

# Risk and Return in High-Frequency Trading

**Matthew Baron (Cornell University)**

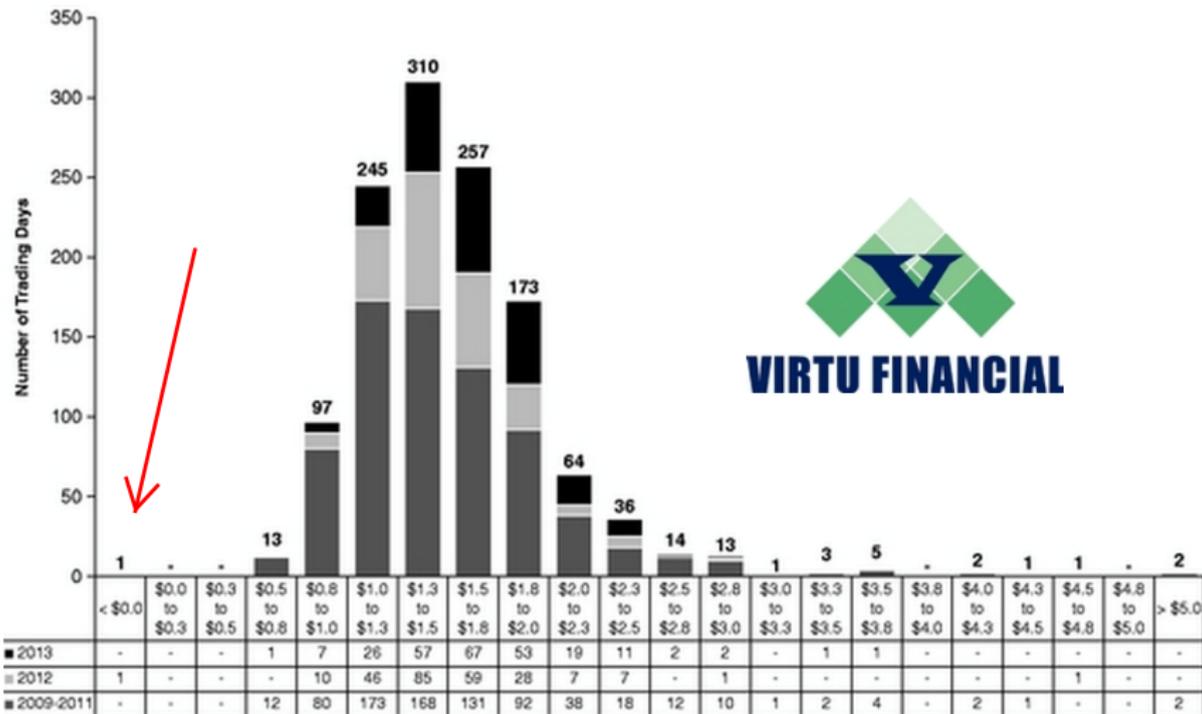
Jonathan Brogaard (University of Washington)

Björn Hagströmer (Stockholm Business School)

Andrei Kirilenko (Imperial College Business School)

**March 1, 2017**

# Virtu's trading record



# Main results

- 1 We find ***large, persistent*** differences in trading performance across HFTs
- 2 Differences in ***relative*** latency account for much of the difference in trading performance across HFTs
  - Better trading performance for HFTs that lower latency after colocation upgrades

# Main results

- ③ Being fastest is important for a **variety** of trading strategies
  - **Short-term information** channel and **risk management** channel
  - **Cross-market arbitrage**: React quicker to changes in futures market
  
- ④ We examine some **implications for market concentration**

# Isn't it obvious that speed is important?

- 1 **Not all HFTs choose co-location upgrades or trade in micro-seconds**
  - But those that do have the best trading performance
  
- 2 **Unclear which is more important for trading performance: *relative* or *nominal latency***
  - *Relative latency* can lead to ([Biais et al., 2015](#); [Budish et al., 2015](#)):
    - high concentration that does not decrease over time
    - over-investment in speed (e.g., microwave transmitters)

# Isn't it obvious that speed is important?

## ④ Unclear through which channels speed is important

- **Short-term informat. advantages** from speed: can reduced market quality
  - [Foucault, Hombert and Roşu \(2016\)](#): fast traders trade aggressively on news, picking off stale quotes.
  - [Chaboud et al. \(2014\)](#), [Foucault, Kazhan, & Tham \(2014\)](#): fast traders better at cross-market arb opportunities.
- **Better risk-management** from speed: can improved market quality
  - [Hoffmann \(2014\)](#): low latency allows liquidity providers to reduce their adverse selection costs
  - [Aït-Sahalia and Saglam \(2014\)](#): fast traders also benefit in terms of reduced inventory costs

## **1. Data & Methodology**

- HFT Identification, HFT trading performance measures

## **2. Relative Latency and Trading Performance**

- Alternative latency measures, Evidence from colocation upgrades

## **3. How do HFTs use latency?**

- Short-term information vs. risk-management channel, Cross-market arbitrage

## **4. Potential implications for market concentration**

- Profitability and concentration over the long-run, Entry and exit

## Sample:

- 25 Swedish large-cap stocks
- January 2010 – December 2014
- All trading venues in Sweden: lit and dark

## Data source:

### Transaction Reporting System

Broker-reported trade proprietary data  
Identifiers for brokers and clients  
Second time stamps

### Thomson Reuters Tick History

Public data feed  
Partial broker identifiers  
Microsecond time stamps

# HFT Identification

We use 25 firms who self-describe as HFTs

- based on the FIA-EPTA membership website

Narrow down to 16 HFTs that “actively trade”

- required to trade  $>10$  MSEK (about 1 M USD) on for  $>50$  days (out of 1,255 trading days)

“Behavior-based” identification based on 1) high trading volume and 2) low intraday & end-of-day inventory *gets nearly identical list*

# HFTs on NASDAQ-OMX (according to public records)

---

Algoengineering  
All Options International  
Citadel Securities  
Flow Traders  
GETCO<sup>a</sup>  
Hardcastle Trading  
IMC Trading  
International Algorithmic Trading (SSW Trading)  
Knight Capital<sup>a</sup>  
Madison Tyler<sup>b</sup>  
MMX Trading  
Optiver  
Spire  
Susquehanna Int. Sec.  
Timber Hill  
WEBB Traders  
Virtu Financial<sup>b</sup>  
Wolverine Trading UK

---

<sup>a</sup> Knight Capital merged with GETCO in July 2013

<sup>b</sup> Madison Tyler merged with Virtu Financial in July 2011

# HFT performance measures

## “Quantity” measures:

$$\text{Revenues} = \sum_{n=1}^N p_n q_n + p_{EOD} q_{EOD}$$

Cash flow for trade  $n$ , where  $q_n$  is the signed quantity

End-of-day position closed at closing price

$$\text{Trading volume} = 10^{-6} \sum_{n=1}^N |p_n q_n|$$

## Risk-adjusted measures:

$$\text{Return} = \frac{\text{Revenues}}{\text{Firm capitalization}}$$

$$\text{Sharpe ratio} = \frac{\text{Mean(Revenues)}}{\text{Sd(Revenues)}} \times \sqrt{252}$$

## “Quality” measure:

$$\text{Revenues per MSEK traded} = \frac{\text{Revenues}}{\text{Trading volume}}$$

# Risk and return in the cross-section of HFTs

	Mean	Std. Dev.	p10	p25	p50	p75	p90
Revenues (SEK)	18,181	29,519	-7,572	-487	6,990	31,968	61,354
Revenues per MSEK Traded	153.25	504.78	-257.94	-43.7	56.45	147.24	472.16
Returns	0.29	0.42	-0.09	0.01	0.09	0.51	0.89
Sharpe Ratio	4.16	6.58	-1.47	0.33	1.61	7.02	11.14
1-factor Alpha	0.29	0.43	-0.08	0.01	0.10	0.51	0.90
3-factor Alpha	0.29	0.43	-0.07	0.01	0.09	0.51	0.94
4-factor Alpha	0.29	0.43	-0.06	0.01	0.09	0.51	0.94
Trading Volume (MSEK)	272.05	378.09	4.20	7.39	63.69	507.67	909.20
Aggressiveness Ratio	0.51	0.26	0.16	0.28	0.56	0.69	0.88
End-of-Day Inventory Ratio	0.23	0.23	0.01	0.02	0.13	0.33	0.63
Max intraday Inventory Ratio	0.28	0.25	0.03	0.07	0.18	0.41	0.70
Average Trade Size (thous SEK)	239.19	697.38	46.17	56.64	72.24	92.18	173.39
Decision Latency (microseconds)	86,859	168,632	42	209	22,522	48,472	508,869

(N = 16 firms)

# Are trading revenues a good proxy for firm profits?

## Public filings of 5 HFTs: comparison of trading revenues with firm net profits

	Virtu				KCG		GETCO				Flow Traders			Jump
	2014	2013	2012	2011	2014	2013	2012*	2011	2010	2009	2014	2013	2012	2010
Trading Revenues (in millions)	685.2	623.7	581.5	449.4	1,274.0	903.8	526.6	896.5	865.1	955.2	240.8	200.5	125.1	511.6
-- % of revenue from proprietary trading	98.5%	98.4%	100%	100%	68.5%	67.0%	89.9%	94.2%			100%	100%	100%	
Trading Costs (% of Trading Revenue)	60.0%	57.8%	72.6%	62.1%	52.4%	59.0%	62.5%	48.5%	48.6%	40.4%	41.6%	43.7%	47.5%	
-- Brokerage, exch. & clearance fees	33.7%	31.3%	34.5%	32.9%	23.9%	27.3%	35.3%	32.2%	35.1%	32.1%	15.7%	15.8%	14.8%	
-- Communication and data processing	10.0%	10.4%	9.5%	10.3%	11.8%	13.7%	17.2%	9.7%	7.1%	4.5%				
-- Equipment rentals, deprec. & amort	4.5%	4.0%	15.7%	11.1%	10.4%	11.0%	9.1%	6.2%	6.2%	3.8%	1.8%	1.9%	2.4%	
-- Net interest (from credit lines, etc.)														
and dividends paid on sec borrowed	8.6%	7.8%	7.1%	6.0%	5.4%	6.5%	1.0%	0.3%	0.1%	0.0%	12.5%	12.8%	12.3%	
-- Other trading costs (e.g., administrative & technical costs)	3.2%	4.4%	5.8%	1.8%	0.8%	0.5%	0.0%	0.0%	0.0%	0.0%	11.5%	13.3%	18.0%	
Trading Profit Margin	40.0%	42.2%	27.4%	37.9%	47.6%	41.0%	37.5%	51.5%	51.4%	59.6%	58.4%	56.3%	52.5%	52.3%
Trading Revenue / (Trading Assets Minus Trading Liabilities)**	228%	196%	184%		96%	60%	62%				118%	119%	103%	237%
Trading Revenue / (Book Equity)	135%	138%	84%		84%	60%	80%				169%	146%	123%	222%

- 1 Profit margins are high (40-60%); do not vary much across firms & time
- 2 Fixed costs are small (15% of the total costs); no obvious relationship between trading profits & fixed costs

We conclude that **HFT trading revenue is a close proxy for HFT profits.**

# Measuring HFT latency



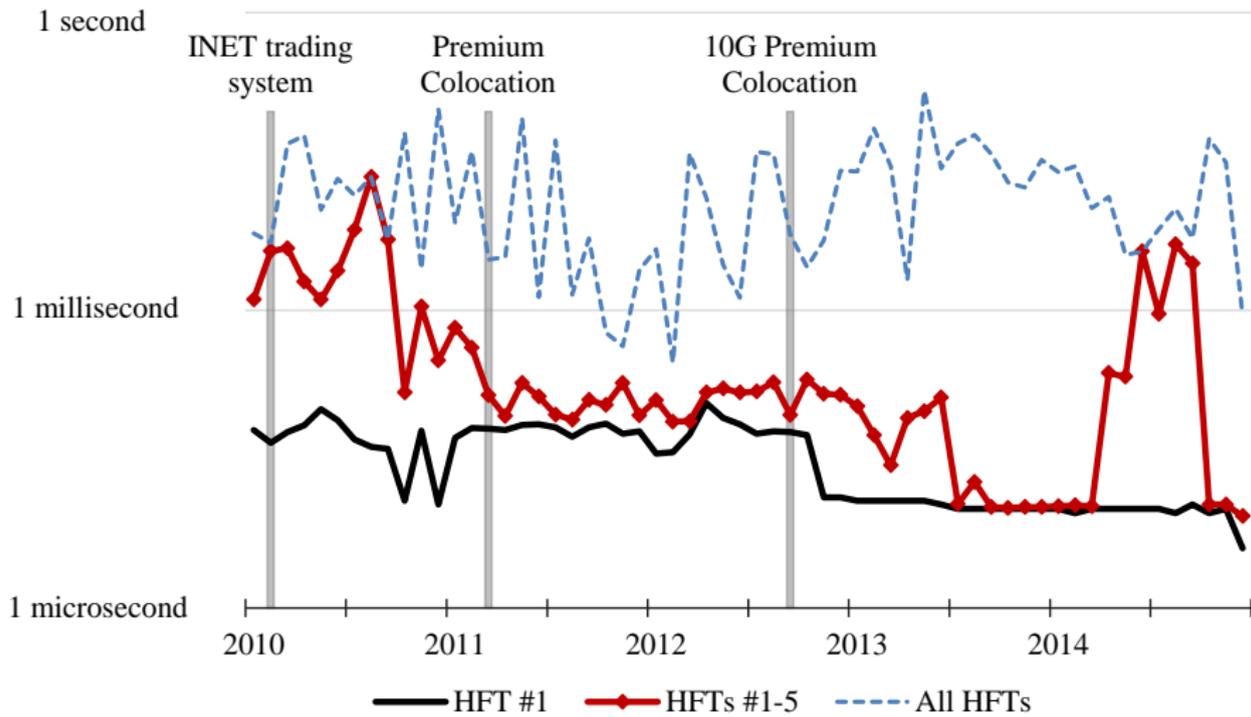
## Main measure: *Decision Latency*

- Aim:
  - Measure how fast HFTs can respond to new information
- Strategy:
  - 1 Measure the time from a passive execution (signal) to a reverse active execution (response) in the same stock and at the same venue (Weller, 2013)
  - 2 Record the 0.1% quantile of the distribution of reactions in each firm-month
    - (Or, alternatively, the mean of this distribution conditional on  $< 1$  millisecond)

## Alternative approaches in this paper:

- ***Queuing Latency***: measures the race to be at the top of the order book (Yao and Ye, 2015; Yueshen, 2014)
- ***Two colocation upgrades***: improve the relative latency of some HFTs, as they jump in rank relative to other HFTs

# HFT latency over time



# HFT latency and trading performance

$$\text{Performance}_{i,t} = \alpha_t + \beta_1 \log(\text{Decision Latency})_{i,t} + \beta_2 \text{Top1}_{i,t} + \beta_3 \text{Top5}_{i,t} \\ + \gamma' \text{Controls}_{i,t} + \text{Month FEs} + \epsilon_{i,t}$$

**Performance measures** = Revenues, Returns, Sharpe Ratio, etc.

$\log(\text{Decision Latency})$  = **nominal speed**

*Top 1* and *Top 5* rank dummies = **relative speed**

**Firm-month controls** = firm's inventory limits, aggressiveness,

**Time FEs** = account for market conditions like volatility and market volume

# HFT latency and trading performance

	Revenues			Returns			Sharpe Ratio		
Log <i>Decision Latency</i>	-14020*** (4311)	-1063 (6358)	9925 (10481)	-.221*** (.0483)	-.059 (.065)	-.00349 (.0852)	-4.38*** (.632)	-1 (1.2)	2.03 (1.46)
Top 1 dummy		29849* (15251)	24639** (12249)		.238* (.134)	.252* (.142)		3.77* (2.21)	4.2* (2.29)
Top 1-5 dummy		24074** (11619)	15451* (8009)		.333** (.155)	.303** (.133)		7.29** (3.24)	5.61** (2.63)
End-of-Day Inv.			2921 (3774)			.0839* (.0494)		2*** (.74)	
Max Intraday Inv.			-21008** (8579)			<i>[omitted]</i>		-3.74*** (1.23)	
Investment Horizon			-5401 (5994)			-.134*** (.0404)		-2.25*** (.726)	
Aggressive Rat.			5481 (3865)			-.0212 (.0558)		-.779 (.823)	
Constant	20278*** (6973)	8466** (4189)	10894** (4885)	.254*** (.0579)	.104* (.0587)	.107** (.0513)	5.1*** (1.26)	1.94 (1.23)	2.26* (1.23)
Month FEs	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
R-squared	0.123	0.168	0.263	0.198	0.233	0.269	0.207	0.254	0.361
N	737	737	737	737	737	737	737	737	737

Log(*Decision Latency*) and Controls in units of standard deviation

Standard errors dually clustered by firm and month

# HFT latency and trading performance

	Revenues			Trading Volume (x 10 <sup>-6</sup> )			Revenues per MSEK Traded		
Log <i>Decision Latency</i>	-14020*** (4311)	-1063 (6358)	9925 (10481)	-247*** (43.7)	-89.7 (59.1)	10.5 (74)	-19.4 (57.5)	-10.7 (69.1)	101** (40.4)
Top 1 dummy		29849* (15251)	24639** (12249)		326*** (97.9)	281*** (104)		6.99 (51.7)	57.6* (32.8)
Top 1-5 dummy		24074** (11619)	15451* (8009)		301** (132)	201** (97.4)		19.4 (93.3)	44.1 (55.9)
End-of-Day Inv.			2921 (3774)			-33.9** (15.9)			326* (168)
Max Intraday Inv.			-21008** (8579)			-183*** (65.3)			-76.3 (127)
Investment Horizon			-5401 (5994)			-76.4 (50.3)			-73.3 (63.8)
Aggressive Rat.			5481 (3865)			41.7 (28.8)			-55.5 (65.8)
Constant	20278*** (6973)	8466** (4189)	10894** (4885)	313*** (75.9)	169*** (57.8)	198*** (56.3)	35.2 (57.3)	27 (80.2)	7.91 (10)
Month FEs	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes	Yes
R-squared	0.123	0.168	0.263	0.294	0.362	0.454	0.080	0.080	0.148
N	737	737	737	737	737	737	737	737	737

Log(*Decision Latency*) and Controls in units of standard deviation

Standard errors dually clustered by firm and month

- ① Results robust: accounting for exchange fees and liquidity rebates, etc.
  - ② More importantly, need for robustness check for **latency measure** that
    - Does not rely on microsecond time stamps, and
    - Captures alternative HFT strategies
- 
- **Queuing latency:**
    - Measures the **race to be at the top of the order book** (Yao and Ye, 2015; Yueshen, 2014).
    - Specifically, when the price changes and a new tick opens up, how often does a given HFT get to the top of the queue?

# Queuing latency and trading performance

	Revenues			Returns	Sharpe Ratio	Trading Volume	Revenues per MSEK Traded
Log ( <i>Queuing Latency</i> + 1)	16761*** (4608)	4150 (4907)	-8265 (11994)	.121* (.0671)	-.283 (1.43)	34.8 (102)	-95 (62.1)
Top 1 dummy	50803*** (18676)	51684*** (15865)	.471** (.218)	12.6*** (2.13)	603*** (119)	35.5 (75.6)	
Top 1-5 dummy	13698* (7180)	9563 (7400)	.127 (.136)	2.39 (1.88)	128 (87.4)	87.7 (66.2)	
End-of-Day Inv.		2718 (3742)	.0858* (.0508)	1.95** (.767)	-32.5* (18.1)	325* (170)	
Max Intraday Inv.		-19571** (8658)	-2.9** (1.25)	-151** (68.5)	-58.6 (129)		
Investment Horizon		-4950 (7114)	-.0917*** (.0333)	-1.96** (.846)	-60.5 (60.8)	-72.8 (62.1)	
Aggressive Rat.		6664 (4551)	.00695 (.0428)	-.442 (.617)	60.3* (33.9)	-51.4 (65.8)	
Constant	20223*** (6685)	10893** (5372)	11153** (4814)	.16*** (.0496)	2.91** (1.19)	203*** (58.2)	-7.21 (45.5)
Month FEs	Yes	Yes	Yes	Yes	Yes	Yes	Yes
R-squared	0.152	0.213	0.302	0.323	0.429	0.547	0.146
N	737	737	737	737	737	737	737

Log(*Decision Latency*) and Controls in units of standard deviation

Standard errors dually clustered by firm and month

# Colocation upgrades

To address endogeneity concerns, we study two colocation upgrades offered by NASDAQ-OMX:

- Disruptive events that cause some HFTs to increase in relative speed
  - March 14, 2011: “Premium Colocation” upgrade
  - September 17, 2012: “10G Colocation” upgrade
    - Previously studied by Brogaard et al. (2015)
    - Only about half of HFTs immediately subscribed to the new connection type
- We compare the change in trading performance for **HFTs that become *relatively faster*** to **HFTs that become *relatively slower***

# Colocation upgrades

HFT latency rank	Revenues				Returns				Sharpe Ratio			
	Before	After	Diff.	(S.E.)	Before	After	Diff.	(S.E.)	Before	After	Diff.	(S.E.)
Faster	9,537	52,770	43,233	(14,841)	0.022	0.158	0.136	(0.055)	1.47	1.68	0.20	(0.46)
Slower	31,557	32,811	1,255	(2,608)	0.777	0.748	-0.030	(0.067)	5.50	4.70	-0.81	(0.99)
Diff-in-diff			<b>41,978***</b>	(6,533)			<b>0.165**</b>	(0.045)			<b>1.01*</b>	(0.50)

HFT latency rank	Trading Volume ( $\times 10^{-6}$ )				Revenues per MSEK Traded			
	Before	After	Diff.	(S.E.)	Before	After	Diff.	(S.E.)
Faster	415.9	537.4	121.5	(101.8)	-21.3	87.7	109.0	(33.9)
Slower	448.6	398.1	-50.5	(43.1)	207.3	282.2	74.9	(65.2)
Diff-in-diff			<b>171.9**</b>	(50.1)			<b>34.1</b>	(40)

# How do HFTs use lower latency?

## Short-lived Information channel

- Theory:
  - Foucault, Hombert and Roşu (2016): fast traders trade aggressively on news, picking off stale quotes.
  - Biais et al. (2015), Chaboud et al. (2014), Foucault, Kazhan, & Tham (2014): fast traders superior ability to react to cross-market arb opportunities.
- We measure: **Active Price Impact** = b.p. change in midpoint from just before a trade initiated by HFT to 10 seconds after

# How do HFTs use lower latency?

## Risk Management channel

- Theory:
  - [Hoffmann \(2014\)](#): low latency allows liquidity providers to reduce their adverse selection costs by revising stale quotes before picked off
  - [Aït-Sahalia and Saglam \(2014\)](#): fast traders also benefit in terms of reduced inventory costs
  
- We measure: **Passive Realized Spread** = b.p. difference between transaction price and midpoint 10 seconds after a trade in which HFT is liquidity provider.
  - Captures the **benefit of earning a wide bid-ask spread**
  - As well as the ability to **avoid supplying liquidity to trades with price impact.**

# How do HFTs use lower latency?

$$Performance_{i,s,t} = \alpha_t + \beta_1 \log(Decision\ Latency)_{i,t} + \beta_2 Top1_{i,t} + \beta_3 Top5_{i,t} + \gamma' Controls_{i,t,s} + Month\ FEs + \epsilon_{i,t}$$

	Price Impact		Realized Spread	
Log decision latency	-.318 (.225)	-.494* (.212)	-.364*** (.0982)	-.384*** (.0958)
Top 1 dummy	.371* (.201)	.337* (.193)	.0214 (.131)	.0599 (.126)
Top 1-5 dummy	.73** (.362)	.645** (.315)	.448*** (.136)	.477*** (.118)
Constant	3.91*** (.182)	3.96*** (.22)	-.0958 (.084)	-.108 (.107)
Month FEs	Yes	Yes	Yes	Yes
Stock FEs	Yes		Yes	
Firm & Stock controls		Yes		Yes
R-squared	0.196	0.016	0.158	0.017
N	11449	11449	11269	11269

Log(*Decision Latency*) and Controls in units of standard deviation  
 Standard errors dually clustered by firm and stock-month

We further examine both channels by focusing on **cross-market trading between the futures market and equities**

- **Cross-Market Short-Lived Information**

- We test if faster HFTs are more likely than slower HFTs to actively trade in equities in quick response to “news” in the futures market
  - “News” is defined to be a price change in the OMXS30 futures above a certain size.

- **Cross-Market Risk Management**

- We test if faster HFTs are less likely than slower HFTs to be adversely selected in a passive trade in equities markets in response to “news” in the futures market.

# Cross-market arbitrage

$$Pr[\text{Fast HFT Trades}] = \Phi[\beta \text{ News} + \gamma' \text{ Controls} + \text{StockFEs}].$$

- This regression captures the **increased probability of a *Fast HFT* trading in equities relative to a *Slow HFT*, in response to “news” in the futures market.**
  - The unit of observation is an equity-markets trade.
  - To capture who is trading quickly in response to “news” in the futures market, we consider equity market trades in the 1-second interval subsequent to a “news” event in the futures market.
    - “News” =  $\pm 1$  (and 0 otherwise) when the return on the OMXS30 futures during a one-second window preceding the stock trade is “large”
  - The dependent variable is 1 when a Fast HFT executes an equities trade in the subsequent one-second and 0 if a Slow HFT does it.
    - Fast HFT = those being Top 1 or Top 1-5 of HFTs by trading speed within a month
    - Slow HFTs are those not among the top 5

# Cross-market arbitrage

	Active trading				Passive trading			
	"Fast" = Top 1 HFT		"Fast" = Top 1-5 HFT		"Fast" = Top 1 HFT		"Fast" = Top 1-5 HFT	
	Probit (1=Fast HFT)	Marginal effects	Probit (1=Fast HFT)	Marginal effects	Probit (1=Fast HFT)	Marginal effects	Probit (1=Fast HFT)	Marginal effects
Constant	1.055*** (0.31)		2.143*** (0.17)		0.551* (0.30)		1.646*** (0.15)	
News	0.139*** (0.04)	0.006	0.199*** (0.03)	0.008	0.001 (0.03)	0.004	-0.097*** (0.02)	-0.015
Lagged Volatility	-0.094 (0.08)	0.000	-0.007*** (0.00)	-0.001	-0.116 (0.12)	0.001	0.008*** (0.00)	0.001
Lagged Volume	-0.005*** (0.00)	0.000	-0.004*** (0.00)	0.000	0.001 (0.00)	0.000	0.000 (0.00)	0.000
Quoted Spread	-0.046*** (0.01)	-0.004	-0.035*** (0.00)	-0.001	-0.024 (0.02)	-0.002	-0.001 (0.00)	0.000
Depth at BBO	0.013 (0.03)	0.006	0.049*** (0.02)	0.001	-0.120** (0.05)	-0.009	-0.015 (0.02)	-0.003
Stock FEs	Yes		Yes		Yes		Yes	
Average N	109684		277044		95268		258409	
Avg. psuedo-R <sup>2</sup>	0.209		0.169		0.204		0.163	

# Implications for market concentration

Competing viewpoints regarding HFT market concentration:

- 1 Traditional models: more competition among market intermediaries → decrease their profits, lower trading costs for other investors
  - Ho and Stoll (1983), Weston (2000)
  
- 2 **Competition on *relative latency*** can lead to a distinct competitive environment
  - Budish, Cramton, Shim (2015), Biais, Foucault, Moinas (2015), Foucault, Kozhan Tham (2015)
  - **Small increases in trading speed lead to large, discontinuous differences in payoffs**
    - As the fastest HFT responds first to profitable trading opportunities, capturing all the gains.
    - Marginally slower HFTs arrive too late.

# Implications for market concentration

**Predictions of this second viewpoint** ([Budish, Cramton, Shim, 2015](#); [Biais, Foucault, Moinas, 2015](#)):

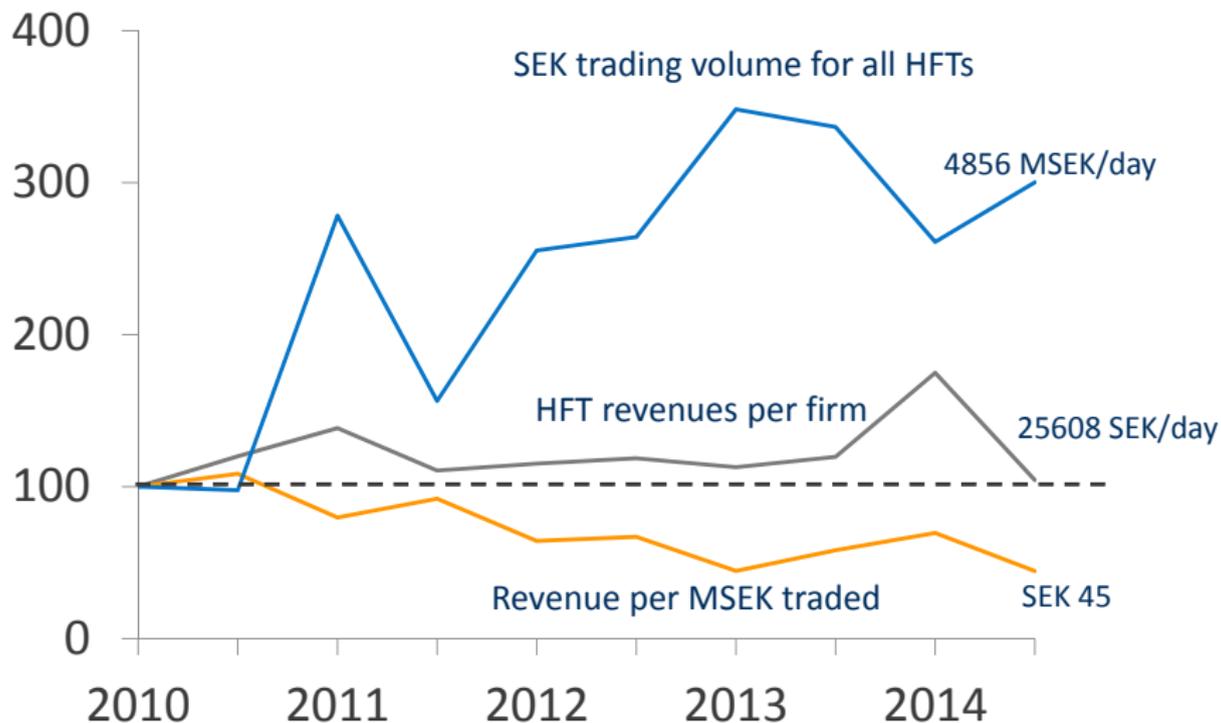
- 1 **Persistence** in performance, both at the firm-level and industry-wide level
- 2 **High concentration** of HFT revenues and trading volume
- 3 **Difficulty of new entry**



# Market concentration over time



# Aggregate Profits and Trading Volume over time



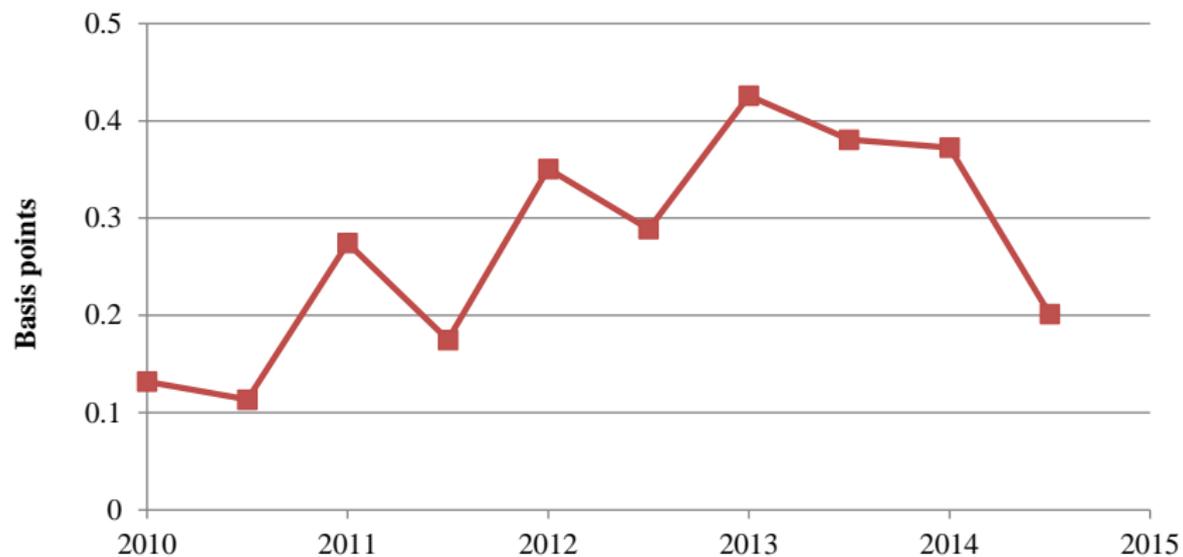
# Entry exit

	Revenues (thous. SEK)	Revenues per MSEK Traded	Returns	Daily Probability of Exit (x 10 <sup>3</sup> )	Decision Latency (in milliseconds, monthly obs.)
One-month dummy	-1.90** (.93)	-97.46 (209.7)	-.032** (.014)	1.455*** (.420)	44.36* (26.73)
Two-month dummy	-3.05** (1.434)	-87.11 (230.3)	-.033*** (.011)	1.486*** (.425)	134.6*** (33.98)
Three-month dummy	-.78 (1.35)	104.7 (196.6)	-.018* (.010)	-.194 (.477)	22.8** (11.19)
Constant	1.43*** (.19)	76.64*** (4.12)	.017*** (.002)	.530*** (.049)	14.87*** (.89)
Day x Stock FEs	Yes	Yes	Yes	Yes	(Month x Stock FEs)
R-squared	0.101	0.129	0.147	0.154	0.432
N	241053	241053	241053	241053	11014

New entrants in a given stock are less profitable, slower, and more likely to exit.

# HFT costs on non-HFTs

**Panel C: Cost of HFT Activities to Non-HFTs**



# Conclusions

- 1 We find *large, persistent* differences in trading performance across HFTs
- 2 Differences in *relative* latency account for much of the difference in trading performance across HFTs
  - Better trading performance for HFTs that lower latency after colocation upgrades
  - Lower latency associated with increased trading opportunities and risk-mitigation
    - No improvements in revenues per trade
- 3 Being fastest is important for a variety of trading strategies
  - **Short-term information** channel and **risk management** channel
  - **Cross-market arbitrage**: React quicker to changes in futures market
- 4 We examine some implications for market concentration