to generate the NBBO, we further restrict our analysis to those quotations that are eligible to establish an exchanges' best offer or best bid (i.e., quotation updates having a condition of A, B, H, O, R, W, or Y). Finally, for our latency analysis in Section 4, we exclude quote or trade records with missing venue

²⁵ To appreciate the importance of location decisions, suppose transaction reports are transmitted at the speed of light. In such a scenario, a trader co-located at Nasdaq's Carteret data center would receive transaction reports from Nasdaq approximately 188 microseconds faster than a trader who also subscribed to Nasdaq's direct feed but is located in Mahwah, 35 miles to the north. If the Mahwah-based trader relied on the Nasdaq SIP, the first trader's speed advantage would increase to nearly 1,000 microseconds after incorporating the SIP-processing delay.

²⁶ The implementation date for Tape C securities was July 27, 2015 and August 3, 2015 for Tape A and Tape B securities. However, the BATS Y exchange did not fully commence using the new timestamps until August 6, 2015.

²⁷ Following the conventions of software, we record time in microseconds as one-millionth of a second, so that 1 microsecond past 9:45am is recorded as 9:45:00.000001.

timestamps or with venue timestamps that are subsequent to the SIP timestamp.²⁸ Imposing these conditions results in a core sample of 385,028,820 trades and 6,212,857,437 quote updates.²⁹

We next use these data to construct two versions of the NBBO that prevailed at the time of each trade in our sample.³⁰ The first version calculates the NBBO using the timestamp showing the time (in microseconds) at which a SIP disseminated a quote update. This version reflects the NBBO that was available from the SIP at the moment of each trade; therefore, we designate it as the “SIP NBBO.” The second version calculates an alternative NBBO using the new “Participant Timestamp,” which shows the time (in microseconds) at which an exchange matching engine reported processing a quote update. This alternative version reflects the NBBO at the moment of each trade in a world with no processing or transmission latencies. Because it is derived directly from exchange data, we designate it the “Direct NBBO.”

Finally, we further exploit the new Participant Timestamps to match each trade to the SIP NBBO and Direct NBBO *that prevailed at the time the trade was executed*. This approach differs from traditional approaches that assign the SIP NBBO to trades using only the SIPs’ timestamp of a trade, which was previously the only timestamp the SIP provided for a transaction. However, the SIP timestamp may not reflect the SIP NBBO that prevailed at the time a venue actually executed the trade due to the transit and processing-related delays associated with the SIP’s dissemination of quotes. For similar reasons, relying on the SIP timestamp of a trade does not permit insight into the Direct NBBO that prevailed at the moment a venue executes a trade. Relying on the Participant Timestamp for trades thus permits a unique insight into how a broker or venue perceived the SIP NBBO and Direct NBBO at the very time they were seeking to price transactions, rather than the time at which the SIP processes the trade report.

To accomplish this matching of trades and NBBOs, we assign to each trade a SIP NBBO and a Direct NBBO based on the microsecond at which the trade was executed. We do so by again

²⁸ This sample selection rule excludes 55,226,095 quote updates (0.9% of all quotes), only 8 of which are due to missing venue timestamps, and 2,811,429 trade records (0.7% of all trades), none of which are due to missing venue timestamps.

²⁹ For all analyses in Section 5, we include in the sample all quote and trade records with venue timestamps that are subsequent to the SIP timestamp, which are excluded in our latency analysis for the reasons set forth in Section 4.

³⁰ For analyses involving the NBBO, we restrict attention to NBBOs as of 10:00am. Because our record of quotes starts at 9:45am, this 15-minute “burn in” phase ensures that our first daily measure of the NBBO reflects the best quotes available across all exchanges.

using the new Participant Timestamp recorded by the SIP for each trade report, which reflects the time (in microseconds) that an exchange or FINRA member reports a trade as occurring on the trading venue. We additionally classify trades as having been buy- or sell-side initiated using the Lee and Ready (1991) algorithm. In so doing, we compare each trade’s execution price to the SIP-NBBO assigned to the trade. This is the logical choice since our research question focuses on whether there is harm to traders on venues that price transactions using the SIP NBBO. For all trades, we retain the SIP-generated timestamp on a trade report to permit analysis of trade reporting latencies, as described in the following section.

Before turning to the results, we want to emphasize that the Direct NBBO is a construct rather than a direct measure. As noted, no trader has access to the Direct NBBO because of the physical distance between exchange matching engines. Nonetheless, the Direct NBBO provides an in-the-limit representation of the advantages of having fast access to exchanges’ trading data. A trader using the fastest direct feeds would have access to market information nearly as current as the Direct NBBO. We note that Ding, Hanna, & Hendershott (2014), while focused only on a subset of exchanges, take advantage of direct measures.

4. Estimating Trade and Quote Reporting Latencies

We define reporting latency as the difference between the timestamp of a transaction reported by a SIP (the “SIP Timestamp”) and the Participant Timestamp, which is the time an exchange matching-engine or broker-dealer records a transaction as having occurred:

$$Latency = Timestamp_{SIP} - Timestamp_{Participant}$$

This definition resembles, but is distinct from, that used by Ding, Hanna, and Hendershott (2014). Those authors analyze the timestamp generated by a proprietary server located at BATS’ trading center that receives transaction reports directly from select exchanges (BATS, Direct Edge, and Nasdaq) as well as from Nasdaq’s SIP feed. Their definition of latency accordingly assesses the delay associated with receiving SIP market data relative to receiving market data from these select exchanges for a trader in Secaucus, New Jersey (the location of the BATS data center).

In contrast, our measure of latency represents the delay between the time a market center processes a transaction and the time when the appropriate SIP disseminates a report for the transaction. As such, it represents the delay created by: (a) the transit time from a market center

to either the NYSE SIP or the Nasdaq SIP, as applicable, and (b) the time it takes for the relevant SIP to process and disseminate the transaction report. In this regard, it can be viewed as the floor latency experienced by all consumers of the SIP data, regardless of their location relative to the SIPs. Our measure also permits analysis of this latency across all market centers and for both NYSE- and Nasdaq-listed securities.

All timestamps are marked in microseconds; therefore, our measure of latency is formally in microseconds. We note, however, that the microsecond timestamps for trades by non-exchange venues are uniformly reflected as having occurred in intervals of 1,000 microseconds (i.e., 1 millisecond). We interpret this pattern as reflecting the fact that most non-exchange venues have continued to record transactions at the level of the millisecond.³¹ Assuming this is the case, our measure of latency will accordingly be biased slightly higher for these trades to the extent the transaction did not occur at precisely the beginning of the reported millisecond. As we discuss below, the delay in transaction reporting for non-exchange trades is so large it could be measured in milliseconds—and hence microsecond precision is not necessary to get an accurate sense of latency for these transactions.

a. Institutional Background on Clock Synchronization.

Because our analysis relies on comparing timestamps imposed by two different data centers, an important preliminary issue to consider is clock synchronization. In particular, if the clock used by a SIP and the clock used by a market participant are not synchronized, our latency measure may be inaccurate. Not surprisingly, addressing similar clock synchronization concerns has also been central to the SEC’s proposed Consolidated Audit Trail (CAT), which is designed to allow the reconstruction of all quote and trade activity across multiple market centers. In this subsection, we provide institutional details regarding why synchronization issues for this study, like the CAT more generally, are unlikely to be material in today’s markets. Readers familiar with these issues from the CAT or otherwise are invited to skip to Section 4(b) where we commence presentation of our empirical results.

³¹ FINRA has required since 2014 that firms report a trade’s execution time in milliseconds when reporting trades to the FINRA facilities if the firm’s system captures time in milliseconds. See FINRA Regulatory Notice 14-21 (May 2014), available at <http://www.finra.org/sites/default/files/NoticeDocument/p506337.pdf>. The new timestamp requirements permit FINRA to convert to microseconds any transaction times submitted in milliseconds by a FINRA member. See NasdaqTrader.com, UTP Vendor Alert #2015 - 7 : New Timestamp Definitions for July 2015 Release, available at <https://www.nasdaqtrader.com/TraderNews.aspx?id=UTP2015-07>. We assume that clocks record transaction time in milliseconds by rounding microseconds to milliseconds.

Considering modern computer clock synchronization protocols, the scope for non-synchronized clocks in recent years is likely small. This is partly because of the demise of manual, mechanical time-stamping of transactions in favor of automated order-entry systems. For instance, Network Timing Protocol (NTP) clients have long been included in servers and personal computers, permitting computer clocks to be synchronized within milliseconds of the US national time standard, or UTC(NIST) (Lombardi, 2000).³² Alternative protocols such as IEEE 1588 Precision Time Protocol (PTP) are also commonly available for more advanced servers, ensuring clock times are within nanoseconds of the UTC(NIST). In releasing the new microsecond time stamp specifications, the CTA, CQS and UTP accordingly required that exchanges use a clock synchronization methodology ensuring timestamp tolerances of 100 microseconds. In releasing its plan for the CAT, the SEC (2016) further reports that these tolerances apply to the two SIPs and that the absolute clock offset on exchanges averages just 36 microseconds.

Clock synchronization is potentially a greater issue for non-exchange venues and broker-dealers. In contrast to the 100 microsecond tolerance used by exchanges, in recent years FINRA required that all computer system clocks and mechanical time stamping devices of FINRA members be synchronized to within one second of the UTC(NIST).³³ In practice, however, brokers responsible for handling the largest share of trading volume appear to utilize clock syncing with much greater precision than this formal requirement. In anticipation of the CAT, for instance, FINRA recently adopted new Rule 4590 which reduces the drift tolerance for computer clocks that record transactions in OTC and NMS equity securities from one second to 50 milliseconds. In adopting the new standard, FINRA noted that firms accounting for 95 percent of reportable transactions to FINRA's Order Trail Audit System (OATS) already report events in milliseconds and comply with the 50 millisecond clock synchronization standard. Likewise, in responding to the proposed rule, dark pool operator IEX noted the standard could be further reduced below 50 milliseconds given the system capabilities of most FINRA firms, citing its own synchronization standard of one millisecond.

³² National time standards are synchronized (essentially, averaged after removing outliers and consistent errors) to yield an international reference called Universal Coordinated Time. In the United States, the National Institute of Standards and Technology (NIST) maintains an atomic clock that serves as the country's primary time standard, or UTC(NIST). It generally tracks the UTC to within 5 nanoseconds.

³³ FINRA Rule 7430 applied through 2015.

Indeed, for many of the most important FINRA members such as dark pool operators and broker-dealer internalizers, the emergence of co-location services has undoubtedly facilitated synchronization tolerances of far less than 50 milliseconds. For instance, firms such as IEX that are hosted by the Equinix NY4 datacenter in Secaucus, New Jersey (which also hosts the matching engines of BATS and Direct Edge) can utilize a service called “High Precision Time” offered through Perseus’ Communications. The service allows synchronization with UTC(NIST) through both NTP and PTP protocols and promises “certified time stamps to subnanosecond accuracy.”³⁴ In December 2014, Nasdaq announced it would be offering the same service to its customers at its U.S. datacenter in Carteret, N.J.³⁵ The 2015 Customer Guide for NYSE Euronext similarly offers four different connection protocols to the UTC(NIST) to ensure timestamp “accuracy on the order of nanoseconds.”³⁶

All of these factors reduce the likelihood that clock synchronization issues materially affect our latency measure, particularly for exchanges, but even for non-exchange venues and broker-dealers. However, as noted by Angel (2014), any non-zero synchronization tolerance and random variation surrounding it will introduce some degree of clock synchronization error when reconstructing market conditions using time-stamped records from multiple market centers. Consistent with these concerns, our data do reveal evidence of such residual noise in the form of transaction reports with negative latency. In particular, approximately 0.88% of quote updates and 0.72% of trade reports had a SIP Timestamp that *preceded* the time reported in the Participant Timestamp. Given that a transaction must be processed by a participant before it is even received by a SIP, these outcomes obviously represent physical impossibilities.

Close inspection of the data reveals that the majority of these reports resulted from clock synchronization issues at the NYSE Arca and the NYSE SIP.³⁷ For instance, between May 16, 2016 and June 6, 2016, more than 75% of the daily quote updates and trade reports emanating from the NYSE Arca had negative reporting latencies. These reports from Arca account for

³⁴ See Equinix Press Release, Equinix is First to Offer Global Access to High Precision Time™ from Perseus Telecom, Sep. 10, 2014, available at <http://www.prnewswire.com/news-releases/equinix-is-first-to-offer-global-access-to-high-precision-time-from-perseus-telecom-274588571.html>.

³⁵ See Nasdaq Press Release, Perseus Selected by Nasdaq for Time Stamping Service at US Data Center, Dec. 3, 2014, available at <http://www.prnewswire.com/news-releases/perseus-selected-by-nasdaq-for-time-stamping-service-at-us-data-center-300004156.html>

³⁶ See Intercontinental Exchange, Infrastructure and Americas User Guide (Feb. 2015), available at www.nyxdata.com/doc/243267.

³⁷ According to NYSE Euronext, “The Arca issue... identified was due to a bug that was fixed.” Personal communication between authors and NYSE officials, dated August 2, 2016.

approximately 65% of quote updates having negative latencies and more than 77% of trade reports with negative latencies. Excluding the negative latencies appearing in these Arca reports for this two week period, negative latencies appeared in 0.31% of quote updates and 0.162% of trade reports. Of these, 99.953% of the quote updates and 99.917% of the trade reports were in Tape A securities and arose across all exchanges trading Tape A securities, indicating occasional clock syncing issues at the NYSE SIP.³⁸

Evidence of these negative latencies within our sample highlight the potential of clock synchronization issues to arise even with the institutional structures described previously. Because we lack a record of the actual UTC(NIST) for each transaction report, we are unable to measure the extent to which clock syncing affects our measure outside of these negative latencies. However, in all analyses in Section 4(b), we exclude from our sample any transaction report having a negative latency. As noted below, our resulting latency estimates generally reflect the institutional structure of the SIP-reporting regime, providing confidence that any residual clock synchronization issues do not materially bias our analyses.

b. SIP Reporting Latencies Across Trading Venues

As we are unaware of any prior work utilizing these new timestamps, we first report in this section some of the basic descriptive patterns of our latency measures. Table 2 presents the mean, standard deviation, median, and 90th percentile measures of latency by trading venue according to where the transaction originated, both for quote updates (Panel A) and trade reports (Panel B).³⁹ Because securities within the Dow Jones are listed on both the NYSE and Nasdaq, we also separate transactions according to whether transaction reports were sent to the NYSE SIP (Tape A securities) or the Nasdaq SIP (Tape C securities).

In both panels, we group exchanges by the location of their matching engines to facilitate analysis of the role of transit time in explaining variation in reporting latencies. All three exchanges controlled by the NYSE (the NYSE, NYSE MKT, and NYSE Arca) are hosted in the

³⁸ Analysis of negative latencies in the transaction reports for all NMS equity securities occurring on two randomly-chosen trading days reveals a clear dependency on the structure in which the NYSE SIP processes transaction reports. Under the technical specifications for the CTA Plan and the CQ Plan, transaction reports in Tape A securities and Tape B securities are processed separately across twenty-six multicast lines with each line processing approximately 250 securities according to its trading symbol. On both trading days, when a negative latency appeared for a security's quote update or trade report, negative latencies also appeared for the quote updates and trade reports of every other security assigned to the same multicast line before ceasing for all securities so assigned.

³⁹ As we show below, our latency measure exhibits a long and thick right-hand tail, implying that the median may be a better estimator of the center of the distribution than the mean. See Appendices A and B, described below.

NYSE's datacenter in Mahwah, New Jersey. The three exchanges owned by Nasdaq (Nasdaq, Nasdaq OMX BX, and Nasdaq OMX PSX) are hosted in Nasdaq's datacenter in Carteret, New Jersey. Five other exchanges are hosted in Equinix's NY4 and NY5 datacenters in Secaucus. This includes the four exchanges owned by BATS Global Markets—BATS Exchange, BATS Y, Direct Edge A, and Direct Edge X—as well as the matching engine of the Chicago Stock Exchange responsible for trades in all Dow Jones equity securities.⁴⁰ The Equinix facility also hosts the trading system for the National Stock Exchange (NSX), which recommenced trading on January 1, 2016 after ceasing operations in early 2015. However, we list NSX separately given its idiosyncratic reporting of Participant Timestamps, as described below.

In Figure 1, we map the approximate location of these three datacenters to illustrate their physical distances from one another, as well as from the two SIPs. Because the NYSE runs the NYSE SIP, transaction reports in Tape A securities transmitted by an NYSE exchange need only travel the distance between the NYSE matching engine and the SIP processor within the Mahwah datacenter. Across all Tape A transaction reports, reporting latencies are accordingly the smallest for transactions occurring on an NYSE-owned venue. For instance, NYSE mean (median) quote update and trade report latencies are 690 (301) microseconds and 356 (298) microseconds, respectively.

In contrast, quote updates and trade reports in Tape A securities occurring on an exchange hosted in the Equinix Secaucus datacenters must travel approximately 21 miles before being processed by the NYSE SIP. Reports in Tape A securities occurring on a Nasdaq exchange must travel even further in light of the approximately 35 miles separating the Nasdaq datacenter from the NYSE's facility. These distances account for the larger latencies for Tape A securities for quote and trade reports arising from transactions on the Nasdaq- and Equinix-based matching engines. With regard to quote updates, for instance, median reporting latencies in Tape A securities for the three Nasdaq exchanges are approximately 900 microseconds, while those for the BATS exchanges range from 491 to 517 microseconds. Median reporting latency for quote updates occurring on the Chicago Stock Exchange (CHX) is slightly higher at 839 microseconds.

⁴⁰ The Chicago Stock Exchange (CHX) also maintains a matching engine in Equinix's CH3 datacenter in Chicago. A "Matching Engine Committee" at CHX determines which of the two matching engines will handle transactions in securities that can be traded on the CHX. At present, only seventy securities are assigned to the Chicago matching engine; all others are matched in New Jersey, including all securities in our sample. See Chicago Stock Exchange, New Jersey Data Center Eligible Symbols (July 18, 2016), available at <http://www.chx.com/market-data/nj-data-center/>.

Given that reporting latencies for the CHX are similarly higher than those of the BATS exchanges in Tape C transactions, we attribute these higher latencies to the superior network performance of the BATS-controlled exchanges. Reporting latencies for Tape C securities display the opposite patterns across the three groups of exchanges, consistent with the fact that the exchanges closest to Nasdaq's datacenter should have the shortest transit times.

The primary exception among exchanges to these geographic-centered patterns appears in the transaction reports for the National Stock Exchange. While median latency for quote updates are just slightly higher than for other Equinix-based exchanges, mean reporting latencies are considerably higher at approximately 18 milliseconds for Tape A securities and over 41 milliseconds for Tape C securities. As suggested by the extraordinarily large standard deviations reported in the table, these very high mean values reflect extreme outliers. Reporting latencies for trade reports are even more out of line with the latencies one would expect given NSX's geographic location relative to the NYSE- and Nasdaq-SIPs. For instance, mean (median) reporting latencies for trade reports at the NSX were nearly 53 milliseconds (52 milliseconds) for Tape A securities and 52 milliseconds (53 milliseconds) for Tape C securities. Notably, transaction reports can traverse the 700 miles from Chicago to the two SIPs in just 9 milliseconds.⁴¹ These reporting latencies would accordingly appear to reflect either an idiosyncratic system for recording Participant Timestamps or extremely slow and inconsistent report processing at the NSX.⁴²

Across exchanges, mean trade latency was generally lower than mean quote latency for securities on both tapes; however, this difference largely reflects the thick right-hand tail of the distribution of quote updates. Among trades on exchanges, for instance, the 90th percentile latency was roughly twice the size of the median; for quote updates, the 90th percentile latency was closer to four times the size of the median quote latency. Focusing on median latencies

⁴¹ In unreported results, we calculated reporting latencies for the approximately seventy securities that continue to match on the CHX's matching engine in Chicago located in Equinix's CH3 datacenter. For Tape B securities, mean (median) reporting latencies were 9,139 (9,005) microseconds for quote updates and 9,402 (9,409) microseconds for trades. For Tape C securities, mean (median) reporting latencies were 8,853 (8,749) microseconds for quote updates and 8,190 (9,207) microseconds for trades. These latencies reflect the fact that both the NYSE SIP and the Nasdaq SIP are approximately 700 miles from the CHX matching engine in Chicago.

⁴² Data in this study were shared with the NSX, which provided the following statement regarding these findings: "The NSX is aware of both the research into the SIP reported quote and trade reporting latencies and the variances reflected with respect to other exchanges regarding the trade reporting latencies. NSX will conduct its own review of the data to better understand the anomalies of the trade reporting latency times and look forward to working with the authors on their continuing market research."

between trade reports and quote updates, the difference between trade and quote latency for exchanges falls considerably (Tape A Quotes=566 microseconds vs. Tape A Trades=604 microseconds; Tape C Quotes=551 microseconds vs. Tape C Trades=555 microseconds).

With the exception of the NSX, the distribution of latencies for both trade reports and quote updates for exchange transactions was unimodal with extreme kurtosis, highlighting both the strong clustering near the median as well as outliers. In Appendix A and B, we present histograms of reporting latencies for all combinations of exchanges and tapes for both quote updates and trade reports.⁴³ As noted previously, the presence of outliers is particularly prominent within quote updates. Given the considerably greater number of quote updates during the trading day, the long right-hand tail for quote update latencies is consistent with concerns that the large volume of quote message traffic can occasionally overwhelm available network capacity (Nanex, 2014; Ye, Yao, and Gai 2012).

A striking exception to the unimodal distribution of reporting latencies appears in the reporting latency of non-exchange trade reports. In Figure 2a and 2b we present histograms of trade reporting latencies for non-exchange trades in Tape A and Tape C securities, respectively.⁴⁴ As shown in both figures, the distribution of latency across the two tapes is both multi-modal and highly-skewed, resulting in mean and median latencies that are considerably higher than latencies for exchange trades. For Tape A securities, the mean (median) reporting latency is approximately 87 milliseconds (7.1 milliseconds); for Tape C, the mean (median) latency is approximately 101 milliseconds (7.0 milliseconds).

Two features of non-exchange trade reporting most likely account for the peculiar shape of these distributions. First, as noted previously, all non-exchange trades must be reported to one of two TRF facilities, thereby aggregating trades executed by automated wholesalers and dark pools as well as by smaller broker-dealers. While many dark pools and retail wholesalers are co-located at exchanges, smaller FINRA members may be located further away and may have slower network connections to the TRFs. Smaller members of FINRA may also have slower trade reporting protocols, particularly given the amount of time brokers are permitted to report

⁴³ With the exception of the histogram for trade reports at the NSX, all histograms presented in Appendix A and B are truncated at latencies of 4 milliseconds (approximately the 95th percentile of the overall distribution of latency) to facilitate visualization of distributional form. We truncate trade reports for the NSX at 100,000 microseconds given the large number of trades having latencies in excess of 4 milliseconds.

⁴⁴ Because this is a data-rich environment, the structure of the density can be inferred from a histogram without resorting to smoothing choices and kernel density techniques (Silverman 1986).

trades under SEC and FINRA rule-making. For instance, Rule 601 of Reg. NMS simply states that brokers must report trades “promptly,” while FINRA requires trades to be reported to a TRF as soon as practicable, but no later than 10 seconds, following trade execution.⁴⁵ The gap between this slow formal requirement and the comparatively rapid actual implementation highlights the extent to which off-exchange reporting is today conducted through automated systems.

Additionally, regardless of the speed with which a broker reports a transaction to a TRF, reporting latencies for non-exchange trades are also increased by the double-legged nature of the TRF-reporting regime. For instance, a broker who chooses to report non-exchange trades to the NYSE TRF will first report trades in Tape C securities to the NYSE TRF in Mahwah, which will then report the trade to the Nasdaq SIP in Cartaret. For a broker based at the Nasdaq facility, such a process guarantees a reporting latency equal to at least the round-trip transit time between the Nasdaq and NYSE facilities before the Nasdaq SIP even begins processing the report.

Finally, our summary results also highlight what appears to be an inconsistency in time-stamping procedures between the NYSE SIP and the Nasdaq SIP. Evidence of this inconsistency appears in comparing the processing latency reported for trade and quote records in Table 1, and our latency measures in Table 2. In Panel A of Table 3, we set forth the median processing latency for quote updates in the second quarter of 2016 for the NYSE and Nasdaq SIPs from Table 1, as well as the median reporting latencies for quote updates in Tape A and Tape C securities for all exchanges other than the NSX.⁴⁶ We also present the difference between these two medians for each exchange, which represents an estimate of the median transit time experienced by a quote update for each exchange. To illustrate how this estimate compares with the theoretical minimum transit time, we present the time it takes for light to travel the same distance in a vacuum. Finally, we present the ratio of our estimated transit latency to this theoretical minimum. Panel B of Table 3 does the same for trade reports.⁴⁷

As shown in both panels, estimated transit times for quote updates and trade reports for Tape A securities range from approximately 2.5 to 8 times the theoretical minimum. These results are to be expected given that message signals travel slower in fiber optic cable than through a

⁴⁵ See FINRA, Regulatory Notice 14-16: Equity Trading Initiatives: OTC Trade Sequencing, available at https://www.finra.org/sites/default/files/notice_doc_file_ref/Notice_Regulatory_14-46.pdf

⁴⁶ We exclude the NSX given that, as noted previously, quote and trade reporting latencies do not appear to reflect the NSX’s geographic location relative to the two SIPs.

⁴⁷ We omit theoretical minimum transit times where an exchange is located in the same facility as the SIP.

vacuum and must navigate additional networking frictions from an exchange matching engine to the NYSE SIP.

Results for Tape C securities, in contrast, reveal transit times that would appear to defy the laws of physics. For instance, the median reporting latency of approximately 523 microseconds for quote updates on BATS would mean that messages traveled the 16 miles from Equinix’s facility in Secaucus to Nasdaq’s datacenter in Carteret in approximately 63 microseconds—an astounding two-thirds the amount of time it would take for light to travel this same distance. For Tape C transactions on one of the three Nasdaq-owned exchanges, Tables 3A and 3B suggest the total time between the moment a transaction occurred and the moment it was processed and disseminated by the Nasdaq SIP was *less* than the time it took the SIP to just process the report.

Given these findings, we sought to document the manner in which the two SIPs calculated message processing times and imposed the SIP timestamp. With regard to processing times, both SIPs define processing latency from the time a message is received from an exchange to the time it takes to place the message on the multicast feed for distribution.⁴⁸ The SIPs are less consistent, however, with respect to when they impose the SIP timestamps. For the NYSE SIP, the technical specifications of the CTS and the CQS were modified in connection with the roll-out of the new timestamps to make clear that the SIP timestamp “indicates the time that processing a message is completed.” With respect to the Nasdaq SIP, technical specifications were also modified at this time; however, the definition of the “SIP Timestamp” was revised to state simply that it provides “the number of microseconds since midnight EST.”

In light of these disclosures and our empirical results, we suspect Nasdaq’s SIP may be placing its timestamp on a transaction report during its message processing routine, rather than at the conclusion of the routine as is done by the NYSE SIP.⁴⁹ Indeed, as shown in the final column of Table 3A and 3B, adding 200 microseconds to each of the median Tape C reporting latencies in Table 3 would bring the estimated Tape C transit times from all exchanges in line with those of the Tape A latencies. In our analyses below, we account for this possibility by using two versions of “SIP Time” for Tape C trade and quote records. In one, we assume the SIP’s timestamp represents the time the Nasdaq SIP placed the message on its multicast line; in the other, we add 200 microseconds to SIP Time.

⁴⁸ See Financial Information Forum, https://www.fif.com/docs/fif_latency__member_input.ppt

⁴⁹ We rule out the possibility that the discrepancy arises from exchange clocks running faster than the Nasdaq SIP’s since fast clocks on an exchange would affect reporting latencies of the NYSE SIP as well as the Nasdaq SIP.

c. Dislocations of the SIP NBBO and Direct NBBO

An inevitable consequence of these SIP reporting latencies is for the SIP NBBO to lag changes in the Direct NBBO. For instance, across all securities in our sample, the NBB from the SIP NBBO and that of the Direct NBBO differed on average 6,839 times per day. These differences—which, following Ding, Hanna, and Hendershott (2014), we refer to as “Dislocations”—ranged from a daily minimum of 206 for General Electric to a maximum of 138,644 for Apple.⁵⁰ However, as one would expect from the latencies set forth in Table 2, the duration of these dislocations was typically short-lived. Across all dislocations of the NBB, for example, the mean (median) duration was 1,001.6 (489) microseconds. A standard deviation of 567,349.5, however, highlights the existence of outliers. In Figure 3a, we present a histogram of the duration of NBB dislocations which illustrates the thick-tailed nature of this distribution.⁵¹

With regard to the size of these dislocations, mean and median dislocations for the NBB were \$0.0109 and \$0.01, with a 99th percentile of \$0.03. Dislocations of the NBO were similarly slight, having a mean, median and 99th percentile measure of \$0.0109, \$0.01, and \$0.03, respectively. These figures are consistent with the fact that securities in the sample often traded at or near penny spreads. Figure 3b shows the histogram of the magnitude of NBB dislocations, which emphasizes how tightly clustered around a penny these dislocations are. Penny dislocations are well over 90 percent of all dislocations. Dislocations of two, three, and four pennies occur, but are rare. Dislocations of a nickel or above occur so infrequently they cannot be discerned in the graph.⁵²

5. Does Pricing Off the SIP NBBO Harm Traders?

We now turn to an empirical investigation of the economic significance of SIP reporting latencies for “slow” traders such as retail traders and institutional investors. These concerns have been at the center of the recent controversies surrounding SIP reporting latencies. The possibility that SIP reporting latencies can be used to harm these investors has gained widespread

⁵⁰ The mean number of dislocations of the NBO was approximately 8,433, ranging from a minimum of 203 for GE to a maximum of 139,997 for Apple. As noted previously, we estimate NBBO dislocations starting at 10:00 am following a 15-minute burn-in phase.

⁵¹ The duration of dislocations for the NBO are similar to those of the NBB. In the interest of space, we present results for the NBB only.

⁵² There are nonetheless some quite rare dislocations that are large in magnitude (e.g., over \$1).

attention since the publication of Michael Lewis' *Flash Boys* in 2014. In general, these concerns are grounded in the fact that trading venues commonly fill marketable orders by reference to the NBBO. For instance, a broker-dealer internalizer such as Citadel will generally fill retail market orders obtained from payment-for-order flow agreements at (or slightly) within the prevailing NBBO. Likewise, both exchanges and dark pools generally permit order types that are "pegged" to the near or far side of the NBBO, and pegged orders are commonly used by institutional investors. When a market order arrives in a venue holding these orders, the order is then priced at the prevailing NBB (for marketable sell orders) or NBO (for marketable buy orders).

In any of these cases, the latency of the SIP-generated NBBO raises the possibility that venues will fill orders at stale prices if they use the SIP NBBO as their benchmark. For example, where a direct feed shows the NBBO changing from \$10.00 x \$10.01 to \$9.99 x \$10.00, a broker-dealer internalizer might fill buy orders by selling at \$10.01 (the stale NBO reflected in the SIP-NBBO) rather than at \$10.00 (the NBO shown in its direct feed). In this fashion, the trader placing the marketable buy order paid \$0.01 more per share than if the broker had priced the transaction using the direct data feed. Conversely, where a dark pool prices trades using the SIP NBBO, an institutional client who provides liquidity to the dark pool by submitting a pegged order to buy at the NBB might have their stale order "picked off" at \$10.00 by an HFT seller who sees the price decline to \$9.99 x \$10.00 through its direct feed. In this situation, our liquidity provider would have sold for \$0.01 less per share than if the venue had priced the transaction at \$9.99 using direct data.

In a similar fashion, SIP latencies can also allow trading venues to misrepresent the quality of their trade execution statistics. Pursuant to Rule 605 of Reg NMS, all trading centers must report how market orders executed by a trading center fared relative to the NBBO to aid investors in their routing decisions. Among other things, these disclosures include the average effective spread paid for market orders, which is defined as twice the difference between the trade price and the midpoint of the NBBO. In the first example above, the existence of two NBBOs would permit our hypothetical broker to report an effective spread of \$0.01 on a trade at \$10.01 (i.e., the effective spread using the SIP midpoint of 10.005) rather than the actual effective spread of \$0.03 (the effective spread using the direct feed's midpoint of 9.995). In May 2016, Reuters (Levinson, 2016) reported that the Justice Department is investigating the market-making firms Citadel and KCG, in part, because of concerns that each firm is using the slower SIP NBBO to

“claim it got the optimal deal for a client based on the prices on the slower data feed, even as the firm knew a better price existed on a faster feed.”

We use the microsecond timestamps to investigate empirically the extent to which traders during our sample period may have been adversely affected by SIP reporting latency, as well as its effect on the reliability of venues’ trade execution statistics.

a. Estimating Investor Trading Losses

1. Liquidity Takers. To estimate investor trading losses for liquidity takers arising from SIP latency arbitrage, we exploit the fact that our dataset includes both the SIP NBBO as well as the Direct NBBO prevailing for every trade in our sample. This basic structure permits us to estimate investor losses in a two-step process. In step one, we identify those trades that match the SIP NBBO by defining an indicator variable “SIP Priced” that equals 1 when the trade price matches either the NBB or NBO as reflected on the SIP NBBO, and equals 0 otherwise.⁵³ Trades that are “SIP priced” represent purchase and sale transactions that place the liquidity taker in the same position as if the venue priced the order using the SIP NBBO. Second, because trades priced at the SIP NBBO represent those trades that are at risk of NBBO arbitrage, we next compare how these SIP-priced trades would have been priced had they been priced at the Direct NBBO. We then measure whether a trade priced at the SIP NBBO rather than the Direct NBBO resulted in a loss or a profit for the trader placing the liquidity taking order.

In Table 4, we illustrate this two part process using 35 trades occurring in Apple, Inc. over a 15 millisecond period on November 13, 2015. The time set forth in the second column is the Participant Timestamp, which is the timestamp reported by the trading venue for when the trade occurred. We use the Participant Timestamp to place trades in chronological order. The Participant Timestamp gives us the ability to sort quotes and trades according to the moment they occurred, conferring knowledge of the actual quoting environment surrounding trades. For comparison, the third column presents the SIP Timestamp. Note that several pairs of trades, such as the fifth and sixth, are in chronological order according to the Participant Timestamp (by design) but not in chronological order according to the SIP Timestamp.

⁵³ As noted in Bartlett & McCrary (2016), trading venues also frequently use the NBBO to price trades at its midpoint. However, because we require trade direction to evaluate a trade’s profitability, we focus only on those trades priced at exactly the NBB or NBO which allows us to assign trading direction using the Lee-Ready (1991) algorithm.

The fourth and fifth columns represent the NBB and NBO in effect at the time of the trade as reflected in the Direct NBBO, while the next two columns reflect the NBB and NBO as reflected in the SIP NBBO. As shown in the table, the trades commenced when the Direct and SIP NBBOs match at \$113.37 x \$113.38. At that time, however, the market data suggest an inter-market sweep order (ISO) to buy approximately 6,000 shares with a limit price of \$113.39 was submitted to all exchanges sitting at the NBO (BATS, Direct Edge A, Nasdaq, and the NYSE Arca).⁵⁴ Evidence of this order can be seen by the manner in which the first 30 trades (each marked with code “F” in column 8 for an ISO) sweep through these four exchanges (column 9), buying all shares on the venues that are offered for less than \$113.40 (column 10).⁵⁵ This order results in the Direct NBBO changing to \$113.39 x \$113.40 by 11:37:47.465000, at which time an apparently unrelated trade occurs in a non-exchange venue. At the time of this latter trade, however, the SIP NBBO now reflects a stale NBBO of \$113.37 x \$113.38. Following this non-exchange trade, the SIP NBBO updates to reflect the true NBBO so that the Direct NBBO and SIP NBBO match one another by the time of the last three trades.

For purposes of analyzing this sequence of trades, we focus on those trades whose price matched the SIP NBBO, identified in the column entitled “SIP Priced.” Were these trades actually priced off the SIP, the SIP’s delay has an economic effect only for the non-exchange trade (Trade #31, highlighted in bold) occurring immediately after the ISO order finished sweeping through the market and inducing a mismatch between the Direct NBBO and the SIP NBBO. Based on the price of this trade, it appears to have been the result of marketable buy order; therefore, the fact that the trade was filled at \$113.38 (the stale NBO) rather than at \$113.40 (the new NBO) allowed the originator of the order to save two pennies per share acquired, or \$2.00 for the total order. The SIP NBBO and the Direct NBBO matched one another for all other trades, so choice of NBBO had no effect on trade profitability for these trades.

⁵⁴ ISO trades are those with an “F” in the trade condition code listed in column 8. An order marketed as an ISO is exempt from the Order Protection Rule of Reg. NMS, which prohibits a venue from filling an in-bound order if superior prices rest at other exchanges. As such, a trading venue receiving an inbound liquidity-taking ISO can fill it without checking other venues for better prices. However, the broker sending the ISO is responsible for sending simultaneous orders that sweep all venues with better prices. As such, ISO orders allow a trader to sweep through multiple levels of a venue’s order book, as occurs in this example.

⁵⁵ Column 8 reports the exchange code used by the SIPs. Codes are: Exchange Z=BATS; Exchange K=Direct Edge A; Exchange Q=Nasdaq; Exchange P=NYSE Arca).

In Table 5, we generalize this type of analysis to our full sample of Dow Jones 30 trades. Panel A summarizes by exchange the percentage of trades that we classify as SIP Priced, weighted by transaction size. As discussed previously, a possibility exists that the Nasdaq SIP printed timestamps on messages approximately 200 microseconds before it actually disseminated a message. Therefore, we present results after conducting the aforementioned analysis with no adjustment to the Nasdaq SIP's timestamp, as well as after adding 200 microseconds to the timestamp for all trade reports and quote updates disseminated by the Nasdaq SIP. In both cases, we find approximately 75% of all shares traded in our sample were traded at prices that exactly match the SIP NBBO. Excluding shares traded in non-exchange venues, this percentage increases to approximately 88%.

Panel B presents by exchange the mean amount of lost profit per share that liquidity takers experienced by having their trades priced at the SIP NBBO rather than at the Direct NBBO. For each exchange, means are size-weighted based on the number of shares traded. Overall, Panel B indicates that liquidity-taking trades priced at the SIP NBBO had average lost profits of approximately -\$0.0002 per share. As indicated in our Apple illustration, lost profits are defined as the difference between the Direct NBB and the SIP NBB for sell orders, and the difference between the Direct NBO and the SIP NBO for buy orders. As such, these negative lost profits suggest that liquidity takers, on average, *benefited* if their trades were priced at the SIP NBBO.

We note, first, that the magnitude of this effect is manifestly small. In terms of dollar value, for instance, the net aggregate dollar value of this benefit for all shares traded in our sample amounted to just \$11.1 million, notwithstanding the fact that the total trading value in our sample exceeded \$4 trillion.⁵⁶ Second, the sign of this effect is the opposite of what would be expected if liquidity takers are systematically receiving inferior pricing due to SIP reporting latencies. This analysis suggests that the widespread concerns about the risk to liquidity takers posed by latency arbitrage of SIP prices are exaggerated and perhaps even misplaced. Moreover, Panel B indicates that this result was generally persistent across trading venues, with the singular exception of trades made on the Chicago Stock Exchange. For trades occurring there, traders were, on average, effectively indifferent between having their trades priced at the SIP NBBO or the Direct NBBO.

⁵⁶ If we exclude offsetting positive values of lost profits, the aggregate dollar value of \$11.1 million rises to \$11.6 million.

To more fully understand this result, we present in Panel C the full distribution of lost profits per share. Given heightened concerns about stale quote arbitrage within dark pools that price off the SIP, we present separately the distribution in exchange and non-exchange venues.⁵⁷ As reflected in the distribution, a trade priced at the SIP NBBO rather than at the Direct NBBO had no economic effect for approximately 97% of shares traded within our sample. As was apparent in our Apple illustration, it is only when the SIP and Direct NBBOs differ that the choice of NBBO matching can affect transaction prices. Accordingly, the high percentage of shares traded with zero lost profits reflects the simple fact that the SIP and Direct NBBO typically match.

For those trades where the use of the SIP NBBO rather than the Direct NBBO produced non-zero lost-profits per share, our unadjusted results show that nearly 90% of the trades (weighted by shares traded) produced better pricing for liquidity takers when the trade's price matched the SIP NBBO rather than the Direct NBBO. Specifically, among trades priced at the SIP NBBO, approximately 2.7% of shares traded on non-exchange venues and 2.2% of shares traded on exchanges had negative lost profits, which were largely unchanged when we used the modified Nasdaq timestamps. Moreover, almost all of these instances cluster at -\$0.01 lost profits per share. We attribute this distribution to the fact that the NBBO will commonly change in response to serial buy (sell) orders so that late-arriving buy (sell) orders benefit from stale SIP quotes. For instance, in the Apple illustration above, the delay in updating the SIP NBBO to reflect the ISO buy order that induced a change in the NBBO allowed the later-arriving non-exchange buy order to benefit by purchasing at the lower, stale NBO.

Reflecting this logic, Panel C highlights the remarkably low likelihood that a marketable order priced at the SIP received poorer pricing than it would have, had it been priced at the Direct NBBO. For non-exchange trades, our unadjusted and adjusted results indicate that, among trades priced at the SIP NBBO, just 0.2% of shares traded had a positive measure of lost profits. This estimate drops to 0.06% for exchange trades. The higher overall incidence of these trades among non-exchange venues relative to exchanges, however, does lend some support to concerns that the likelihood of this form of latency arbitrage—while low overall—might be somewhat higher for trades executed in non-exchange venues.

⁵⁷ For ease of presentation, we exclude from the distribution trades having a lost profit per share of more than \$0.10 and less than -\$0.10, which in the aggregate comprise less than 0.0004% of shares traded in our sample.

2. *Liquidity Providers.* Of course, since there are two sides to every trade, while the foregoing results suggest liquidity takers generally benefit when trades are priced at the SIP NBBO, the reverse conclusion applies to liquidity providers. In our Apple illustration, for example, the buy order completed at the stale SIP NBO of \$113.38 rather than at the new NBO of \$113.40 meant the seller in the non-exchange venue who had posted the resting liquidity lost \$0.02 per share by selling at the stale SIP NBO rather than at the Direct NBO. The mean measure of lost profits of approximately -\$0.0002 per share accordingly highlights that to the extent SIP pricing adversely affects traders, these costs are more likely to be borne by liquidity providers than by liquidity takers.

Depending on the identity of the trader taking liquidity in these trades, this latter finding may point to the presence of an alternative form of SIP latency arbitrage occurring in the market. In particular, our results have largely assumed marketable orders reflect uninformed order flow, such as orders submitted by retail investors. Our basic finding that liquidity takers benefit from being priced at the SIP NBBO, however, is in principle also consistent with HFT firms using marketable orders to exploit dislocations between the SIP- and Direct-NBBO.

The sequence of Apple trades in Table 4 provides an example of how such a strategy might work in practice. After having secured a “buy” trade at \$113.38 (the stale NBO) rather than at \$113.40 (the new NBO), the trader submitting the buy order need only sell at the new NBB of \$113.39 to realize an immediate profit of \$0.01 per share (excluding trading fees). Because the ensuing four trades each reflected “sell” transactions at this price, our Apple example—and the results in Table 5 more generally—may simply reflect the strategic use of marketable orders by HFT firms to “pick off” stale limit orders posted in venues that use the SIP NBBO as the benchmark for pricing orders that are pegged to the NBBO.

To explore this possibility, we leverage the new participant timestamp data and the fact that an HFT firm following such a strategy would need to make a pair of trades. To see how this works, consider trades immediately subsequent to those trades where trading at the SIP NBBO yielded more favorable pricing for the liquidity-taking order than trading at the Direct NBBO—that is, where the trade produced a negative measure of lost profits. For each of these potential first-leg trades, suppose the trade originated from an HFT firm submitting to a venue an immediate-or-cancel buy or sell order after having observed a momentary dislocation between

the SIP-NBBO and the Direct-NBBO. The success of this HFT strategy requires an off-setting second-leg trade—which one should be able to see in the data.

To execute this pairing strategy, we sort trades based on the Participant Timestamp and identify each potential first-leg trade based on whether it produced a negative measure of lost profits. We then search forward for a matching second-leg trade until a window of 1,000 microseconds from the first-leg trade timestamp has been exhausted. For a trade to match the first-leg trade, it must match both on the direction of the trade and the trade price. In particular, for first-leg buy orders, we require a matching second-leg trade to be a sell order at a price that is higher than the first-leg purchase price; conversely, for first-leg sell orders, we require a second-leg buy order at a price that is less than the first-leg sales price. We impose a 1,000 microsecond trading window following each first-leg trade to ensure there is sufficient time for a trader to receive a trade confirmation on the first-leg trade before placing the second-leg trade at either an exchange or non-exchange venue.⁵⁸

Before presenting our results, it is worth emphasizing that this simple empirical strategy almost certainly over-estimates—potentially by a wide margin—the actual incidence of second-leg matches. Among other things, for instance, our strategy disregards order size and transaction fees and focuses purely on identifying subsequent transactions that are priced higher (lower) than first-leg buy (sell) orders. Moreover, our approach seeks to identify matching second-leg trades independently for each first-leg trade, creating the possibility that the same second-leg trade can be matched to two different first-leg trades. Finally, our strategy also permits second-leg trades to occur in non-exchange venues, even though the absence of displayed liquidity in these venues makes such an approach for executing second-leg trades extraordinarily risky from an ex ante

⁵⁸ We suspect a 1,000 microsecond trading window is most likely too generous for first-leg transactions occurring on stock exchanges. For instance, a trader subscribing to exchanges' fastest fiber optic data feeds and co-located at Nasdaq would receive a trade confirmation of a first-leg trade occurring at the NYSE (the furthest exchange from Nasdaq) in approximately 200 microseconds based on Table 3, allowing it to execute a second-leg trade even at the NYSE in approximately 400 microseconds from the time of the first-leg trade. Trade confirmations for transactions occurring at a BATS exchange or on Nasdaq would require even less transit time for such a trader. Our choice of a 1,000 microsecond trading window is driven instead by the possibility that first-leg transactions occur on non-exchange venues. Given the execution risk assumed by a trader executing a first-leg trade, we assume an HFT firm choosing to use a non-exchange venue for the first-leg of this strategy would focus on those automated venues based in Figure 1's "Equity Triangle" that are capable of providing a trade confirmation with latencies comparable to those of the primary exchanges. However, as noted previously, participant timestamps are recorded in milliseconds (rather than microseconds); therefore, imposing a 1,000 microsecond trading window for these trades has the effect of imposing a maximum window of between 501 microseconds (for a non-exchange trade that actually occurs at the 499th microsecond of a second) and 1,500 microseconds (for a non-exchange trade that actually occurs at the 500th microsecond of a second).

perspective.⁵⁹ This analysis is predicated on measuring the scope for fast traders to engage in *risk-free* profitable liquidity-taking, a type of trading strategy commonly ascribed to HFT firms in contemporary debates.

Even with this bias in favor of finding second-leg trades, the results of this analysis, which we present in Table 6, reveal an extremely low incidence of matches between first- and second-legs of this type of trading pairs. Consequently, the proper interpretation of the entries in Table 6 is an estimate of an upper bound on the prevalence of this form of latency arbitrage within our sample.

Table 6 stratifies the analysis between exchange and non-exchange venues because non-exchange venues may be more likely to price orders using the SIP NBBO. The results show that only 1.4% of all first-leg trades occurring in non-exchange venues can be matched to a second-leg trade within 1,000 microseconds. For first-leg trades occurring on exchanges, this percentage falls to less than 1%. Given the strong bias our empirical strategy creates in favor of finding a second-leg matching trade, we interpret these results as confirming that the pricing advantage for liquidity-takers of having orders priced at the SIP NBBO are unlikely to be the result of HFT firms seeking to exploit liquidity providers in venues that price transactions using the SIP NBBO. Our estimated upper bounds demonstrate that although anecdotal evidence may establish that these trading strategies exist, they are unlikely to be allocatively important in recent years for the Dow Jones 30. After all, if all lost profits amount to only \$11 million, and roughly 0.8% of those trades might be part of an arbitrage play, arbitrage profits for the Dow Jones 30 are at most \$110,000 over our sample period.

A natural question is whether we can approximate how much money might be at stake for the entire market over the course of a year, as opposed to just our sample period for just our sample stocks. Our sample period covers the period August 6, 2015 to June 30, 2016, or 228 trading days, so the first order of business is to multiply the upper bound on arbitrage profits by 253/228. Doing so yields an upper bound on annual arbitrage profits for the Dow Jones 30 of \$122,061. A second consideration is that the Dow Jones 30 obviously cover only 30 stocks. To expand our computation to the entire market is a computational challenge, because there are nearly 8,000 stocks observed over this time period. To spare computation, we first examined how many

⁵⁹ Even for second-leg trades aimed at hitting an exchange's displayed liquidity, this strategy would appear to involve non-trivial execution risk. As emphasized in Fox, Glosten & Rauterberg (2015), an HFT firm attempting to profit from "slow-market" arbitrage "must be able to transact against the new best quote before anyone else can."

stocks corresponded to what fraction of the overall shares traded during our study period.⁶⁰ The top 257 stocks correspond to half the trading volume during our sample period, and the top 887 stocks correspond to three-quarters of the trading volume during our sample period.

Consequently, we examine how our \$11 million figure might increase, were we to focus on the top 257 and 887 stocks instead, and then multiply those figures by 2 and 4/3, respectively, to obtain an estimate for the overall market.

For the top 257 stocks, the total amount of lost profits is roughly \$83 million. Multiplying by $2 \times 253/228 \times 0.008$ yields our best estimate of the upper bound on annual market-wide arbitrage profits of \$1.5 million. For the top 887 stocks, our best estimate of the upper bound on annual market-wide arbitrage profits, following the same methodology outlined above, is about \$3 million. Neither of these is large enough to sustain an industry with many competing firms and large annual operating costs.

b. Trade Execution Costs.

As noted previously, a secondary concern with the availability of direct data feeds relates to the possibility that market centers might misreport their trade execution statistics using the SIP NBBO. In general, these concerns are typically coupled with concerns that retail investors are receiving inferior pricing at the SIP NBBO, as indicated in the example provided at the beginning of Section 4. To the extent retail investors send orders to venues based on these statistics, such misreporting might therefore compound the risk that liquidity-taking orders will be harmed by receiving SIP-priced trades.

At their most general level, these claims find little support in our finding that marketable orders, on average, benefit from pricing at the SIP NBBO. Yet even if a marketable order benefits from SIP pricing, any divergence between the SIP and Direct NBBOs nevertheless creates the possibility for conflicting trade execution measures. In the prior example using trades in Apple, for instance, the fact that a dark venue priced a buy order at the stale SIP NBO of \$113.38 rather than the current NBO of \$113.40 created two possible measures of price improvement. Using the SIP NBBO as the benchmark, the trade received zero price

⁶⁰ We ignore stocks with any suffixes on their trading symbol. Stocks with no suffixes correspond to 98.5% of trading volume over our study period, so this is of little consequence for the calculations we report here. Ignoring suffixes is computationally advantageous because of the database index structure.

improvement—it was priced exactly at the SIP NBO of \$113.38. However, using the Direct NBBO of 113.4, the trade would have received 2 cents of price improvement.

The challenge of dueling trade execution statistics is even more extreme for effective spreads, which all trading centers must disclose in their Rule 605 reports and which are routinely used as “the industry’s acid-test quality measure” to rank trading venues (Alpert, 2015).

Returning to the Apple example of Table 4, a venue that benchmarked trade execution to the SIP NBBO would record an effective spread of 0.01 for Trade #31; however, using the Direct NBBO, effective spreads for that trade would be 0.03.⁶¹ Thus, even though the buyer paid two cents less than the Direct NBO of \$113.40, the effective spread of 0.03 would suggest the trader received an *inferior* trade execution than if she had simply transacted at the Direct NBO.

This counterintuitive result stems from the basic arithmetic for calculating effective spreads, which seeks to infer price improvement based on the difference between a trade’s price and the midpoint of the benchmark NBBO. This emphasis on a trade’s distance from the midpoint of the benchmark NBBO can cause effective spreads to increase when a venue calculates them using an NBBO other than the one used to price trades. This is especially true when an exchange handles orders that are to be priced by reference to the NBBO, such as orders pegged to the near, far, or midpoint of the NBBO.

To illustrate, consider a situation where the Direct NBBO is \$10.00 x \$10.01, but the SIP NBBO is \$10.01 x \$10.02. A venue that priced off the Direct NBBO and filled pegged orders at \$10.00 (the NBB), \$10.01 (the NBO), and \$10.005 (the midpoint) would record effective spreads on these trades of 0.01, 0.01, and 0, respectively, if it used the Direct NBBO as its benchmark. If it used the SIP NBBO, these measures would be 0.03, 0.01, and 0.02. Conversely, a venue that priced pegged orders off the SIP NBBO and filled orders at \$10.01 (the SIP NBB), \$10.02 (the SIP NBO), and \$10.015 (the SIP midpoint) would record effective spreads for these trades of 0.01, 0.01 and 0, respectively if it used the SIP NBBO as its benchmark and 0.03, 0.01 and 0.02 if it used the Direct NBBO instead. Situations could also arise in which effective spreads improve when a venue used a different NBBO benchmark than the one used to price trades,

⁶¹ As noted previously, effective spreads are calculated as twice the difference between the trade price and the midpoint of the benchmark NBBO.

which simply underscores the potential for divergent NBBOs to create conflicting measures of trade execution for the same trade.⁶²

At the same time, however, the extent to which rival NBBO benchmarks actually affect a trading venues' aggregate trade performance disclosures should be mitigated by the fact that dislocations between the SIP-generated NBBO and the NBBO generated by direct feeds are—as discussed above—typically infrequent and, when they do occur, short-lived. As noted previously, for instance, pricing a trade at the SIP NBBO rather than at the Direct NBBO had no economic effect for approximately 97% of the trades within our sample, while the mean (median) duration of NBBO dislocations was 1,001.6 (489) microseconds.

To estimate empirically how much the choice of NBBO benchmark affects venues' trade execution statistics, we calculate effective spreads for each trade in our sample using as our benchmark both the SIP NBBO and the Direct NBBO. Specifically, for each trade, we first calculate effective spreads using the prevailing SIP NBBO for the trade as our NBBO benchmark followed by using the prevailing Direct NBBO as our benchmark. Since we are interested in understanding the pricing of marketable orders at the NBBO, we exclude ISOs given that ISOs can be filled at prices worse than the NBBO. In all cases, we calculate effective spreads as a percentage of the quoted NBBO spread—generally known as the effective/quoted spread ratio (E/Q)—to account for variation in the size of the quoted spread for our sample securities.⁶³

Table 7 presents the results of this examination. In the first three rows, we present separately the analysis for the NYSE MKT, the Chicago Stock Exchange (CHX), and the NSX. We distinguish these three exchanges for institutional reasons: each disclosed using the SIP NBBO

⁶² These latter situations can occur, for example, if a venue attempts to fill a trade at the NBB or NBO of its benchmark NBBO, which happens to be the midpoint of the alternative NBBO. For instance, if the SIP NBBO stands at \$10.00 x \$10.01 and the Direct NBBO stands at \$10.00 x \$10.02, a venue that tries to fill a “buy” order at the SIP NBO of \$10.01 would record an effective spread of 0.005 on the trade using the SIP NBBO as its benchmark for calculating effective spreads. However, using the Direct NBBO as the benchmark for calculating effective spreads would yield an effective spread for a trade at \$10.01 of 0 since it happened to occur at the midpoint of the Direct NBBO.

⁶³ We base our calculation of the E/Q ratio on the methodology described by BATS Global Markets. See Execution Quality Definitions, available at https://batstrading.com/market_data/execution_quality/definitions/. In summary, this method restricts attention to trades that (a) occur when markets are neither locked nor crossed, and (b) that are within 10% of the NBBO. Because these conditions would imply analyzing a slightly different subsample of trades when the benchmark NBBO changes, we first construct a sample of trades meeting the above criteria using the SIP NBBO, second construct an analogous sample using the Direct NBBO, and finally use as our analysis sample the set of trades that are in both the first and second set.

to price all un-priced, pegged orders during our sample period.⁶⁴ As discussed previously, this institutional choice can often favor the use of the SIP NBBO as the relevant E/Q benchmark to the extent these venues process a material volume of these un-priced, pegged orders.

Consistent with this prediction, using the SIP NBBO to calculate the E/Q ratio produces a more favorable trade execution measure for the NYSE MKT. Specifically, the E/Q ratio for non-ISO trades on the NYSE MKT was approximately 84.82% when calculated using the SIP NBBO as the benchmark, and 85.06% when calculated with the Direct NBBO. Results for the Chicago Stock Exchange and the National Stock Exchange, in contrast, were inconsistent with this prediction, most likely reflecting the lower volume of un-priced, pegged orders processed on these exchanges. For instance, in unreported results, we find that trades priced at the midpoint of the SIP NBBO constitute 6.5% of non-ISO trades on the NYSE MKT, but only 0.66% on the Chicago Stock Exchange and 0% on the National Stock Exchange. Because these trades reflect the filling of midpoint peg orders (Bartlett & McCrary, 2016), this evidence would suggest these latter two venues process a lower volume of trades pegged to the NBBO.⁶⁵

The subsequent nine rows present results for the remaining exchanges, which all disclose using direct feeds to price trades, as opposed to the SIP NBBO.⁶⁶ For these venues, Table 7 indicates that using the SIP NBBO as the benchmark generally results in a worse E/Q ratio, while using the Direct NBBO produces a more favorable measure of trade execution costs.⁶⁷ For all exchanges showing a statistically significant difference in E/Q ratios, however, the difference between using the SIP NBBO and the Direct NBBO as a performance benchmark changes the measure by a relatively small amount. The effect ranges from a low of 0.01 percentage points for the Nasdaq PSX to a high of 1.85 percentage points for BATS X. These figures highlight the

⁶⁴ All U.S. stock exchanges have disclosed the market data sources used to price and route trades since 2015. These disclosures were prompted in a June 5, 2014 speech by SEC Chair Mary Jo White, where she requested equity exchanges to file with the Commission the data feeds used for purposes of order handling, order execution, and order routing.

⁶⁵ The trading rules for the CHX and the NXS also suggest these venues do not ordinarily rely on the NBBO to price trades. For instance, while the CHX permits “midpoint cross” orders, it does not support other pegged order types. See Chicago Stock Exchange, CHX Order Types Primer, available at <http://www.chx.com/trading-information/order-types/>. The NSX supports orders that are pegged to the near, far, and midpoint of the NBBO, however, all such orders are non-displayed. See National Stock Exchange, Select NSX Order Types and Modifiers, available at http://www.nsx.com/images/documents/publications/NSX_Order_Types_v3_0_1.pdf.

⁶⁶ We include the NYSE within this group, notwithstanding the fact that its SEC filings indicate that it uses the SIP to obtain top-of-the-book quote updates from other exchanges when pricing pegged orders. Given that the NYSE trades in only NYSE-listed securities, the fact that it also uses order data obtained directly from its own matching engine has the practical effect of giving it a direct feed to a critical source of quote updates for Tape A securities.

⁶⁷ The single exception is for Nasdaq PSX, which accounts for less than 1% of all trades.

fact that, while the choice of NBBO benchmark affects effective spread calculations, the degree to which it does so is likely to be small in magnitude.

While we lack data on how individual non-exchange venues calculate effective spreads, the final row in Table 7 provides an analysis analogous to that above for all non-ISO FINRA trades within our sample. Calculating the E/Q ratio using the SIP NBBO produces a ratio of approximately 73.71%—modestly lower than the 74.09% obtained using the Direct NBBO as a benchmark. Given that these venues are likely to price a large number of orders by reference to the NBBO (Bartlett & McCrary, 2016), this finding is consistent with claims that a substantial portion of non-exchange venues continue to price trades using the SIP NBBO. At the same time, the extraordinarily small difference between the two calculations further underscores the conclusion that the short-lived nature of dislocations between the SIP and Direct NBBOs greatly diminishes the potential for a venue’s choice of NBBO to have a meaningful effect on its published effective spreads.⁶⁸

6. Conclusion

In his 2014 book *Flash Boys*, Michael Lewis captured international attention through his depiction of an equity market that systematically favors high frequency traders over slower traders such as retail and institutional investors. Central to his critique was the sale to HFT firms of fast access to exchange quotation data, which enables them to predict changes in the SIP-generated NBBO that trading venues have historically used to price both market and limit orders. For retail market-making firms such as Citadel and KCG, this speed advantage means the possibility of filling in-bound market orders at NBBO prices they know to be stale. For other

⁶⁸ That NBBO dislocations can matter at all, however, nevertheless underscores the limitations of the prevailing system governing order execution disclosures. Initially implemented as Rule 11Ac1-5 in 2001, Rule 605 makes no mention of *which* NBBO to utilize as a performance benchmark when calculating order execution statistics; however, subsequent SEC guidance suggests market centers should utilize data from the SIP when complying with the rule. Given the large number of venues using direct feeds to price transactions, we believe any such endorsement of the SIP NBBO in Rule 605 reporting has the potential to bias trade performance metrics, as shown in Table 7. At the same time, permitting venues to choose their NBBO benchmark (as appears to be tolerated by the SEC) complicates interpretation of a venue’s order execution information without disclosure of this information. For instance, certain venues have expressly declined to follow the SEC’s guidance to use the SIP NBBO in calculating their Rule 605 reports, opting instead to calculate trade performance statistics using the same market data used to price transactions. See, e.g., IEX ATS Rule 605 Disclosure of Order Execution Information, available at <http://50.116.60.129/regulation/605/>. Requiring venues to disclose the NBBO benchmark used for calculating their performance metrics would represent a logical modification of Rule 605 given the divergent use of market data among trading centers.

HFT firms, it means the possibility of picking-off mispriced limit orders pegged to a SIP NBBO that has yet to reflect the prices these fast traders can foresee.

Using recently released data from the two SIPs, we present novel evidence regarding the merits of these claims in the current trading environment. Due in large part to the political fallout from Lewis' narrative, these data now include the precise time at which a quote update or trade report was processed by the relevant SIP along with the time it actually occurred on a trading venue. As we show, the availability of this latter timestamp is especially important as it enables for the first time the reconstruction of the real-time sequencing of quote updates and trades across the entire market and, critically, how they relate to one another and to the SIP NBBO.

Exploiting these new data, we show that since the release of these timestamps in August 2015, liquidity-taking orders gain on average \$0.0002 per share when priced at the SIP-reported NBBO rather than the NBBO calculated using exchanges' direct data feeds. In all likelihood, we suspect this finding reflects the simple fact that dislocations between the SIP NBBO and Direct NBBO can occur in response to serial buy and sell orders, allowing late-arriving market orders to benefit if they are priced at an NBBO that has yet to reflect the new trading interest. To the extent this is the case, concerns about trading at the slower SIP NBBO would accordingly seem more relevant for traders providing liquidity in venues that price limit orders pegged to the NBBO using the slower SIP data. Yet while these concerns are consistent with claims that HFT firms pick-off mispriced limit orders in these venues, we find virtually no evidence of this strategic behavior using the new Participant Timestamp data.

In short, our findings reveal that pricing at the SIP-NBBO can benefit liquidity takers to the detriment of liquidity providers. However, the incidence of these gains and losses between these two forms of trading interest appears to be primarily a product of chance rather than of HFT design. Because our data commence in August 2015, we emphasize that these findings may very well reflect a new market environment in which the HFT strategies depicted in *Flash Boys* are less prevalent than in the past. Among other things, for instance, the increasing processing speed of the SIPs shown in Table 1, enhanced regulatory scrutiny of HFT, and the emergence of venues such as IEX that shield traders from HFT trading may have simply made these SIP-oriented arbitrage strategies increasingly infeasible.

Finally, while our findings are consistent with the incentive of liquidity providers to invest in fast access to trading data to avoid trading at stale NBBO prices, our results suggest these incentives play, at most, a subsidiary role in promoting the socially costly arms-race for trading speed described in Budish, Cramton & Shim (2015). Although our sample includes over \$4 trillion of trades, liquidity providers trading at the SIP NBBO could have saved just \$11 million in lost profits had they transacted at the Direct NBBO instead. To the extent traders participate in this arms race, the primary incentives today would accordingly appear to rest outside a desire to avoid the costs of trading at stale SIP prices.

References

- Alpert, Bill, Who Makes Money on Your Stocks, *Barrons*, Feb. 27, 2015.
- Angel, James, 2014, When Finance Meets Physics: The Impact of the Speed of Light on Financial Markets and their Regulation, *Financial Review* 49, 271-281.
- Angrist, Joshua D. and Jörn-Steffen Pischke, 2009, *Mostly Harmless Econometrics: An Empiricist's Companion*, Princeton University Press.
- Barreca, Alan I., Jason M. Lindo, and Glen R. Waddell, 2011, Heaping-Induced Bias in Regression-Discontinuity Designs, NBER Working Paper No. 17408.
- Bartlett, Robert, and Justin McCrary, 2016, Dark Trading at the Midpoint: Pricing Rules, SEC Enforcement Policy, and Latency Arbitrage, UC Berkeley School of Law Working Paper.
- Budish, Eric B., Peter Cramton, and John J. Shim, 2015, The High-Frequency Trading Arms Race: Frequent Batch Auctions as a Market Design Response, *Quarterly Journal of Economics* 130, 1547-1621.
- Ding, Shengwei, John Hanna, and Terrence Hendershott, 2014, How Slow is the NBBO? A Comparison with Direct Exchange Feeds, *Financial Review* 49, 313–332.
- Fox, Merritt, Lawrence R. Glosten, and Gabriel Rauter, 2015, The New Stock Market: Sense and Nonsense, *Duke Law Journal* 65, 191-277.
- Lee, Charles M.C. and Mark J. Ready, 1991, Inferring Trade Direction from Intraday Data, *Journal of Finance*, 46, 733-746.
- Lewis, Michael, 2014, *Flash Boys: A Wall Street Revolt*, W.W. Norton & Co.
- Levinson, Charles, U.S. investigates market-making operations of Citadel, KCG, Reuters, May 10, 2016.
- Lombardi Michael, 2000, Computer time synchronization, National Institute of Standards and Technology working paper, available at: <http://www.tf.nist.gov/service/pdf/computertime.pdf>
- Nanex, The Quote Stuffing Trading Strategy, Aug. 15, 2014, available at <http://www.nanex.net/aqck2/4670.html>
- O'Hara, Maureen and Mao Ye, 2011, Is market fragmentation harming market quality? *Journal of Financial Economics* 100, 459-474.
- Securities and Exchange Commission, 2016, Joint Industry Plan; Notice of Filing of the National Market System Plan Governing the Consolidated Audit Trail, available at <https://www.sec.gov/rules/sro/nms/2016/34-77724.pdf>.

Silverman, Bernard, 1986, *Density Estimation for Statistics and Data Analysis*. Boca Raton: Chapman & Hall/CRC Press.

Tabb, Larry, Latency Arbitrage and the Problem With the SIP, *Tabb Forum*, July 19, 2016, available at <http://tabbforum.com/opinions/latency-arbitrage-and-the-problem-with-the-sip>. Accessed July 30, 2016.

Wah, Elaine, and Michael Wellman, 2013, *Latency arbitrage, market fragmentation, and efficiency: A two-market model*, Proceedings of the fourteenth ACM Conference.

White, Mary Jo, 2015, Testimony on “Examining the SEC’s Agenda, Operations and FY 2016 Budget Request”, Before the United States House of Representatives Committee on Financial Services, available at <https://www.sec.gov/news/testimony/2015-ts032415mjw.html>.

Ye, Mao, Chen Yao, and Jiading Gai, 2012, The externalities of high frequency trading, Working paper, University of Illinois at Urbana-Champaign

Figure 1

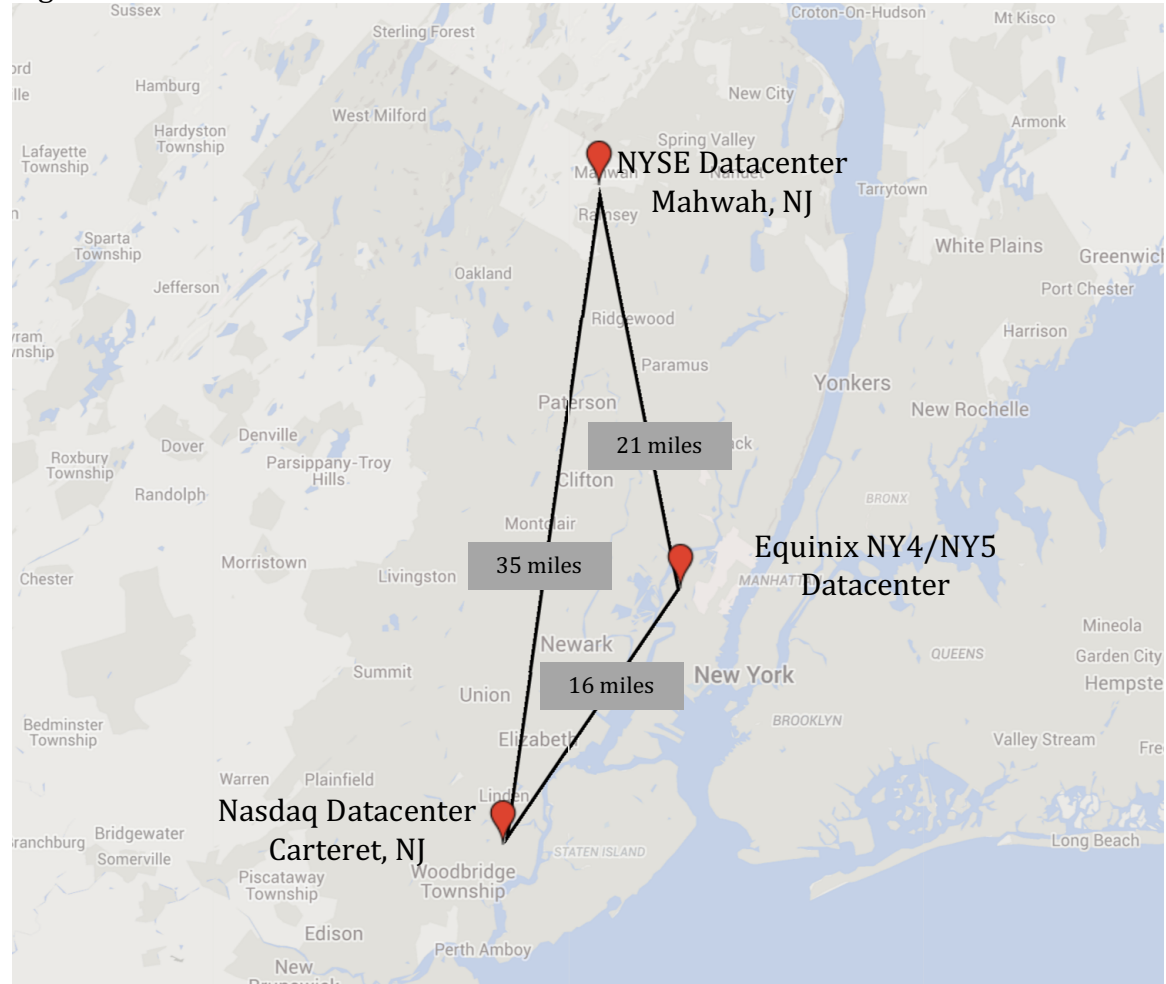


Figure 2a

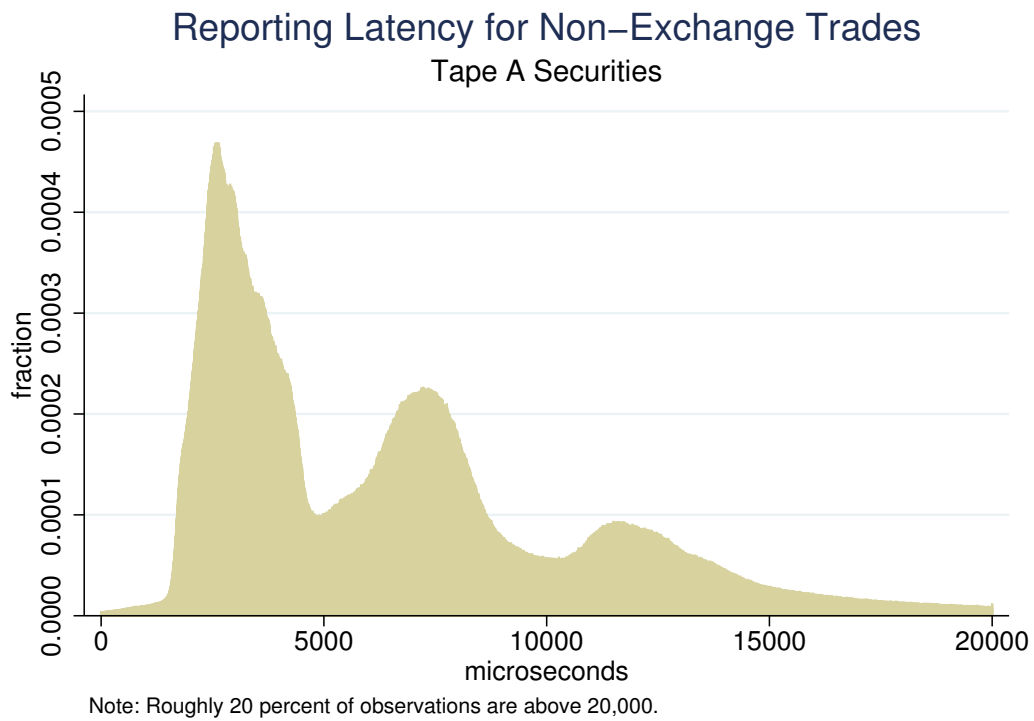


Figure 2b

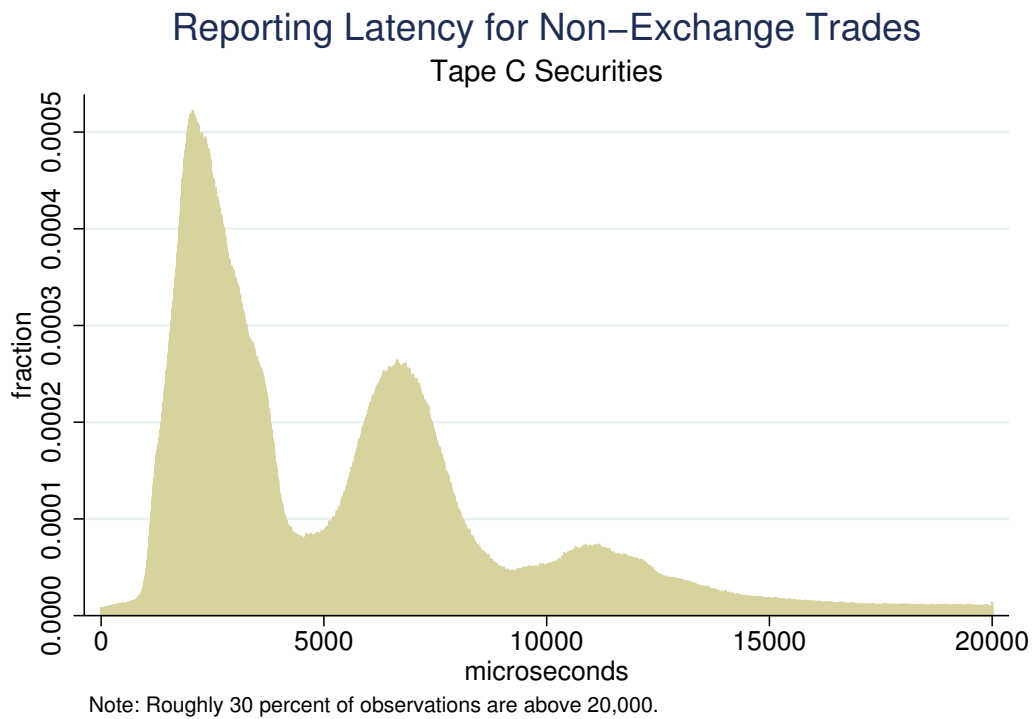


Figure 3a

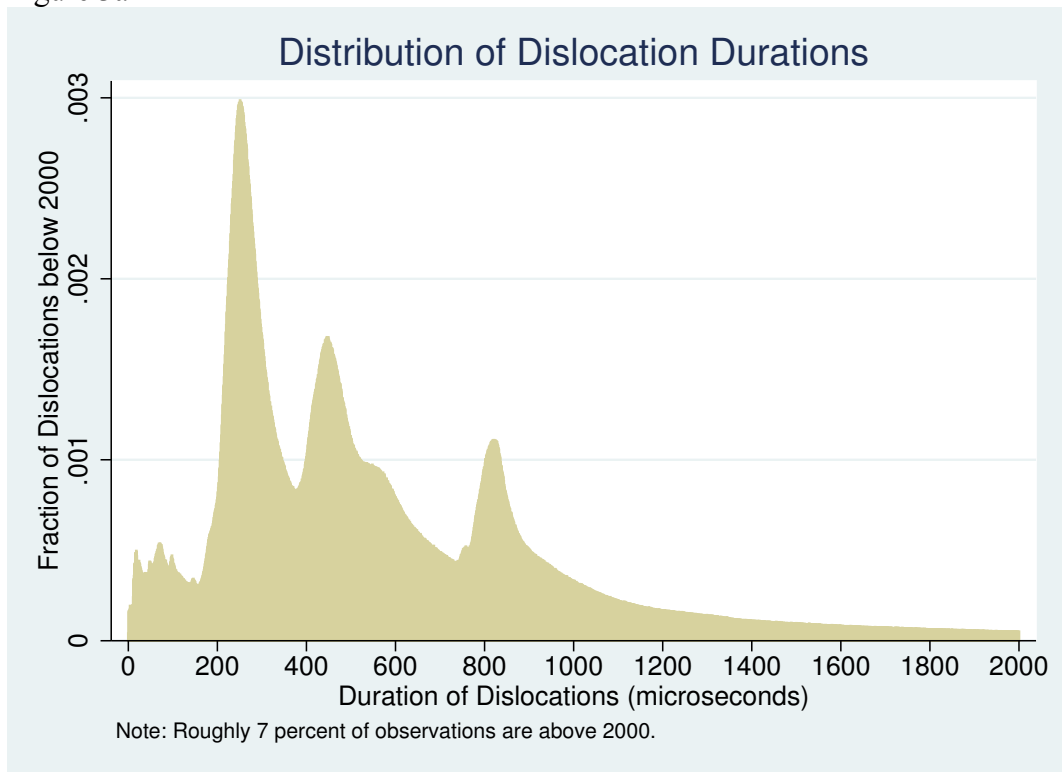


Figure 3b

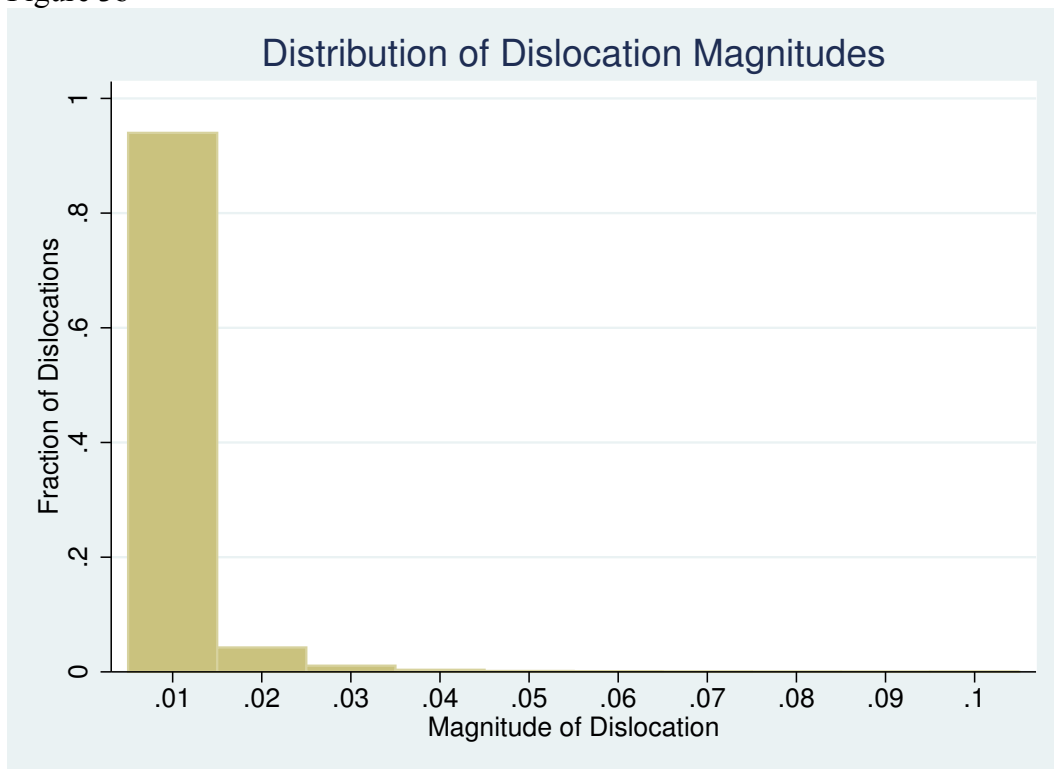


Table 1: SIP Processing Times

Panel A: SIP Processing Time for Trades

	Tape A&B Trade metrics						Tape C Trade metrics					
	Peak Messages per 100 Milliseconds (thousands)	Capacity Messages per 100 Milliseconds (thousands)	Capacity vs Peak Ratio	Average Latency	Median Latency	90th percentile latency	Peak Messages per 100 Milliseconds (thousands)	Capacity Messages per 100 Milliseconds (thousands)	Capacity vs Peak Ratio	Average Latency	Median Latency	90th percentile latency
1q14	21.80	60.00	2.75	0.51	n/a	0.71	19.30	39.40	2.04	1.32	1.25	1.67
2q14	23.50	60.00	2.55	0.51	n/a	0.66	20.50	39.40	1.92	0.82	0.54	0.74
3q14	22.70	65.00	2.86	0.51	n/a	0.66	17.60	48.50	2.76	0.59	0.49	0.68
4q14	24.20	65.00	2.69	0.45	n/a	0.60	19.40	48.50	2.50	0.59	0.49	0.67
1q15	22.10	70.00	3.17	0.45	n/a	0.59	20.10	68.70	3.42	0.53	0.45	0.60
2q15	31.80	70.00	2.20	0.34	n/a	0.43	22.80	132.80	5.82	0.54	0.46	0.62
3q15	27.10	75.00	2.77	0.32	0.24	0.41	16.10	132.80	8.25	0.58	0.47	0.64
4q15	43.70	75.00	1.72	0.31	0.24	0.41	18.60	132.80	7.14	0.62	0.47	0.66
1q16	42.40	86.00	2.03	0.33	0.25	0.43	19.40	132.80	6.85	0.77	0.49	0.76
2q16	37.40	96.00	2.57	0.34	0.24	0.45	28.20	132.80	4.71	0.63	0.48	0.68
mean	29.67	72.20	2.53	0.41	0.24	0.54	20.20	90.85	4.54	0.70	0.56	0.77

Panel B: SIP Processing Time for Quotes

	Tape A&B Trade metrics						Tape C Trade metrics					
	Peak Messages per 100 Milliseconds (thousands)	Capacity Messages per 100 Milliseconds (thousands)	Capacity vs Peak Ratio	Average Latency	Median Latency	90th percentile latency	Peak Messages per 100 Milliseconds (thousands)	Capacity Messages per 100 Milliseconds (thousands)	Capacity vs Peak Ratio	Average Latency	Median Latency	90th percentile latency
1q14	121.10	300.00	2.48	0.45	n/a	0.90	51.50	70.70	1.37	1.20	1.08	1.62
2q14	131.70	300.00	2.28	0.44	n/a	0.76	51.20	70.70	1.38	0.69	0.48	0.70
3q14	121.10	325.00	2.68	0.45	n/a	0.88	49.80	83.80	1.68	0.59	0.43	0.79
4q14	141.80	325.00	2.29	0.41	n/a	0.75	95.40	83.80	0.88	0.55	0.43	0.66
1q15	146.40	350.00	2.39	0.39	n/a	0.68	85.50	166.90	1.95	0.50	0.44	0.62
2q15	142.60	350.00	2.45	0.46	n/a	1.02	48.00	215.00	4.48	0.65	0.44	0.69
3q15	158.40	375.00	2.37	0.51	0.23	1.13	37.10	215.00	5.80	0.80	0.45	0.79
4q15	162.30	375.00	2.31	0.44	0.21	0.93	41.00	215.00	5.24	0.81	0.45	0.81
1q16	163.30	392.00	2.40	0.49	0.22	1.08	60.10	215.00	3.58	0.92	0.47	1.04
2q16	168.40	400.00	2.38	0.49	0.22	1.09	83.00	215.00	2.59	0.80	0.46	0.93
mean	145.71	349.20	2.40	0.45	0.22	0.92	60.26	155.09	2.90	0.75	0.51	0.87

Source: CTA and UTP disclosures.

Table 2: Quote and Trade Latencies

Panel A: Quote Updates										
<i>Venue</i>	Tape A Securities					Tape C Securities				
	<i>N</i>	<i>mean</i>	<i>sd</i>	<i>median</i>	<i>90p</i>	<i>N</i>	<i>mean</i>	<i>sd</i>	<i>median</i>	<i>90p</i>
NYSE	1,362,744,432	690	1,579	301	1,297	-	-	-	-	-
NYSE MKT	-	-	-	-	-	27,706,440	1,304	2,554	937	1,476
NYSE Arca	403,251,799	783	8,394	302	1,513	294,642,117	1,547	3,207	977	2,245
Nasdaq OMX BX	221,620,468	1,799	10,441	877	2,745	62,733,627	762	2,926	325	1,018
NASDAQ OMX PSX	258,795,046	1,972	12,730	923	3,297	87,393,187	886	3,335	367	1,246
NASDAQ	793,107,717	1,587	10,066	933	2,551	419,195,751	1,194	10,353	404	2,017
BATS	590,111,028	1,255	2,679	507	2,630	242,481,473	999	3,305	523	1,251
BATS Y	355,567,830	916	2,029	486	1,609	100,514,420	974	3,354	510	1,202
Direct Edge A	223,325,479	829	1,776	491	1,406	86,102,843	1,065	3,620	529	1,384
Direct Edge X	442,063,443	1,147	2,627	517	2,238	239,827,518	1,017	3,324	526	1,274
Chicago Stock Exchange	827,450	1,019	5,418	839	1,120	209,724	849	2,130	722	994
National Stock Exchange	529,478	18,073	3,657,389	1,228	2,073	106,167	41,176	6,307,078	962	1,992
All:	4,651,944,170	1,116	39,536	566	2,015	1,560,913,267	1,152	52,368	551	1,697

Panel B: Trades

<i>Venue</i>	Tape A Securities					Tape C Securities				
	<i>N</i>	<i>mean</i>	<i>sd</i>	<i>median</i>	<i>90p</i>	<i>N</i>	<i>mean</i>	<i>sd</i>	<i>median</i>	<i>90p</i>
NYSE	41,035,340	356	352	298	410	-	-	-	-	-
NYSE MKT	-	-	-	-	-	590,311	1,166	3,066	954	1,161
NYSE Arca	24,039,351	573	7,136	368	846	15,719,357	1,309	4,940	992	1,408
Nasdaq OMX BX	7,913,775	1,153	9,425	893	1,479	3,432,209	495	1,886	334	692
NASDAQ OMX PSX	3,478,058	1,125	6,511	903	1,575	2,656,520	503	1,869	345	746
NASDAQ	53,492,822	1,218	5,148	957	1,508	23,557,237	639	7,744	375	943
BATS	26,033,986	773	898	585	1,077	14,344,829	1,154	5,178	559	1,114
BATS Y	15,244,556	704	680	565	902	6,994,690	788	2,583	528	953
Direct Edge A	10,100,225	682	455	576	871	5,053,145	712	1,862	547	881
Direct Edge X	22,368,427	728	587	590	984	15,308,697	855	2,733	575	1,029
Chicago Stock Exchange	10,811	1,262	627	1,168	1,470	3,676	1,101	917	1,010	1,255
National Stock Exchange	11,886	52,824	29,633	52,447	94,009	4,705	52,473	29,350	53,255	92,826
FINRA TRF	64,940,748	86,979	2,311,033	7,149	115,260	28,693,459	101,277	371,501	6,982	176,022
All	268,669,985	21,624	1,136,806	849	8,344	116,358,835	25,648	189,542	717	10,226
All (excluding FINRA)	203,729,237	791	4,203	604	1,155	87,665,376	894	5,243	555	1,157

Table 3

Panel A: Quotes

<i>Exchange</i>	Tape A Securities					Tape C Securities					
	<i>Median SIP Processing Time</i>	<i>Median Reporting Latency</i>	<i>Difference (Estimated Transit Time)</i>	<i>Theoretical Minimum Transit Time</i>	<i>Estimated Transit Time / Theoretical Minimum</i>	<i>Median SIP Processing Time</i>	<i>Median Reporting Latency</i>	<i>Difference (Estimated Transit Time)</i>	<i>Theoretical Minimum Transit Time</i>	<i>Estimated Transit Time / Theoretical Minimum</i>	<i>Adjusted Ratio (+200 ms Reporting Latency)</i>
NYSE	220	301	81	-	-	460	-	-	-	-	-
NYSE MKT	220	-	-	-	-	460	937	477	188	2.5	3.6
NYSE Arca	220	302	82	-	-	460	977	517	188	2.8	3.8
Nasdaq OMX	220	877	657	188	3.5	460	325	-135	-	-	-
NASDAQ OMX	220	923	703	188	3.7	460	367	-93	-	-	-
NASDAQ	220	933	713	188	3.8	460	404	-56	-	-	-
BATS	220	507	287	113	2.5	460	523	63	86	0.7	3.1
BATS Y	220	486	266	113	2.4	460	510	50	86	0.6	2.9
Direct Edge A	220	491	271	113	2.4	460	529	69	86	0.8	3.1
Direct Edge X	220	517	297	113	2.6	460	526	66	86	0.8	3.1
Chicago	220	839	619	113	5.5	460	722	262	86	3.0	5.4

Panel B: Trades

<i>Exchange</i>	Tape A Securities					Tape C Securities					
	<i>Median SIP Processing Time</i>	<i>Median Reporting Latency</i>	<i>Difference (Estimated Transit Time)</i>	<i>Theoretical Minimum Transit Time</i>	<i>Estimated Transit Time / Theoretical Minimum</i>	<i>Median SIP Processing Time</i>	<i>Median Reporting Latency</i>	<i>Difference (Estimated Transit Time)</i>	<i>Theoretical Minimum Transit Time</i>	<i>Estimated Transit Time / Theoretical Minimum</i>	<i>Adjusted Ratio (+200 ms Reporting Latency)</i>
NYSE	240	298	58	-	-	480	-	-	-	-	-
NYSE MKT	240	-	-	-	-	480	954	474	188	2.5	3.6
NYSE Arca	240	368	128	-	-	480	992	512	188	2.7	3.8
Nasdaq OMX	240	893	653	188	3.5	480	334	-146	-	-	-
NASDAQ OMX	240	903	663	188	3.5	480	345	-135	-	-	-
NASDAQ	240	957	717	188	3.8	480	375	-105	-	-	-
BATS	240	585	345	113	3.1	480	559	79	86	0.9	3.2
BATS Y	240	565	325	113	2.9	480	528	48	86	0.6	2.9
Direct Edge A	240	576	336	113	3.0	480	547	67	86	0.8	3.1
Direct Edge X	240	590	350	113	3.1	480	575	95	86	1.1	3.4
Chicago	240	1168	928	113	8.2	480	1010	530	86	6.2	8.5

Table 4: Apple Trades Ordered by Participant Timestamp, November 13, 2015

(1) Trade No.	(2) Participant Timestamp	(3) SIP Timestamp	(4) NBB Direct	(5) NBO Direct	(6) NBB SIP	(7) NBO SIP	(8) Trade Cond.	(9) Exch.	(10) Trade Price	(11) Trade Size	(12) Buy Order	(13) SIP Priced	(14) Lost Profits
1	11:37:47.464119	11:37:47.464616	113.37	113.38	113.37	113.38	@F	Z	113.38	2500	1	1	0
2	11:37:47.464119	11:37:47.464706	113.37	113.38	113.37	113.38	@F	Z	113.38	100	1	1	0
3	11:37:47.464119	11:37:47.464762	113.37	113.38	113.37	113.38	@F	Z	113.39	100	1	0	
4	11:37:47.464119	11:37:47.464792	113.37	113.38	113.37	113.38	@F	Z	113.39	100	1	0	
5	11:37:47.464119	11:37:47.464848	113.37	113.38	113.37	113.38	@F	Z	113.39	200	1	0	
6	11:37:47.464135	11:37:47.464743	113.37	113.38	113.37	113.38	@F	K	113.38	100	1	1	0
7	11:37:47.464135	11:37:47.464820	113.37	113.38	113.37	113.38	@F	K	113.38	200	1	1	0
8	11:37:47.464135	11:37:47.464861	113.37	113.38	113.37	113.38	@F	K	113.39	100	1	0	
9	11:37:47.464135	11:37:47.464889	113.37	113.38	113.37	113.38	@F	K	113.39	100	1	0	
10	11:37:47.464135	11:37:47.464916	113.37	113.38	113.37	113.38	@F	K	113.39	100	1	0	
11	11:37:47.464298	11:37:47.464673	113.37	113.38	113.37	113.38	@F	Q	113.38	100	1	1	0
12	11:37:47.464298	11:37:47.464727	113.37	113.38	113.37	113.38	@F	Q	113.38	100	1	1	0
13	11:37:47.464298	11:37:47.464777	113.37	113.38	113.37	113.38	@F	Q	113.38	100	1	1	0
14	11:37:47.464298	11:37:47.464806	113.37	113.38	113.37	113.38	@F	Q	113.38	100	1	1	0
15	11:37:47.464315	11:37:47.464834	113.37	113.38	113.37	113.38	@F	Q	113.38	200	1	1	0
16	11:37:47.464315	11:37:47.464875	113.37	113.38	113.37	113.38	@F	Q	113.39	100	1	0	
17	11:37:47.464315	11:37:47.464903	113.37	113.38	113.37	113.38	@F	Q	113.39	100	1	0	
18	11:37:47.464315	11:37:47.464929	113.37	113.38	113.37	113.38	@F	Q	113.39	100	1	0	
19	11:37:47.464315	11:37:47.464943	113.37	113.38	113.37	113.38	@F	Q	113.39	100	1	0	
20	11:37:47.464360	11:37:47.465298	113.37	113.38	113.37	113.38	@F	P	113.38	100	1	1	0
21	11:37:47.464360	11:37:47.465320	113.37	113.38	113.37	113.38	@F	P	113.38	100	1	1	0
22	11:37:47.464360	11:37:47.465337	113.37	113.38	113.37	113.38	@F I	P	113.38	73	1	1	0
23	11:37:47.464360	11:37:47.465352	113.37	113.38	113.37	113.38	@F	P	113.38	200	1	1	0
24	11:37:47.464397	11:37:47.465380	113.37	113.39	113.37	113.38	@F	P	113.39	500	1	0	
25	11:37:47.464397	11:37:47.465423	113.37	113.39	113.37	113.38	@F	P	113.39	100	1	0	
26	11:37:47.464397	11:37:47.465441	113.37	113.39	113.37	113.38	@F	P	113.39	100	1	0	
27	11:37:47.464397	11:37:47.465456	113.37	113.39	113.37	113.38	@F	P	113.39	100	1	0	
28	11:37:47.464397	11:37:47.465472	113.37	113.39	113.37	113.38	@F	P	113.39	100	1	0	
29	11:37:47.464397	11:37:47.465487	113.37	113.39	113.37	113.38	@F	P	113.39	100	1	0	
30	11:37:47.464397	11:37:47.465502	113.37	113.39	113.37	113.38	@F I	P	113.39	72	1	0	
31	11:37:47.465000	11:37:47.467422	113.39	113.40	113.37	113.38		D	113.38	100	1	1	-0.02
32	11:37:47.466000	11:37:47.511814	113.39	113.40	113.39	113.40		D	113.39	100	0	1	0
33	11:37:47.466018	11:37:47.466459	113.39	113.40	113.39	113.40		Z	113.39	100	0	1	0
34	11:37:47.475000	11:37:47.478795	113.39	113.40	113.39	113.40		D	113.40	245	1	1	0
35	11:37:47.479000	11:37:47.482618	113.39	113.40	113.39	113.40		D	113.40	805	1	1	0

Note: Table illustrates trades matched to the prevailing SIP NBBO and Direct NBBO. *Participant Timestamp* is the time in microseconds at which a venue reports executing a trade. *SIP Timestamp* is the time the SIP placed the trade report on its multicast line for dissemination, which incorporates transit and SIP-processing latencies. The Direct NBBO is calculated using the Participant Timestamp for quote updates, which reflects the time an exchange matching engine processed a quote. The SIP NBBO is calculated using the traditional SIP Timestamp assigned to quotes, which reflects the time a SIP disseminated a quote update. The Direct NBBO is matched to each trade based on the Participant Timestamp of the trade and the Participant Timestamp of the Direct NBBO. The SIP NBBO is matched to each trade based on the Participant Timestamp of a trade and the SIP Timestamp of the SIP NBBO. See Sections 3 and 5 for additional details.

Table 5

Panel A				
<i>Exchange</i>	<i>% of Trades Matching SIP NBBO (Unadjusted)</i>	<i>Transaction Value (Unadjusted)</i>	<i>% of Trades Matching SIP NBBO (Adjusted)</i>	<i>Transaction Value (Adjusted)</i>
NYSE	92.66%	\$560,605,000,000	92.62%	\$559,738,000,000
NYSE MKT	72.10%	\$3,521,261,553	72.10%	\$3,521,713,446
NYSE Arca	90.15%	\$329,627,000,000	90.10%	\$329,286,000,000
Nasdaq OMX BX	88.85%	\$72,776,300,000	88.84%	\$72,725,600,000
NASDAQ OMX PSX	93.47%	\$49,616,300,000	93.47%	\$49,511,800,000
NASDAQ	90.17%	\$611,460,000,000	90.13%	\$610,534,000,000
BATS	88.06%	\$290,765,000,000	88.01%	\$290,312,000,000
BATS Y	90.28%	\$131,344,000,000	90.25%	\$131,229,000,000
Direct Edge A	93.03%	\$93,885,300,000	93.01%	\$93,800,300,000
Direct Edge X	92.54%	\$330,709,000,000	92.48%	\$330,312,000,000
Chicago Stock Exchange	10.07%	\$75,388,700,000	10.07%	\$75,388,500,000
National Stock Exchange	95.47%	\$55,348,258	95.39%	\$55,328,300
FINRA TRF	51.43%	\$1,407,850,000,000	51.44%	\$1,407,010,000,000
All venues:	75.33%	\$3,957,600,000,000	75.30%	\$3,953,400,000,000
All Exchanges:	88.53%	\$2,549,800,000,000	88.48%	\$2,546,400,000,000

Panel B		
<i>Exchange</i>	Lost Profit Per Share (Unadjusted)	Lost Profit Per Share (Adjusted)
NYSE	-0.0002*** (0.00001)	-0.0002*** (0.00001)
NYSE MKT	-0.0002*** (0.00002)	-0.0002*** (0.00002)
NYSE Arca	-0.0002*** (0.00001)	-0.0002*** (0.00002)
Nasdaq OMX BX	-0.0002*** (0.00001)	-0.0002*** (0.00001)
NASDAQ OMX PSX	-0.0003*** (0.00002)	-0.0003*** (0.00002)
NASDAQ	-0.0003*** (0.00002)	-0.0003*** (0.00003)
BATS	-0.0003*** (0.00003)	-0.0003*** (0.00003)
BATS Y	-0.0001*** (0.00001)	-0.0001*** (0.00001)
Direct Edge A	-0.0002*** (0.00001)	-0.0002*** (0.00001)
Direct Edge X	-0.0002*** (0.00002)	-0.0002*** (0.00002)
Chicago	0.0000 (0.00000)	0.0000 (0.00000)
NSX	-0.0002*** (0.00004)	-0.0002*** (0.00004)
FINRA TRF	-0.0003*** (0.00002)	-0.0003*** (0.00002)
All Venues	-0.0002*** (0.00001)	-0.0002*** (0.00001)
All Exchanges	-0.0002*** (0.00001)	-0.0002*** (-0.00023)

Note: Estimates reflect the mean amount of lost profit per share that liquidity takers experienced by having their trades priced at the SIP NBBO rather than at the Direct NBBO. Robust standard errors are in parentheses.

*** p<.01, ** p<.05, * p<.1

Panel C

Lost Profit Per Share Traded	Unadjusted		Adjusted	
	Non-Exchange Venues	Exchange Venues	Non-Exchange Venues	Exchange Venues
-0.1	0.000%	0.000%	0.000%	0.000%
-0.09	0.000%	0.000%	0.000%	0.000%
-0.08	0.000%	0.000%	0.000%	0.000%
-0.07	0.001%	0.001%	0.000%	0.001%
-0.06	0.001%	0.001%	0.001%	0.001%
-0.05	0.003%	0.003%	0.003%	0.003%
-0.04	0.005%	0.007%	0.005%	0.007%
-0.03	0.015%	0.020%	0.016%	0.020%
-0.02	0.065%	0.078%	0.068%	0.080%
-0.01	2.578%	2.094%	2.635%	2.126%
0	97.129%	97.732%	97.061%	97.698%
0.01	0.198%	0.060%	0.204%	0.060%
0.02	0.004%	0.002%	0.004%	0.002%
0.03	0.001%	0.000%	0.001%	0.000%
0.04	0.000%	0.000%	0.000%	0.000%
0.05	0.000%	0.000%	0.000%	0.000%
0.06	0.000%	0.000%	0.000%	0.000%
0.07	0.000%	0.000%	0.000%	0.000%
0.08	0.000%	0.000%	0.000%	0.000%
0.09	0.000%	0.000%	0.000%	0.000%
0.1	0.000%	0.000%	0.000%	0.000%

Table 6

	Frequency of First-Leg Trades Having a Second-Leg Match	Std. Dev.	N
All Exchanges	0.007	0.002	9,201,335
Non-Exchanges	0.014	0.004	1,824,470
Combined	0.008	0.004	11,025,805

Table 7

Venue:	E/Q Ratio SIP NBBO As Benchmark	E/Q Ratio Direct NBBO As Benchmark	Difference	N
NYSE MKT	0.8482	0.8506	-0.0024***	137,825
Chicago Stock Exchange	15.6757	15.6728	0.0029	11,515
National Stock Exchange	0.9986	0.9943	0.0043***	1,337
NYSE	0.9115	0.8985	0.0130***	17,379,603
NYSE Arca	0.8933	0.8860	0.0073***	14,579,292
Nasdaq	0.9114	0.8957	0.0157***	34,388,463
Nasdaq BSX	0.9324	0.9300	0.0024***	6,727,836
Nasdaq PSX	0.9103	0.9104	-0.0001***	2,197,674
BATS X	0.8825	0.8639	0.0185***	17,470,674
BATSY	0.9691	0.9572	0.0119***	14,079,200
DirectEdge A	0.9714	0.9613	0.0101***	9,377,659
DirectEdge J	0.9536	0.9379	0.0157***	16,640,596
FINRA	0.7371	0.7409	-0.0038***	92,262,303

*** p<.01, ** p<.05, * p<.1