# Ultra-high-frequency lead-lag relationship and information arrival

**Abstract**

To our knowledge, this paper is the first study on the effect of information arrival on the lead-lag relationship among related spot instruments. Based on a large dataset of ultra-high-frequency transaction prices time-stamped to the millisecond of the S&P500 index and its two most liquid tracking ETFs, we find that their lead-lag relationship is affected by the rate of information arrival whose proxy is the unexpected trading volume of these instruments. Specifically, when information arrives, the leadership of the leading instrument may strengthen or weaken depending on whether the leading or lagging instrument responds to that information. An increase in the unexpected volume of the leader strengthens its leadership whereas an increase in the unexpected volume of the lagger weakens this leadership. In addition to the strength of leadership, an increase in the unexpected volume in response to information arrival may also have opposite effects on the lead-lag correlation coefficient depending on whether that volume increase belongs to the leader or the lagger. Finally, we find that sophisticated investors have a more significant effect on the lead-lag relationship than non-sophisticated ones.

## 1. Introduction

The efficient market hypothesis (EMH) suggests that return predictability and arbitrage opportunities should not exist in financial markets. Accordingly, the returns of related instruments (e.g. an equity index and its futures) should show contemporaneous correlations in efficient and frictionless markets (Stoll and Whaley, 1990, Brooks et al., 1999). However the lead-lag effect, a phenomenon where a security follows the movements of another with a time delay (Huth and Abergel, 2014), is often found in the literature. Yet according to the EMH, even if returns of security A are correlated with past returns of security B, it should still be impossible to use price changes of B to forecast and make abnormal profits from price changes of A due to practical constraints.

Motivated by the literature on lead-lag effect in returns (e.g. Kawaller et al., 1987, Fleming et al., 1996, Chen and Gau, 2009, Alsayed and McGroarty, 2014, Curme et al., 2015), our paper is the first to investigate the effect of information arrival on the lead-lag relationship between related stocks. To examine related stocks, we focus on ETFs because the purpose of ETFs is to track the performance of some index or asset and ETFs tracking the same index or asset can be considered very much related. Moreover, the ETF market has been growing rapidly (Shin and Soydemir, 2010, Kearney et al., 2014) since their introduction with many investors choosing ETFs as their investment vehicle. ETFs track different asset classes (e.g. stock, bond, commodity) and we choose equity ETFs for our analysis because they are easily accessible to investors, very liquid (Marshall et al., 2013) and should be representative of the US economy. Equity ETFs, which track stock indices or economic sectors, are more liquid than single stocks (Ruan and Ma, 2012).

Our study on the lead-lag effect in relation to information is also motivated by the fact that the lead-lag relationship exists because some financial instruments reflect information faster than others. In general, information plays an important role in financial markets. Hanousek and Podpiera (2003) state that informed trading affects the bid-ask spread because market makers set the spread to compensate for the risk of adverse selection which they face when trading with informed traders. Similarly, Gregoriou et al. (2005) argue that market makers mitigate their informational disadvantage compared to informed traders by increasing the bid-ask spread. On the other hand, the flow of information to the market also helps explain the ARCH and GARCH effects in daily stock returns (Lamoureux and Lastrapes, 1990, Sharma et al., 1996, Aragó and Nieto, 2005). In addition to prices and volatility, Frank and Kenneth (2005) suggest

that the relative trade size (i.e. trade size scaled by market depth) is also affected by information, showing that informed traders prefer to trade larger volume.

We find that the lead-lag relationship is indeed influenced by the rate of information arrival. When information arrives, the leadership of the leading instrument may strengthen or weaken depending on whether the leading or lagging instrument responds to that information. An increase in the unexpected trading volume of the leader strengthens its leadership whereas an increase in the unexpected volume of the lagger weakens this leadership. In addition to the strength of leadership, an increase in the unexpected volume in response to information arrival may also have opposite effects on the lead-lag correlation coefficient depending on whether that volume increase belongs to the leader or the lagger. We also find that sophisticated investors have a more significant effect on the lead-lag relationship than non-sophisticated ones.

Our research is conducted in the high-frequency context which has become more and more important in recent times. In current financial markets, speed is considered such an essential competitive edge that many market participants are willing to make significant technological investments to increase their speed of analysis and execution (even by only a small amount), in both absolute and relative terms (i.e. trying to be faster than their competitors). This competition for speed has pushed the boundary to extreme levels; specifically, Hasbrouck and Saar (2013) find that high-frequency traders can operate with a latency of only a few milliseconds while an eye blink takes a few hundred milliseconds. Even more extreme, Goldstein et al. (2014) report that it is possible to trade in the microsecond environment. In any case, high-frequency traders are among the most active and important market participants and thus we use the high-frequency setting to examine the lead-lag relationship. This setting is appropriate for our analysis because the lead-lag relationship is a common phenomenon in high-frequency data (Huth and Abergel, 2014).

Moreover, using high-frequency data to analyse the lead-lag relationship is suitable since the increasing electronification of financial markets and high-frequency trading activities have reduced the lead-lag time dramatically, to the point where data sampled at regular intervals can no longer capture this time delay (Huth and Abergel, 2014). In other words, it is not suitable to use regularly sampled data to measure high-frequency correlation (Zhang, 2011), especially when one security is traded more often than the other (Lo and MacKinlay, 1990). However, in order to measure correlation, using high-frequency data requires an approach different from

3

that used for regularly sampled data. Following Alsayed and McGroarty (2014), we apply the model of Hayashi and Yoshida (2005) to calculate the contemporaneous correlation, and its extension by Hoffmann et al. (2013) to include leads and lags. This model uses the original tick data and does not require any modification such as interpolation or resampling at regular intervals (Huth and Abergel, 2014).
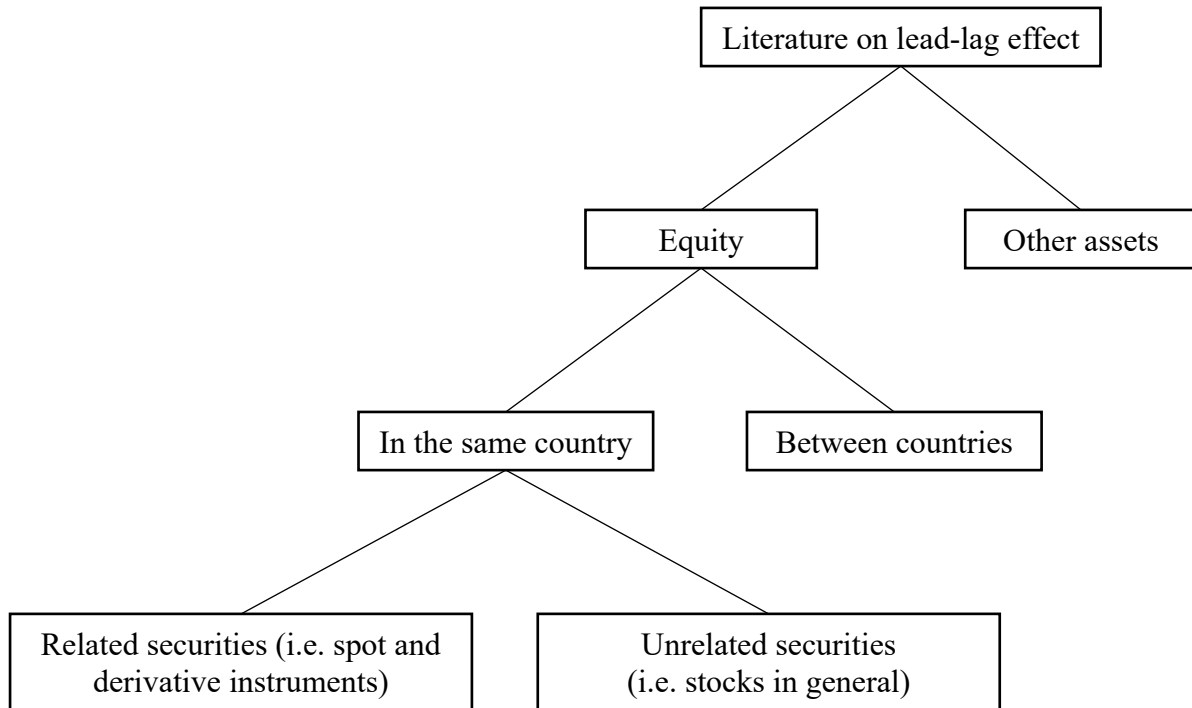
We contribute to the literature on lead-lag effects in returns by analysing the role of information in the lead-lag relationship among related spot instruments (i.e. equity index and ETFs) which, to our knowledge, has not been studied before. We examine, at the tick level, the most liquid ETFs that track the S&P500 index thus representative of the performance of the US economy. We show that information arrival has a large part of play in the lead-lag effect among these instruments. Moreover, we conduct our analysis in the high-frequency context, which is increasingly important in financial markets, using a large dataset and a novel approach proposed by Hayashi and Yoshida (2005) and Hoffmann et al. (2013).

The structure of this paper is as follows. Section 2 reviews the literature on lead-lag effects in returns. Section 3 presents our dataset. Section 4 describes our estimation of the lead-lag relationship and our analysis of the effect of information arrival on this relationship. Section 5 and 6 provide the results and conclusions respectively.

## 2. Literature review
The literature on lead-lag effects is plentiful and we classify the literature in Figure 1.

**Figure 1.** Classification of the literature on lead-lag effects in returns. Related securities refer to securities which have the same underlying asset (i.e. stocks and their derivatives). Unrelated securities refer to stocks in general.

```
                    ┌──────────────────────────┐
                    │ Literature on lead-lag effect │
                    └──────────────────────────┘
                       /                    \
            ┌──────────────┐          ┌──────────────┐
            │    Equity    │          │ Other assets │
            └──────────────┘          └──────────────┘
               /           \
  ┌────────────────────┐  ┌──────────────────┐
  │ In the same country │  │ Between countries │
  └────────────────────┘  └──────────────────┘
        /          \
┌──────────────────────┐  ┌──────────────────────┐
│ Related securities   │  │ Unrelated securities │
│ (i.e. spot and       │  │ (i.e. stocks in      │
│ derivative instruments) │ │ general)            │
└──────────────────────┘  └──────────────────────┘
```

Regarding Figure 1, a large part of the literature on lead-lag effects focuses on equity markets so we divide the literature into studies on equity markets and studies on other assets. Studies on equity markets can be divided further into studies on instruments in the same country and studies on instruments between countries. Finally, studies on instruments in the same country include studies on related securities (i.e. spot and derivative instruments) and studies on unrelated securities (i.e. stocks in general). Because our paper investigates related securities (i.e. ETFs tracking the same index), the following literature review will focus on the lead-lag effect between related securities. We mainly look at the spot – futures and spot – options relationship since there are few studies on other relationships.

The spot – futures relationship is the most extensively studied relationship in the literature on lead-lag effects of related securities. Many papers find that the lead-lag effect is bi-directional (i.e. futures lead the index and vice versa), although the futures' lead is stronger and longer than the index's lead (e.g. Chiang and Fong, 2001, Nam et al., 2006, Ergün, 2009). While futures' lead can be up to 45 minutes (Kawaller et al., 1987), the index's lead does not exceed 15 minutes (Chan, 1992). Regarding US indices, many papers find a bi-directional spot –

futures lead-lag relationship (e.g. Chan, 1992, Pizzi et al., 1998, Ergün, 2009) and only a few find a uni-directional effect where futures lead the index (e.g. Fleming et al., 1996). Regarding non-US indices, it is common to find both a bi-directional effect (e.g. Brooks et al., 1999, Frino and West, 1999) and a uni-directional effect (e.g. Najand and Min, 1999, Frino et al., 2000). However, Brooks et al. (2001) find that although the lead-lag effect can be used to produce accurate forecasts, trading these forecasts does not outperform the benchmark after considering transaction costs.

Some researchers attribute the lead-lag relationship to infrequent and nonsynchronous trading (Shyy and Vijayraghavan, 1996, Brooks et al., 1999) while others still find the lead-lag effect after considering infrequent and nonsynchronous trading (Stoll and Whaley, 1990, Grünbichler et al., 1994, Martikainen and Perttunen, 1995, Fleming et al., 1996). Another explanation for the lead-lag effect, especially the futures' lead, is the trading cost hypothesis (Nam et al., 2008). Specifically, because trading the index is cheaper in the derivative markets than in the spot market, new information should be updated in the derivative markets before the spot market (Martikainen and Perttunen, 1995, Fleming et al., 1996). Consistent with the trading cost hypothesis, Chen and Gau (2009) find that when the bid-ask spread in the spot market decreases (due to a decrease in the minimum tick size), the spot index's contribution to price discovery becomes more significant. In addition to nonsynchronous trading and trading costs, the trading mechanism also affects the lead-lag relationship. When the futures are screen-traded and the index is floor-traded, the futures' lead is longer than when both are floor-traded, since screen trading increases the price discovery speed (Grünbichler et al., 1994).

In terms of time variation, the lead-lag effect is regime-switching and non-linear (Chung et al., 2011). The futures' lead has weakened and the spot – futures integration has strengthened over time (Frino and West, 1999, Brooks et al., 1999). Lien et al. (2003) use daily data and even find that in more recent time, information flows from the spot market to the futures market, which means the index leads the futures. However, Nam et al. (2008) warn that using low-frequency data may lead to information loss and incorrect results, which is why we use tick data in this paper.

In addition to the spot – futures relationship, the spot – options relationship has often been studied and the findings are mixed. Some authors find no lead-lag effect (Panton, 1976, Chan et al., 1993) while others find a uni-directional effect (Stephan and Whaley, 1990, Fleming et

al., 1996) or a bi-directional effect (Chiang and Fong, 2001, Nam et al., 2006). However, the spot market tends to have a longer lead than the options. The spot market can lead the options by up to 20 minutes whereas the options lead the spot market by up to only 10 minutes (Chiang and Fong, 2001). Interestingly, bi-directional effects are usually found in non-US markets (e.g. Chiang and Fong, 2001, Nam et al., 2006).

Although option prices contain information not reflected in stock prices (Manaster and Rendleman, 1982), this information is not lucrative enough to cover transaction costs and search costs (Bhattacharya, 1987). On the other hand, the spot market's lead over options might be due to the infrequent trading and illiquidity of options (Chan et al., 1993, Fleming et al., 1996, Chiang and Fong, 2001). Lead-lag effects can also be explained by trading costs. Generally, information will be updated faster where it is cheaper to trade. Fleming et al. (1996) find that stocks lead options because trading stocks in the spot market is cheaper than in the option market and that futures lead options since trading costs in the future market are lower than in the option market. Ryu (2015) also finds that in the futures – options relationship, futures play a more significant role in price discovery than options. However, regarding the options – options relationship, there is no lead-lag effect between call and put options because of their similar trading costs (Fleming et al., 1996).

In summary, our literature review has focused on the lead-lag effect in returns of related securities (i.e. stocks, futures and options). The findings range from no lead-lag effect to a uni-directional effect to a bi-directional effect, with liquid and cheap instruments often leading illiquid and expensive ones. The lead-lag relationship may be affected by the trading mechanism of securities (i.e. screen-traded or floor-traded). However, this relationship might not be exploited profitably after considering transaction costs. Finally, there is also evidence of a weakening lead-lag effect and strengthening integration between markets over time.

## 3. Data

Our dataset includes the S&P 500 index and the two most liquid equity ETFs in the US market which are constructed to reflect the performance of the S&P 500 index. They are SPDR S&P 500 ETF Trust (ticker symbol SPY) and iShares Core S&P 500 ETF (ticker symbol IVV), provided by State Street Global Advisors and BlackRock respectively. Using Thomson Reuters Tick History database, we collect index value and ETF transaction price data (time-stamped to milliseconds) between August 2014 and July 2015. Following Marshall et al. (2013), in every

trading day, only the main trading session from 9:30am to 4pm is considered to ensure maximum liquidity.

Following the data cleaning procedure employed by Marshall et al. (2013) to remove potential data errors, we exclude observations where (i) the logarithmic return of price is higher than 25% or lower than -25% [1], or (ii) the time-stamp is within the first or last five minutes of the trading session. Table 1 shows the summary statistics of the data after the cleaning process. The index has the highest number of observations because it consists of a large number of stocks and moves whenever one or more component stocks move. The mean returns are small due to the large number of observations while the median returns are zero since there are many instances where consecutive transactions occur at the same price, resulting in zero returns. The returns range from -1.5% to 1.5%, with SPY and IVV showing the lowest and highest standard deviation respectively. Because of the positive skewness and leptokurtosis, the non-normality of the data is confirmed by the Jarque-Bera statistic, which is statistically significant at 1%.

**Table 1.** Summary statistics of logarithmic returns. The returns are in percentage terms. *** superscript denotes statistical significance at 1% level.

|  | S&P 500 | SPY | IVV |
|---|---|---|---|
| Mean | 4.93E-07 | 2.41E-06 | 3.1E-05 |
| Median | 0 | 0 | 0 |
| Maximum | 1.48 | 1.31 | 1.30 |
| Minimum | -1.36 | -1.44 | -1.47 |
| Standard deviation | 4.75E-03 | 2.60E-03 | 0.01 |
| Skewness | 30.04 | 10.03 | 2.65 |
| Kurtosis | 30349.33 | 41274.72 | 2271.93 |
| Jarque-Bera normality | 1.35E+14 *** | 1.25E+15 *** | 2.10E+11 *** |
| Number of observations | 17745864 | 4117085 | 281734 |

## 4. Methodology

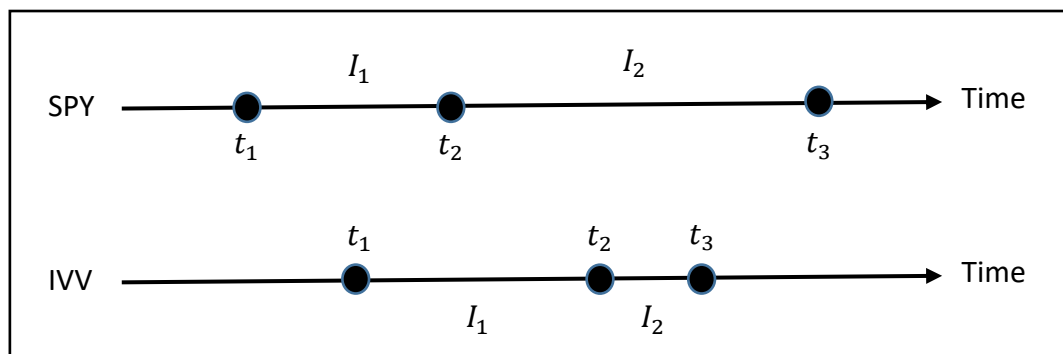### 4.1. Estimation of the lead-lag relationship

We analyse the lead-lag relationship between the S&P500 index and each ETF as well as between the two ETFs. To examine the lead-lag relationship between two series of non-synchronous tick data, we use the method of Hayashi and Yoshida (2005) and Hoffmann et al. (2013). Our purpose is to calculate the correlation coefficients between one series and

---

[1] In addition to the 25% threshold, we have used alternative cut-off points (i.e. 5%, 15%, 35% and 45%) and still got the same results.

timestamp-adjusted versions of the other to find the time adjustment which maximises their correlation. We describe the specific steps below using the SPY – IVV pair as an example. These steps are equally applicable to the S&P500 – SPY and S&P500 – IVV pairs.

The first step is to estimate the contemporaneous correlation coefficient between SPY and IVV (i.e. correlation where the timestamps of both series are kept unchanged). Because of the non-synchronicity of the data (illustrated in Figure 2), we use the method of Hayashi and Yoshida (2005). Their method does not require data synchronisation (e.g. through interpolation) and thus can avoid potential biases.

**Figure 2.** Non-synchronicity of the data. Each dot is a data point; t is the arrival time of observations; I is the interval between two consecutive observations.



Letting R denote the return and I denote the interval between two consecutive observations, the covariance C between SPY and IVV is given by

$$C = \sum_{i,j} R_{SPY}^{I_i} R_{IVV}^{I_j} \; \mathbb{I} \tag{1}$$

$$\text{where} \quad \mathbb{I} = \begin{cases} 1 & if \; I_i \cap I_j \neq \emptyset \\ 0 & if \; I_i \cap I_j = \emptyset \end{cases}$$

Equation (1) shows that the covariance is calculated by summing the products of every SPY interval return and its overlapping IVV interval returns. For example, in Figure 2, the covariance is calculated by summing the products of the following pairs of returns: $(R_{SPY}^{I_1}, R_{IVV}^{I_1})$, $(R_{SPY}^{I_2}, R_{IVV}^{I_1})$ and $(R_{SPY}^{I_2}, R_{IVV}^{I_2})$.

The standard deviation $\sigma$ of SPY and IVV is given by

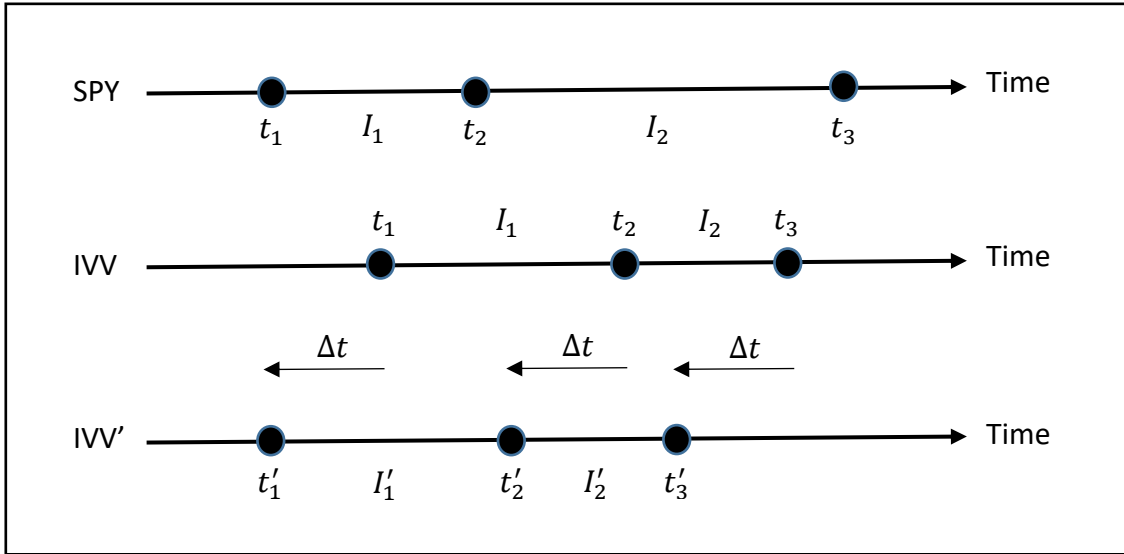$$\sigma_{SPY} = \sqrt{\Sigma_i (R_{SPY}^{I_i})^2} \qquad (2)$$

$$\sigma_{IVV} = \sqrt{\Sigma_j (R_{IVV}^{I_j})^2} \qquad (3)$$

The correlation coefficient $\rho$ of SPY and IVV is given by

$$\rho = \frac{C}{\sigma_{SPY} \, \sigma_{IVV}} \qquad (4)$$

After estimating the contemporaneous correlation by equation (4), the next step is to adjust the timestamp of either SPY or IVV to allow for leads and lags, and re-estimate their correlation as suggested by Hoffmann et al. (2013). We choose to fix the timestamp of SPY and adjust that of IVV. Regarding the S&P500 – SPY and S&P500 – IVV pairs, we choose to fix the timestamp of the ETFs and adjust that of the index. Figure 3 illustrates this process.

**Figure 3.** Example of time adjustment and correlation re-estimation. Each dot is a data point; t is the arrival time of observations; I is the interval between two consecutive observations. IVV' is created by moving every IVV observation backward in time by the same amount $\Delta t$. Then the correlation is re-estimated between SPY and IVV'.



Letting IVV' denote the timestamp-adjusted IVV series, I' denote the interval between two consecutive IVV' observations and C' denote the covariance between SPY and IVV'; the correlation coefficient $\rho$' between SPY and IVV' is given by

$$\rho' = \frac{C'}{\sigma_{SPY} \, \sigma_{IVV'}} = \frac{\Sigma_{i,j} \, R_{SPY}^{I_i} \, R_{IVV'}^{I'_j} \, \mathbb{I}'}{\sqrt{\Sigma_i (R_{SPY}^{I_i})^2 \, \Sigma_j (R_{IVV'}^{I'_j})^2}} \qquad (5)$$

$$\text{where} \quad \mathbb{I}' = \begin{cases} 1 & if \quad I_i \cap I_j' \neq \emptyset \\ 0 & if \quad I_i \cap I_j' = \emptyset \end{cases}$$

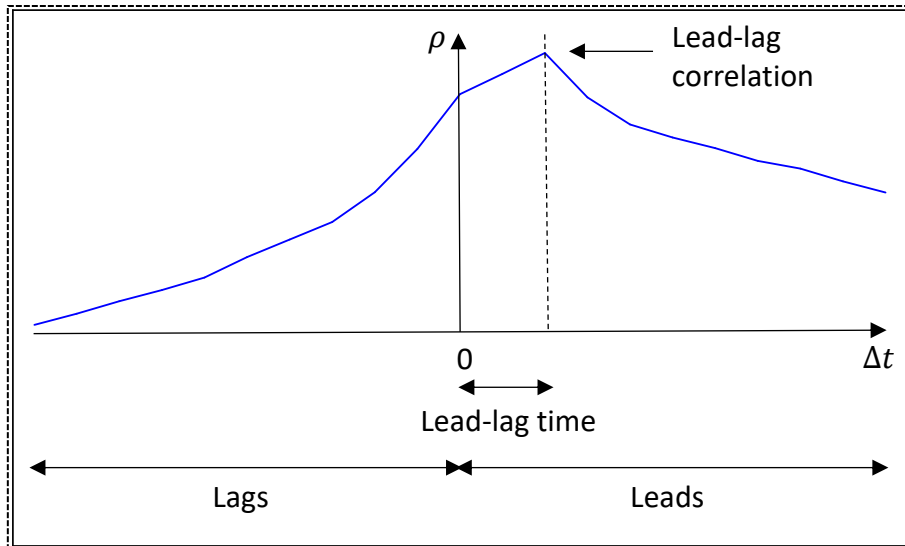Calculating $\rho$' with different time adjustments $\Delta t$ of IVV produces the correlation curve which shows the correlation coefficient between SPY and IVV at different leads and lags of IVV. To capture the ultra-high-frequency lead-lag relationship, we consider $\Delta t$ from -100 milliseconds (i.e. moving IVV backward) to 100 milliseconds (i.e. moving IVV forward) with 10-millisecond increments.

After producing the correlation curve, the final step is to find the $\Delta t$ which maximises the correlation. This $\Delta t$ shows the temporal relationship between SPY and IVV. If it is zero, there is no lead-lag relationship; if it is negative, IVV lags SPY by $\Delta t$; if it is positive, IVV leads SPY by $\Delta t$. Also, for ease of reference, we refer to the maximum correlation on the correlation curve as the lead-lag correlation coefficient hereafter. In addition, following Huth and Abergel (2014), we calculate the lead-lag ratio (LLR) as follows.

$$LLR = \frac{\sum_i (\rho_{(\Delta t)_i})^2}{\sum_i (\rho_{-(\Delta t)_i})^2} \quad (\Delta t > 0) \tag{6}$$

The numerator of LLR is the sum of squared correlation coefficients at all leads of IVV while the denominator is the sum of squared correlation coefficients at all lags of IVV. LLR measures the relative strength of leadership (i.e. if LLR is higher than one, IVV tends to lead SPY more than lag and vice versa). Figure 4 shows an example of the correlation curve.

**Figure 4.** Example of the correlation curve. This curve is obtained by calculating the correlation between two instruments while fixing the timestamp of one and adjusting that of the other. The horizontal axis shows the time adjustment $\Delta t$. The vertical axis shows the correlation coefficient $\rho$. The peak of the curve is the lead-lag correlation and its corresponding $\Delta t$ is the lead-lag time.
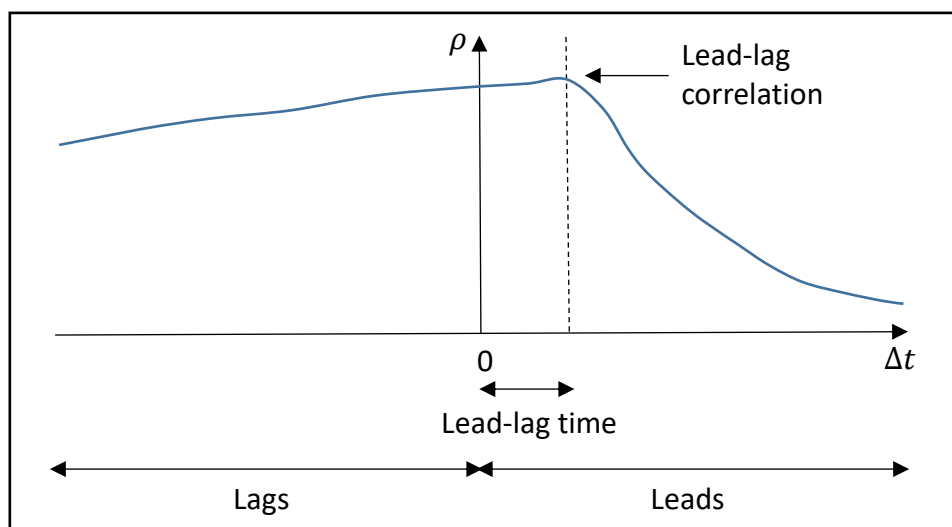
## 4.2. The effect of information arrival on the lead-lag relationship

### 4.2.1. Dependent variables

To study the effect of information arrival on the lead-lag relationship, we use regression analysis. The dependent variables are variables which represent the lead-lag relationship, namely the lead-lag correlation coefficient, the lead-lag time and the strength of leadership (measured by the lead-lag ratio). Although previous studies focus on the lead-lag correlation and the lead-lag time (e.g. Fleming et al., 1996, Nam et al., 2006, Ergün, 2009), the lead-lag ratio (calculated by Equation 6) is necessary to provide a more comprehensive analysis of the lead-lag relationship. Let us consider the following example.

**Figure 5.** Example of a correlation curve. This curve is obtained by calculating the correlation between two instruments while fixing the timestamp of one and adjusting that of the other. The horizontal axis shows the time adjustment $\Delta t$. The vertical axis shows the correlation coefficient $\rho$. The peak of the curve is the lead-lag correlation and its corresponding $\Delta t$ is the lead-lag time.



If we focus on the lead-lag correlation and the lead-lag time, we may conclude that the time-adjusted instrument leads the time-fixed one because the peak of the curve is on the 'leads' side. However, the lead-lag ratio, which examines not only the peak but also a range of leads and lags, results in a different conclusion. Since the correlation is generally higher on the 'lags' side than on the 'leads' side, the lead-lag ratio suggests that the time-adjusted instrument tends to lag the time-fixed one. To cover this type of situation, it is important to consider the lead-lag ratio in addition to the lead-lag correlation and the lead-lag time. For each pairwise combination of the S&P 500 index and the two ETFs, we follow the steps in section 4.1 to obtain these lead-lag quantities for every trading day in the sample period. As a result, each

pair of instruments has three daily series corresponding to the three lead-lag variables which are used as the dependent variables in our regression analysis.

### 4.2.2. Independent variables

Previous research has found that the lead-lag relationship is affected by factors such as the trading cost and the trading mechanism of the instruments. Regarding the trading cost, information is generally updated faster where it is cheaper to trade. For example, because trading the index is cheaper in the derivative markets than in the spot market, new information should be updated in the derivative markets before the spot market (Martikainen and Perttunen, 1995, Fleming et al., 1996, Nam et al., 2008). However, as for the derivative markets, Fleming et al. (1996) show that there is no lead-lag effect between call and put options because of their similar trading costs. Regarding the trading mechanism, Grünbichler et al. (1994) find that when the leading instrument changes from being floor-traded to being screen-traded, its leadership strengthens since screen trading increases the price discovery speed. In our study, these factors should not contribute to the lead-lag effect because we only examine electronically traded spot instruments and no derivative.

Unlike the trading cost and trading mechanism, the information flow to the market may influence the lead-lag effect in our study. Motivated by (i) the fact that the lead-lag relationship exists because some instruments reflect information faster than others and (ii) the importance of information in financial markets (e.g. Hanousek and Podpiera, 2003, Gregoriou et al., 2005, Frank and Kenneth, 2005), we hypothesise that changes in the information flow have an impact on the lead-lag relationship. Therefore, our independent variables are variables which represent information arrival. A common proxy for the rate of information arrival is trading volume (e.g. Lamoureux and Lastrapes, 1990, Sharma et al., 1996, Aragó and Nieto, 2005) so our independent variables are daily trading volume of the S&P500 index and the two ETFs. Table 2 shows the summary statistics of the trading volume. The volume is highest for the index and lowest for the IVV ETF. The index volume is platykurtic while the ETFs' volume is leptokurtic. The volume of all instruments is positively skew and non-normal, as shown by the Jarque-Bera statistic.

**Table 2.** Summary statistics of the daily trading volume. The volume is in million shares. *** superscript denotes statistical significance at 1% level.

|  | S&P 500 | SPY | IVV |
|---|---|---|---|
| Mean | 525.074 | 29.944 | 0.770 |
| Median | 509.568 | 26.997 | 0.684 |
| Maximum | 927.539 | 100.688 | 2.673 |
| Minimum | 224.945 | 10.861 | 0.175 |
| Standard deviation | 101.001 | 13.250 | 0.432 |
| Skewness | 0.742 | 1.719 | 1.715 |
| Kurtosis | 2.286 | 4.274 | 3.763 |
| Jarque-Bera normality | 27.887 *** | 138.382 *** | 127.118 *** |

Because the information content of trading activities is not uniform among different types of investors, we distinguish sophisticated investors from non-sophisticated ones. We divide the total trading volume of each instrument into the 'block trades' and 'non-block trades' component to proxy for sophisticated and non-sophisticated investors respectively since Madhavan and Sofianos (1998) show that block trades (i.e. trades of a large number or value of shares) are typically initiated by institutional traders. After obtaining the block volume data, the non-block volume is the difference between the total volume and the block volume. Letting V, BV and NBV denote the total, block and non-block volume respectively, we have the following equation.

$$NBV_t = V_t - BV_t \tag{7}$$

To reflect the information flow more accurately, we divide both the block and non-block volume into the expected and unexpected component, as suggested by Bessembinder and Seguin (1993) and Aragó and Nieto (2005). Aragó and Nieto (2005) point out that the expected and unexpected volume capture the normal level of market activity and the arrival of new information respectively. Moreover, Bessembinder and Seguin (1993) find that the expected component in a given day is equal to yesterday's level and the unexpected component shows the change during the day. Letting EBV, UBV, ENBV and UNBV denote the expected block, unexpected block, expected non-block and unexpected non-block volume respectively, we have the following equations.

$$EBV_t = BV_{t-1} \tag{8}$$

$$UBV_t = BV_t - EBV_t \tag{9}$$

$$ENBV_t = NBV_{t-1} \tag{10}$$

$$UNBV_t = NBV_t - ENBV_t \tag{11}$$

14

In summary, we divide the total volume of each instrument into four components based on two dimensions, namely (i) block and non-block and (ii) expected and unexpected. As a result, each of the three instruments has four daily volume series which are used as the independent variables.

### 4.2.3. Regression analysis

We hypothesise the following.

1. The information flow (i.e. unexpected volume) affects the lead-lag relationship (i.e. lead-lag correlation coefficient, lead-lag time and lead-lag ratio).
2. Sophisticated investors (i.e. block volume) have a more significant impact on the lead-lag relationship than non-sophisticated investors (i.e. non-block volume).

Regarding the regressions, for each index – ETF pair, the independent variables are the trading volume of each instrument in that pair. However, for the ETF – ETF pair, the independent variables include not only the volume of each ETF but also the index volume, because both ETFs track the index so the index volume may be relevant to the ETFs. For each pair, we estimate a separate regression for each of the three lead-lag variables (i.e. lead-lag correlation coefficient, lead-lag time or lead-lag ratio). Letting Y denote the lead-lag variable, we estimate the following regression for each index – ETF pair.

$$Y_t = \alpha + \beta_1 EBV_{t,SP500} + \beta_2 UBV_{t,SP500} + \beta_3 ENBV_{t,SP500} + \beta_4 UNBV_{t,SP500} \\ + \gamma_1 EBV_{t,ETF} + \gamma_2 UBV_{t,ETF} + \gamma_3 ENBV_{t,ETF} + \gamma_4 UNBV_{t,ETF} + \varepsilon_t \quad (12)$$

We estimate the following regression for the ETF – ETF pair.

$$Y_t = \alpha + \beta_1 EBV_{t,SP500} + \beta_2 UBV_{t,SP500} + \beta_3 ENBV_{t,SP500} + \beta_4 UNBV_{t,SP500} \\ + \gamma_1 EBV_{t,IVV} + \gamma_2 UBV_{t,IVV} + \gamma_3 ENBV_{t,IVV} + \gamma_4 UNBV_{t,IVV} \\ + \delta_1 EBV_{t,SPY} + \delta_2 UBV_{t,SPY} + \delta_3 ENBV_{t,SPY} + \delta_4 UNBV_{t,SPY} + \varepsilon_t \quad (13)$$

## 5. Results

### 5.1. Overall lead-lag relationship

Table 3 shows the overall lead-lag relationships among the S&P 500 index and its tracking ETFs, namely SPY and IVV. These results are based on the whole sample, as opposed to the daily estimation used for the regression analysis. The lead-lag relationship of each pairwise combination of the three instruments is represented by three lead-lag quantities. After obtaining
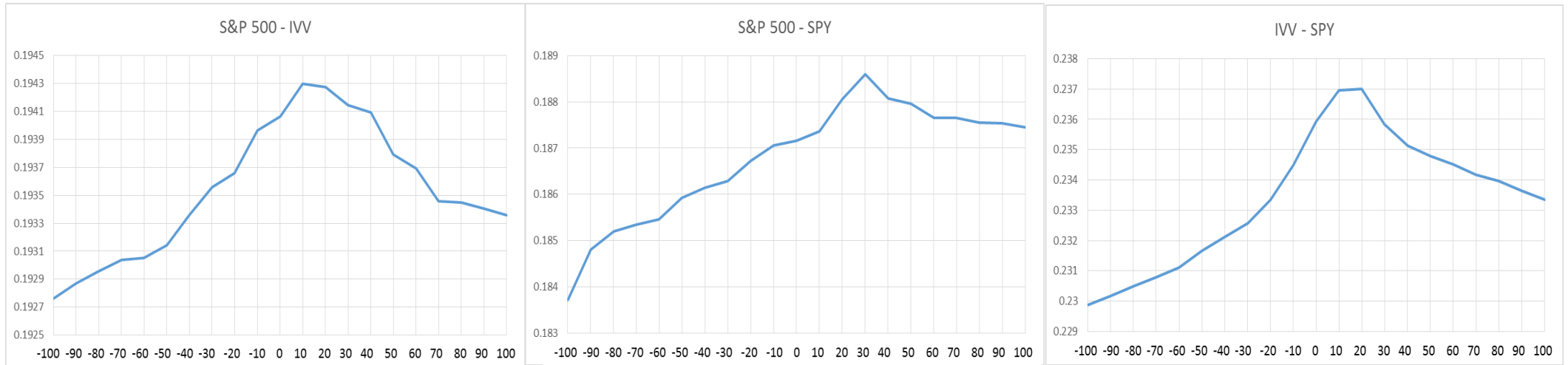
the correlation curve using Equation 5, the lead-lag correlation and the lead-lag time are measured at the peak of the curve while the lead-lag ratio is calculated by Equation 6.

**Table 3.** Lead-lag relationship among the S&P 500 index and its tracking ETFs. The first name in each pair is the leader. The lead-lag time is in milliseconds. For example, that the lead-lag time of the S&P500 – IVV pair is 10 means that the S&P 500 index leads the IVV ETF by 10 milliseconds.

|  | Lead-lag correlation | Lead-lag time (ms) | Lead-lag ratio |
|---|---|---|---|
| S&P500 – IVV | 0.1943 | 10 | 1.0058 |
| S&P500 – SPY | 0.1886 | 30 | 1.0230 |
| IVV – SPY | 0.2370 | 20 | 1.0284 |

Although the S&P 500 index and its tracking ETFs are highly correlated using daily data (i.e. their daily correlation coefficient is 99.9%), Table 3 shows that they are only moderately correlated (i.e. around 20% correlation) using tick data due to the Epps effect (i.e. an effect documented by Epps (1979) where financial instruments become less correlated at higher sampling frequencies). The ETF – ETF pair is more correlated than the two index – ETF pairs. The lead-lag time is relatively short, ranging from 10 to 30 milliseconds. The index leads both ETFs, so price discovery starts from the index and the ETFs follow. Regarding the lead-lag ratio, it is slightly higher than one and highest for the ETF – ETF pair. Figure 5 shows the correlation curves of the three pairs of instruments.

**Figure 5.** Correlation curves. This figure shows the correlation curves of the three pairs of instruments. The vertical axis shows the correlation coefficient. The horizontal axis shows the time adjustment in milliseconds (we adjust the timestamp of the first instrument in each pair). Because the correlation range of each pair is small, if we plot all the three curves on the same chart, they look like straight lines. Therefore, we plot each curve on a separate chart.

According to Figure 5, the peaks of all three curves correspond to a positive time adjustment on the horizontal axis, which means that the instrument whose timestamp is adjusted (i.e. the first instrument in each pair) leads the instrument whose timestamp is fixed (i.e. the second instrument in each pair). When we move away from the peak on both sides, the correlation coefficient decreases gradually. The ETF – ETF curve is smoother than the two index – ETF curves and the two IVV curves are more symmetrical on both sides of the peak than the S&P500 – SPY curve. Because Figure 5 shows that the lead-lag relationship exists among the index and ETFs, we attempt to exploit this relationship by applying the trading strategy of Alsayed and McGroarty (2014) where we trade the lagger in the direction of the leader's previous price change. However, the strategy is not profitable, which is consistent with Brooks et al. (2001), who find that although the lead-lag effect can be used to produce accurate forecasts, trading these forecasts does not yield favourable returns.

## 5.2. The effect of information arrival on the lead-lag relationship

Table 4 shows the results of regression (12) and (13), which estimate the effect of information arrival on the daily lead-lag relationship among the S&P500 index and ETFs. The daily trading volume of each instrument (i.e. the independent variable, in million shares) is divided into four components based on two dimensions, namely (i) block and non-block and (ii) expected and unexpected. The expected and unexpected volume are used as a proxy for the normal market activity and information arrival respectively while the block and non-block volume are used as a proxy for sophisticated and non-sophisticated investors respectively. The lead-lag relationship (i.e. the dependent variable) is represented by three lead-lag quantities. After obtaining the correlation curve using Equation 5, the lead-lag correlation (panel A) and the lead-lag time (panel B) are measured at the peak of the curve while the lead-lag ratio (panel C) is calculated by Equation 6.

**Table 4.** Effect of information arrival on the lead-lag relationship. This table shows the estimated coefficients for regression (12) and (13). The first instrument in each pair is the leader. The lead-lag time is in milliseconds. The numbers in brackets are standard errors of the coefficients. The $^{***}$, $^{**}$ and $^{*}$ superscripts denote statistical significance at 1%, 5% and 10% level respectively.

| | Index - IVV | | Index - SPY | | IVV - SPY | |
|---|---|---|---|---|---|---|
| *Panel A: lead-lag correlation coefficient* | | | | | | |
| Intercept | 0.15634$^{***}$ | (0.02066) | 0.20387$^{***}$ | (0.01383) | 0.26919$^{***}$ | (0.01640) |
| Expected S&P500 block volume | -0.00007 | (0.00072) | 0.00031 | (0.00023) | -0.00020 | (0.00027) |
| Unexpected S&P500 block volume | -0.00073$^{**}$ | (0.00032) | 0.00012 | (0.00048) | -0.00016 | (0.00058) |
| Expected S&P500 non-block volume | -0.00013 | (0.00009) | -0.00007 | (0.00007) | -0.00001 | (0.00008) |
| Unexpected S&P500 non-block volume | -0.00001 | (0.00018) | 0.00001 | (0.00012) | -0.00003 | (0.00014) |
| Expected IVV block volume | 0.03171$^{*}$ | (0.01647) | - | | 0.01968 | (0.01464) |
| Unexpected IVV block volume | 0.05245$^{**}$ | (0.02096) | - | | 0.04518$^{**}$ | (0.01991) |
| Expected IVV non-block volume | 0.03043$^{*}$ | (0.01593) | - | | 0.02738$^{*}$ | (0.01435) |
| Unexpected IVV non-block volume | 0.03540$^{*}$ | (0.01856) | - | | 0.02530 | (0.01735) |
| Expected SPY block volume | - | | -0.00006 | (0.00039) | -0.00127$^{**}$ | (0.00058) |
| Unexpected SPY block volume | - | | -0.00032 | (0.00043) | -0.00118$^{**}$ | (0.00059) |
| Expected SPY non-block volume | - | | -0.00077$^{*}$ | (0.00043) | -0.00110$^{*}$ | (0.00057) |
| Unexpected SPY non-block volume | - | | -0.00086$^{**}$ | (0.00043) | -0.00112$^{*}$ | (0.00059) |
| Adjusted R$^2$ | 0.27455 | | 0.20854 | | 0.27274 | |
| *Panel B: lead-lag time (ms)* | | | | | | |
| Intercept | 9.01535$^{***}$ | (2.38260) | 23.85950$^{***}$ | (5.76255) | 15.42479$^{***}$ | (5.55488) |
| Expected S&P500 block volume | 0.04093 | (0.03675) | -0.00957 | (0.09441) | -0.07944 | (0.09076) |
| Unexpected S&P500 block volume | 0.11544 | (0.08324) | 0.29438 | (0.20158) | 0.12359 | (0.19666) |
| Expected S&P500 non-block volume | -0.00904 | (0.01068) | 0.00303 | (0.02821) | 0.02119 | (0.02726) |
| Unexpected S&P500 non-block volume | -0.03131 | (0.02030) | -0.07552 | (0.04950) | -0.02889 | (0.04838) |
| Expected IVV block volume | 0.92662 | (2.41662) | - | | 3.32228 | (6.74262) |
| Unexpected IVV block volume | 2.38254 | (1.89976) | - | | 1.14084 | (4.95906) |
| Expected IVV non-block volume | 1.67334 | (2.14030) | - | | -5.70446 | (5.87593) |
| Unexpected IVV non-block volume | 2.85017 | (1.83670) | - | | -6.78641 | (4.86146) |
| Expected SPY block volume | - | | 0.06611 | (0.16393) | -0.01539 | (0.19880) |
| Unexpected SPY block volume | - | | 0.11651 | (0.17921) | 0.01657 | (0.19765) |
| Expected SPY non-block volume | - | | 0.02900 | (0.18099) | 0.01984 | (0.20063) |
| Unexpected SPY non-block volume | - | | 0.08662 | (0.17958) | 0.05234 | (0.19309) |
| Adjusted R$^2$ | 0.15143 | | 0.11322 | | 0.14527 | |
| *Panel C: lead-lag ratio* | | | | | | |
| Intercept | 1.03604$^{***}$ | (0.01999) | 1.04818$^{***}$ | (0.00984) | 1.02697$^{***}$ | (0.02650) |
| Expected S&P500 block volume | -0.00036 | (0.00031) | 0.00016 | (0.00016) | -0.00079$^{*}$ | (0.00043) |
| Unexpected S&P500 block volume | 0.00196$^{***}$ | (0.00070) | -0.00027 | (0.00034) | -0.00029 | (0.00094) |
| Expected S&P500 non-block volume | 0.00009 | (0.00009) | -0.00007 | (0.00005) | 0.00016 | (0.00013) |
| Unexpected S&P500 non-block volume | 0.00045$^{***}$ | (0.00017) | 0.00006 | (0.00008) | 0.00007 | (0.00023) |
| Expected IVV block volume | -0.01787 | (0.02027) | - | | 0.03732 | (0.03216) |
| Unexpected IVV block volume | -0.00847 | (0.01594) | - | | 0.03306 | (0.02365) |
| Expected IVV non-block volume | -0.02726 | (0.01796) | - | | 0.00242 | (0.02803) |
| Unexpected IVV non-block volume | -0.02190 | (0.01541) | - | | 0.03576 | (0.02319) |
| Expected SPY block volume | - | | 0.00006 | (0.00028) | -0.00098 | (0.00095) |
| Unexpected SPY block volume | - | | -0.00022 | (0.00031) | -0.00246$^{***}$ | (0.00094) |
| Expected SPY non-block volume | - | | -0.00002 | (0.00031) | -0.00085 | (0.00096) |
| Unexpected SPY non-block volume | - | | -0.00013 | (0.00031) | -0.00227$^{**}$ | (0.00092) |
| Adjusted R$^2$ | 0.20238 | | 0.13908 | | 0.25208 | |

According to Table 4, there is evidence that the lead-lag relationship among the S&P500 index and its tracking ETFs is influenced by the rate of information arrival which is captured by the unexpected trading volume of these instruments. In particular, the impact of information is significant on the lead-lag correlation coefficient and the lead-lag ratio. Specifically, when more information comes to the market, the resultant increase in unexpected volume of the leader and the lagger has opposite effects on the lead-lag correlation coefficient. For the S&P500 – IVV pair, increased trading of the leader (lagger) leads to a lower (higher) correlation coefficient while for the other pairs, increased trading of the leader (lagger) leads to a higher (lower) correlation coefficient. Compared to the ETFs, the unexpected index volume has a less pronounced effect on the lead-lag correlation (statistically significant in only the index – IVV pair). Regarding the strength of leadership (measured by the lead-lag ratio), although an increase in information intensity may lead to an increase in the unexpected trading volume of both the leader and the lagger, their changes have opposite effects on the strength of leadership. If the unexpected volume of the leader increases, its leadership strengthens while if the unexpected volume of the lagger increases, this leadership weakens. This effect is significant at 1% level for the index – IVV pair and the IVV – SPY pair. Finally, sophisticated investors (i.e. block volume) have a more significant effect on the lead-lag relationship than non-sophisticated ones (i.e. non-block volume).

Table 5 summarises the magnitude of statistically significant results, namely the lead-lag correlation and the lead-lag ratio. It shows the changes in these lead-lag variables caused by a given increase in the trading volume of the leader and the lagger. Because each instrument has a different level of trading activities, we report the changes caused by a volume increase of one standard deviation to make the results more comparable. The changes in the lead-lag variables are reported in absolute and relative terms (i.e. as a percentage of the overall levels in Table 3).

**Table 5.** Magnitude of statistically significant results. This table shows the changes in the lead-lag correlation and the lead-lag ratio (in absolute and relative terms) caused by an increase of one standard deviation in the trading volume of the leader and the lagger. The first instrument in each pair is the leader.

|  | Index - IVV | | Index - SPY | | IVV - SPY | |
|---|---|---|---|---|---|---|
|  | Absolute | Relative | Absolute | Relative | Absolute | Relative |
| *Lead-lag correlation* |  |  |  |  |  |  |
| Leader | -0.0477 | -24.57% | 0.0065 | 3.43% | 0.0429 | 18.08% |
| Lagger | 0.0545 | 28.03% | -0.0241 | -12.77% | -0.0553 | -23.35% |
| *Lead-lag ratio* |  |  |  |  |  |  |
| Leader | 0.1453 | 14.45% | -0.0036 | -0.36% | 0.0395 | 3.84% |
| Lagger | -0.0267 | -2.66% | -0.0039 | -0.38% | -0.0800 | -7.78% |

The SPY volume has the largest impact on the lead-lag correlation in absolute terms (-0.0553 for the IVV – SPY pair) whereas the IVV volume has the largest impact in relative terms (28.03% for the index – IVV pair). For all pairs, the laggers have a larger influence on the lead-lag correlation than the leaders. On the other hand, the index volume has the largest effect on the lead-lag ratio in both absolute and relative terms (0.1453 and 14.45% respectively for the index – IVV pair).

## 6. Conclusion

To our knowledge, this paper is the first study on the effect of information arrival on the lead-lag relationship among related spot instruments. Based on a large dataset of ultra-high-frequency transaction prices time-stamped to the millisecond and a novel approach proposed by Hayashi and Yoshida (2005) and Hoffmann et al. (2013), we find lead-lag effects among the S&P500 index and its two most liquid tracking ETFs. The lead-lag correlation coefficients are relatively low (i.e. around 20%) due to the Epps effect (Epps, 1979) where financial instruments become less correlated at higher sampling frequencies. The index leads both ETFs by 10 and 30 milliseconds respectively so price discovery starts from the index.

Using daily unexpected trading volume of the index and ETFs to proxy for the arrival rate of information, we find that information intensity affects the lead-lag relationship, especially the lead-lag correlation coefficient and the strength of leadership. Regarding the strength of leadership, the leadership of the leader may strengthen or weaken depending on whether the leader or lagger responds to the arrival of information. This leadership (i) strengthens when the unexpected volume of the leader increases and (ii) weakens when the unexpected volume of the lagger increases. Regarding the lead-lag correlation coefficient, although the unexpected volume of both the leader and the lagger may increase in response to information arrival, their changes may also have opposite effects on the correlation coefficient. Moreover, we find that sophisticated investors have a more significant effect on the lead-lag relationship than non-sophisticated ones. All in all, our study suggests that future research on the lead-lag relationship, especially in high-frequency data, should pay attention to the influence of information arrival on this relationship.

# References

ALSAYED, H. & MCGROARTY, F. 2014. Ultra-high-Frequency Algorithmic Arbitrage across International Index Futures. *Journal of Forecasting,* 33**,** 391-408.

ARAGÓ, V. & NIETO, L. 2005. Heteroskedasticity in the returns of the main world stock exchange indices: volume versus GARCH effects. *Journal of International Financial Markets, Institutions & Money,* 15**,** 271-284.

BESSEMBINDER, H. & SEGUIN, P. J. 1993. Price Volatility, Trading Volume, and Market Depth: Evidence from Futures Markets. *Journal of Financial and Quantitative Analysis,* 28**,** 21-39.

BHATTACHARYA, M. 1987. Price Changes of Related Securities: The Case of Call Options and Stocks. *Journal of Financial & Quantitative Analysis,* 22**,** 1-15.

BROOKS, C., GARRETT, I. & HINICH, M. J. 1999. An alternative approach to investigating lead-lag relationships between stock and stock index futures markets. *Applied Financial Economics,* 9**,** 605-613.

BROOKS, C., REW, A. G. & RITSON, S. 2001. A trading strategy based on the lead–lag relationship between the spot index and futures contract for the FTSE 100. *International Journal of Forecasting,* 17**,** 31-44.

CHAN, K. 1992. A Further Analysis of the Lead--Lag Relationship Between the Cash Market and Stock Index Futures Market. *Review of Financial Studies,* 5**,** 123-152.

CHAN, K., CHUNG, Y. P. & JOHNSON, H. 1993. Why Option Prices Lag Stock Prices: A Trading-based Explanation. *Journal of Finance,* 48**,** 1957-1967.

CHEN, Y. L. & GAU, Y. F. 2009. Tick Sizes and Relative Rates of Price Discovery in Stock, Futures and Options Markets: Evidence from the Taiwan Stock Exchange. *Journal of Futures Markets,* 29**,** 74-93.

CHIANG, R. & FONG, W.-M. 2001. Relative informational efficiency of cash, futures, and options markets: The case of an emerging market. *Journal of Banking and Finance,* 25**,** 355-375.

CHUNG, H.-L., CHAN, W.-S. & BATTEN, J. A. 2011. Threshold non-linear dynamics between Hang Seng stock index and futures returns. *European Journal of Finance,* 17**,** 471-486.

CURME, C., STANLEY, H. E., KENETT, D. Y., TUMMINELLO, M. & MANTEGNA, R. N. 2015. Emergence of statistically validated financial intraday lead-lag relationships. *Quantitative Finance,* 15**,** 1375-1386.

EPPS, T. W. 1979. Comovements in Stock Prices in the Very Short Run. *Journal of the American Statistical Association,* 74**,** 291-298.

ERGÜN, A. T. 2009. NYSE Rule 80A restrictions on index arbitrage and market linkage. *Applied Financial Economics,* 19**,** 1675-1685.

FLEMING, J., OSTDIEK, B. & WHALEY, R. E. 1996. Trading Costs and the Relative Rates of Price Discovery in Stock, Futures, and Option Markets. *Journal of Futures Markets,* 16**,** 353-387.

FRANK, H. & KENNETH, W. S. 2005. Trade Size And Informed Trading: Which Trades Are "Big"? *Journal of Financial Research,* 28**,** 133-163.

FRINO, A., WALTER, T. & WEST, A. 2000. The Lead-Lag Relationship between Equities and Stock Index Futures Markets around Information Releases. *Journal of Futures Markets,* 20**,** 467-487.

FRINO, A. & WEST, A. 1999. The Lead-Lag Relationship Between Stock Indices and Stock Index Futures Contracts: Further Australian Evidence. *Abacus,* 35**,** 333-341.

GOLDSTEIN, M. A., KUMAR, P. & GRAVES, F. C. 2014. Computerized and High-Frequency Trading. *Financial Review,* 49**,** 177-202.

GREGORIOU, A., IOANNIDIS, C. & SKERRATT, L. 2005. Information Asymmetry and the Bid-Ask Spread: Evidence from the UK. *Journal of Business Finance and Accounting,* 32**,** 1801-1826.

GRÜNBICHLER, A., LONGSTAFF, F. A. & SCHWARTZ, E. S. 1994. Regular Article: Electronic Screen Trading and the Transmission of Information: An Empirical Examination. *Journal of Financial Intermediation,* 3**,** 166-187.

HANOUSEK, J. & PODPIERA, R. 2003. Informed trading and the bid–ask spread: evidence from an emerging market. *Journal of Comparative Economics,* 31**,** 275-296.

HASBROUCK, J. & SAAR, G. 2013. Low-latency trading. *Journal of Financial Markets,* 16**,** 646-679.

HAYASHI, T. & YOSHIDA, N. 2005. On Covariance Estimation of Non-Synchronously Observed Diffusion Processes. *Bernoulli,* 11**,** 359-379.

HOFFMANN, M., ROSENBAUM, M. & YOSHIDA, N. 2013. Estimation of the lead-lag parameter from non-synchronous data. *Bernoulli,* 19**,** 426-461.

HUTH, N. & ABERGEL, F. 2014. High frequency lead/lag relationships — Empirical facts. *Journal of Empirical Finance,* 26**,** 41-58.

KAWALLER, I. G., KOCH, P. D. & KOCH, T. W. 1987. The Temporal Price Relationship Between S&P 500 Futures and the S&P 500 Index. *Journal of Finance,* 42**,** 1309-1329.

KEARNEY, F., CUMMINS, M. & MURPHY, F. 2014. Outperformance in Exchange-Traded Fund Pricing Deviations: Generalized Control of Data Snooping Bias. *Journal of Financial Markets,* 19**,** 86-109.

LAMOUREUX, C. G. & LASTRAPES, W. D. 1990. Heteroskedasticity in Stock Return Data: Volume versus GARCH Effects. *Journal of Finance,* 45**,** 221-229.

LIEN, D., TSE, Y. K. & XIBIN, Z. 2003. Structural change and lead-lag relationship between the Nikkei spot index and futures price: a genetic programming approach. *Quantitative Finance,* 3**,** 136-144.

LO, A. W. & MACKINLAY, A. C. 1990. An Econometric Analysis of Nonsynchronous Trading. *Journal of Econometrics,* 45**,** 181-211.

MADHAVAN, A. & SOFIANOS, G. 1998. An Empirical Analysis of NYSE Specialist Trading. *Journal of Financial Economics,* 48**,** 189-210.

MANASTER, S. & RENDLEMAN, J. R. J. 1982. Option Prices as Predictors of Equilibrium Stock Prices. *Journal of Finance,* 37**,** 1043-1057.

MARSHALL, B., NGUYEN, N. H. & VISALTANACHOTI, N. 2013. ETF arbitrage: Intraday evidence. *Journal of Banking & Finance,* 37**,** 3486-3498.

MARTIKAINEN, T. & PERTTUNEN, J. 1995. On The Dynamics of Stock Index Futures and Individual Stock Returns. *Journal of Business Finance & Accounting,* 22**,** 87-100.

NAJAND, M. & MIN, J. H. 1999. A further investigation of the lead-lag relationship between the spot market and stock index futures: Early evidence from Korea. *Journal of Futures Markets,* 19**,** 217-232.

NAM, S. O., OH, S. & KIM, H. K. 2008. The time difference effect of a measurement unit in the lead–lag relationship analysis of Korean financial market. *International Review of Financial Analysis,* 17**,** 259-273.

NAM, S. O., OH, S., KIM, H. K. & KIM, B. C. 2006. An empirical analysis of the price discovery and the pricing bias in the KOSPI 200 stock index derivatives markets. *International Review of Financial Analysis,* 15**,** 398-414.

PANTON, D. 1976. Chicago Board Call Options As Predictors of Common Stock Price Changes. *Journal of Econometrics,* 4**,** 101-113.

PIZZI, M. A., ECONOMOPOULOS, A. J. & O'NEILL, H. M. 1998. An Examination of the Relationship between Stock Index Cash and Futures Markets: A Cointegration Approach. *Journal of Futures Markets,* 18**,** 297-305.

RUAN, J. & MA, T. 2012. Ex-Dividend Day Price Behavior of Exchange-Traded Funds. *Journal of Financial Research,* 35**,** 29-53.

RYU, D. 2015. The information content of trades: An analysis of KOSPI 200 index derivatives. *Journal of Futures Markets,* 35**,** 201-221.

SHARMA, J. L., MOUGOUE, M. & KAMATH, R. 1996. Heteroscedasticity in stock market indicator return data: volume versus GARCH effects. *Applied Financial Economics,* 6**,** 337-342.

SHIN, S. & SOYDEMIR, G. 2010. Exchange-traded funds, persistence in tracking errors and information dissemination. *Journal of Multinational Financial Management,* 20**,** 214-234.

SHYY, G. & VIJAYRAGHAVAN, V. 1996. A Further Investigation of the Lead-Lag Relationship between the Cash Market and Stock Index Futures Market with the Use of Bid/Ask Quotes: the Case of France. *Journal of Futures Markets,* 16**,** 405-420.

STEPHAN, J. A. & WHALEY, R. E. 1990. Intraday Price Change and Trading Volume Relations in the Stock and Stock Option Markets. *Journal of Finance,* 45**,** 191-220.

STOLL, H. R. & WHALEY, R. E. 1990. The Dynamics of Stock Index and Stock Index Futures Returns. *Journal of Financial & Quantitative Analysis,* 25**,** 441-468.

ZHANG, L. 2011. Estimating covariation: Epps effect, microstructure noise. *Journal of Econometrics,* 160**,** 33-47.