# Detecting Wash Trade in Financial Market Using Digraphs and Dynamic Programming

Yi Cao, Yuhua Li, *Senior Member, IEEE*, Sonya Coleman, Member, *IEEE*, Ammar Belatreche, *Member*, *IEEE*, and Thomas Martin McGinnity, *Senior Member, IEEE*

*Abstract*— **Wash trade refers to the illegal activities of traders who utilise carefully designed limit orders to manually increase the trading volumes for creating a false impression of an active market. As one of the primary formats of market abuse, wash trade can be extremely damaging to the proper functioning and integrity of capital markets. Existing work focuses on collusive clique detections based on certain assumptions of trading behaviours. Effective approaches for analysing and detecting wash trade in a real-life market have yet to be developed. This paper analyses and conceptualises the basic structures of the trading collusion in a wash trade by using a directed graph of traders. A novel method is then proposed to detect the potential wash trade activities involved in a financial instrument by first recognizing the suspiciously matched orders and then further identifying the collusions among the traders who submit such orders. Both steps are formulated as a simplified form of the Knapsack problem, which can be solved by dynamic programming approaches. The proposed approach is evaluated on seven stock datasets from NASDAQ and the London Stock Exchange. Experimental results show that the proposed approach can effectively detect all primary wash trade scenarios across the selected datasets.**

*Index Terms*—**Market Abuse, Directed Graph, Dynamic Programming, Wash Trade.**

## I. Introduction

Surveillance of a financial exchange market for preventing market abuse activities has been attracting significant academic and industrial attention after the financial crisis in 2008 and especially since the flash crash in 2010. The abuse of financial markets can occur in a variety of ways, all of which can be extremely damaging to the proper functioning and integrity of the market. Trade-based manipulation, where the manipulation tactic is carried out only by simply buying and selling (Franklin & Douglas, Stock Price Manipulation, 1992), is one of the primary forms. Price and volume are usually two major objects to be manipulated, and the former format, price

manipulation, is thoroughly studied in a previous work by the authors (Cao, Li, Coleman, Belatreche, & T.M.McGinnity, Adaptive Hidden Markov Model with Anomaly States for Price Manipulation Detection, 2015) (Cao, Li, Coleman, Belatreche, & McGinnity, A Hidden Markov Model with Abnormal States for Detecting Stock Price Manipulation, Oct. 2013) (Cao, Li, Coleman, Belatreche, & McGinnity, Detecting price manipulation in the financial market, Mar. 2014) (Cao, Li, Coleman, Belatreche, & McGinnity, Detecting wash trade in the financial market, Mar. 2014) (Zhai & Cao, Mar. 2014). Another format of trade-based abuse is volume manipulation, the manipulation actions intending to increase the transaction volume for the purpose of giving a false impression of high trading volume on the market (Franklin & Douglas, Stock Price Manipulation, 1992) (Franklin, Litov, & Mei, Large investors, price manipulation, and limits to arbitrage: An anatomy of market corners, 2006). The major form of volume manipulation is wash trade, which occurs when the same individuals or a group of collusive clients are on both sell and buy sides of a financial instrument (i.e. stock) trading. While there is no beneficial change in ownership, wash trading has the effect of creating a misleading appearance of an active interest in the stock (Douglas, Feng, & Aitken, 2012).

Wash trade usually does not contain any illegal actions such as financial rumour spreading and market resource squeezing but is carried out only by legitimate trading activities. With carefully designed buy and sell order sequences, manipulators can make the transaction follow their expectation. In the wash trade tactics, a series of orders is often submitted as a number of order pairs. The monitoring of any single leg of one pair or part of a pair would not be concluded as collusive trading. Most of the existing related literature studies the collusive cliques according to the "activity similarity", which is defined under certain assumptions. Very few address quantitative analysis of the features of different wash trade scenarios and the corresponding detection approaches. This paper follows on from our previous work on trade-based manipulation (Cao, Li, Coleman, Belatreche, & T.M.McGinnity, Adaptive Hidden Markov Model with Anomaly States for Price Manipulation Detection, 2015) and proposes a detection approach that considers a complete spectrum of the wash trade detection. The main contributions of the work are as follows: the problem of wash trade is thoroughly discussed including the analysis of all possible scenarios, from which the key features are extracted and quantified. This provides a clear problem formulation and explains the significance of exploring the conceptual models. To the best of our knowledge, this is the first theoretical study of wash trade market manipulation. A two-step algorithm is proposed to detect wash trade activities. The proposed two

S. Coleman, and A. Belatreche are with the Intelligent Systems Research Centre, University of Ulster, Londonderry BT48 7JL, U.K. (e-mail:sa.coleman@ulster.ac.uk;a.belatreche@ulster.ac.uk).

Yuhua Li is with the School of Computing, Science and Engineering, University of Salford, Salford M5 4WT, UK. (Y.Li@salford.ac.uk ).

T. M. McGinnity is with the Intelligent Systems Research Centre, University of Ulster, Londonderry BT48 7JL, U.K. and also with the School of Science and Technology, Nottingham Trent University, Nottingham, NG1 4BR, U.K. (e-mail: martin.mcginnity@ntu.ac.uk).

Y. Cao is with the Institute of Management Science and Engineering, Henan University, Kaifeng, China. (jason.caoyi@gmail.com)

steps, which consist of discovering the matching orders and further recognizing the collusions, are both formulated as a combinatorial optimisation problem and solved by one unified algorithm. Extensive experiments have been conducted on real data from both USA and UK markets for testing the practicability of the proposed wash trade detection method in real-life.

The remainder of the paper is organised as follows. Section II provides a review of wash trade manipulation and the corresponding detection methods. The features of all types of wash trade scenarios as well as the proposed detection approach are analysed, formulated and characterised in Section III. Performance evaluation of the proposed approach is provided in Section IV. Finally, Section V concludes the paper and discusses potential improvements and future work.

## II. WASH TRADE AND ITS DETECTION

### A. Wash Trade

In capital markets, limit orders, indicate the trading intention of the trader to buy or sell volumes of a specific equity at a specific price or better (SEC, 2011) ("Better price" refers to higher selling prices or lower buying prices). The transaction occurs when eligible orders meet order-matching rules. The outstanding unmatched limit orders are recorded in order books of the exchange market, in which the highest buying price decides the best bid price while the lowest selling price is the best ask price. The gap between the best bid and ask price is defined as bid-ask spread (Kojo & Paudyal, 2000). In most of the exchange markets, the matching rule selects the earliest order with the matched price for execution. In the following examples in Table 1, three limit orders, #01, #02 and #03 are submitted in sequence to the exchange market. According to the matching rule, order #03 is firstly executed by 300 shares with #01, which has the same price but is earlier than order #02, and then the remaining 100 shares are executed with #02.

Table 1 Limit order sequences

| Order # | Trader | Time | Buy/Sell | Price | Volume |
|---|---|---|---|---|---|
| 01 | *A* | 09:00:000 | Buy | 125 | 300 |
| 02 | *A* | 09:05:000 | Buy | 125 | 300 |
| 03 | *B* | 09:06:100 | Sell | 125 | 400 |

Wash trades follow the same matching rules as legitimate transactions with the special feature defined by the Financial Conduct Authority (FCA) as "no change in beneficial interest or market risk", or "the transfer of beneficial interest or market risk only between parties acting in concert or collision, other than for legitimate reasons" (FSA, 2006). The Committee of European Securities Regulators (CESR) further indicates that a wash trade is the "deliberate arrangement in concert or collusion" (EU, Market Abuse Directive, 2014). On August 28 2014, Chicago Mercantile Exchange (CME) released a new rule (adopted by U.S. Commodity Futures Trading Commission CFTC), termed as "Rule 575" (CME, 2014). Rule 575 clearly states that "no person shall enter messages to the market as pre-arranged collusion (wash trade) with intent to mislead other participants". The definition in Rule 575 in the

U.S. shows the consistent regulation to CESR in Europe that the pre-arranged collusive trading is wash trade and shall be strictly prohibited. Although clearly defined the wash trade activity, the regulators (FCA, CESR, and CFTC) do not provide any quantitative approach on detecting such activities.

As illustrated by the example in Table 2, the simplest format of wash trade is the simultaneous submission of two opposite limit orders with identical price (125 in Table 2) and similar volume (495 in Table 2) from one trader *A*. By the matching rules, order #01 and #02 match and 495 shares are executed immediately after the submission. Additionally, the wash trade actions can be also carried out by multiple orders and traders as the example formats shown in Table 3 and Table 4. In Table 3, order #03 is matched and executed with #01 and #02 sequentially so that a transaction of 490 shares can be artificially created by trader *A*. In Table 4, two transactions are created by four matched orders between traders *A* and *B*. After the transactions (450 matched volumes), there is almost no effective transfer of beneficial interest among the two traders.

Table 2 Basic format of wash trade

| Order # | Trader | Time | Buy/Sell | Price | Volume |
|---|---|---|---|---|---|
| 01 | *A* | 09:00:000 | Buy | 125 | 500 |
| 02 | *A* | 09:00:001 | Sell | 125 | 495 |

Summarising the typical formats in Table 2 and Table 4 as well as the definitions from the regulators, we obtain three features of a successful execution of a wash trade manipulation:
1. Tight submission intervals between the matched buy and sell orders (to minimize the risk of the orders being unintentionally picked up by other traders);
2. Executable prices (to make the orders an immediate execution);
3. Mostly matched volumes (to minimize the risk of loss from the unmatched volumes executed with other traders).

Table 3 Wash trade with multiple orders

| Order # | Trader | Time | Buy/Sell | Price | Volume |
|---|---|---|---|---|---|
| 01 | *A* | 09:00:000 | Buy | 125 | 250 |
| 02 | *A* | 09:00:001 | Buy | 125 | 250 |
| 03 | *A* | 09:00:002 | Sell | 125 | 490 |

Table 4 Wash trade with multiple traders

| Order # | Trader | Time | Buy/Sell | Price | Volume |
|---|---|---|---|---|---|
| 01 | *A* | 09:00:000 | Buy | 125 | 500 |
| 02 | *B* | 09:00:001 | Sell | 124.2 | 490 |
| 03 | *B* | 09:10:000 | Buy | 125.5 | 490 |
| 04 | *A* | 09:10:001 | Sell | 125 | 500 |

Perfect matching orders, which have the same price, volume and submission time according to the summarised features, guarantee the execution but are obviously easy to be suspected as market abuse trade by the regulators. Therefore, to avoid being easily detected, "smart manipulators" design the wash trade orders to be "mostly matched" such as the examples in Table 2 and Table 4, where around 99% volumes are executed, respectively. Similarly, due to the matching rules in most exchange markets (Bowen, 2013), that is buy (sell) limit order

matching sell (buy) limit orders with the same price or lower (higher), the limit prices in the examples in Table 4, which are different but executable, are also deliberately designed to avoid inspection. In Table 4, order #02 can be executed with order #01 at price 125 and order #04 can be executed with order #03 at price 125.5. The 125 and 125.5 are the execution prices of the two possible transactions; we refer to such prices as transaction prices.

*B. Wash Trade Detection*

To the best of our knowledge, there is no related work on the detection of wash trade activities in capital markets. The only analogous research is work on the detection of collusive cliques based on certain similar trading behaviours, which are defined as the buy/sell activities of equities in a similar way. A spectral clustering based approach was developed (Markus., Hoser, & Schröder, 2007), where a trading-behavioural network is generated and any behaviour that deviates from the network is reported as an irregularity. The assumption of this work is the strong consistency between a trader's current behaviours and his/her previous trading network. A graph clustering algorithm for detecting a set of collusive traders has been proposed in (Palshikar & Apte, 2008). The relationship between traders is constructed as a stock flow graph, and those with "heavy trading" within their network are clustered as a collusion set.

A new trading collusion detection approach, the correlation matrix of one trading day, was presented in recent work (Wang, Zhou, & Guan, 2012), where trader behaviour was represented by an aggregated time series of signed volumes of submitted orders. The similarities of behaviours among multiple traders are measured by Pearson's product-moment coefficient and the cliques with a coefficient higher than a user-specified threshold were considered as suspicious collusions. The experiments of this study evaluated the real order data of futures traded in the Shanghai Futures Exchange. The "signed order volume" is constructed by volumes and directions (buy/sell) of the order. The order price information is ignored according to the assumption that order prices are not related to the trader's behaviours (Wang, Zhou, & Guan, 2012). However, the market impact measure shows that order price significantly impacts the market (Hautsch & Huang, The market impact of a limit order, 2012) so that the market moves caused by the traders' own actions (orders) become the principal part of the transaction costs (ITG, 2013). It is therefore unacceptable to ignore the order price information, which not only distinguishes traders' intention, but is a key feature of wash trade manipulation tactics.

A technique developed by the Chicago Mercantile Exchange (CME) to prevent wash trades at the "engine level" was rolled out in the middle of 2011 (Patterson, Strasburg, & Trindle, 2013) and updated in the summer of 2013 (Bowen, 2013). However, it only monitored the same-priced buy/sell orders from trading accounts with the same beneficial ownership (Bowen, 2013) (example in Table 2). The lack of the surveillance mechanisms for wash trades with multiple orders or traders (example illustrations in Table 3 and

Table 4) left it possible for collusive parties to create a number of transactions that give a false appearance of large trading volumes.

In December 2012, a wash trade case was manually inspected and documented by the Securities and Exchange Commission of Pakistan (Jamal, 2012). In March 2013, the US regulators started to investigate traders acting as both buyer and seller in the same transactions and reported that several hundred potential wash trades occur each day on CME and Intercontinental Exchange (ICE) (Patterson, Strasburg, & Trindle, 2013). In June 2012, the Hong Kong financial regulator claimed that the attempts of entering wash trade or matched trade were financial manipulation crimes whether or not the wash trade or matched trade in fact has, or is likely to have, the effect of misleading appearance (Loh & Cumming, 2012). This ruling was also accepted by "Rule 575" (CME, 2014) and the Market Abuse Directive II (EU, European Commission, 2015). This rule provided an aggressive restriction: any attempts of wash trade or matched trade are financial crimes.

To date, academic research has mainly focused on detecting the overall trading collusions according to defined analogous behaviours. The detection of mass market behaviours can hardly reach a precise and determinable manipulation detection result but can show a collective correlation of trading activities among different trader clusters. Industry techniques merely covered the simple format of wash trade scenarios. A slightly improved manipulation tactic can bypass the wash trade monitoring. However, no efforts appear to have been made in the analysis of wash trade strategic behaviour or the design of a detection approach identifying any tactics of attempts of wash trade. Given the gap in the field, it is this aspect of market manipulation that this paper seeks to address. This paper proposes a wash trade detection algorithm that monitors all incoming limit orders that can possibly attempt to compose a wash trade. Recognising such attempts helps the regulators to prevent market abuse by a strict regulation.

### III. WASH TRADE DETECTION METHODOLOGY

*A. Analysis Terminologies*

To analyse the wash trade strategic behaviours, the definitions and terminologies in (Tsang, Olsen, & Masry, 2013) are adopted and revised to formalise the trading properties and market changes. The effect of wash trade can be represented by the position of the whole trading collusion, where "position" is the amount of equities held by a trader. As the wash trade is merely fraudulent activities rather than true trading actions, each participated trader tends to maintain his own positions unchanged for minimising the unnecessary financial loss, and therefore the position of the whole wash trade collusive group is also not changed. During the wash trade process, the position change is caused by a number of orders from the trader in the collusive group and can be defined as:

$$\text{Position} + \text{Orders} \rightarrow \text{Position},$$

Position is comprised of a sequence of orders:

$$\text{Position} = \{(\text{Order}_1), (\text{Order}_2)\dots (\text{Order}_n)\},$$

where each order is defined as:

$$\text{Order} = (\text{Trader\_ID}, \text{Type}, \text{Price}, \text{Volume}),$$

where Type = buy | sell. Representing the order Type buy and sell by positive and negative signs respectively and affixing the sign to the Trader_ID and Volume, a sell order can be represented as:

$$\text{Order} = (\text{-Trader\_ID, Price, -Volume}). \qquad (1)$$

By this, the orders in Table 4 can be illustrated as:

Position={(A, 125, 500), (-B, 124.2, -450),
(B, 125.5, 450), (-A, 125, -500)  }.

The buy/sell orders having matched prices can be merged as:

Position={ ( A-B, 125, 500-450=50 ),
( B-A, 125.5, 450-500=-50 ) },

As discussed in Section II.A, prices 125 and 125.5 are represented as transaction prices. The difference between the executable limit prices is calculated as the margins of the transaction prices. In this case, the transaction price 125 has the margin 125-124.2=0.8 and the transaction price 125.5 has the margin 125.5-125=0.5. We merge the potential transactions who price margins are overlapped, i.e. 125+0.8 and 125.5+0.5 are overlapped. After the merge, we re-represent the positions, i.e. the margin between 124.2 and 125.5 is represented as: 124.85±0.65.

Position={ A-B+B-A, 124.85±0.65, 50-50 },
={ "0", 124.85±0.65, 0  },

where the Trader_ID calculation is carried out as a symbolic operation and 0.65 is represented as the transaction margin $\delta^T$ and 124.85 is the transaction price $P^T$. The zero-valued "signed trader ID" implies that each collusive trader transact at both sides (buy and sell) of the market and the zero "signed volume" indicate the total amounts of the transactions in both sides are zero: no equity is really bought or sold. Therefore, the unchanged position, represented through zero-valued "signed trader ID" and "signed volume", indicate the wash trade activities in certain collusion.

*B.  Wash Trade among multiple traders*

As the FCA and CESR pointed out in their consultation reports (FSA, 2006) (EU, Market Abuse Directive, 2014), it is difficult to distinguish a wash trade because the format of trading collusions varies and the collusive transactions can be buried in mass numbers of normal trading activities, such as the complex network reported by Nanex on 31 May 2013 (NANEX, Chicago PMI, 2013), where vertices illustrate traders and directional connections among vertices represent the transaction between traders. We utilise this idea in (NANEX, Chicago PMI, 2013) and represent submitted limit orders (from a number of traders) by a graph, where vertices represent traders, the short arrows affixed to the vertex represent the orders submitted by the trader (buying and selling orders are represented by arrows pointing inward and outward, respectively) and the dotted arrow lines represent the possible executed orders according to the matching rule discussed in Section II.A. An example of wash trade action mixed up with legitimate trading orders is shown in Table 5 and illustrated by the graph in Fig. 1. Among the 14 orders submitted by six traders in this example, four pairs (#1-#4 in Table 5) of wash trade orders are deliberately submitted by four traders with tight

submission intervals, executable prices, and mostly matched volumes so that orders in each pair are suspiciously easy to match and execute. In Fig. 1, the possible executions of the orders are illustrated by four dotted arrow lines: each dotted arrow line connecting one pair of matched orders and the arrowhead indicating the transaction direction of the financial equity, i.e. *A* pointing to *B* means trader *A* sells shares of equity to trader *B*. From the illustration in Fig. 1, when participating wash trade activities, traders (*A*, *B*, *C*, and *D*) connect as a closed simple cycle (dotted arrow lines) and continuous transactions among the traders flow throughout the cycle in one single direction (either "clockwise" or "counter clockwise") with each trader along the pathway "passing the parcel" (Aitken, Harris, & Ji, Nov. 2009). After a complete transaction loop, the beneficial interest has been transferred across the collusive group and no traders in the group have an actual position change.

Table 5 Example of wash trade in a sequence of limit orders from a number of traders

| # | Trader | Time | Buy/Sell | Price | Volume | Pairs |
|---|--------|------|----------|-------|--------|-------|
| 01 | *A* | 9:00:000 | Sell | 125.00 | 1450 | # 1 |
| 02 | *B* | 9:00:001 | Buy | 125.01 | 1500 | |
| 03 | *B* | 9:05:000 | Sell | 124.95 | 1500 | # 2 |
| 04 | *C* | 9:05:001 | Buy | 125.01 | 1450 | |
| 05 | *E* | 9:16:000 | Sell | 124.90 | 200 | |
| 06 | *C* | 9:20:000 | Buy | 124.90 | 235 | |
| 07 | *C* | 9:30:001 | Sell | 125.00 | 1450 | # 3 |
| 08 | *D* | 9:30:002 | Buy | 125.01 | 1500 | |
| 09 | *C* | 9:45:000 | Sell | 124.80 | 250 | |
| 10 | *F* | 10:05:000 | Buy | 124.70 | 350 | |
| 11 | *D* | 10:50:000 | Sell | 125.01 | 1450 | # 4 |
| 12 | *A* | 10:50:001 | Buy | 125.01 | 1450 | |
| 13 | *F* | 11:35:000 | Sell | 124.80 | 200 | |
| 14 | *E* | 11:50:000 | Buy | 124.50 | 550 | |



*A, B, C, D, E* and *F*: traders; — ➤ : possible execution;
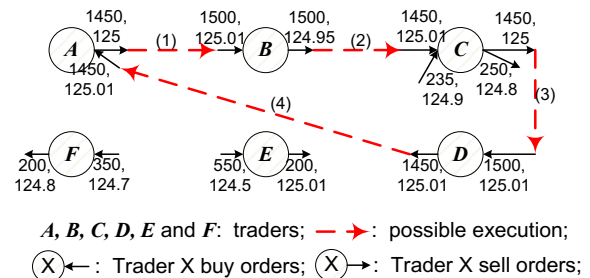Ⓧ◄— : Trader X buy orders; Ⓧ—➤ : Trader X sell orders;

Fig. 1 Closed connection cycle of traders and the possible execution flow along the cycle in wash trade action (14 orders in Table 5 are mapped to the graph)

The "no beneficial interest change" of all collusive traders in wash trade activities can also be calculated by the terminologies defined in Section III.A as the equation (2).

Equation (2) shows the possible execution (the dotted arrow (1) in Fig. 1) of two orders in pair #1 in Table 5 due to the matching rule, execution occurring on earliest orders with matched prices, discussed in Section II.A. Similarly, the executions of matched pairs #2-4 in Table 5 (and the dotted arrows (2)-(4) in Fig. 1) are represented by equation (2). The aggregated results of those executions are calculated in

equation (2), where 50 shares of volumes are remained due to the "mostly matched volumes" tactic between any two "smart manipulator neighbours" to avoid regulatory inspections (Aitken, Harris, & Ji, Nov. 2009). The unmatched volumes (for example 2%) can then be defined as the *matching margin* $(\delta_v)$. Similarly, the differences between the limit order prices and transaction prices can be defined as limit price margin $(\delta_p)$ and the transaction margin $(\delta_p^T)$ respectively. In the following case, $\delta_p^T = 0.005$.

$$
\begin{aligned}
\text{Position} = \{ \quad & (-A, 125.00, -1450), (B, 125.01, 1500), \\
& (-B, 124.95, -1500), (C, 125.01, 1450), \\
& (-C, 125.00, -1450), (D, 125.01, 1500), \\
& (-D, 125.01, -1450), (A, 125.01, 1450) \quad \} \\
= \{ \quad & (-A+B, 125.00+0.01, +50), \\
& (-B+C, 124.95+0.06, -50), \quad\quad (2) \\
& (-C+D, 125.00+0.01, +50), \\
& (-D+A, 125.01+0, 0) \quad \} \\
= \{ \quad & \text{`-}A+B\text{-}B+C\text{-}C+D\text{-}D+A\text{'}, \\
& 124.95+0.06, 50\text{-}50+50+0 \quad \} \\
= \{ \quad & \text{`0'}, 125.005\pm0.005, +50 \quad \}
\end{aligned}
$$

Furthermore, as shown in Table 5, the time intervals between different pairs can vary as random events occurred in one single trading day. To avoid being detected as suspiciously trading action, in practice, "smart manipulators" tactically place the pairs at separated time points as the examples in Table 5, where the time differences among any two pairs are completely different and random. To achieve this, manipulators carefully design each pair of matched orders to minimise the possible financial loss from price changes in the time period (i.e. from 9:00 to 10:50 in Table 5) and to maintain the positions of their whole collusive group at zero. The separated arrangement of the matched pairs increases the complexity of detecting a wash trade under a mixture environment of both "normal" and "manipulative" trades.

Additional to the example in Table 5 and Fig. 1, the matched pairs among any two manipulators can also be constructed by a number of limit orders as illustrated in Table 3, rather than simply matched one-to-one sell and buy orders (as the pairs in Table 5). For example, the matched pair #1 in Table 5 can be constituted by four selling orders and one buying orders as shown in Table 6 and illustrated in Fig. 2.

Table 6 Example of matched pair composed of multiple orders in wash trade activity

| # | Trader | Time | Buy/Sell | Price | Volume | Pairs |
|---|--------|------|----------|-------|--------|-------|
| 01 | A | 9:00:000 | Sell | 124.99 | 450 | |
| 02 | A | 9:00:000 | Sell | 124.98 | 450 | |
| 03 | A | 9:00:000 | Sell | 124.97 | 350 | #1 |
| 04 | A | 9:00:000 | Sell | 124.96 | 200 | |
| 05 | B | 9:00:001 | Buy | 125.01 | 1500 | |

In the examples, the submission of four sell orders is followed tightly by one large buy order, which matches, potentially executes and removes all (or most) volumes of previous four sell orders. The graph of the traders and the transaction flow are revised in Fig. 2, where the #1 matched pair between *A* and *B* is illustrated by four short outward arrows affixed to *A* connecting with one short inward arrows affixed to *B* through the dotted arrow and other parts of the structure of

the whole closed cycle of the traders is remained. In the example in Table 6, since the buy order #05 is submitted later than the sell orders, it will be executed at the prices of four sell orders, i.e., order #05 will be firstly executed as 450 shares at 124.99 with order #01 and then another 450 shares executed at 124.98 with order #02 and so on.

### C. Wash trade features

From the discussion in Section III.A and III.B, the strategy that constructs a wash trade activity has the following two key features:

Feature 1: Matched orders - as the first step of wash trade manipulation, traders deliberately submit the matched orders to the market in tiny time intervals to guarantee the execution; those orders can be one-to-one (examples in Table 5) or one-to-many matched (example in Table 6); this feature refers to dotted arrow lines in Fig. 1 and Fig. 2;

Feature 2: Closed transaction cycle - any single execution of the matched orders does not refer to wash trade manipulation unless those executions constitute a closed cycle as illustrated in the examples shown in Fig. 1 and Fig. 2; this feature refers to closed cycle of dotted arrows among the traders in Fig. 1 and Fig. 2.
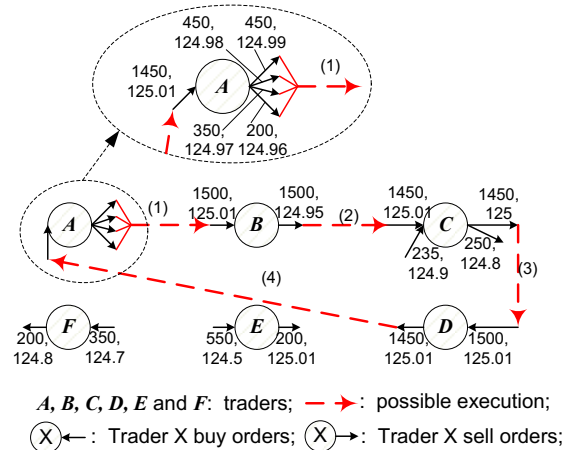


*A, B, C, D, E* and *F*: traders; — ▸ : possible execution;
(X)◂— : Trader X buy orders; (X)—▸: Trader X sell orders;

Fig. 2 Multiple matched orders between two manipulators in wash trade action (14 orders in Table 5 and five orders in Table 6 are mapped into the graph)

Considering the example in Table 5, the manipulators set up the matched orders from the #01 order at time 9:00:000, but the wash trade is not completely constructed until the submission of the #12 order at time 10:50:001, which closes the transaction cycle. Therefore, wash trade can be detected through detecting the matched orders and closed cycle in two steps:

Step 1: Detect the suspiciously matched order pairs $S$ according to the matching rule and wash trade features, tight submission intervals, executable prices and mostly matched volumes:

$$
\text{Order Pair} = \{\textstyle\sum Orders\} = \{+T_m - T_n, P^T \pm \delta_p^T, \pm\delta_v\},
$$

where $+T_m - T_n$ represents trader $T_n$ selling shares of equity to $T_m$ and $\delta_v$ and $\delta_p^T$ represent the *matching margin* of volume and transaction price $P^T$;

Step 2: Among $S$, find the order pairs whose transaction price margins are overlapped, in those pairs, if some pairs fulfil the condition:

$$\text{Position} = \{\sum_{k \in S} \text{Order Pair}_k\} = \{"0", P^T \pm \delta_p^T, \pm\delta_v\},$$

a wash trade alert is triggered.

To further formulate those features, we define the #k order $L$ submitted by trader $T_n$ at time $t_k$ as:

$$L_k = (t_k, \pm T_n, P_k, \pm V_k)$$

where $P_k$ and $V_k$ are #k order price and volume respectively and the positive and negative sign $\pm$ represent buy and sell operation. The *matching margin* $\delta$ is defined as a vector $\delta = [\delta_p, \delta_t, \delta_v]$ with three small positive values for price, time and volume respectively. If buy order #**K** is matched with **K-1** sell orders from #**1** to #**K-1**, their features have 1) tiny time interval:

$$|t_1 - t_K| < \delta_t, \tag{3}$$

2) executable tiny price difference:

$$P_K - \min(P_1, \dots, P_{K-1}) < \delta_p, \tag{4}$$

3) and mostly matched volume:

$$|\sum_{k=1}^{K-1} V_k - V_K| < \delta_v. \tag{5}$$

If **K** orders among $N$ traders construct wash trade action, their features meet the following condition, where $r_{nk}$ is the indicator that if order #k from trader $T_n$ is a sell order, then $r_{nk} = -1$, and $r_{nk} = +1$ for buy order.

$$\begin{aligned}\text{Position} &= \{\sum_{n=1}^N r_{nk} T_n, P^T \pm \delta_p^T, \pm\delta_v\} \\ &= \{"0", P^T \pm \delta_p^T, \pm\delta_v\}\end{aligned} \tag{6}$$

The features in equations (2)-(5) are detected in Step 1 and the feature in equation (6) is detected in Step 2.

*D. Problem formulation*

To discover the wash trade before it completely occurs (fulfilling the recent regulations on preventing the attempts of wash trade), the detection approach is applied to the limit order streams instead of the trade records. The order stream is the sequence of limit orders received by the trading platform from numerous traders. The stream is updated by the "order event", which could be submission, modification, cancellation or execution. As shown in Table 1, an order includes ID, trader ID, time, buy/sell sign, price and volume. In this study, we assume that the orders in the stream are on one specific stock. Thus the stock information in the stream can be ignored once the specific stock is determined. This assumption, on one hand, narrows the scope of this study specifically on the underlying problem, on the other hand conforms the practical trading platform environment, where the algorithm can be easily applied to selected equity.

The Step 1, detecting the suspiciously matched order pairs according to equations (3)-(5), is termed as coarse detection while Step 2, recognising the closed cycle based on equation (6), is termed as fine detection. The limit order stream is then required to be pre-organised to commence with those two tasks. A physical time sliding window sized $\theta_T$ is specified and the trading order stream can be split into two queues of consecutive orders: *buy order queue, $Q_b$* and *sell order queue, $Q_s$* each of which maintains a size $\theta_T$. That is, if a new order $L_k$ is a buy order, push it into $Q_b$; otherwise push it into $Q_s$. If the length of the updated queue is larger than $\theta_T$, pop the earliest orders to maintain the length of the sliding window. The algorithm is described in Algorithm 1. Since the order stream is measured in "order event time", $\theta_T$ is maintained by calculating the difference between the physical time stamps of the first and the last orders in the queue. Hence the number of orders in each queue ultimately depends on the underlying frequency of order activities and differs across time. (Algorithm 1 is named as WASH_TRADE_DETECT because it will involves all detection sub-functions, which are discussed in follow up sections.)

The intention of the wash trade, increasing transaction volume, indicates that wash trades are usually associated with large-sized orders. Consequently, the orders with volumes smaller than a predefined threshold $\theta_v$ are ignored, where the threshold can be setup according to the requirements of the detection solidness. Given the limit order queues $Q_b$ and $Q_s$, the *coarse detection* can then be formulated as follows: for a large incoming order, examine in the opposite order queue for one or multiple potential matching orders which are characterised by equation (3)-(5). The result of the *coarse detection* comprises all order combinations matched with the incoming order. Collusions may exist among those combinations.

Algorithm 1. Wash Trade Detection – Pre-organisation

```
WASH_TRADE_DETECT(L_k )
1   Q_s = ∅; Q_b = ∅
2     while L_k is a valid limit order
3        if L_k is buy
4           Push L_k into Q_b
5           while Q_b length > θ_T
6              Pop Q_{b,1} to maintain θ_T
7        else
8           Push L_k into Q_s
9           while Q_s length > θ_T
10             Pop Q_{s,1} to maintain θ_T
```

Similarly, the *fine detection* can be formulated as follows: given the matched order pairs, find certain sets of pairs in which the sum of "signed trader ID" and "signed volume" have zero values as the illustrations in equation (6). Defining *coarse detection* and *fine detection* as the function **COARSE_DETECT** and **FINE_DETECT** respectively, the wash trade detection is further designed in following section.

*E. Coarse Detection - Matching Search*

The matching relationship of wash trade order pairs is summarised in equations (3)-(5). In the *Coarse Detection* process, three conditions are sequentially checked to identify the potential matching.

The time matching margin $\delta_t$ in equation (3) shows the tiny interval between the orders in a pair. Setting the length of the order queue $\theta_T$ in Algorithm 2 and Algorithm 1 equivalent to $\delta_t$, the *coarse detection* is designed as the illustration in Fig. 3: given the incoming order $L_k$, examining the opposite orders in previous $\delta_t$ ($\theta_T$) period for potential matched orders, which are determined by price and volume margin, $\delta_p$ and $\delta_v$. Algorithm 1 is then revised as Algorithm 2 which includes both the **COARSE_DETECT** and **FINE_DETECT** functions, where

the $\{MP\}$ is the detected matched pairs of **COARSE_DETECT**.

In financial markets, only the orders following executable price rules (Bowen, 2013) match and execute. Therefore the price margin $\delta_p$ in equation (4) is constrained by the following rules:

Rule 1. Sell order matches buy orders with equal or higher prices;

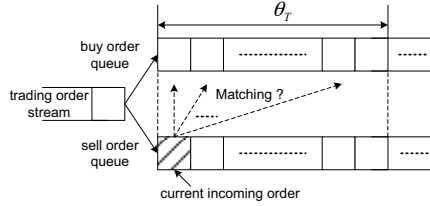Rule 2. Buy order matches sell orders with equal or lower prices.



Fig. 3 Coarse Detection Scheme

The example in Table 6, where the #5 buy order price is slightly higher than all previous sell orders, shows the Rule 2 of price margin $\delta_p$. Considering the price margin $\delta_p$, the *coarse detection* is designed as follows: given the incoming buy (sell) order $L_k$, among all executable orders (in terms of the executable limit prices) in the previous $\delta_t$ ($\theta_T$) period, find the order pairs having the best matching volumes.

Algorithm 2. Wash Trade Detection Algorithm

| WASH_TRADE_DETECT($L_k$) |
|---|
| 1  $Q_s = \emptyset$; $Q_b = \emptyset$; {Matched Pairs}= $\emptyset$; |
| 2  while $L_k$ is a valid limit order |
| 3      if $L_k$ is buy |
| 4          Push $L_k$ into $Q_b$ |
| 5          while $Q_b$ length $> \theta_T$ |
| 6              Pop $Q_{b,1}$; |
| 7          {$MP$}= **COARSE_DETECT**($Q_s$, $L_k$); |
| 8      else |
| 9          Push $L_k$ into $Q_s$ |
| 10          while $Q_s$ length $> \theta_T$ |
| 11              Pop $Q_{s,1}$; |
| 12          { $MP$ }= **COARSE_DETECT**($Q_b$, $L_k$); |
| 13      if { $MP$ } $\neq \emptyset$ |
| 14          **FINE_DETECT**({$MP$}); |

The volume matching can be defined as a function **VOL_MATCH**($Q^{t,p}, L_k$), where $Q^{t,p}$ is a set of orders after being filtered by $\delta_t$ and $\delta_p$. Given this, **COARSE_DETECT**($Q$, $L_k$) is defined in Algorithm 3, where $Q$ contains all opposite orders in the previous $\delta_t$ periods and $L_k$ is the incoming order. Based on the above discussions and the constraints in equation (5), the function **VOL_MATCH**($Q^{t,p}, L_k$) is defined as follows: given incoming order $L_k$ and a set of matched orders $Q^{t,p}$, find subsets $S$ of the order pairs from $Q^{t,p}$ such that

$$|(\textstyle\sum_{i \in S} V_i) - V_k| \leq \delta_v.$$

The number of limit orders in subset $S$ is $n_s$ ($n_s$ is smaller than the size of $Q^{t,p}$). In essence, the problem of **VOL_MATCH** is a practical case of a more general problem called the *Knapsack Problem* (Rumen, Vincent, & Sanjay, 2000) (Vincent, Yanev, & Andonov, 2009) (Zukerman, Jia, Neame, & Woeginger, 2001). The name *Knapsack* refers to the problem of filling a knapsack of capacity $W$ using a subset of $m$ items $\{1, \dots, m\}$, each of which has a mass and a value, such as the total weight of the selected items is less than or equal W and their total value is maximised. The volume matching problem can be viewed as a simplified form of the *Knapsack Problem*: given a capacity $V_k$ (the knapsack size) and a set $Q^{t,p}$ of items, each having non-negative size $V_i$, find all possible subsets $S$ of items to eventually make

$$|(\textstyle\sum_{i \in S} V_i) - V_k| \leq \delta_v.$$

Due to the similarity of the two problems, the widely used approach solving the *Knapsack Problem,* dynamic programming, is employed in **VOL_MATCH**($Q^{t,p}, L_k$). The main principles of dynamic programming are that we have to come up with a number of sub-problems so that each sub-problem can be solved easily from "smaller" sub-problems, and the solution of the original problem can be obtained easily once we know the solutions to all the sub-problems (Kleinberg & Tardos, 2005). Dynamic programming has been studied thoroughly in optimization problems in (Zhen, He, Wen, & Xu, 2013) (Jiang & Jiang, 2014).

Algorithm 3. Coarse Detection

| COARSE_DETECT ($Q$, $L_k$) |
|---|
| 1      $Q^{t,p} = \emptyset$; |
| 2      if $L_k$ is a valid buy order |
| 3          for each order $L_i$ in Q |
| 4              if $P_i <= P_k$ |
| 5                  push $L_i$ into $Q^{t,p}$ |
| 6      else if $L_k$ is a sell order |
| 7          for each order $L_i$ in Q |
| 8              if $P_i >= P_k$ |
| 9                  push $L_i$ into $Q^{t,p}$ |
| 10      if $Q^{t,p} \neq \emptyset$ |
| 11          $S$ = **VOL_MATCH**($Q^{t,p}, L_k$) |

To solve the special form of the *Knapsack Problem* under $N$ limit orders and volume $V_k$, denoting the final subset of orders in an optimum solution for the original problem as $S_N$, we then use the notation **OPT**($N, V_k$) to denote the sum of the order volumes of the first $N$ orders in the subset $S$ under the constraint $|\textbf{OPT}(N, V_k) - V_k| \leq \delta_v$. The sum in the first *N-1*, *N-2,…,1* orders can then be represented as **OPT**($N-1, V_k$), **OPT**($N-2, V_k$), …, **OPT**($1, V_k$). To determine **OPT**($N, V_k$), we not only need the solution of **OPT**($N-1, V_k$), but also need to know **OPT**($N-1, V_k - V_N$), the best solution for the first *N-1* orders with the remaining capacity $V_k - V_N$, which constructs the constraint as $|\textbf{OPT}(N-1, V_k) - (V_k - V_N)| \leq \delta_v$. The recursion can then be summarised as follows: if $L_N$ is not one of the orders in the final subset $S_N$, we can ignore the order $N$ and determine **OPT**($N-1, V_k$); however if $L_N$ is one of the orders, we need to seek an optimal solution for the remaining orders, *1,.., N-1*, which is **OPT**($N-1, V_k - V_N$). Using this set of sub-problems, we are able to express the **OPT**($N, V_k$) as a simple expression in terms of values from "smaller" problems. Therefore, the recursion is summarised as two conditions:

1. If $L_i \notin S_N$, then OPT($N, V_k$)= OPT($N-1, V_k$);

2. If $L_i \in S_N$, then OPT$(N, V_k) = V_N +$ OPT$(N - 1, V_k - V_k)$;

This recursive process is re-organised based on the above two conditions to give Algorithm 4. This recursive algorithm can be used by invoking OPT $(N, V_k)$ for $N$ limit orders and the capacity $V_k$.

### F. Fine detection – collusion search

$S_N$, orders from $Q$ matched with the incoming order $L_k$, is the result of the *coarse detection*. To further detect the potential closed cycle of transactions, the orders in $S_N$ are represented by equation (1), where the trader ID and volumes are affixed with trading direction signs. After the conversion, $S_N$ is defined as $S_N^c$, the input of the fine detection algorithm FINE_DETECT. As discussed in Section III.A and equation (2), the order pairs with potential transaction prices with overlapped price margins are grouped together for potential collusion detection.

Detecting trader collusion is treated as discovering the combinations $C$ from $S_N$ such that the sum of the signed trader equals zero as illustrated in equation (6). This process can be considered equivalent to a special case of the previously defined volume matching problem: given a capacity W=0 (the knapsack size) and a set of signed trader pairs, each having a value (e.g. {+A} and {−A}), select all possible subsets C of signed trader pairs to make $\sum_{n \in C} r_{nk} T_n = W$, where the symbolic computation of the signed trader pairs are defined in Section A and can be implemented by operator overloading. The subset $C$ is considered as trading collusion in a wash trade.

Algorithm 5, derived from Algorithm 4, provides the recursive solution for FINE_DETECT($S_N^c$).

**Algorithm 4 Volume Matching Detection by recursion**

| 1 | VOL_MATCH($Q^{t,p}, L_k$) | // original limit order set $Q^{t,p}$; |
|---|---|---|
| 2 | $S_N = \emptyset$; | // solution subset, initialized to empty; |
| 3 | $N$ = length($Q^{t,p}$); | // $N$: size of $Q^{t,p}$; |
| 4 | OPT($N, L_k$) | // $N$ decreases on each recursion step; |
| 5 | if $N < 1$ or $|L_k| \leq \delta_v$ | // if $N$ reaches the last one or $\delta_v$ |
| 6 | return; | condition is satisfied; |
| 7 | if $|V_N - L_k| \leq \delta_v$ | // if condition is satisfied, then orders |
| 8 | output $S_N$; | in $S_N$ is one solution; |
| 9 | push $L_N$ into $S_N$ | // assume $L_N \in S_N$; |
| 10 | OPT$(N - 1, V_N - L_{N,v})$; | // recursively find solution by condition 2; |
| 11 | Discard $L_N$ from $S_N$ | // assume $L_N \notin S_N$; |
| 12 | OPT$(N - 1, L_v^n)$; | // recursively find solution by condition 1; |
| 13 | end of OPT | |
| 14 | return $S_N$; | |

**Algorithm 5 Collusion Search by recursion**

| 1 | FINE_DETECT ($S_N^c$); | // original signed trader set $S_N^c$ |
|---|---|---|
| 2 | $\vec{C} = \emptyset$; | // solution subset, initialized to empty; |
| 3 | $N$ = length($S_N^c$); sum = 0 ; | // $N$: size of $S_N^c$; |
| 4 | OPT($N$, sum) | // $N$ decreases on each recursion step; |
| 5 | if $N < 1$ | // if N reaches the last one, done |
| 6 | return; | |
| 7 | if rT + sum = 0 | // if the sum of signed trader including current one is zero |
| 8 | output $\vec{C}$; | // the signed trader in $S_N^c$ is a solution; |
| 9 | push rT into $\vec{C}$; | // assume $\vec{t_N} \in \vec{C}$; |
| 10 | OPT$(N - 1$, sum + rT); | // recursively find solution by |
| | | condition 2; |
| 11 | Discard rT from $\vec{C}$; | // assume $\vec{t_N} \notin \vec{C}$; |
| 12 | OPT($N - 1$, sum); | // recursively find solution by condition 1; |
| 13 | end of OPT | |
| 14 | return $\vec{C}$; | |

## IV. EXPERIMENTS AND EVALUATION

Evaluating a detection model usually relies on real data of both "normal" and "abuse" cases. However, due to the limited reports on wash trade manipulation and regulatory rules prohibiting the disclosure of illegitimate market data, the availability of the examples of wash trade behaviours in capital markets is far less than the availability of routine normal trading records. Therefore to evaluate the proposed detection model, it is acceptable to the financial industry that all the characteristic patterns of wash trade examples are reproduced then injected into original trading records to generate a mixed dataset of normal and abuse cases (NANEX, Exploratory Trading in the eMini, 2013). Randomly synthesised exploratory manipulation cases can mimic any possibility of wash trade scenarios, i.e. we can generate the matched order at any time with any volume size as well as matching margins. Synthetic exploratory financial data are also accepted in academia for evaluating the proposed model when real market data are hard to collect (Palshikar & Apte, 2008) (Ou, Cao, Luo, & Zhang, Dec. 2008) (Markus., Hoser, & Schröder, 2007). In this paper, the experimental evaluation is composed of two parts:

Part 1: experimental evaluation using original trading datasets from the market;

Part 2: experimental evaluation using original trading datasets injected with synthetically generated wash trade scenarios following the analysis in Section II.A.

### A. Experiment Setup

The experimental data used in this work involve real market data (trading orders) of seven stocks Google (GOOG), Microsoft (MSFT) and Apple (AAPL) from NASDAQ and First Quantum Minerals (FQM), Yamana Gold (YAU), Gazprom (OGZD), and Vodafone (VOD) from London Stock Exchange (LSE). The selection of these datasets is due to their active trading activities, relatively high trading volumes and more volatile price fluctuation, the factors that might increase the likelihood of market abuse across the exchanges (Douglas, Feng, & Aitken, 2012) (Lee, Eom, & Park, 2013). The datasets from NASDAQ cover messages over five trading days from $11^{th} - 15^{th}$ June 2012, and consist of more than 400,000 trading orders in total for each stock. The datasets from LSE cover $23^{rd} - 27^{th}$ May 2011, and consist of more than 100,000 orders in total for each stock. Table 1 shows an excerpt of the trading records used in this study. The wash trade detection algorithms are evaluated on the original seven datasets for detecting any transactions which are suspiciously similar to wash trade manipulation. Additionally, typical wash trade activities are reproduced according to the discussions and examples in Table 5 Table 6 and injected into those seven datasets for further experimental evaluations.

## B. Determining the marginal parameters

As discussed in Section II.A, the submissions of the matched orders in a wash trade are usually within tiny time intervals $\delta_t$ so that the manipulated execution can compete against the action of normal traders who may pick the orders unintentionally (FSA, 2006) (Bowen, 2013). Consequently, the normal execution time shows a reasonable reference to the time interval $\delta_t$, which otherwise is not available because of the lack of the statistical studies of the real wash trade cases.

Usually, the execution time of a limit order is strongly associated with its volume (Douglas, Feng, & Aitken, 2012) (Hautsch & Huang, The market impact of a limit order, 2012) (Hautsch & Huang, Limit Order Flow, Market Impact and Optimal Order Sizes: Evidence from NASDAQ TotalView-ITCH Data, 2011). Therefore, a more reasonable measure of the average execution time of normal limit orders can be given by volume-weighted average execution time (VWAT), defined as

$$T_{VWAT} = \frac{\sum_j (T_j * v_j)}{\sum_j v_j}, \tag{7}$$

where $T_{VWAT}$ is volume-weighted average execution time; $T_j$ is the execution time of order $j$; $v_j$ is the volume of order $j$; and $j$ is each individual order (Hautsch & Huang, Limit Order Flow, Market Impact and Optimal Order Sizes: Evidence from NASDAQ TotalView-ITCH Data, 2011). In practice, if the wash trade orders are submitted with time intervals larger than $T_{VWAT}$, they are apparently easy to pick by other legitimate traders. Accordingly, by setting $\delta_t = T_{VWAT}$, this approach covers a time period for all possible wash trade activities. The order execution time $T_j$ and the $T_{VWAT}$ across the seven stocks in the test dataset are calculated and summarized in Table 7.

Theoretically, the wash trade can be carried out by a large number of small orders. However, in practice, the wash trade orders are usually larger than the average volume of the normal trading orders because a large number of orders can significantly increase the uncertainty of the order executions, which may bring a risk of loss if it does not follow the expected arrangements. Therefore, the average order volume of each stock is selected as the threshold $\theta_V$ for the order volume filtering discussed in Section D. The average volume across seven stocks is also calculated and summarized in Table 7.

Table 7 Volume weighted average execution Time and average volume

| | $T_j$ (sec.) | $T_{VWAT}$ (sec.) | Avg. Vol (share) |
|---|---|---|---|
| GOOG | 2.77 | 118.79 | 635.57 |
| MSFT | 3.07 | 107.68 | 530.70 |
| AAPL | 5.92 | 87.04 | 900.04 |
| FQM | 10.19 | 83.87 | 163.20 |
| YAU | 14.35 | 104.25 | 878.46 |
| OGZD | 6.04 | 52.35 | 796.30 |
| VOD | 12.97 | 71.15 | 661.16 |

In addition, the volume matching margin $\delta_v$ is selected as percentages: 0%, 1%, 2%, 3%, 4% and 5% indicating the ratio of not matching (1% refers identifying orders with 99% matching volumes). In the example in Table 6, the #5 buy order volume (1500 shares) is around 96.7% matched with all previous sell orders (1450 shares). The price margin $\delta_p$ is unconstrained in the detection so that any orders following the price matching rules Rule 1 and Rule 2 are scanned for possible matching pairs under the condition in equation (5).

Under the configurations of $\delta_t$, $\theta_V$, $\delta_v$, and $\delta_p$, Algorithm 4 reflects the fact that given an order $L_k$, among all executable priced orders (unconstrained $\delta_p$ but following Rule 1 and 2) with volume not smaller than $\theta_V$ in a previous $\delta_t$ time period, find the matched orders that executed at least $(1-\delta_v)\%$ volumes of $L_k$.

## C. Part 1: Experiments on original datasets

In Part 1 experiment, the wash trade detection algorithm is evaluated on the original seven datasets using the parameters in Section IV.B. The evaluation shows the applicability of the proposed algorithms to real transaction data and also examines the legitimacy of the transactions in original dataset. Since the original datasets do not contain any reported wash trade manipulation activities, it is assumed to only contain legitimate transactions. Thus the evaluation measure is based on *false negative rate*, $FNR = \frac{FN}{FN+TP}$, which is based on *false negative*, *FN*, defined as normal cases detected as a wash trade, and *true positive*, *TP*, defined as normal cases detected as normal.

The results of the experiments (max FNR values on each stock dataset are highlighted) are shown in Table 8. It is clear that in each dataset, some transactions are detected as suspicious wash trade actions and the numbers of the detected actions increase across the increases of volume margins. Most of the datasets do not contain any suspicious actions when the volume margin is set to 0% and the Apple stock shows the highest FNR rate (1.263%) at the 5% volume margin.

With careful inspection and consultation with the financial industry experts, we determined that the detected false negative cases show very similar features to the wash trade actions although not reported by the regulators. The detected false negative cases fall into two formats as shown in Table 9.

Table 8 Experiment results (FNR) across original datasets of seven stocks

| stock | Volume Margins | | | | | |
|---|---|---|---|---|---|---|
| | 0% | 1% | 2% | 3% | 4% | 5% |
| GOOG | 0.000% | 0.046% | 0.059% | 0.073% | 0.093% | **0.096%** |
| MSFT | 0.000% | 0.030% | 0.176% | 0.275% | 0.519% | **0.530%** |
| AAPL | 0.000% | 0.000% | 0.166% | 0.576% | 1.153% | **1.263%** |
| FQM | 0.389% | 0.499% | 0.526% | 0.553% | 0.673% | **0.926%** |
| YAU | 0.000% | 0.669% | 0.761% | 0.780% | 0.853% | **1.186%** |
| OGZD | 0.000% | 0.346% | 0.519% | 0.680% | 0.693% | **0.853%** |
| VOD | 0.000% | 0.953% | 1.048% | 1.066% | 1.143% | **1.219%** |

In case #1, the trader Client12 sold 6600 shares to Client3 at price 58.0 and bought 6606 back 2 seconds later at the same price. The 99.9% matched transacted volumes, 100% matched prices and the closed cycle of the transaction directions between Client12 and Client3 make this case extremely suspicious and potentially be a wash trade action according to the regulation (SEC, 2011) although not reported yet. Detecting such suspicious activities shows effectiveness of the proposed

algorithms, although recognising the real intention behind such cases requires more inspections from the regulators, which is out of the scope of our work. In case #2, Client1 sold 15000 shares to Client5 at price 58.56 at market closing time and bought them all back at slightly higher price at market opening time on the next day. Those transactions also fulfill the conditions of a wash trade action except the trading dates. According to the suggestions from financial experts, case #2 refers to pre-arranged trading, which is defined as "a sell is coupled with a buy back at the same or pre-arranged price that limits the risks" (SEC, 2011). The only difference between pre-arranged trading and wash trade is that the former is usually among merely two parties and may occur in different days and the latter can involve a number of collusive traders and usually occurs as intra-day trading. When only targeting wash trade, the proposed algorithms can be applied on intra-day transactions to avoid picking up the pre-arranged trading as case #2, although the pre-arranged trading is also illegal (SEC, 2011) and needs to be monitored and banned from the capital markets.

Table 9 False Negative cases of stock AAPL

| case# | time | volume | price | seller | buyer |
|---|---|---|---|---|---|
| 1 | 21/06/2012 15:18:48.768 | 6600 | 58.00 | Client12 | Client3 |
| | 21/06/2012 15:20:11.811 | 6606 | 58.00 | Client3 | Client12 |
| 2 | 21/06/2012 16:28:40.629 | 15000 | 58.56 | Client1 | Client5 |
| | 22/06/2012 10:00:45.187 | 15000 | 58.60 | Client5 | Client1 |

*D. Part 2: Experiments on datasets with injected wash trade*

Testing with synthetic data can mimic any possible wash trade cases and also can evaluate the robustness of the proposed algorithms under any wash trade scenarios, i.e. random combinations of one or multiple traders wash trade activities.

*1) Wash Trade Case Generation*

The typical wash trade activities are reproduced and injected in each stock dataset. The activities are reproduced in two format groups:

Group 1: one order matched with single opposite order, termed as "single-matching";

Group 2: one order matched with multiple opposite orders, termed as "multi-matching".

Each group contains three different sets according to trader numbers in the wash trade collusion: set #1 has examples with one trader in a trading collusion; set #2 and #3 has two and four traders in a trading collusion respectively. To ensure a comprehensive assessment of the approach, in each set, volume matching margin $\delta_v$ is selected as a percentage 0%, 1%, 2%, 3%, 4% and 5% indicating the ratio of not matching (1% indicating the orders from two sides are 99% matching). There are 10 examples for each combination of the above parameters.

Table 10 Generated single matched
Wash Trade cases ($\delta_v$=5%)

| Case | Trader | Time | Buy/Sell | Price | Volume | pairs |
|---|---|---|---|---|---|---|
| #1 | A | 9:00:000 | Sell | 58.00 | 5000 | 1 |
| | B | 9:00:001 | Buy | 58.01 | 4750 | |
| | A | 9:15:000 | Buy | 58.01 | 5000 | 2 |
| | B | 9:15:001 | Sell | 58.00 | 4750 | |
| #2 | A | 12:16:000 | Sell | 58.00 | 5000 | 1 |
| | B | 12:16:100 | Buy | 58.05 | 4750 | |
| | C | 13:00:001 | Buy | 58.05 | 4750 | 2 |

| | B | 13:00:002 | Sell | 58.00 | 5000 | |
| | C | 13:20:001 | Sell | 58.00 | 5000 | |
| | D | 13:20:002 | Buy | 58.05 | 4750 | 3 |
| | A | 14:20:001 | Buy | 58.05 | 4750 | 4 |
| | D | 14:20:002 | Sell | 58.00 | 5000 | |

Table 11 Generated multiple matched
Wash Trade cases ($\delta_v$=5%)

| Case | Trader | Time | Buy/Sell | Price | Volume | pairs |
|---|---|---|---|---|---|---|
| #3 | A | 9:00:100 | Sell | 58.00 | 1100 | |
| | A | 9:00:100 | Sell | 58.01 | 1200 | |
| | A | 9:00:100 | Sell | 58.02 | 1000 | 1 |
| | A | 9:00:100 | Sell | 58.03 | 1400 | |
| | B | 9:00:101 | Buy | 58.05 | 5000 | |
| | B | 10:10:000 | Sell | 58.00 | 1000 | |
| | B | 10:10:000 | Sell | 58.01 | 1300 | |
| | B | 10:10:000 | Sell | 58.02 | 1200 | 2 |
| | B | 10:10:000 | Sell | 58.03 | 1250 | |
| | A | 10:10:100 | Buy | 58.05 | 5000 | |

The examples in Table 10 and Table 11 show an excerpt of the generated wash trade cases: case #1: two traders with 5% single matched volumes; case #2: four traders with 5% single matched volumes; case #3: two traders with 5% multiple matched volumes. The volume, time and matching margin of the synthetic orders are all randomly generated. For example, in case 3 in Table 11, buy order volume $v_b$ in pair 1 is randomly generated (under condition: $v_b \geq \theta_V$ ) and all sell orders in pair 1 are also randomly generated under the condition that volume sum $V_s$ of all sell orders satisfies: $v_b * (1 - \delta_v) \leq V_s \leq v_b$. The time of orders in pair 2 are also randomly generated as long as they are much later than the time of pair 1. Similarly, the prices of order in each pairs are randomly generated following the price matching rules discussed in Section III.E. Similar to the examples in Table 6, two order pairs in Table 11 have different transaction prices. The buy order in pair #1 in Table 11 will be executed with the previous four sell orders at 58, 58.01, 58.02 and 58.03 respectively. Therefore, the generated examples have different transaction prices within transaction margins.

Such random generation of synthetic cases provides the possibility of thorough evaluation of the proposed algorithms using any possible wash trade cases.

As discussed before, the models are tested on seven real stocks, each of which contains two groups of injected examples. Each group has three sets (1, 2 and 4 traders) and each set contains six margin configurations. Under each configuration, there are 10 examples. There are overall $7\times2\times3\times6\times10 = 2520$ different experiments carried out as a robust evaluation plan for the proposed detection model.

The generated wash trade orders are then injected into the data of corresponding stocks making the test data a mixture of both "normal" and "abuse" patterns. The time intervals between different pairs are selected randomly as examples in Table 10 and Table 11. For example, in case #1 in Table 10, time of pair 2 is randomly selected after the pair 1 occurs. In addition, the generated orders in each pair are separated by several normal orders in original datasets to mimic the practical case in the markets. This is a practical approach to simulate how these wash trade scenarios occur in the real world (Cao,

Ou, & Yu, 2012).

*2)  Performance evaluation metrics*

The performance evaluation of the proposed model is based on two popular statistical measures: sensitivity (SEN) and specificity (SPE). Both of them are based on the confusion matrix, where a false positive (FP) is defined as a wash trade case detected as normal; a true negative (TN) is defined as a wash trade case detected as wash trade, and a false negative (FN) and a true positive (TP) which are defined in Section IV.C. The sensitivity, defined as $SEN = TP/(TP + FN)$ , represents the rate of correctly detecting normal trading orders (a.k.a. the true positive rate) while the specificity, defined as $SPE = TN/(FP + TN)$ , refers to the rate of correctly detecting wash trade cases (a.k.a the true negative rate).

*3)  Experimental Results*

The experimental evaluations across seven stocks are summarized in Fig. 4, where the average SEN and SPE values across different numbers of traders are illustrated against the margin values.

From Fig. 4, the SPE values for single-matching show that the algorithm completely detects the single matching cases, which is the simplest wash trade format and is apparently easy to detect. The SPE values for multi-matching vary across the margins and the different stocks as the illustrations in Fig. 4. The SPE values increase with the increase of the margins and approach 100% when the margin is higher than 5%. The result conforms to the design expectation of the detection approach: more possible collusions will be detected under bigger matching margins. As discussed in Section II.A, "mostly matched" (for example 98%) orders might be built by "smart manipulators" for standing aside from the inspections. A big marginal value compensates this "smart tactic" and the configurability of the margin increases the practicability of the model in a real trading context.
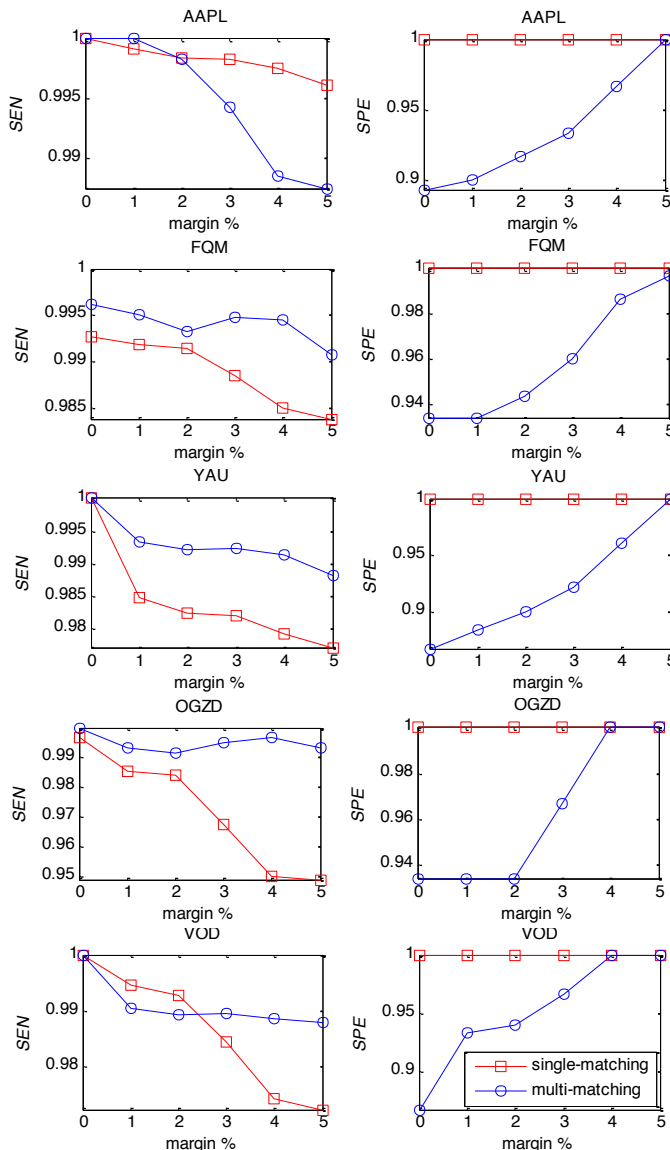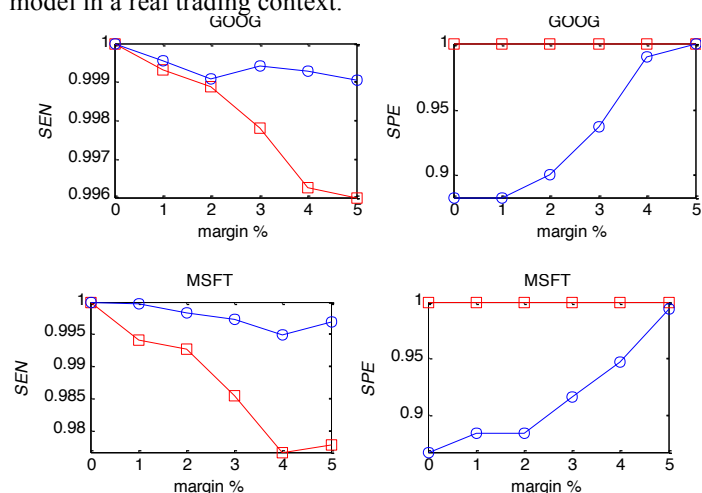


Fig. 4 Experiment results across seven stock dataset.

The SEN values show more volatile results across the margins. In most experiments, the sensitivity values reduce as the margin increases indicating more normal activities incorrectly detected as wash trade cases. On the contrary, the highest SEN value appears at the zero margin value.

From the experimental results, it can be concluded that the proposed approach detects the primary wash trade scenarios effectively and consistently across the selected stocks with SEN values in a range of 97% to 100%.

## V.  CONCLUSION AND FUTURE WORK

A wash trade activity detection approach is proposed after thoroughly studying the various scenarios of wash trade behaviours. The analysis of the collusive activities in wash trades through a graph of traders with transactions represented by the directed connections among the vertexes shows the basic structure of the collusion among multiple traders following a closed cycle of the transactions among certain traders. Further studies also show that the limit orders in wash trades are usually

submitted fast with mutually executed prices and matched volumes. According to the analysed features, the proposed method is then split into steps defined separately in Algorithm 4 and Algorithm 5.

There are two major innovations in the proposed method:

1. Graph theory has been used to represent and model the collusive relationships of the traders in wash trade activities. The concluded fundamental structure of the closed cycle structure within a trader graph simplifies the detection from the complexity of the collusive networks;

2. The wash trade order detection has been approached as a Knapsack problem which can be solved in two steps by traditional dynamic programming approaches.

Instead of only detecting the same-priced buy/sell orders in the "engine level" detection mechanism in CME, the proposed method determines the wash trade activities by considering the suspicious matched orders as well as the collusive groups, which are according to the trading activities in a certain time period rather than a tiny time interval in real-time detection. Therefore, the proposed approach best suits over-night detection in real financial world. However, the rapidly growing trading frequency challenges detection mechanisms and hence implementing the proposed approach in real-time in a computationally efficient way will be the focus of future work.

**Reference**Aitken, M. J., Harris, F. H., & Ji, S. (Nov. 2009). Trade-based manipulation and market efficiency: a cross-market comparison. *Proceedings of the 22nd Australasian Finance and Banking Conference*, (p. Vol 18). Sydney.

Bowen, C. (2013, Jul. 9). *Market Regulation Advisory Notice.* Retrieved from US Commodity and Future Trading Commission: http://www.cftc.gov/stellent/groups/public/@rulesandproducts /documents/ifdocs/rul070913cmecbotnymexcomandkc1.pdf

Cao, L., Ou, Y., & Yu, P. (2012). Coupled Behavior Analysis with Applications. *IEEE Transaction on Knowledge and Data Engeering , 24* (8), 1378-1392.

Cao, Y., Li, Y., Coleman, S., Belatreche, A., & McGinnity, T. (Oct. 2013). A Hidden Markov Model with Abnormal States for Detecting Stock Price Manipulation. *Proceedings of the 2013 IEEE International Conference on Systems, Man, and Cybernetics (SMC)*, (pp. 3014 - 3019). Manchester.

Cao, Y., Li, Y., Coleman, S., Belatreche, A., & McGinnity, T. (Mar. 2014). Detecting price manipulation in the financial market. *Proceedings of the 2014 IEEE Conference on Computational Intelligence for Financial Engineering & Economics (CIFEr)*, (pp. 77 - 84). London.

Cao, Y., Li, Y., Coleman, S., Belatreche, A., & McGinnity, T. (Mar. 2014). Detecting wash trade in the financial market. *Proceedings of the 2014 IEEE Conference on Computational Intelligence for Financial Engineering & Economics (CIFEr)*, (pp. 85 - 91). London.

Cao, Y., Li, Y., Coleman, S., Belatreche, A., & T.M.McGinnity. (2015). Adaptive Hidden Markov Model with Anomaly States for Price Manipulation Detection. *IEEE Transactions on Neural Networks and Learning Systems , 26* (2), 318 - 330.

CME. (2014, Aug.). Retrieved from U.S. Commodity Futures Trading Commision: http://www.cftc.gov/filings/orgrules/rule082814cmedcm001.p df

Douglas, C., Feng, Z., & Aitken, M. J. (2012, Sep. 12). *High Frequency Trading and End-of-Day Manipulation.* Retrieved from SSRN: http://dx.doi.org/10.2139/ssrn.2145565

EU. (2015). Retrieved from European Commission: http://ec.europa.eu/finance/securities/prospectus/index_en.htm

EU. (2014, Jun.). *Market Abuse Directive.* Retrieved from European Commission: http://ec.europa.eu/finance/securities/abuse/index_en.htm

Franklin, A., & Douglas, G. (1992). Stock Price Manipulation. *The Review of Financial Studies , 5* (3), 503-529.

Franklin, A., Litov, L., & Mei, J. (2006). Large investors, price manipulation, and limits to arbitrage: An anatomy of market corners. *Review of Finance , 10* (4), 645-693.

FSA. (2006, Mar.). *The Code of Market Conduct.* Retrieved from http://www.fsa.gov.uk/pubs/hb-releases/rel52/rel52mar.pdf

Hautsch, N., & Huang, R. (2011, Aug.). *Limit Order Flow, Market Impact and Optimal Order Sizes: Evidence from NASDAQ TotalView-ITCH Data.* Retrieved from SSRN: http://dx.doi.org/10.2139/ssrn.1914293

Hautsch, N., & Huang, R. (2012). The market impact of a limit order. *Journal of Economic Dynamics and Control , 36* (4), 501 - 522.

ITG. (2013, Jul.). *Transaction Cost Analysis.* Retrieved from ITG: http://www.itg.com/marketing/ITG_GlobalCostReview_Q120 13_20130725.pdf

Jamal, N. (2012, Dec.). *LSE broker fined for 'wash trade'.* Retrieved from DAWN: http://dawn.com/news/771335/lse-broker-fined-for-wash-trade

Jiang, Y., & Jiang, Z. P. (2014). Robust Adaptive Dynamic Programming and Feedback Stabilization of Nonlinear Systems. *IEEE Transactions on Neural Networks and Learning Systems , 25* (5), 882 - 893.

Kleinberg, J., & Tardos, E. (2005). *Algorithm Design* (1 edition ed.). Pearson.

Kojo, M., & Paudyal, K. (2000). The components of bid-ask spreads on the London Stock Exchange. *Journal of Banking & Finance , 24* (11), 1767-1785.

Lee, E. J., Eom, K. S., & Park, K. S. (2013). Microstructure-based manipulation: Strategic behavior and performance of spoofing traders. *Journal of Financial Markets , 16* (2), 227 - 252.

Loh, T., & Cumming, G. (2012, Jun.). *Market Manipulation: Safe Harbour for wash trades and matched orders upheld.* Retrieved from http://www.timothyloh.com/publications/120606_market_man ipulation_cfa.html

Markus., F., Hoser, B., & Schröder, J. (2007). On the Analysis of Irregular Stock Market Trading Behavior. *Proceedings of the Data Analysis, Machine Learning and Applications*, (pp. 355-362). Freiburg.

NANEX. (2013, May 31). *Chicago PMI.* Retrieved from http://www.nanex.net/aqck2/4304.html

NANEX. (2013, Jul. 10). *Exploratory Trading in the eMini*. (NANEX) Retrieved from http://www.nanex.net/aqck2/4136.html

Ou, Y., Cao, L., Luo, C., & Zhang, C. (Dec. 2008). Domain-Driven Local Exceptional Pattern Mining for Detecting Stock Price Manipulation. *Proceedings of the PRICAI 2008: Trends in Artificial Intelligence*, (pp. 849-858). Hanoi.

Palshikar, G. K., & Apte, M. M. (2008). Collusion set detection using graph clustering. *Data Mining and Knowledge Discovery , 16* (2), 135-164.

Patterson, S., Strasburg, J., & Trindle, J. (2013, Mar.). *'Wash Trades' Scrutinized*. (Wall Street Journal) Retrieved from The Wall Street Journal: http://online.wsj.com/article/SB1000142412788732363960457836649149707020 4.html

Rumen, A., Vincent, P., & Sanjay, R. (2000). Unbounded knapsack problem: Dynamic programming revisited. *European Journal of Operational Research , 123* (2), 394–407.

SEC. (2011, 10). *US Securities & Exchange Commision*. Retrieved from http://www.sec.gov/answers/limit.htm

Tsang, E., Olsen, R., & Masry, S. (2013). A formalization of double auction market dynamics. *Quantitative Finance , 13* (7), 981-988.

Vincent, P., Yanev, N., & Andonov, R. (2009). A hybrid algorithm for the unbounded knapsack problem. *Discrete Optimization , 6* (1), 110–124.

Wang, J., Zhou, S., & Guan, J. (2012). Detecting potential collusive cliques in futures markets based on trading behaviors from real data. *Neurocomputing , 92*, 44 - 53.

Zhai, J., & Cao, Y. (Mar. 2014). On the calibration of stochastic volatility models: A comparison study. *Proceedings of the 2104 IEEE Conference on Computational Intelligence for Financial Engineering & Economics (CIFEr)*, (pp. 303 - 309). London.

Zhen, N., He, H., Wen, J., & Xu, X. (2013). Goal Representation Heuristic Dynamic Programming on Maze Navigation. *IEEE Transactions on Neural Networks and Learning Systems , 24* (12), 2038 - 2050.

Zukerman, M., Jia, L., Neame, T., & Woeginger, G. J. (2001). A polynomially solvable special case of the unbounded knapsack problem. *Operations Research Letters , 29* (1), 13-16.

**Yi Cao** received the B.Eng. degree in navigation and control in aeronautics from the Beihang University, Beijing, China, the M.S. degree in computer science from Florida International University, Miami, FL, USA, and Ph.D. degree in financial machine learning in University of Ulster, Londonderry, UK in 2002, 2005 and 2015, respectively. He is currently an associate professor with Institute of Management Science and Engineering, Henan University, Kaifeng, Henan. He was an Integrated Circuit and Hardware Engineer at Vimicro, Beijing, and Conexant Beijing Design Centre, Beijing, respectively, from 2005 to 2008. From 2008 to 2011, he was a Senior System Engineer at ERICSSON, Beijing.

**Yuhua Li** (SM'11) received the Ph.D. degree in general engineering from the University of Leicester, Leicester, U.K.

He was with Manchester Metropolitan University, Manchester, U.K., and then the University of Manchester, Manchester, from 2000 to 2005, as a Senior Research Fellow and a Research Associate, respectively. He was a Lecturer with the School of Computing and Intelligent Systems, University of Ulster, U.K. from 2005 to 2014. He is currently a lecturer with the School of Computing, Science and Engineering, University of Salford, U.K. His current research interests include pattern recognition, machine learning, data science, knowledge-based systems, and condition monitoring and fault diagnosis.

**Sonya Coleman** (M'11) received the B.Sc (Hons.) degree in mathematics, statistics, and computing, and the Ph.D. degree in mathematical image processing from the University of Ulster, Londonderry, U.K.

She is a Reader in the University of Ulster. She has authored and co-authored more than 80 research papers in image processing, robotics, and computational neuroscience. She has experience managing research grants (with respect to technical aspects and personnel) as both a principal and co-investigator. In addition, she is a co-investigator on the EU FP7 funded projects RUBICON and VISUALISE.

**Ammar Belatreche** (M'09) received the Ph.D. degree in computer science from the University of Ulster, Londonderry, U.K.

He was a Research Assistant with the Intelligent Systems Engineering Laboratory and is currently a Lecturer with the School of Computing and Intelligent Systems, University of Ulster. His current research interests include bioinspired adaptive systems, machine learning, pattern recognition, and image processing and understanding. Dr. Belatreche is a fellow of the Higher Education Academy, a member of IEEE Computational Intelligence Society (CIS), and the Northern Ireland representative of the IEEE CIS. He is an Associate Editor of Neurocomputing and has served as a PC Member and reviewer for several international conferences and journals.

**Thomas Martin McGinnity** (SM'09) received the B.Sc. (Hons.) degree in physics and the Ph.D. degree from the University of Durham, Durham, U.K., in 1975 and 1979, respectivly. He is currently a Dean of Science and Technology with Nottingham Trent University and was formerly a Professor of Intelligent Systems Engineering and Director of the Intelligent Systems Research Centre with the Faculty of Computing and Engineering, University of Ulster. He was a Director with the University of Ulster's technology transfer company, Innovation Ulster, and a spin out company Flex Language Services. He has authored and co-authored approximately 300 research papers and has attracted over £24 million in research funding.

Prof. McGinnity is a fellow of the Institution of Engineering and Technology and a Chartered Engineer.