



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

How to Make Money in Microseconds

Citation for published version:

MacKenzie, D 2011, 'How to Make Money in Microseconds', *London Review of Books*, vol. 33, no. 10, pp. 16-18. <<http://www.lrb.co.uk/v33/n10/donald-mackenzie/how-to-make-money-in-microseconds>>

Link:

[Link to publication record in Edinburgh Research Explorer](#)

Document Version:

Publisher's PDF, also known as Version of record

Published In:

London Review of Books

Publisher Rights Statement:

© MacKenzie, D. (2011). How to Make Money in Microseconds. London Review of Books, 16-18.

General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact openaccess@ed.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.



London Review of Books

How to Make Money in Microseconds

Donald MacKenzie

What goes on in stock markets appears quite different when viewed on different timescales. Look at a whole day's trading, and market participants can usually tell you a plausible story about how the arrival of news has changed traders' perceptions of the prospects for a company or the entire economy and pushed share prices up or down. Look at trading activity on a scale of milliseconds, however, and things seem quite different.

When two American financial economists, Joel Hasbrouck and Gideon Saar, did this a couple of years ago, they found strange periodicities and spasms. The most striking periodicity involves large peaks of activity separated by almost exactly 1000 milliseconds: they occur 10-30 milliseconds after the 'tick' of each second. The spasms, in contrast, seem to be governed not directly by clock time but by an event: the execution of a buy or sell order, the cancellation of an order, or the arrival of a new order. Average activity levels in the first millisecond after such an event are around 300 times higher than normal. There are lengthy periods – lengthy, that's to say, on a scale measured in milliseconds – in which little or nothing happens, punctuated by spasms of thousands of orders for a corporation's shares and cancellations of orders. These spasms seem to begin abruptly, last a minute or two, then end just as abruptly.

Little of this has to do directly with human action. None of us can react to an event in a millisecond: the fastest we can achieve is around 140 milliseconds, and that's only for the simplest stimulus, a sudden sound. The periodicities and spasms found by Hasbrouck and Saar are the traces of an epochal shift. As recently as 20 years ago, the heart of most financial markets was a trading floor on which human beings did deals with each other face to face. The 'open outcry' trading pits at the Chicago Mercantile Exchange, for example, were often a mêlée of hundreds of sweating, shouting, gesticulating bodies. Now, the heart of many markets (at least in standard products such as shares) is an air-conditioned warehouse full of computers supervised by only a handful of maintenance staff.

The deals that used to be struck on trading floors now take place via 'matching engines', computer systems that process buy and sell orders and execute a trade if they find a buy order and a sell order that match. The matching engines of the New York Stock Exchange, for example, aren't in the exchange's century-old Broad Street headquarters with its Corinthian columns and sculptures, but in a giant new 400,000-square-foot plain-brick data centre in Mahwah, New Jersey, 30 miles from downtown Manhattan. Nobody minds you taking photos of the Broad Street building's striking neoclassical façade, but try photographing the Mahwah data centre and you'll find the police quickly taking an interest: it's classed as part of the critical infrastructure of the United States.

Human beings can, and still do, send orders from their computers to the matching engines, but this accounts for less than half of all US share trading. The remainder is algorithmic: it results from share-trading computer programs. Some of these programs are used by big

institutions such as mutual funds, pension funds and insurance companies, or by brokers acting on their behalf. The drawback of being big is that when you try to buy or sell a large block of shares, the order typically can't be executed straightaway (if it's a large order to buy, for example, it will usually exceed the number of sell orders in the matching engine that are close to the current market price), and if traders spot a large order that has been only partly executed they will change their own orders and their price quotes in order to exploit the knowledge. The result is what market participants call 'slippage': prices rise as you try to buy, and fall as you try to sell.

In an attempt to get around this problem, big institutions often use 'execution algorithms', which take large orders, break them up into smaller slices, and choose the size of those slices and the times at which they send them to the market in such a way as to minimise slippage. For example, 'volume participation' algorithms calculate the number of a company's shares bought and sold in a given period – the previous minute, say – and then send in a slice of the institution's overall order whose size is proportional to that number, the rationale being that there will be less slippage when markets are busy than when they are quiet. The most common execution algorithm, known as a volume-weighted average price or VWAP algorithm (it's pronounced 'veewap'), does its slicing in a slightly different way, using statistical data on the volumes of shares that have traded in the equivalent time periods on previous days. The clock-time periodicities found by Hasbrouck and Saar almost certainly result from the way VWAPs and other execution algorithms chop up time into intervals of fixed length.

The goal of execution algorithms is to avoid losing money while trading. The other major classes of algorithm are designed to *make* money by trading, and it is their operation that gives rise to the spasms found by Hasbrouck and Saar. 'Electronic market-making' algorithms replicate what human market makers have always tried to do – continuously post a price at which they will sell a corporation's shares and a lower price at which they will buy them, in the hope of earning the 'spread' between the two prices – but they revise prices as market conditions change far faster than any human being can. Their doing so is almost certainly the main component of the flood of orders and cancellations that follows even minor changes in supply and demand.

'Statistical arbitrage' algorithms search for transient disturbances in price patterns from which to profit. For example, the price of a corporation's shares often seems to fluctuate around a relatively slow-moving average. A big order to buy will cause a short-term increase in price, and a sell order will lead to a temporary fall. Some statistical arbitrage algorithms simply calculate a moving average price; they buy if prices are more than a certain amount below it and sell if they are above it, thus betting on prices reverting to the average. More complicated algorithms search for disturbances in price patterns involving more than one company's shares. One example of such a pattern, explained to me by a former statistical arbitrageur, involved the shares of Southwest Airlines, Delta and ExxonMobil. A rise in the price of oil would benefit Exxon's shares and hurt Delta's, while having little effect on Southwest's (because market participants knew that, unlike Delta, Southwest entered into hedging trades to offset its exposure to changes in the price of oil). In consequence, there was normally what was in effect a rough equation among relative changes in the three corporations' stock prices: $\Delta + \text{ExxonMobil} = \text{Southwest Airlines}$. If that equation temporarily broke down, statistical arbitrageurs would dive in and bet (usually successfully) on its reasserting itself.

No one in the markets contests the legitimacy of electronic market making or statistical arbitrage. Far more controversial are algorithms that effectively prey on other algorithms. Some algorithms, for example, can detect the electronic signature of a big VWAP, a process called 'algo-sniffing'. This can earn its owner substantial sums: if the VWAP is programmed to buy a particular corporation's shares, the algo-sniffing program will buy those shares faster than the VWAP, then sell them to it at a profit. Algo-sniffing often makes users of VWAPs and other execution algorithms furious: they condemn it as unfair, and there is a growing business in adding 'anti-gaming' features to execution algorithms to make it harder to detect and exploit them. However, a New York broker I spoke to last October defended algo-sniffing:

I don't look at it as in any way evil ... I don't think the guy who's trying to hide the supply-demand imbalance [by using an execution algorithm] is any better a human being than the person trying to discover the true supply-demand. I don't know why ... someone who runs an algo-sniffing strategy is bad ... he's trying to discover the guy who has a million shares [to sell] and the price then should readjust to the fact that there's a million shares to buy.

Whatever view one takes on its ethics, algo-sniffing is indisputably legal. More dubious in that respect is a set of strategies that seek deliberately to fool other algorithms. An example is 'layering' or 'spoofing'. A spoofer might, for instance, buy a block of shares and then issue a large number of buy orders for the same shares at prices just fractions below the current market price. Other algorithms and human traders would then see far more orders to buy the shares in question than orders to sell them, and be likely to conclude that their price was going to rise. They might then buy the shares themselves, causing the price to rise. When it did so, the spoofer would cancel its buy orders and sell the shares it held at a profit. It's very hard to determine just how much of this kind of thing goes on, but it certainly happens. In October 2008, for example, the London Stock Exchange imposed a £35,000 penalty on a firm (its name has not been disclosed) for spoofing.

Some, but not all, automated trading strategies require ultra-fast 'high-frequency trading'. Electronic market making is the clearest example. The 'spread' between the price at which a market-making program will buy shares and the price at which it will sell them is now often as little as one cent, so market-making algorithms need to change the quotes they post very quickly as prices and the pattern of orders shift. An algo-sniffer or statistical arbitrageur may have a little more time: I've been told, for example, that statistical arbitrage programs may hold a position for as long as a day (and in some cases even longer) before liquidating it. Even in those cases, however, an opportunity will vanish very quickly indeed if another algorithm spots it first.

Speeds are increasing all the time. In Hasbrouck and Saar's data, which come from 2007 and 2008, the salient unit of trading time was still the millisecond, but that's now beginning to seem almost leisurely: time is often now measured in microseconds (millionths of a second). The London Stock Exchange, for example, says that its Turquoise trading platform can now process an order in as little as 124 microseconds. Some market participants are already talking in terms of nanoseconds (billionths of a second), though that's currently more marketing hype than technological reality.

Because the timescales of trading have changed, the significance of space has also altered. A few years ago, it was common to proclaim the 'end of geography' in financial markets, and it's certainly true that if one is thinking in terms of hour-by-hour or even minute-by-minute

market movements, it doesn't really matter whether a trader is based in London, New York, Tokyo, Singapore or São Paulo. However, that's not the case in high-frequency trading. Imagine, for example, that your office is in Chicago, the second largest financial centre in the US, and you want to trade on the New York Stock Exchange. You are around 800 miles away from the matching engines in Mahwah, and sending a message that distance, using the fastest fibre-optic route between Chicago and New Jersey that I know of, takes around 16 milliseconds. That's a huge delay: you might as well be on the moon. Technical improvements in the amplifiers needed to boost signal strength and in other aspects of fibre-optic transmission will reduce the delay somewhat, as would straightening the route (fibre-optic cables still tend to follow railway lines because it's easy to negotiate rights of way there, but railways don't usually run in straight lines for long distances, instead going via centres of population). Ultimately, however, the speed of light is an insuperable barrier. If Einstein is right, no message is ever going to get from Chicago to Mahwah in less than four milliseconds.

The solution is what's called 'colocation': placing the computer systems on which your algorithms run next to the matching engines in data centres such as Mahwah. Colocation isn't cheap – a single rack on which to place your server can cost you \$10,000 a month, and it has become a big earner for exchanges and other electronic trading venues – but it's utterly essential to high-frequency trading. Even the precise whereabouts of your computers within data centres is a matter of some sensitivity: you hear tales (possibly apocryphal) of traders gaining entry to centres and trying to have holes drilled in walls so that the route from their server to the matching engine is shorter. The New York Stock Exchange has put quite a lot of effort into ensuring that no one spot within the Mahwah facility is better than any other in terms of speed of access to the matching engines.

Tales of computers out of control are a well-worn fictional theme, so it's important to emphasise that it is not at all clear that automated trading is any more dangerous than the human trading it is replacing. If the danger had increased, one way it would manifest itself is in higher volatility of the prices of shares traded algorithmically. The evidence on that is not conclusive – like-for-like comparison is obviously hard, and the academic literature on automated trading is still small – but data we do have suggest, if anything, that automated trading reduces volatility. For example, statistical arbitrage algorithms that buy when prices fall and sell when they rise can normally be expected to dampen volatility.

The bulk of the research also suggests that automated trading makes the buying and selling of shares cheaper and usually easier. Renting rack space in a data centre may be expensive, but not nearly as expensive as employing dozens of well-paid human traders. Twenty years ago the 'spread' between the price at which a human market maker would buy and sell a share was sometimes as much as 25 cents; the fact that it is now often as little as one cent means substantial savings for mutual funds, pension funds and other large institutions, almost certainly outweighing by far their losses to algo-sniffers. When assessed on criteria such as the cost of trading, the effects of automation are probably beneficial nearly all of the time.

What needs weighing against this, however, are the implications of one strange and disturbing episode that lasted a mere 20 minutes on the afternoon of 6 May 2010, beginning around 2.40 p.m. The overall prices of US shares, and of the index futures contracts that are bets on those prices, fell by about 6 per cent in around five minutes, a fall of almost unprecedented rapidity (it's typical for broad market indices to change by a maximum of between 1 and 2 per cent in an entire day). Overall prices then recovered almost as quickly,

but gigantic price fluctuations took place in some individual shares. Shares in the global consultancy Accenture, for example, had been trading at around \$40.50, but dropped to a single cent. Sotheby's, which had been trading at around \$34, suddenly jumped to \$99,999.99. The market was already nervous that day because of the Eurozone debt crisis (in particular the dire situation of Greece), but no 'new news' arrived during the critical 20 minutes that could account for the huge sudden drop and recovery, and nothing had been learned about Accenture to explain its shares losing almost all their value.

Sotheby's price of \$99,999.99 is, of course, the giveaway. What happened between 2.40 and 3 p.m. – the 'flash crash' as it's called – was primarily an 'internal' crisis of the financial markets, not a response to external events. For five months, large teams from the Securities and Exchange Commission (SEC) and Commodity Futures Trading Commission (CFTC) researched what had gone wrong in great detail, ploughing through terabytes of data. While some market participants disagree with specific aspects of the analysis they published last September, most seem to feel that it is broadly correct.

The trigger was indeed an algorithm, but not one of the sophisticated ultra-fast high-frequency trading programs. It was a simple 'volume participation' algorithm, and while the official investigation does not name the firm that deployed it, market participants seem convinced that it was the Kansas City investment managers Waddell & Reed. The firm's goal was to protect the value of a large position in the stock market against further declines, and it did this by programming the algorithm to sell 75,000 index future contracts. (These contracts track the S&P 500 stock-market index, and each contract was equivalent to shares worth a total of around \$55,000. The seller of index futures makes money if the underlying index falls; the buyer gains if it rises.) The volume participation algorithm calculated the number of index futures contracts that had been traded over the previous minute, sold 9 per cent of that volume, and kept going until the full 75,000 had been sold. The total sell order, worth around \$4.1 billion, was unusually large, though not unprecedented: the SEC/CFTC investigators found two efforts in the previous year to sell the same or larger quantities of futures in a single day. But the pace of the sales on 6 May was very fast.

On both those previous occasions, the market had been able to absorb the sales without crashing. In the first few minutes after the volume participation algorithm was launched, at 2.32 p.m. on 6 May, it looked as if the market would be able to do so again. Electronic market-making algorithms bought the futures that the volume participation algorithm was selling, as did index-arbitrage algorithms. (These programs exploit discrepancies between the price of index futures and the price of the underlying shares. A large sell order in the index futures market will often create just such a discrepancy, which can be profited from by buying index futures and selling the underlying shares.) Algorithmic trading was still in the benign zone that it occupies most of the time: electronic market makers and arbitrageurs were 'providing liquidity', as market participants put it, making it possible for the volume participation algorithm to do its intended large-scale selling.

However, high-frequency traders usually program their algorithms to be 'market neutral', in other words to insulate their trading positions from fluctuations in overall market levels. From around 2.41 p.m., therefore, those algorithms started to sell index futures to counterbalance their purchases, and the electronic index futures market entered a spasm of the kind identified by Hasbrouck and Saar. One algorithm would sell futures to another algorithm, which in its turn would try to sell them again, in a pattern that the SEC/CFTC investigators call 'hot potato' trading. In the 14-second period following 2.45 and 13 seconds,

more than 27,000 futures contracts were bought and sold by high-frequency algorithms, but their aggregate net purchases amounted to only around 200 contracts. By 2.45 and 27 seconds, the price of index futures had declined by more than 5 per cent from its level four and a half minutes earlier. The market had entered a potentially catastrophic self-feeding downward spiral.

Fortunately, though, the electronic trading platform on which these index futures were being bought and sold – the Chicago Mercantile Exchange’s Globex system – is programmed to detect just such a spiral. Its ‘Stop Logic Functionality’ is designed to interrupt self-feeding crashes and upward price spikes. A ‘stop’ is an order that is triggered automatically when prices reach a preset adverse level. Buyers of index futures, for example, will sometimes try to protect themselves from catastrophic losses by placing stop orders that will sell those futures if their prices fall below a given level. However, these sales can potentially begin a cascade, causing further price falls which in turn trigger further stop orders. The goal of the Stop Logic Functionality is to halt this process by giving human traders time to assess what is happening, step in and pick up bargains.

At 2.45 and 28 seconds, the price falls triggered Globex’s Stop Logic Functionality, and it imposed a five-second pause in trading. It worked. As Alison Crosthwait of Instinet (one of the oldest electronic trading venues) told the readers of an internet discussion forum hosted by the TABB Group, the five-second pause ‘provided ample time for market participants to consider their positions and return to the market or not, depending on the conclusions they reached ... [It] allowed market participants to regain confidence.’ Their purchases stopped the downward spiral of the price of index futures when trading restarted five seconds later.

But the crisis was not yet over. Index arbitrage and other mechanisms tie the index futures market intimately to the underlying stock market, and by 2.45 p.m. the latter was largely paralysed. High-frequency trading systems are often programmed to cease operating if unusually large price movements occur, and other systems are monitored by human beings who have what is in effect a large red stop button on their screens. Throughout the United States the automated systems stopped and the red buttons were pushed. Some market participants told the SEC/CFTC investigators they had been scared that the price falls meant some catastrophe had occurred, but that somehow they hadn’t heard about it. Others seem simply to have worried that there was a technical fault, such as corruption of the incoming data feeds that carry price information. Orders were cancelled on a massive scale and no replacements posted. In the case of some corporations’ shares, the market effectively ceased to exist.

The world of human trading that algorithmic trading has largely replaced had at its heart a subtle compromise. To be a market maker on the steps of Chicago’s raucous open-outcry trading pits or in the marble-walled main trading room of the New York Stock Exchange conferred certain privileges. Unlike other market participants, a New York Stock Exchange ‘specialist’ (as the exchange’s official market makers were called) could see the ‘book’ of buy and sell orders that had not yet been executed. In return for this considerable advantage, specialists were required to keep trading going, even in the event of a considerable imbalance between buy and sell orders, by using their own capital to fill the gap, all the while adjusting prices until the imbalance disappeared. Market makers in both Chicago and New York sometimes overstepped the line, opportunistically exploiting their privileged positions. However, as the sociologist Mitchel Abolafia documented in *Making Markets* (1997), in general such opportunism was held in check, not just by the formal rules but by the presence

of informal norms among people who interacted with each other face to face day in, day out, year after year.

These delicate social ecosystems have not survived the transition to fully electronic trading. Market makers' privileges have largely vanished: for instance, you don't now need to be a market maker to get fast access to the New York Stock Exchange's 'book', you need only pay for what is called 'level two' access and to rent a few racks at Mahwah to ensure that the data arrive with minimal delay. Their obligations have been reduced commensurately, though some traces still linger: for example, official market makers are still obliged to quote a price at which they would buy and a price at which they would sell the shares in which they are making a market.

Pushing the red button on an official market maker's system, therefore, did not entirely remove the bids to buy and offers to sell, but reduced the bids to the lowest possible price that could be entered into electronic trading systems (one cent), and increased the offers to the maximum possible price (\$99,999.99). These 'stub quotes' allow market makers to fulfil their formal obligations, while being so hopelessly unattractive that under normal circumstances no one would ever want to take a market maker up on them. In the case of several stocks, however, the evaporation of the market by around 2.45 p.m. was so complete that stub quotes were the only ones left. In consequence, 'market orders' (orders simply to buy or to sell at the best available price) were executed against stub quotes, hence Accenture's price of a cent and Sotheby's of \$99,999.99.

The recovery was gradual, although largely complete by 3 p.m. It seems to have been led by futures prices on Globex stabilising and then rebounding after the five-second pause. Traders began to spot what appeared to be extraordinary opportunities, although they were later often to be disappointed when exchanges cancelled sales at a cent and purchases at \$99,999.99 on the grounds that they were 'clearly erroneous'. In fact what had happened was largely that trading had effectively stopped, so the sums of money lost (and made) were only modest, and wider financial damage limited. It isn't even likely that Waddell & Reed – if indeed it was their volume participation algorithm that sparked the flash crash – lost overwhelming amounts. The algorithm simply kept going through the turmoil – algorithms, after all, don't panic – and finally completed its 75,000 sales at 2.51 p.m. By then, futures prices were already well on the way back up, thus limiting the losses caused by the algorithm selling at temporarily very low price levels.

Despite the modesty of the losses incurred, many market participants and regulators found the flash crash deeply unnerving, and I think they were right to do so. What troubles me most about the episode is not something that happened, nor even something that was said, but something that was not said. Alison Crosthwait's posting elicited only five comments from other TABB forum members, and none disagreed with her judgment that five seconds was 'ample time for market participants to consider their positions'. She was certainly right to identify the triggering of the Stop Logic Functionality as the turning point, and the stabilisation of futures prices after the five-second pause shows that she was correct: five seconds *was* enough time. But bear in mind that she was talking about human beings coming to decisions and not computer systems recalibrating themselves: we don't ordinarily talk of computers 'considering' things and 'regaining confidence'. This is a situation that in the terminology of the organisational sociologist Charles Perrow is one of 'tight coupling': there is very little 'slack', 'give' or 'buffer', and decisions need to be taken in what is, on any

ordinary human scale, a very limited period of time. It takes me five seconds to blow my nose.

With the rise of electronic trading, the stock market (especially in the US) has become a system, and it is one of at least moderate complexity. True, there is nothing too dreadfully complicated about trading on any one exchange. The programs controlling Globex were, in consequence, perfectly able to detect a dangerous condition and pause trading accordingly. However, while the Chicago Mercantile Exchange has a dominant position in the trading of 'derivatives' such as index futures, the traditional stock exchanges such as those in New York and London have been losing business rapidly to other electronic trading venues. There are now some 50 such venues on which US shares are traded, and they don't operate in isolation. They are tied together by algorithms exploiting discrepancies in prices circulating among them, and also by the rules imposed by the Securities and Exchange Commission, which for decades has been trying to fuse the diverse exchanges of the US into a National Market System. The SEC requires, for example, that brokers don't simply execute their customers' orders at their preferred venue, but look for the most favourable prices (the 'national best bid and offer', as those prices are called). As Steve Wunsch, one of the pioneers of electronic exchanges, put it in another TABB forum discussion, US share trading 'is now so complex as a system that no one can predict what will happen when something new is added to it, no matter how much vetting is done.' If Wunsch is correct, there is a risk that attempts to make the system safer – by trying to find mechanisms that would prevent a repetition of last May's events, for example – may have unforeseen and unintended consequences.

Systems that are both tightly coupled and highly complex, Perrow argues in *Normal Accidents* (1984), are inherently dangerous. Crudely put, high complexity in a system means that if something goes wrong it takes time to work out what has happened and to act appropriately. Tight coupling means that one doesn't have that time. Moreover, he suggests, a tightly coupled system needs centralised management, but a highly complex system can't be managed effectively in a centralised way because we simply don't understand it well enough; therefore its organisation must be decentralised. Systems that combine tight coupling with high complexity are an organisational contradiction, Perrow argues: they are 'a kind of Pushmepullyou out of the Doctor Dolittle stories (a beast with heads at both ends that wanted to go in both directions at once)'.

Perrow's theory is just that, a theory. It has never been tested very systematically, and certainly never proved conclusively, but it points us in a necessary direction. When thinking about automated trading, it's easy to focus too narrowly, either pointing complacently to its undoubted benefits or invoking a sometimes exaggerated fear of out of control computers. Instead, we have to think about financial systems as a whole, desperately hard though that kind of thinking may be. The credit system that failed so spectacularly in 2007-8 is slowly recovering, but governments have not dealt with the systemic flaws that led to the crisis, such as the combination of banks that are too big to be allowed to fail and 'shadow banks' (institutions that perform bank-like functions but aren't banks) that are regulated too weakly. Share trading is another such system: it is less tightly interconnected in Europe than in the United States, but it is drifting in that direction here as well. There has been no full-blown stock-market crisis since October 1987: last May's events were not on that scale.^[*] But as yet we have done little to ensure that there won't be another.

[*] Donald MacKenzie wrote about the 1987 crash in the *LRB* of 4 August 2005.

[Vol. 33 No. 10 · 19 May 2011](#) » [Donald MacKenzie](#) » [How to Make Money in Microseconds](#)
pages 16-18 | 5183 words