

# Dependent Microstructure Noise and Integrated Volatility

## Estimation from High-Frequency Data

Z. Merrick Li\*

Faculty of Economics  
University of Cambridge

Roger J. A. Laeven†

Amsterdam School of Economics  
University of Amsterdam, EURANDOM  
and CentER

Michel H. Vellekoop‡

Amsterdam School of Economics  
University of Amsterdam

October 10, 2019

### Abstract

In this paper, we develop econometric tools to analyze the integrated volatility (IV) of the efficient price and the dynamic properties of microstructure noise in high-frequency data under general dependent noise. We first develop consistent estimators of the variance and autocovariances of noise using a variant of realized volatility. Next, we employ these estimators to adapt the pre-averaging method and derive consistent estimators of the IV, which converge stably to a mixed Gaussian distribution at the optimal rate  $n^{1/4}$ . To improve the finite sample performance, we propose a multi-step approach that corrects the finite sample bias, which turns out to be crucial in applications. Our extensive simulation studies demonstrate the excellent performance of our multi-step estimators. In an empirical study, we analyze the dependence structures of microstructure noise and provide intuitive economic interpretations; we also illustrate the importance of accounting for both the serial dependence in noise and the finite sample bias when estimating IV.

*Keywords:* Dependent microstructure noise, realized volatility, bias correction, integrated volatility, mixing sequences, pre-averaging method.

*JEL classification:* C13, C14, C55, C58.

---

\*Corresponding author. Faculty of Economics, University of Cambridge, Austin Robinson Building, Sidgwick Avenue, Cambridge, CB3 9DD, United Kingdom. Email: [Z.Merrick.Li@gmail.com](mailto:Z.Merrick.Li@gmail.com). Phone: +44 (0) 1223 335288.

†University of Amsterdam, Amsterdam School of Economics, PO Box 15867, 1001 NJ Amsterdam, The Netherlands. Email: [R.J.A.Laeven@uva.nl](mailto:R.J.A.Laeven@uva.nl). Phone: +31 (0)20 5254219.

‡University of Amsterdam, Amsterdam School of Economics, PO Box 15867, 1001 NJ Amsterdam, The Netherlands. Email: [M.H.Vellekoop@uva.nl](mailto:M.H.Vellekoop@uva.nl). Phone: +31 (0)20 5254210.

# 1 Introduction

Over the past decade and a half, high-frequency financial data have become increasingly available. In tandem, the development of econometric tools to study the dynamic properties of high-frequency data has become an important subject area in economics and statistics. A major challenge is provided by the accumulation of market microstructure noise at higher frequencies, which can be attributed to various market microstructure effects including, for example, information asymmetries (see [Glosten and Milgrom \(1985\)](#)), inventory controls (see [Ho and Stoll \(1981\)](#)), discreteness of the data (see [Harris \(1990\)](#)), and transaction costs (see [Garman \(1976\)](#)).

It has been well-established (see, e.g., [Black \(1986\)](#)) that the observed transaction price<sup>1</sup>  $Y$  can be decomposed into the unobservable “efficient price” (or “frictionless equilibrium price”)  $X$  plus a noise component  $U$  that captures market microstructure effects. That is, it is natural to assume that

$$Y_t = X_t + U_t, \tag{1}$$

where further assumptions on  $X$  and  $U$  need to be stipulated. While estimating the IV of the efficient price is a canonical problem in high-frequency financial econometrics (see, for example, [Aït-Sahalia and Jacod \(2014\)](#)), the study of microstructure noise, e.g., its magnitude, dynamic properties, etc., is the main focus of the market microstructure literature (see, for example, [Hasbrouck \(2007\)](#)). A common challenge, however, is that the two components of the observed price  $Y$  in (1) are latent. Therefore, distributional features of one component, say, of the microstructure noise, will affect the estimation of characteristics of the other, such as the IV of the efficient price.<sup>2</sup>

While the semimartingale framework provides the natural class to model the efficient price (see, e.g., [Duffie \(2010\)](#)), the statistical assumptions on noise induced by microeconomic financial models range from simple to very complex, depending on which phenomena the model aims to capture. For example, the classic Roll model (see [Roll \(1984\)](#)) postulates an i.i.d. bid-ask bounce resulting from uncorrelated order flows; [Hasbrouck and Ho \(1987\)](#), [Choi et al. \(1988\)](#), and [Stoll \(1989\)](#) introduce autocorrelated order flows, yielding autoregressive microstructure noise; and [Gross-Kluschmann and Hautsch \(2013\)](#) model microstructure noise with long-memory properties. Therefore, being able to account for the potentially complex statistical behavior of microstructure noise that contaminates our observations of the semimartingale efficient price dynamics, would be an appealing property of any method that aims at disentangling the efficient price and microstructure noise.

---

<sup>1</sup>In this paper, “price” always refers to the “logarithmic price”.

<sup>2</sup>Indeed, while high-frequency data in principle facilitate the asymptotic and empirical analysis of volatility estimators, the pronounced presence of microstructure noise at high frequency subverts the desirable properties of traditional estimators such as realized volatility.

To estimate the IV of the efficient price, several de-noise methods have been developed, mostly assuming i.i.d. microstructure noise. Examples include the two-scale and multi-scale realized volatility estimators developed in Zhang et al. (2005) and Zhang (2006), the likelihood approach initiated by Aït-Sahalia et al. (2005) and Xiu (2010), the realized kernel methods developed in Barndorff-Nielsen et al. (2008), and the pre-averaging method developed in a series of papers by Podolskij and Vetter (2009b) and Jacod et al. (2009, 2010), see also Podolskij and Vetter (2009a). The variance of noise is usually obtained as a by-product.

In this paper, we allow the microstructure noise to be serially dependent in a general setting, nesting many special cases (including independence). We do not impose any parametric restrictions on the distribution of the noise, except for some rather general mixing conditions that guarantee the existence of limit distributions, hence our approach is essentially nonparametric. In this setting, we first derive the stochastic limit of the realized volatility of observed prices after  $j$  lags. Using this limit result, we develop consistent estimators of the variance and covariances of noise. The aim of estimating the second moments of noise is twofold. On the one hand, we would like to explore the dynamic properties of microstructure noise. In particular, we would like to compare these properties to those induced by various parametric models of microstructure noise based on leading microstructure theory, and obtain corresponding economic interpretations to achieve a better understanding of the microstructure effects in high-frequency data. On the other hand, the second moments of noise become nuisance parameters when estimating the IV, which is a prime objective in the analysis of high-frequency financial data.

To estimate the IV, we next adapt the pre-averaging approach (PAV) to allow for serially dependent noise in our general setting, first based on non-overlapping sampling blocks and next based on overlapping sampling blocks, in both cases using general weight functions (i.e., general kernels). We find that the stochastic limits of the adapted PAV estimators are functions of the volatility and the variance and covariances of noise, and the latter, constituting an *asymptotic bias*, can be consistently estimated by our realized volatility estimator. Hence, we can correct the asymptotic bias, resulting in centered estimators of the IV, for which we establish the associated central limit theorems.

A key interest in this paper is to unravel the interplay between asymptotic and finite sample biases when estimating the IV. In a formal finite sample analysis, we find that the realized volatility estimator has a finite sample bias that is proportional to the IV. This bias term becomes significant when the number of lags (in computing the variant of realized volatility) is large, or the noise-to-signal ratio<sup>3</sup> is small. Therefore, we are in a situation in which the IV generates a *finite sample bias* to the estimators of the second moments of noise, while the latter introduce an *asymptotic bias* when estimating the former. This “feedback effect” in the bias corrections motivates us to develop *multi-step estimators*.

---

<sup>3</sup>That is, the ratio of the variance of noise and the IV.

First, we simply ignore the dependence in noise and proceed with the pre-averaging method to obtain an estimator of the IV. Next, we use this estimator to obtain *finite sample bias* corrected estimators of the second moments of noise, which can then be used to correct the asymptotic bias, yielding the second-step estimator of the IV. Repeating this process leads to three-step estimators (and beyond). Figure 1 gives a simple graphical illustration of the implementation of the multi-step estimators. We establish consistency and a central limit theorem for our multi-step estimators.

We conduct extensive Monte Carlo experiments to examine the performance of our estimators, which proves to be excellent. We demonstrate in particular that they can accommodate both serially dependent and independent noise and perform well in finite samples with realistic data frequencies and sample sizes. The experiments reveal the importance of a unified treatment of asymptotic and finite sample biases when estimating IV.

Empirically, we apply our new estimators to a sample of Citigroup transaction data. We find that the associated microstructure noise tends to be positively autocorrelated. This is in line with earlier findings in the microstructure literature, see [Hasbrouck and Ho \(1987\)](#), [Choi et al. \(1988\)](#), and [Huang and Stoll \(1997\)](#). When we attribute this positive autocorrelation to order flow continuation, the estimated probability that a buy (or sell) order follows another buy (or sell) order is found to be 0.87. Furthermore, microstructure noise turns out to be negatively autocorrelated under tick time sampling. This is consistent with inventory models, in which dealers alternate quotes to maintain their inventory position. We obtain an estimate of the probability of reversed orders equal to 0.84. Turning to the estimators of IV, we find that with positively autocorrelated noise the commonly adopted methods that hinge on the i.i.d. assumption of noise tend to overestimate the IV. Under two alternative (sub)sampling schemes our estimators also appear to work well. This testifies to the critical relevance of the bias corrections embedded in our multi-step estimators.

In earlier literature, [Aït-Sahalia et al. \(2011\)](#) show that the two-scale and multi-scale realized volatility estimators are robust to exponentially decaying dependence in noise. In this paper, we provide explicit estimators of the second moments of noise and analyze their asymptotic behavior, develop bias-corrected estimators of the IV based on these moments of noise, and empirically assess the noise characteristics. Furthermore, [Hautsch and Podolskij \(2013\)](#) study  $q$ -dependent microstructure noise, develop consistent estimators of the first  $q$  autocovariances of microstructure noise and define the associated pre-averaging estimators. An appealing feature of their approach is that their autocovariance-type estimators of  $q$ -dependent noise consider non-overlapping increments which avoids finite sample bias. We allow for more general assumptions on the dependence structure of microstructure noise. Owing to its generality our setting incorporates many microstructure models as special cases. We therefore do not need to advocate any particular model of microstructure noise.

In two contemporaneous works, [Jacod et al. \(2017, 2019\)](#) also study dependent noise in high-frequency data. In [Jacod et al. \(2017\)](#), they develop a novel local averaging method to “recover” the noise and can, in principle, estimate any finite (joint) moments of noise with diurnal features. Moreover, they also allow observation times to be random. Empirically, they find some interesting statistical properties of noise. In particular, they find that noise is strongly serially dependent, with polynomially decaying autocorrelations. Employing this local averaging method, [Jacod et al. \(2019\)](#) develop an estimator of IV that allows for dependent noise. The local averaging method differs from, and allows to analyze more general noise characteristics than, the simpler realized volatility method developed here. The key difference is our explicit treatment of the feedback effect between the asymptotic and finite sample biases: we show that in a finite sample, the IV and second moments of microstructure noise should be estimated in a unified way, since they induce biases in each other. We design novel and easily implementable multi-step estimators to correct for the intricate biases. Our multi-step estimators of the IV, which are designed to allow for dependent noise, also perform well in the special case of independent noise, and in a sample of reasonable size as encountered in practice. This robustness to (mis)specification of noise and to sampling frequencies is an important advantage of our multi-step estimators. Our unified treatment of the asymptotic and finite sample biases may help explain why the empirical studies in [Jacod et al. \(2017\)](#) render the strong dependence in noise they find (and question themselves); see our empirical analysis in [Section 7](#).

In another independent paper, [Da and Xiu \(2019\)](#) introduce a novel quasi maximum likelihood approach to estimate both the volatility and the autocovariances of moving-average microstructure noise. They also extend their estimators to general settings that allow for irregular observation times, intraday patterns of noise and jumps in asset prices. Their approach treats “large” and “small” microstructure noise in a uniform way which leads to a potential improvement in the convergence rate. Our approach is essentially of a nonparametric nature and provides unified estimators of a class of volatility functionals (see [Theorem 4.1](#)) including the asymptotic variance, which account for the feedback between finite sample and asymptotic biases. Our empirical study also has a different focus. Our investigation is not as extensive as in [Da and Xiu \(2019\)](#),<sup>4</sup> but we explicitly consider different sampling frequencies,<sup>5</sup> analyzing the autocovariance patterns of noise in connection to microstructure noise models and their impact on IV estimation.

The remainder of this paper is organized as follows. In [Section 2](#), we introduce the basic setting and notation. In [Section 3](#), we analyze realized volatility with dependent noise and develop consistent estimators of the second moments of noise. The pre-averaging method with dependent noise is studied

---

<sup>4</sup>Da and Xiu maintain a website to provide up-to-date daily annualized volatility estimates for all S&P 1500 index constituents, see <http://dachxiu.chicagobooth.edu/#risklab>.

<sup>5</sup>In their empirical studies, [Da and Xiu \(2019\)](#) only consider tick time sampling.

in Section 4. Section 5 introduces our multi-step estimators. Section 6 reports extensive simulation studies. Our empirical study is presented in Section 7. Section 8 concludes the paper. All proofs and some additional Monte Carlo simulation and empirical results are collected in an online appendix, see Li et al. (2019).

## 2 Framework and Assumptions

We state the following assumption regarding the efficient log-price process:

**Assumption 2.1** (Efficient log-price). *The efficient log-price process  $X$  follows a continuous Itô semimartingale defined on a filtered probability space  $(\Omega, \mathcal{F}, (\mathcal{F}_t)_{t \geq 0}, \mathbb{P})$ :*

$$X_t = X_0 + \int_0^t b_s ds + \int_0^t \sigma_s dW_s, \quad (2)$$

where  $W$  is a standard Brownian motion, the drift process  $b_s$  is optional and locally bounded, and the volatility process  $\sigma_s$  is adapted with càdlàg paths.

We assume that all price observations are collected in the fixed time interval  $[0, T]$ , where without losing generality we let  $T = 1$ . We let  $n + 1$  be the number of observations and denote  $\Delta_n = 1/n$ . The observation times are given by  $t_i^n = i\Delta_n, i = 0, \dots, n$ . We make the following assumption regarding the market microstructure noise:

**Assumption 2.2** (Market microstructure noise). *The noise process  $(U_i)_{i \in \mathbb{N}}$  is defined on the probability space  $(\Omega^{(0)}, \mathcal{G}, \mathbb{P}^{(0)})$ , which has discrete filtrations  $\mathcal{G}_i = \sigma(U_k : k \leq i)$ ,  $\mathcal{G}^i = \sigma(U_k : k \geq i)$  that satisfy  $\mathcal{G} = \mathcal{G}^\infty = \mathcal{G}_\infty$ . Moreover, we assume*

1.  *$U$  is stationary and  $\rho$ -mixing and the mixing coefficients<sup>6</sup>  $\{\rho_h\}_{h=1}^\infty$  decay at a polynomial rate, i.e., there exist some constants  $C > 0, v > 0$  such that*

$$\rho_h \leq \frac{C}{h^v}. \quad (3)$$

2.  *$v > 1$ ,  $\mathbb{E}(U) = 0$ , and all moments of noise exist.*

The mixing conditions in Assumption 2.2 item (1.) ensure that the noise process evaluated at different time instances, say,  $i$  and  $i + h$ , is increasingly limited in dependence when the lag  $h$  increases. In

<sup>6</sup>The mixing coefficients constitute a sequence satisfying

$$\rho_h = \sup \left\{ |\mathbb{E}(V_k V_{k+h})| : \mathbb{E}(V_k) = \mathbb{E}(V_{k+h}) = 0, \|V_k\|_2 \leq 1, \|V_{k+h}\|_2 \leq 1, V_k \in \mathcal{G}_i, V_{k+h} \in \mathcal{G}^{i+h} \right\}.$$

We refer to Bradley (2007) or Chapter VIII of Jacod and Shiryaev (2003) for further details on and properties of mixing sequences.

particular, that assumption implies that there exists some  $C' > 0$  such that

$$|\gamma(h)| \leq \frac{C'}{h^v}, \quad (4)$$

where  $\gamma(h) = \mathbf{Cov}(U_i, U_{i+h})$  is the autocovariance function of  $U$ . We assume all moments of noise exist because this is required for the validity of Theorem 4.1 below for any even integer  $r \geq 2$ .

At stage  $n$ , we will denote  $U_i$  by  $U_i^n$ ,  $\forall i \leq n$ . The  $i$ -th observed price is thus given by

$$Y_i^n = X_i^n + U_i^n, \quad (5)$$

where  $X_i^n = X_{i\Delta_n}$ .

**Remark 2.1** (Microstructure noise in discrete time). *We allow the noise process  $U$  to generate dependencies in sampling time. Hence, our noise process essentially constitutes a discrete-time model — it does not depend explicitly on the time between successive observations. Aït-Sahalia et al. (2005), Hansen and Lunde (2006), and Hansen et al. (2008) study various continuous-time models of dependent microstructure noise. In these continuous-time models, the noise component of a log-return over a time interval  $\Delta$  is of order  $O_p(\sqrt{\Delta})$ , the same order as the logarithmic return of the efficient price. Our theory focuses primarily on sampling in calendar time.<sup>7</sup> In our simulations and empirical work, we also analyze sampling in transaction time,<sup>8</sup> and tick time.<sup>9</sup>*

**Remark 2.2** (General dynamic properties of microstructure noise). *Our assumptions on the dependence of noise are quite general, nesting many models as special cases including, for example, i.i.d. noise,  $q$ -dependent noise (i.e.,  $\gamma(h) = 0$ ,  $\forall h > q$ ), ARMA( $p, q$ ) noise (see Mokkadem (1988)) and some long-memory processes (see Tsay (2005)). We note that AR(1) and AR(2) noise are studied in Barndorff-Nielsen et al. (2008) and Hendershott et al. (2013) respectively,  $q$ -dependent noise is considered by Hansen et al. (2008) and Hautsch and Podolskij (2013), while Gross-Kluschmann and Hautsch (2013) study long-memory bid-ask spreads.*

*Another recent strand of the literature explores the variety of microstructure data including observable information, seeking to parameterize the microstructure noise; see Li et al. (2016), Chaker (2017), Clinet and Potiron (2017) and Clinet and Potiron (2019). The parametrization allows for rich dynamics of the microstructure noise and at the same time improves the convergence rates of associated volatility*

<sup>7</sup>Under this sampling scheme,  $Y_i^n$  (resp.  $X_i^n, U_i^n$ ) is the observed log-price (resp. efficient log-price, microstructure noise) at regular time  $i\Delta_n$ , with  $\Delta_n = 1/n$  in the main text.

<sup>8</sup>Under this sampling scheme,  $Y_i^n$  (resp.  $X_i^n, U_i^n$ ) is the observed log-price (resp. efficient log-price, microstructure noise) associated with the  $i$ -th trade. The observation times  $(t_i^n)_{0 \leq i \leq n}$  can, in general, be deterministic or random, and regular or irregular.

<sup>9</sup>Tick time sampling removes all zero returns; see Aït-Sahalia et al. (2011) and Griffin and Oomen (2008). Hence,  $Y_i^n$  is by definition different from  $Y_{i-1}^n$  and  $Y_{i+1}^n$  under this sampling scheme.

estimators. Specifically, the noise component in these models can be serially correlated. The correlation is attributed to persistent observable quantities, e.g., trading volume and trading directions, that constitute the “observable part” of the microstructure noise. By contrast, we introduce an essentially nonparametric model of microstructure noise, without singling out the sources of the noise.

### 3 Estimation of the Variance and Covariances of Noise

In this section, we develop consistent estimators of the second moments of noise under Assumptions 2.1 and 2.2. These estimators will later serve as important inputs to adapt the pre-averaging method. We also analyze our estimators’ finite sample properties.

#### 3.1 Realized volatility with dependent noise

We start with the following preliminary result:

**Proposition 3.1.** *Assume that the efficient log-price satisfies Assumption 2.1, the observations follow (5), the noise process satisfies Assumption 2.2, and that  $\mathcal{G}$  is independent of  $\mathcal{F}$ . Let  $j > 0$  be a fixed integer and assume the sequence of integers  $j_n$  satisfies  $j_n \rightarrow \infty$ ,  $j_n \Delta_n \rightarrow 0$ . Then we have the following convergences in probability as  $n \rightarrow \infty$ :*

$$\widehat{\langle Y, Y \rangle}(j)_n := \frac{\sum_{i=0}^{n-j} (Y_{i+j}^n - Y_i^n)^2}{2(n-j+1)} \xrightarrow{\mathbb{P}} \gamma(0) - \gamma(j), \quad (6)$$

$$\widehat{\gamma(0)}_n := \frac{\sum_{i=0}^{n-j_n} (Y_{i+j_n}^n - Y_i^n)^2}{2(n-j_n+1)} \xrightarrow{\mathbb{P}} \gamma(0), \quad (7)$$

$$\widehat{\gamma(j)}_n := \widehat{\gamma(0)}_n - \widehat{\langle Y, Y \rangle}(j)_n \xrightarrow{\mathbb{P}} \gamma(j). \quad (8)$$

The special case of (6) that occurs when  $j = 1$  appears in Ait-Sahalia et al. (2011) assuming exponential decay. We also note that in the most recent version of Jacod et al. (2017) similar estimators as  $\widehat{\langle Y, Y \rangle}(j)_n$  are mentioned but without formal analysis of their limiting behavior. To our best knowledge, our paper is the first to estimate the variance and covariances of noise using realized volatility under a general dependent noise setting.

#### 3.2 Finite sample bias correction

The theoretical validity of our realized volatility estimators in (6)–(8) hinges on the increasing availability of observations in a fixed time interval, the so-called *infill asymptotics*. In general, an estimator derived from asymptotic results can, however, behave very differently in finite samples. Our realized volatility



estimators of the second moments of noise are an example for which the asymptotic theory provides a poor representation of the estimators' finite sample behavior.<sup>10</sup>

Intuitively, the finite sample bias stems from the diffusion component, when computing the realized volatility  $\widehat{\langle Y, Y \rangle}(j)_n$  over large lags  $j$  in a finite sample, and we will explain later (e.g., in Remark 3.3) why it is critically relevant to account for it in real applications. In the remainder of this section, we assume the drift  $b_t$  in (2) to be zero. As shown by, for example, Bandi and Russell (2008) and Lee and Mykland (2012) this is not restrictive in high-frequency analysis. This will be confirmed in our Monte Carlo simulation studies in Section 6 and Appendix B.

**Proposition 3.2.** *Assume that the efficient log-price follows (2) with  $b_s = 0 \forall s$ , and assume there is some  $\delta > 0$  so that  $\sigma_t$  is bounded for all  $t \in [0, \delta] \cup [1 - \delta, 1]$ . Furthermore, assume the observations follow (5), the noise process satisfies Assumption 2.2 and  $\mathcal{G}$  is independent of  $\mathcal{F}$ . Then,*

$$\mathbb{E}_\sigma \left( \widehat{\langle Y, Y \rangle}(j)_n \right) = \frac{j\text{IV}}{2(n-j+1)} + \gamma(0) - \gamma(j) + O_p(j^2/n^2), \quad (9)$$

where  $\text{IV} := \int_0^1 \sigma_t^2 dt$  is the integrated volatility. Here,  $\mathbb{E}_\sigma(\cdot)$  denotes the expectation conditional on the volatility path.

**Remark 3.1.** *If  $\sigma_t$  is locally bounded, then the assumptions on  $\sigma_t$  required for Proposition 3.2 will hold. The regularity conditions with respect to  $\sigma_t$  in Proposition 3.2 trivially hold if the volatility process is assumed to be continuous. (Volatility is usually assumed to be continuous when making finite sample bias corrections.)*

**Remark 3.2.** *Let  $j = 1$ . In that special case the result in Proposition 3.2 bears similarities to Theorem 1 in Hansen and Lunde (2006). Contrary to Hansen and Lunde (2006) we assume that the efficient log-price  $X$  is independent of the noise  $U$ . Therefore, any correlations between the two drop out.*

Proposition 3.2 reveals that  $\widehat{\langle Y, Y \rangle}(j)_n - \frac{j\text{IV}}{2(n-j+1)}$  will be a better estimator of  $\gamma(0) - \gamma(j)$  in finite samples than  $\widehat{\langle Y, Y \rangle}(j)_n$ , and this motivates the following finite sample bias corrected estimators:

$$\widehat{\langle Y, Y \rangle}^{(\text{adj})}(j)_n := \widehat{\langle Y, Y \rangle}(j)_n - \frac{\hat{\sigma}^2 j}{2(n-j+1)}; \quad (10)$$

$$\widehat{\gamma(0)}_n^{(\text{adj})} := \widehat{\gamma(0)}_n - \frac{\hat{\sigma}^2 j_n}{2(n-j_n+1)}; \quad (11)$$

$$\widehat{\gamma(j)}_n^{(\text{adj})} := \widehat{\gamma(0)}_n^{(\text{adj})} - \widehat{\langle Y, Y \rangle}^{(\text{adj})}(j)_n; \quad (12)$$

where  $\hat{\sigma}^2$  is an estimator of IV. We note that the bias corrected estimators are still consistent, as the

<sup>10</sup>This applies to the local averaging estimators developed in Jacod et al. (2017) as well; see Footnote 13 for further details.

fraction  $\frac{j}{n-j+1}$  is negligible when  $j$  is much smaller than  $n$ .

**Remark 3.3** (Why the finite sample bias matters). *We now explain why the finite sample bias correction is crucial in applications. We first rewrite (9):*

$$\begin{aligned}\mathbb{E}_\sigma\left(\widehat{\langle Y, Y \rangle}(j)_n\right) &= \frac{j\text{IV}}{2(n-j+1)} + \gamma(0) - \gamma(j) + O_p(j^2/n^2) \\ &= (\gamma(0) - \gamma(j)) \left(1 + \frac{\frac{j}{2(n-j+1)}}{\frac{\gamma(0) - \gamma(j)}{\text{IV}}}\right) + O_p(j^2/n^2).\end{aligned}\tag{13}$$

Observe that the finite sample bias is determined by the ratio of the two terms  $\frac{j}{2(n-j+1)}$  and  $\frac{\gamma(0) - \gamma(j)}{\text{IV}}$ . The first term,  $\frac{j}{2(n-j+1)}$ , depends on the data frequency ( $n$ ) and “target parameters” ( $j$ ); the second term,  $\frac{\gamma(0) - \gamma(j)}{\text{IV}}$ , is the (latent) noise-to-signal ratio. If the second term is “relatively larger (smaller)” than the first one, then the finite sample bias will be small (large). In other words, the finite sample bias is not only determined by the data frequency and target parameters, but also by other properties of the underlying efficient price and noise processes.

In high-frequency financial data, the noise-to-signal ratio  $\frac{\gamma(0)}{\text{IV}}$  is typically small, but it can vary from  $O(10^{-2})$  (see [Bandi and Russell \(2006\)](#)) to  $O(10^{-6})$  (see [Christensen et al. \(2014\)](#)) in empirical studies. The ratio  $\frac{j}{2(n-j+1)}$ , while typically small as well, can still be relatively large, depending on the specific situation. Consider the following two scenarios:

- 1) We have ultra high-frequency data with  $n = O(10^5)$  (recall that the number of seconds in a business day is 23,400), and we select  $j_n = 20$ . Then, the ratio  $\frac{j_n}{2(n-j_n+1)} = O(10^{-4})$ .
- 2) We have i.i.d. noise and we would like to estimate the variance of noise by  $\widehat{\langle Y, Y \rangle}(1)_n$  using high-frequency data with average duration of 20 seconds (thus  $n \approx 10^3$ ); see, e.g., [Bandi and Russell \(2006\)](#). Hence,  $\frac{j}{2(n-j+1)} = O(10^{-3})$ .

In both scenarios, the ratio of  $\frac{j}{2(n-j+1)}$  and  $\frac{\gamma(0) - \gamma(j)}{\text{IV}}$  can vary widely, depending on the magnitude of the latent noise-to-signal ratio. It is then clear from the first line of (13) that the finite sample bias term, which is proportional to the IV, may well wipe out the variance of noise, depending on the specific situation.

**Remark 3.4.** Note that increasing the sample size by extending the time horizon to  $[0, T]$  with large  $T$  will not remove the finite sample bias. Hence, the finite sample bias may be viewed as a low frequency bias.

Throughout the remainder of this paper, we assume the following conditions hold:<sup>11</sup>

$$v > 3, \quad j_n \asymp \Delta_n^{-\delta}, \quad \ell_n \asymp \Delta_n^{-\kappa}, \quad \delta \in \left(\frac{5}{36}, \frac{1}{6}\right), \quad \kappa \in \left(\frac{1}{8}, \frac{1}{6}\right), \quad (14)$$

with  $\ell_n$  another sequence of integers. The following proposition provides an estimator of the “long-run variance” of microstructure noise. As we shall see later, the long-run variance of noise appears as an asymptotic bias in the de-noise method developed in this paper.

**Proposition 3.3.** *Assume that the efficient log-price satisfies Assumption 2.1, the observations follow (5), the noise process satisfies Assumption 2.2 and  $\mathcal{G}$  is independent of  $\mathcal{F}$ . Define*

$$\widehat{\Sigma}_{U_n} := \widehat{\gamma(0)}_n + 2 \sum_{j=1}^{\ell_n} \widehat{\gamma(j)}_n, \quad (15)$$

where  $\widehat{\gamma(0)}_n$  and  $\widehat{\gamma(j)}_n$  are defined in (7) and (8). Then,

$$\widehat{\Sigma}_{U_n} \xrightarrow{\mathbb{P}} \Sigma_U, \quad (16)$$

where

$$\Sigma_U = \gamma(0) + 2 \sum_{j=1}^{\infty} \gamma(j). \quad (17)$$

For i.i.d. noise,  $\Sigma_U$  reduces to  $\gamma(0)$ , and it is known (see Zhang et al. (2005) and Bandi and Russell (2008)) that the variance of noise can then be consistently estimated by the standardized realized volatility of observed returns. However, when noise is dependent we face a much more complex situation: all variance and covariance terms constitute  $\Sigma_U$ . Nevertheless, Proposition 3.3 above provides a consistent estimator of  $\Sigma_U$ .

## 4 The Pre-Averaging Method with Dependent Noise

In this section, we adapt a popular “de-noise” method — the pre-averaging method — to allow for serially dependent noise in our general setting. The pre-averaging method was originally introduced by Podolskij and Vetter (2009b) (see also Jacod et al. (2009), Jacod et al. (2010), Podolskij and Vetter (2009a), Hautsch and Podolskij (2013), and the textbook treatment in Aït-Sahalia and Jacod (2014)). We first construct our pre-averaged statistics based on non-overlapping sampling blocks and next based

<sup>11</sup>Some results, e.g., Proposition 3.3, hold already under weaker conditions. The conditions (14) are, however, needed to establish our main theorems in the next sections.

on overlapping sampling blocks, in both cases using general weight functions.

#### 4.1 Pre-averaging based on non-overlapping intervals

Let  $k_n$  be a sequence of integers, with  $k_n \rightarrow \infty$  and  $k_n \Delta_n \rightarrow 0$  as  $n \rightarrow \infty$ , satisfying

$$\sqrt{\Delta_n} k_n = \theta + o(\Delta_n^{1/4}), \quad (18)$$

where  $\theta > 0$  is a constant. Furthermore, let  $g$  be a general kernel (i.e., weight function). We assume  $g$  is continuous, piecewise  $C^1$  with a piecewise Lipschitz derivative  $g'$ , and satisfies  $g(s) = 0, \forall s \notin (0, 1)$ , and  $\int_0^1 g^2(s) ds > 0$ , as in [Jacod et al. \(2009\)](#). We introduce the following notation associated with  $g$ :

$$\begin{cases} g_i^n = g(i/k_n); & \bar{g}_i^n = g_{i+1}^n - g_i^n; \\ \phi_0^n = \frac{1}{k_n} \sum_{i \in \mathbb{Z}} (g_i^n)^2; & \phi_1^n(j) = k_n \sum_{i \in \mathbb{Z}} \bar{g}_i^n \bar{g}_{i-j}^n; \\ \phi_0(s) = \int_s^1 g(u)g(u-s)du; & \phi_1(s) = \int_s^1 g'(u)g'(u-s)du; \\ \Phi_{ij} = \int_0^1 \phi_i(s)\phi_j(s)ds, & \psi_i = \phi_i(0), \quad i, j \in \{0, 1\}. \end{cases}$$

**Example 4.1** (Triangular kernel). *A simple canonical example of  $g$  is given by the triangular kernel  $g(x) = x \wedge (1 - x)$ . Then,*

$$\psi_0 = 1/12, \quad \psi_1 = 1, \quad \Phi_{00} = 151/80640, \quad \Phi_{01} = 1/96, \quad \Phi_{11} = 1/6.$$

For any sequence  $\{Z_i^n\}_{i=0}^n$ , denote  $\Delta_i^n Z = Z_i^n - Z_{i-1}^n, i = 1, 2, \dots$ , and let its pre-averaged value be given by

$$\bar{Z}_i^n := \sum_{j=1}^{k_n-1} g_j^n \Delta_{i+j}^n Z = - \sum_{j=0}^{k_n-1} \bar{g}_j^n Z_{i+j}^n, \quad i = 0, 1, \dots \quad (19)$$

Furthermore, let  $M_n = \lfloor \frac{n}{k_n} \rfloor$ , where  $\lfloor \cdot \rfloor$  is the floor function. For any real  $r \geq 2$ , the pre-averaged statistics of the log-price process  $Y$  based on *non-overlapping intervals* are defined as follows:

$$\text{PAV}(Y, r)_n := n^{\frac{r-2}{4}} \sum_{m=0}^{M_n-1} \left| \bar{Y}_{mk_n}^n \right|^r, \quad r \geq 2. \quad (20)$$

Under our general setting of dependent noise, we establish in the following results first a consistency theorem for the general functional form of the pre-averaged statistics, based on which we derive a consistent estimator of the IV, and next a central limit theorem providing the associated limit distribution, with a consistent estimator of the asymptotic variance.

**Theorem 4.1.** *Assume that the efficient log-price satisfies Assumption 2.1, the observations follow (5), and the noise process satisfies Assumption 2.2. Furthermore, assume  $\mathcal{G}$  and  $\mathcal{F}$  are independent. Then, for any even integer  $r \geq 2$ ,*

$$\text{PAV}(Y, r)_n \xrightarrow{\mathbb{P}} \text{PAV}(Y, r) := \frac{\mu_r}{\theta} \int_0^1 \left( \theta \psi_0 \sigma_s^2 + \frac{\psi_1}{\theta} \Sigma_U \right)^{\frac{r}{2}} ds, \quad (21)$$

where  $\Sigma_U$  is defined in (17) and  $\mu_r = \mathbb{E}(Z^r)$  for a standard normal random variable  $Z$ .

Aided by this result, we obtain consistent estimators of the IV and the integrated quarticity  $\text{IQ} := \int_0^1 \sigma_s^4 ds$ , as follows:

**Corollary 4.1.** *Under the assumptions of Theorem 4.1, we have the following consistency result for the IV and the IQ:*

$$\widehat{\text{IV}}_n := \frac{\text{PAV}(Y, 2)_n}{\psi_0} - \frac{\psi_1 \widehat{\Sigma}_{U_n}}{\psi_0 \theta^2} \xrightarrow{\mathbb{P}} \text{IV}, \quad (22)$$

$$\widehat{\text{IQ}}_n := \frac{\text{PAV}(Y, 4)_n}{3\psi_0^2 \theta} - \frac{2\psi_1 \widehat{\Sigma}_{U_n} \widehat{\text{IV}}_n}{\psi_0 \theta^2} - \frac{\psi_1^2 (\widehat{\Sigma}_{U_n})^2}{\theta^4 \psi_0^2} \xrightarrow{\mathbb{P}} \text{IQ}, \quad (23)$$

where  $\widehat{\Sigma}_{U_n}$  is defined in (15).

**Theorem 4.2.** *Assume all conditions in Theorem 4.1 hold. Furthermore, assume that the process  $\sigma$  is a continuous Itô semimartingale. Then,*

$$\Delta_n^{-\frac{1}{4}} \left( \widehat{\text{IV}}_n - \text{IV} \right) \xrightarrow{\mathcal{L}^{-s}} \sqrt{\frac{2}{\theta \psi_0^2}} \int_0^1 \left( \theta \psi_0 \sigma_s^2 + \frac{\psi_1}{\theta} \Sigma_U \right) dW'_s, \quad (24)$$

where  $\xrightarrow{\mathcal{L}^{-s}}$  denotes stable convergence in law and where  $W'$  is a standard Wiener process independent of  $\mathcal{F}$ . Moreover, letting  $\widehat{\Sigma}_n := 2\text{PAV}(Y, 4)_n / 3\psi_0^2$ , we have that  $\Delta_n^{-\frac{1}{4}} \left( \widehat{\text{IV}}_n - \text{IV} \right) / \sqrt{\widehat{\Sigma}_n}$  converges stably in law to a standard normal random variable, which is independent of  $\mathcal{F}$ .

A main advantage of the pre-averaging approach and the associated estimators introduced in this section is their simplicity. In fact, we obtain from Theorem 4.1 a class of consistent estimators of  $\int_0^1 \sigma_s^r ds$  with arbitrary even integer  $r$ , since we have a consistent estimator of  $\Sigma_U$ . When only estimation of the IV is concerned, this leads to a simple estimator of the asymptotic variance of the IV estimator.<sup>12</sup>

<sup>12</sup>Our simulation experiments presented later show that, compared to the pre-averaging estimators based on overlapping intervals introduced in the next subsection, the pre-averaging estimators based on non-overlapping intervals often deliver a somewhat smaller bias, although their standard deviations are typically somewhat larger.

## 4.2 Pre-averaging based on overlapping intervals

Now we extend our pre-averaging estimator of the IV in two directions. First, we allow for *overlapping intervals* to conduct pre-averaging; second, we accommodate more general stochastic volatility processes when deriving the respective limit distribution. (We recall that we assumed the process  $\sigma$  to be a continuous Itô semimartingale in Theorem 4.2.)

In particular, we propose the following estimator of the IV, with  $\widehat{\Sigma}_{U_n}$  as introduced in (15):

$$\widetilde{\text{IV}}_n := \frac{\sqrt{\Delta_n}}{\theta\psi_0} \sum_{i=0}^{n-k_n+1} \left(\overline{Y}_i^n\right)^2 - \frac{\psi_1 \widehat{\Sigma}_{U_n}}{\theta^2 \psi_0}. \quad (25)$$

**Theorem 4.3.** *Assume that the efficient log-price satisfies Assumption 2.1, the observations follow (5), the noise process satisfies Assumption 2.2, and  $\mathcal{G}$  is independent of  $\mathcal{F}$ . Then,*

$$\Delta_n^{-\frac{1}{4}} \left(\widetilde{\text{IV}}_n - \text{IV}\right) \xrightarrow{\mathcal{L}-\xi} \Upsilon_1, \quad (26)$$

with  $\Upsilon_t = \int_0^t V_s dW'_s$ , where  $W'$  is a standard Wiener process independent of  $\mathcal{F}$ , and where  $V_t$  satisfies

$$V_t^2 := \frac{4}{\psi_0^2} \left( \Phi_{00} \theta \sigma_t^4 + 2\Phi_{01} \frac{\sigma_t^2 \Sigma_U}{\theta} + \frac{\Phi_{11} \Sigma_U^2}{\theta^3} \right). \quad (27)$$

**Remark 4.1.** *The tuning parameter  $\theta$  (recall (18)) can be chosen such that it minimizes the asymptotic variance, which will improve the efficiency of our estimators. The optimal  $\theta$  is given by*

$$\theta^* = \left( \frac{\sqrt{\Phi_{01}^2 \text{IV}^2 \Sigma_U^2 + 3\Phi_{00} \Phi_{11} \Sigma_U^2 \text{IQ}} + \Phi_{01} \text{IV} \Sigma_U}{\Phi_{00} \text{IQ}} \right)^{1/2}. \quad (28)$$

The optimal choice of  $\theta$  requires an estimate of IQ. Therefore, we provide a consistent estimator, as follows:

$$\widetilde{\text{IQ}}_n := \frac{\sum_{i=0}^{n-k_n+1} \left(\overline{Y}_i^n\right)^4}{3\theta^2 \psi_0^2} - \frac{2\psi_1 \widehat{\Sigma}_{U_n} \widetilde{\text{IV}}_n}{\psi_0 \theta^2} - \frac{\psi_1^2 \left(\widehat{\Sigma}_{U_n}\right)^2}{\theta^4 \psi_0^2} \xrightarrow{\mathbb{P}} \text{IQ}. \quad (29)$$

Note that that  $\widetilde{\text{IQ}}_n$  is analogous to  $\widehat{\text{IQ}}_n$  introduced in (23).

To apply the limit distribution result in Theorem 4.3 above to construct confidence intervals, we need a consistent estimator of the asymptotic variance  $\int_0^1 V_t^2 dt$ . Among other possibilities, we propose the

following:

$$\tilde{\Sigma}_n := \frac{4\Phi_{00}}{3\theta\psi_0^4} \sum_{i=0}^{n-k_n+1} (\bar{Y}_i^n)^4 + \frac{8\widehat{\Sigma}_{U_n}\widetilde{V}_n}{\theta\psi_0^2} \left( \Phi_{01} - \frac{\psi_1\Phi_{00}}{\psi_0} \right) + \frac{4(\widehat{\Sigma}_{U_n})^2}{\theta^3\psi_0^2} \left( \Phi_{11} - \frac{\psi_1^2\Phi_{00}}{\psi_0^2} \right). \quad (30)$$

**Corollary 4.2.** *Under the assumptions of Theorem 4.3, we have*

$$\tilde{\Sigma}_n \xrightarrow{\mathbb{P}} \int_0^1 V_t^2 dt. \quad (31)$$

Therefore, the sequence  $\Delta_n^{-\frac{1}{4}} (\widetilde{V}_n - \text{IV}) / \sqrt{\tilde{\Sigma}_n}$  converges stably in law to a standard normal random variable, which is independent of  $\mathcal{F}$ .

**Remark 4.2** (Irregular Observation Schemes). *We note that, following similar arguments as in Jacod and Mykland (2015), our results, in particular those in Theorem 4.3, extend to (i.e., are robust to) mildly irregular observation schemes, as follows. Let  $\mathcal{T}$  be a function with strictly positive Lipschitz derivative. Assume  $\mathcal{T}(0) = 0$  and  $\mathcal{T}(1) = 1$ . Now let  $\tilde{t}_i^n := \mathcal{T}(i\Delta_n)$ . Such irregular observation schemes have been considered e.g., by Barndorff-Nielsen et al. (2008) and Mykland and Zhang (2012).*

First, we note that such a time transformation theoretically does not affect the microstructure noise process, as the noise is a discrete-time process that does not depend on the time between successive observations. Thus, under the new observation scheme, we have that

$$Y_{\tilde{t}_i^n} = X_{\tilde{t}_i^n} + U_i^n. \quad (32)$$

Denote the time-transformed efficient price process by  $X_{\mathcal{T},t} := X_{\mathcal{T}(t)}$  with  $b_{\mathcal{T},t} := b_{\mathcal{T}(t)}\mathcal{T}'(t)$  and  $\sigma_{\mathcal{T},t} := \sigma_{\mathcal{T}(t)}\sqrt{\mathcal{T}'(t)}$ .

Several conclusions are immediate. First, the new process  $X_{\mathcal{T}}$  satisfies Assumption 2.1; second, under the transformation  $\mathcal{T}$ , the irregular observation scheme becomes regular in the sense that  $X_{\tilde{t}_i^n} = X_{\mathcal{T},i\Delta_n}$ ; third, the integrated volatility is unchanged due to the properties of  $\mathcal{T}$ , upon a change of variable; finally, the probabilistic and statistical behavior of the noise is unchanged, in particular,  $\Sigma_U$  is unchanged and its consistent estimator remains valid.

Thus, upon replacing  $i\Delta_n$  by  $\tilde{t}_i^n$ , we can apply our Theorem 4.3 to observed noisy prices  $Y_{\tilde{t}_i^n} = X_{\mathcal{T},i\Delta_n} + U_i^n$ , which agrees exactly with (32). The limit distribution remains valid but the limit has a slightly different asymptotic variance:

$$V_{\mathcal{T},t}^2 := \frac{4}{\psi_0^2} \left( \Phi_{00}\theta\sigma_t^4\mathcal{T}'(\mathcal{T}^{-1}(t)) + 2\Phi_{01}\frac{\sigma_t^2\Sigma_U}{\theta} + \frac{\Phi_{11}\Sigma_U^2}{\theta^3} \right). \quad (33)$$

**Remark 4.3** (Jumps in the Efficient Price). *Assumption 2.1 does not allow for jumps in the efficient price process  $X$ . (Jumps in the stochastic volatility process are allowed.) We note from the proof of Proposition 3.1 that jumps in the efficient price will not affect the convergences of the realized volatility estimators of the second moments of noise, as the noise has larger asymptotic orders. For the pre-averaging estimators, we conjecture that under suitable conditions both  $\widehat{\text{IV}}_n$  and  $\widetilde{\text{IV}}_n$  will converge to the quadratic variation of  $X$  instead of to the IV. One can apply the truncation method (Mancini (2001)) to eliminate the jumps. But this is beyond the scope of this paper. In this context, it is worth mentioning an extensive empirical study by Christensen et al. (2014), in which the authors show that, as far as IV estimation is concerned, the jump component of the efficient price process in (very) high-frequency data typically only accounts for a small portion of the total price variation.*

### 4.3 Efficiency

It is well-known that estimators of volatility from noisy observations can achieve efficiency when the volatility is a constant,  $c_\sigma$ , i.e., the integrated volatility over  $[0, t]$  equals  $tc_\sigma$ , and the noise takes the form of Gaussian white (i.e., i.i.d.) noise with variance  $\mathbf{Var}(U)$ ; see Gloter and Jacod (2001a) and Gloter and Jacod (2001b) for a detailed account. In this case, an efficient estimator of the IV will converge at rate  $\Delta_n^{-\frac{1}{4}}$  with an asymptotic variance equal to  $\Sigma_t^{\text{opt}} = 8tc_\sigma^{3/2}\sqrt{\mathbf{Var}(U)}$ . This result has been extended to time-varying but non-random volatility processes plus Gaussian additive noise; see Reiß (2011). When the assumption of constant volatility is maintained but the noise is serially dependent, the optimal asymptotic variance becomes  $\overline{\Sigma}_t^{\text{opt}} = 8tc_\sigma^{3/2}\sqrt{\Sigma_U}$ , with the variance of noise replaced by the long-run variance of noise; see Da and Xiu (2019). We can show that the asymptotic variance of our estimator  $\widetilde{\text{IV}}_n$ , with the optimally selected  $\theta$  (recall Remark 4.1) and using the triangular kernel, is quite close to  $\overline{\Sigma}_t^{\text{opt}}$  under constant volatility:

$$\frac{\int_0^t V_s^2 ds}{\overline{\Sigma}_t^{\text{opt}}} \approx 1.07. \quad (34)$$

With stochastic volatility, it is still possible to achieve (34) asymptotically using local estimation — divide all observations into  $B$  blocks and perform estimation on each block and then aggregate the block estimates; see, e.g., Jacod and Mykland (2015), Clinet and Potiron (2018) and Da and Xiu (2019). Our simulation experience (not reported here) shows that in finite samples our estimators often do, but need not always, improve when following such a local estimation procedure. In those cases in which there is a lack of improvement, this may be partially due to a relatively worse estimation of the optimal  $\theta$  in a smaller sample.

Any proper estimation of  $\theta$ , whether local or global, requires accurate estimates of characteristics of the efficient price and noise processes. We will show through our extensive simulations and empirical



studies that model (mis)specification and finite sample biases play first-order roles in the estimation of such characteristics, and that our multi-step method introduced in the next section provides a robust approach. In our analyses presented later, we don't pursue local estimation, but focus on illustrating the robustness of our multi-step approach to model misspecification and to finite sample biases.

## 5 Multi-Step Estimators

In this section, we introduce our multi-step estimators of the IV and the second moments of noise based on both our asymptotic theory and finite sample analysis.

We observe from Theorem 4.3 that the second moments of noise contribute to an *asymptotic bias* in the estimation of the IV. Our finite sample analysis indicates, however, that we need an estimator of the IV to correct the *finite sample bias* when estimating the second moments of noise. Our multi-step estimators are specifically designed for the purpose of correcting the “interlocked” bias.

In the first step, we ignore the dependence in noise and estimate the variance of noise by realized volatility. Hence, our first-step estimators of the second moments of noise are given by

$$\widetilde{\gamma(0)}_n^{(1)} := \widehat{\langle Y, Y \rangle}(1)_n; \quad \widetilde{\gamma(j)}_n^{(1)} := 0, \quad j \neq 0; \quad \widetilde{\Sigma}_{U_n}^{(1)} := \widetilde{\gamma(0)}_n^{(1)}. \quad (35)$$

Next, we proceed with the pre-averaging method to obtain the first-step estimator of the IV (cf. (25)):

$$\widetilde{IV}_n^{(1)} = \frac{\sqrt{\Delta_n}}{\theta\psi_0} \sum_{i=0}^{n-k_n+1} \left( \overline{Y}_i^n \right)^2 - \frac{\psi_1 \widetilde{\Sigma}_{U_n}^{(1)}}{\theta^2 \psi_0}. \quad (36)$$

To initiate the second step, we first replace  $\hat{\sigma}^2$  by  $\widetilde{IV}_n^{(1)}$  in (10) and (11) and obtain the second-step estimators of the variance and covariances of noise. They will in turn be used to correct the asymptotic bias in the estimation of the IV, to eventually obtain the second-step estimator of the IV. Upon iterating this procedure, one obtains multi-step estimators. Specifically, for any  $k \geq 2$ , we define the  $k$ -step estimators recursively as follows:

$$\widetilde{\langle Y, Y \rangle}(j)_n^{(k)} := \widehat{\langle Y, Y \rangle}(j)_n - \frac{j \widetilde{IV}_n^{(k-1)}}{2(n-j+1)}; \quad (37)$$

$$\widetilde{\gamma(0)}_n^{(k)} := \widetilde{\gamma(0)}_n - \frac{j_n \widetilde{IV}_n^{(k-1)}}{2(n-j_n+1)}; \quad (38)$$

$$\widetilde{\gamma(j)}_n^{(k)} := \widetilde{\gamma(0)}_n^{(k)} - \widetilde{\langle Y, Y \rangle}(j)_n^{(k)}; \quad (39)$$

$$\widetilde{\Sigma}_{U_n}^{(k)} := \widetilde{\gamma(0)}_n^{(k)} + 2 \sum_{j=1}^{\ell_n} \widetilde{\gamma(j)}_n^{(k)}; \quad (40)$$

$$\widetilde{IV}_n^{(k)} := \frac{\sqrt{\Delta_n}}{\theta\psi_0} \sum_{i=0}^{n-k_n+1} \left(\overline{Y}_i^n\right)^2 - \frac{\psi_1 \widetilde{\Sigma}_{U_n}^{(k)}}{\theta^2 \psi_0}; \quad (41)$$

$$\widetilde{\Sigma}_{IV_n}^{(k)} := \frac{4\Phi_{00}}{3\theta\psi_0^4} \sum_{i=0}^{n-k_n+1} \left(\overline{Y}_i^n\right)^4 + \frac{8\widetilde{\Sigma}_{U_n}^{(k)}\widetilde{IV}_n^{(k)}}{\theta\psi_0^2} \left(\Phi_{01} - \frac{\psi_1\Phi_{00}}{\psi_0}\right) + \frac{4\left(\widetilde{\Sigma}_{U_n}^{(k)}\right)^2}{\theta^3\psi_0^2} \left(\Phi_{11} - \frac{\psi_1^2\Phi_{00}}{\psi_0^2}\right). \quad (42)$$

We state the following theorem:

**Theorem 5.1.** *Under the assumptions of Theorem 4.3, for any fixed  $K \in \mathbb{N}^*$ , we have*

$$\widetilde{\Sigma}_{U_n}^{(K)} \xrightarrow{\mathbb{P}} \Sigma_U, \quad (43)$$

and the sequence  $\Delta_n^{-\frac{1}{4}} \left( \widetilde{IV}_n^{(K)} - IV \right) / \sqrt{\widetilde{\Sigma}_{IV_n}^{(K)}}$  converges stably in law to a standard normal random variable, which is independent of  $\mathcal{F}$ .

We note that, for brevity, our multi-step estimators introduced above are based only on the pre-averaging estimators using overlapping intervals. Of course, we can adopt the same approach and develop, by analogy, consistent and asymptotically normal multi-step estimators from the pre-averaging estimators using non-overlapping intervals as well. They will henceforth be denoted by  $\widehat{IV}_n^{(k)}$  and will be analyzed alongside  $\widetilde{IV}_n^{(k)}$  later.

**Remark 5.1.** *As the simulation results in the next section will reveal, our multi-step estimators introduced above perform well. An advantage of our multi-step estimators is that they are quite robust to the choice of the tuning parameter  $\theta$ . To offer some insight into this issue, we briefly analyze the relationship between the choice of  $\theta$  and the theoretical finite sample bias of two estimators: our  $\widetilde{IV}_n$  and the benchmark estimator  $\widetilde{IV}_n^{\text{JLZ}}$  recently proposed by [Jacod et al. \(2019\)](#), which employs the local averaging (LA) method to correct the asymptotic bias of pre-averaging estimators. A simple calculation shows that the finite sample errors of  $\widetilde{IV}_n$  and  $\widetilde{IV}_n^{\text{JLZ}}$  (as a percentage) are approximately given by*

$$\text{Err}_{\text{RV}} \approx \frac{(2\ell_n + 1)j_n + \sum_{|\ell| \leq \ell_n} |\ell| \phi_1^r(\ell)}{2n\theta^2\psi_0}, \quad \text{Err}_{\text{JLZ}} \approx \frac{4K_n \sum_{|\ell| \leq \ell_n} \phi_1^r(\ell)}{3n\theta^2\psi_0}, \quad (44)$$

respectively, where  $K_n$  is the tuning parameter of the LA method. While these errors can be significant for both estimators, and moreover a small change in  $\theta$  can lead to sharp changes in the errors, our multi-step estimators are specifically designed to remove this error. Consequently, they are much more robust to changes in  $\theta$  than estimators without unified bias corrections.

## 6 Simulation Study

### 6.1 Simulation design

Motivated by the empirical studies in [Aït-Sahalia et al. \(2011\)](#), we consider an ARMA(1,1) noise process  $U$  given by the following dynamics:

$$U_t = e_t + \epsilon_t, \quad (45)$$

where  $e$  is centered i.i.d. Gaussian and  $\epsilon$  is an AR(1) process with first-order coefficient  $\iota$ ,  $|\iota| < 1$ . We will examine the performance of our estimators for different values of this coefficient:  $\iota \in \{-0.7, -0.3, 0, 0.3, 0.7\}$ . The processes  $e$  and  $\epsilon$  are assumed to be statistically independent. We set  $\mathbb{E}(e_t^2) = 1.9 \times 10^{-7}$ , and  $\mathbb{E}(\epsilon_t^2) = 1.3 \times 10^{-7}$ . These values are chosen to mimic the results of our empirical studies.

We assume that the efficient price is generated by the following dynamics:

$$\begin{cases} dX_t = -\kappa_1(X_t - \mu_1)dt + \sigma_t dW_t, \\ d\sigma_t^2 = \kappa_2(\mu_2 - \sigma_t^2)dt + \kappa_3\sigma_t dB_t + \xi_t dN_t, \end{cases}$$

where  $B$  and  $W$  are standard Brownian motions with quadratic covariation  $\langle B, W \rangle_t = \varrho t$ ,  $N$  is a Poisson process with parameter  $\lambda$ , and  $\xi_t$  is an independent jump size following an exponential distribution with parameter  $\kappa_3$ . We set the parameters as follows:  $\kappa_1 = 0.5$ ,  $\mu_1 = 1.6$ ,  $\kappa_2 = 5/252$ ,  $\mu_2 = 0.04/252$ ,  $\kappa_3 = 0.05/252$ ,  $\lambda = 3$ , and  $\varrho = -0.5$ . We assume the processes  $X$  and  $U$  to be mutually independent. We simulate each sample path within a fixed time interval  $[0, 1]$  that represents one trading day.

### 6.2 Realized volatility estimators of the second moments of noise

To get a first impression of the properties of our estimator  $\widehat{\langle Y, Y \rangle}(j)_n$  defined in (6), we plot  $\widehat{\langle Y, Y \rangle}(j)_n$  against the number of lags  $j$  in Figure 2. In addition to  $\widehat{\langle Y, Y \rangle}(j)_n$ , we also plot the bias adjusted version  $\widehat{\langle Y, Y \rangle}^{(\text{adj})}(j)_n$  defined in (10), in which we employ three ‘‘approximations’’ to the IV that  $\widehat{\langle Y, Y \rangle}^{(\text{adj})}(j)_n$  depends on:  $\hat{\sigma}_H^2 = 1.2\text{IV}$ ,  $\hat{\sigma}_M^2 = \text{IV}$ , and  $\hat{\sigma}_L^2 = 0.8\text{IV}$ . Figure 2 shows that a prominent feature of our realized volatility estimator  $\widehat{\langle Y, Y \rangle}(j)_n$  is that it deviates from its stochastic limit  $\gamma(0) - \gamma(j)$  almost linearly in the number of lags  $j$ , as predicted by Proposition 3.2. The deviation, induced by the finite sample bias, can be largely corrected when only rough ‘‘estimates’’ of the IV are available. In the ideal but infeasible situation that we know exactly the true volatility ( $\hat{\sigma}_M^2 = \text{IV}$ ), the bias corrected estimators almost perfectly match the underlying true values.

Next, we estimate the second moments of noise by our realized volatility estimators (RV) and, for

comparison purposes, by the local averaging estimators (LA) proposed by [Jacod et al. \(2017\)](#). We demonstrate the importance of the finite sample bias correction to obtain accurate estimates, and this applies to both estimators.<sup>13</sup> In [Figure 3](#), we plot the means of the autocorrelations of noise estimated by RV and LA based on 1,000 simulations. In the top panel we plot the estimators without finite sample bias correction and we plot the estimators with finite sample bias correction in the bottom panel, in which we use the true IV instead of any approximation/estimator to make the bias correction. We will analyze the case in which IV is estimated in the next subsection.

We observe that both estimators (RV and LA) perform poorly without finite sample bias correction. In particular, the noise autocorrelations estimated by the LA estimators decay slowly and hover above 0 up to 20 lags, from which we might conclude that the noise exhibits strong and long-memory dependence, while the underlying noise is, in fact, only weakly dependent. However, both estimators perform well after the finite sample bias correction. In [Figure 3](#), we also plot the 95% simulated confidence intervals of the two bias corrected estimators. In terms of mean squared errors, both estimators, after bias correction, yield accurate estimates.

[Figures 2-3](#) reveal that the finite sample bias correction is crucial to obtain reliable estimates of noise moments. The key ingredient of this correction, however, is (an estimate of) the IV. Yet, to obtain an estimate of the IV, we need to estimate the second moments of noise first — whence the feedback loop of bias corrections. This is why we introduced our multi-step estimators, which allow successive bias corrections in estimates for both the IV and noise autocorrelations.

### 6.3 Multi-step estimators of IV

In this subsection, we examine the performance of different estimators of the IV. We compare the estimator  $\widehat{IV}_n$  in [\(22\)](#) which is generated by the pre-averaging method using non-overlapping intervals, with the estimator  $\widetilde{IV}_n$  defined in [\(25\)](#) using overlapping intervals. We then assess the improvement in accuracy from our unified treatment of asymptotic and finite sample biases that can be achieved by using the  $K$ -step estimators  $\widehat{IV}_n^{(K)}$  and  $\widetilde{IV}_n^{(K)}$  introduced in [\(41\)](#). We also compare  $\widehat{IV}_n^{(1)}$  and  $\widetilde{IV}_n^{(1)}$  to  $\widehat{IV}_n^{(2)}$  and  $\widetilde{IV}_n^{(2)}$  to assess the gained accuracy by dropping the possibly misspecified assumption of independent noise.

In [Table 1](#), we report the centered means of our estimators and the standard deviations (between

---

<sup>13</sup>The finite sample bias corrected local averaging estimators of the noise covariances are given by

$$\widehat{R}(j)_n = \frac{1}{n} U((0, j))_n - \frac{K_n}{n} \left( \frac{4}{3} \widehat{\sigma}^2 \right),$$

where  $U((0, j))_n/n$  is the local averaging estimator of the  $j$ -th covariance without bias correction and  $\widehat{\sigma}^2$  is an estimator of the IV; see [Jacod et al. \(2017\)](#) for more details. While [Jacod et al. \(2017\)](#) provide a finite sample bias correction when developing their local averaging estimators of noise covariances, they don't consider the feedback between, and unified treatment of, asymptotic and finite sample biases, which is a key interest in this paper.

parentheses), based on 1,000 simulations.<sup>14</sup> Throughout this subsection, the tuning parameter  $j_n$  is fixed at 20, we take  $\ell_n = 10$  and  $\theta = 0.4$ , and use the triangular kernel. When comparing the estimators  $\widehat{IV}_n$  and  $\widehat{IV}_n^{(2)}$  in the first and the third rows of Table 1, we observe the important advantage of our multi-step estimators over the pre-averaging method that ignores the finite sample bias, since our estimators yield strongly improved accuracy. Furthermore, a comparison to the results for  $\widehat{IV}_n^{(1)}$  and  $\widehat{IV}_n^{(2)}$  in the second and third rows leads to the striking conclusion that ignoring the finite sample bias yields even more inaccuracy than ignoring the dependence in noise. This shows that one should be cautious when applying estimators without appropriate bias corrections even with data on a 1-sec time scale (i.e., 23,400 observations in a day of 6.5 trading hours). The “cost” of applying our multi-step estimators  $\widehat{IV}_n^{(K)}$  is the slightly larger standard deviations they induce. This increased uncertainty is introduced by correcting the “interlocked” bias. However, the reduction in bias strictly dominates the slight increase in standard deviations when noise is dependent. Therefore, the multi-step estimators have smaller mean-squared errors than their counterparts in the first two rows of Table 1. These standard deviations can be reduced by the use of overlapping intervals, as can be observed when we compare the standard deviations of  $\widehat{IV}_n^{(K)}$  with those of  $\widetilde{IV}_n^{(K)}$  (i.e., the first four rows in Table 1 and the next four rows). Although the centered means of the estimators become slightly worse when we adopt overlapping pre-averaging intervals, the significant reduction in the standard deviations implies a better overall performance under a mean-squared error criterion.

The estimator  $\widetilde{IV}_n^{\text{JLZ}}$  recently proposed in [Jacod et al. \(2019\)](#), which corrects the asymptotic bias of pre-averaging estimators by local averaging but does not include a unified treatment of asymptotic and finite sample biases, performs better than the estimators  $\widehat{IV}_n$  and  $\widetilde{IV}_n$ , but worse than all estimators with finite sample bias corrections. The method proposed in [Da and Xiu \(2019\)](#) generates an estimator  $\widehat{IV}_n^{\text{QMLE}}$  which outperforms our method when the autocorrelation in the noise is small, but its performance deteriorates when the noise autocorrelation parameter  $\iota$  is closer to  $-1$  or  $1$ .

In Table 2, we replicate the results of Table 1 but now with a higher data frequency, which corresponds to sampling every 0.2 seconds (i.e., 117,000 observations in a day of 6.5 trading hours). We observe that, with such very high-frequency data, the multi-step estimators still perform much better than their counterparts in rows 1 and 2, and 5 and 6, of Table 2 — with much smaller biases and only slightly larger standard deviations. Indeed, both the errors caused by ignoring the finite sample bias and the inconsistencies caused by a potential misspecification of the dependence structure in noise when using the first-step estimators remain clearly visible. The biases in the estimates typically reduce further when we replace  $\widehat{IV}_n^{(2)}$  by  $\widehat{IV}_n^{(3)}$ , but not in all cases where  $\widetilde{IV}_n^{(2)}$  is replaced by  $\widetilde{IV}_n^{(3)}$ . We also observe that increasing  $K$  in our multi-step estimators  $\widehat{IV}_n^{(K)}$  and  $\widetilde{IV}_n^{(K)}$  gives only a slight increase in the estimators’

---

<sup>14</sup>The numbers are multiplied by  $10^5$ .

standard deviations. As before, the standard deviations of  $\widetilde{IV}_n^{(2)}$  and  $\widetilde{IV}_n^{(3)}$  are substantially smaller than for  $\widehat{IV}_n^{(2)}$  and  $\widehat{IV}_n^{(3)}$ , and for  $\widehat{IV}_n^{\text{JLZ}}$  and  $\widehat{IV}_n^{\text{QMLE}}$ . In terms of CPU, the QMLE-estimator is relatively more time-consuming to compute. Indeed, in the setting of Table 2, 0.1% of the total computing time was spent on our four estimators based on non-overlapping intervals; 3.1% was spent on our four estimators based on overlapping intervals; 7.2% was spent on  $\widehat{IV}_n^{\text{JLZ}}$ ; and 89.6% was spent on  $\widehat{IV}_n^{\text{QMLE}}$ .

To numerically “verify” the central limit theorem established in Theorem 5.1, we plot the quantiles of the normalized estimators  $\Delta_n^{-\frac{1}{4}} \left( \widehat{IV}_n^{(2)} - \text{IV} \right) / \sqrt{\widehat{\Sigma}_{IV_n}^{(2)}}$  and  $\Delta_n^{-\frac{1}{4}} \left( \widetilde{IV}_n^{(2)} - \text{IV} \right) / \sqrt{\widetilde{\Sigma}_{IV_n}^{(2)}}$  against standard normal quantiles in Figure 4. We observe that the established limit distributions are clearly verified.

**Remark 6.1** (Dependence between  $X$  and  $U$ ). *The theoretical results in this paper assume independence between  $X$  and  $U$ . In practice, the efficient price and the microstructure noise processes may be correlated. In Appendix B, we provide additional Monte Carlo simulation results that assess the effects of price discreteness and correlation between  $X$  and  $U$ . (Price discreteness renders dependence between  $X$  and  $U$ .) Our results show that the presence of minimal ticks has relatively little impact on the estimation of the moments of noise and the IV. Furthermore, our results show that in the situation when  $X$  and  $U$  are correlated our multi-step estimators still appear to be performing well.*

## 7 Empirical Study

### 7.1 Data description

We analyze the NYSE TAQ transaction prices of Citigroup (trading symbol: C) over the month January 2011. We discard all transactions before 9:30 and after 16:00. We retain a total of 4,933,059 transactions over 20 trading days, thus on average 10.5 observations per second. The estimation is first performed on the full sample, and then on subsamples obtained by different sampling schemes. We demonstrate how the sampling methods affect the properties of the noise, and thus affect the estimation of the IV. We employ pre-averaging on overlapping intervals, and use the triangular kernel. Throughout this section, the tuning parameter of the RV estimator is fixed at  $j_n = 30$  and  $\theta$  is selected according to the optimal rule (28).

### 7.2 Estimating the second moments of noise

We estimate the  $j$ -th autocovariance and autocorrelation of microstructure noise with  $j = 0, 1, \dots, 30$  by three estimators: our realized volatility (RV) estimators in (7) and (8), the local averaging (LA) estimators proposed by Jacod et al. (2017), and the bias corrected realized volatility (BCRV) estimators in (38) and (39). We perform the estimation over each trading day and end up with 20 estimates (of

the 30 lags of autocovariances or autocorrelations) for each estimator. In Figure 5 we plot the average of the 20 estimates (over the month) as well as the approximated confidence intervals that are two sample standard deviations away from the mean.

We observe that the three estimators yield quite close estimates by virtue of the high data frequency. Noise in this sample tends to be positively autocorrelated — with the BCRV estimators yielding the fastest decay. Empirically this positive autocorrelation is consistent with the finding that the arrivals of buy and sell orders are positively autocorrelated; see Hasbrouck and Ho (1987). This corresponds to the trading practice that informed traders split their orders over (a short period of) time and trade on one side of the market, rendering continuation in their orders.

We emphasize that the finite sample bias can be much more pronounced than what we observe in Figure 5, even if we conduct estimation on a full transaction data sample. Indeed, in Appendix C, we analyze a sample of transaction prices of General Electric (GE) and show that, when the data frequency is very high, the finite sample bias correction is particularly essential when the noise-to-signal ratio is very small (recall Remark 3.3).

### 7.3 Estimating the IV

Turning to the estimation of the IV, we study four estimators of the pre-averaging class:  $\widetilde{IV}_n$ ,  $\widetilde{IV}_n^{(1)}$ ,  $\widetilde{IV}_n^{(2)}$ , and  $\widetilde{IV}_n^{\text{JLZ}}$ . In the top panel of Figure 6, we plot the four estimators of the IV for each trading day. We note that the three estimators  $\widetilde{IV}_n$ ,  $\widetilde{IV}_n^{\text{JLZ}}$ , and  $\widetilde{IV}_n^{(2)}$  yield quite close results. This is expected, as the three methods, RV, LA, and BCRV, provide close estimates of the second moments of noise. However, the estimator  $\widetilde{IV}_n^{(1)}$ , which ignores the dependence in noise, yields very different estimates, and the differences are persistent —  $\widetilde{IV}_n^{(1)}$  yields higher estimates over each trading day. Moreover, the differences are statistically significant by virtue of Theorem 5.1 — all the 20 estimates fall outside of the 95% confidence intervals, as the bottom panel of Figure 6 reveals.

### 7.4 Decaying rate of autocorrelation

Figure 5 shows that the positive autocorrelations of noise drop to zero rapidly. To assess the rate of decay, we perform a logarithmic transformation of the autocorrelations estimated by BCRV.<sup>15</sup> In Figure 7, we plot the logarithmic autocorrelations for each trading day (top panel) and the mean logarithmic autocorrelations over all days (bottom panel). The plots indicate that the logarithmic autocorrelation is approximately a linear function of the number of lags, i.e., the autocorrelation function is decaying at

<sup>15</sup>We restrict attention to the lags up to  $j = 10$ . The logarithmic autocorrelations at higher lags are very volatile since the autocorrelations are close to zero.

an exponential rate.<sup>16</sup>

## 7.5 Robustness check — estimation under other sampling schemes

It is interesting to analyze how our estimators perform when the data is sampled at different time scales. In this section, we consider two alternative (sub)sampling schemes: regular time sampling and tick time sampling (recall Remark 2.1 for details on the sampling schemes).

### 7.5.1 Regular time sampling

The prices in this sample are recorded on a 1-second time scale. If there were multiple prices in a second, we select the first one; and we do not record a price if there is no transaction in a second. We end up with 21,691 observations on average per trading day. Figure 8 is analogous to Figure 5. The three estimators, RV, LA, and BCRV, now produce very different patterns. Both the RV and LA estimators suggest that the noise is strongly autocorrelated in this subsample, even stronger than in the original full sample. This would be counterintuitive since we eliminate more than 90% of the full sample in a fairly random way — the elimination should have weakened the serial dependence of noise in the remaining sample. However, the estimates by BCRV reveal that in fact the noise is approximately uncorrelated — it is the finite sample bias that makes the autocorrelations of noise seem strong and persistent if not taken into account.

If the noise is close to being independent, then  $\widetilde{IV}_n^{(1)}$ , which assumes i.i.d. noise, would be a sound estimator of the IV. An alternative estimator, e.g.,  $\widetilde{IV}_n^{(2)}$ ,  $\widetilde{IV}_n$ , or  $\widetilde{IV}_n^{\text{JLZ}}$  would then be robust if it would deliver similar estimates. In the top panel of Figure 9, we observe that  $\widetilde{IV}_n^{(1)}$  and  $\widetilde{IV}_n^{(2)}$  yield virtually identical estimates. The other two estimators,  $\widetilde{IV}_n$  and  $\widetilde{IV}_n^{\text{JLZ}}$  which don't apply finite sample bias corrections, however, yield even negative estimates. It is interesting to briefly elaborate on the performance of  $\widetilde{IV}_n$  and  $\widetilde{IV}_n^{\text{JLZ}}$  in this scenario. Using the triangular kernel, with the selected tuning parameters  $j_n = 30, \ell_n = 4, K_n = 7$  and a reasonable choice of  $\theta = 0.2$ , we have by (44) that  $\text{Err}_{\text{JLZ}} = 103.69\%$ ,  $\text{Err}_{\text{RV}} = 175.64\%$ . Therefore,  $\widetilde{IV}_n^{\text{JLZ}}$  and  $\widetilde{IV}_n$  are in fact estimating  $-3.69\%$  and  $-75.64\%$  of the IV, and this is in line with the estimates in the top panel of Figure 9. We conclude that Figures 6 and 9 jointly confirm the importance of our multi-step approach. Indeed,  $\widetilde{IV}_n^{(1)}$ , which assumes i.i.d. noise, exhibits unreliable behavior in Figure 6, while  $\widetilde{IV}_n$ , which doesn't apply finite sample bias corrections, shows unreliable behavior in Figure 9.

<sup>16</sup>The autocorrelation decay rate would be slower without unified treatment of the bias corrections, which may explain the relatively strong polynomial dependence in noise found in Jacod et al. (2017) and questioned by these authors themselves.



### 7.5.2 Tick time sampling

In a tick time sample, prices are collected with each price change, i.e., all zero returns are suppressed, see, e.g., Zhou (1996), Griffin and Oomen (2008), Ait-Sahalia et al. (2011), Kalnina (2011) and Da and Xiu (2019). For the Citigroup transaction data, 70% of the returns are zero. The corresponding average number of prices per second in our tick time sample is 3.2. Figure 10 shows that the microstructure noise has a different dependence pattern in the tick time sample — its autocorrelation function is alternating. Masked by alternating noise, the observed returns at tick time have a similar pattern; see Ait-Sahalia et al. (2011) and Griffin and Oomen (2008). This dependence structure of noise is perceived to be due to the discreteness of price changes, irrespective of the distributional features of noise in the original transactions or quotes data.

Interestingly, the bottom panel of Figure 9 shows that the two estimators of the IV,  $\widetilde{IV}_n^{(1)}$  and  $\widetilde{IV}_n^{(2)}$ , remain close. It is not immediate why  $\widetilde{IV}_n^{(1)}$  and  $\widetilde{IV}_n^{(2)}$  deliver almost identical estimates, given the fact that the dependence of noise in this tick time sample is drastically different from i.i.d. noise. However, a clue is provided by the observation that negatively autocorrelated noise has less impact on the estimation of the IV, as the high-order alternating autocovariances partially cancel out, and thus contribute less to the asymptotic bias  $\sigma_U^2$ .<sup>17</sup>  $\widetilde{IV}_n$  and  $\widetilde{IV}_n^{\text{JLZ}}$  are still significantly underestimating the IV due to the finite sample bias.

## 7.6 Economic interpretation and empirical implication

The dependence structure of microstructure noise depends on the sampling frequency and scheme. In an original transaction data sample, noise is likely to be positively autocorrelated as a result of various trading practices that entail continuation in order flows. The dependence of noise can be reduced by sampling sparsely, say, every few (or more) seconds as we show in Section 7.5.1; noise is close to independent in such sparse subsamples. If, however, we remove all zero returns, thus sample in tick time, noise typically exhibits an alternating autocorrelogram.

Microstructure theories can provide some intuitive economic interpretations of the dynamic properties of microstructure noise recovered in this paper. The positive autocorrelation function displayed in Figure 5 is consistent with the findings in Hasbrouck and Ho (1987), Choi et al. (1988) and Huang and Stoll (1997) that explicitly model the probability of order reversal  $\pi$  (or order continuation by

<sup>17</sup>For a tractable analysis, one may consider AR(1) noise processes. Let  $\iota \in (0, 1)$  be the absolute value of the AR(1) coefficient. When the noise is positively autocorrelated, the asymptotic bias  $\sigma_U^2$  corrected by  $\widetilde{IV}_n^{(1)}$  and  $\widetilde{IV}_n^{(2)}$  is  $(1 - \iota)\gamma(0)$  and  $\frac{1+\iota}{1-\iota}\gamma(0)$ , respectively; when the noise is negatively autocorrelated, it is  $(1 + \iota)\gamma(0)$  and  $\frac{1-\iota}{1+\iota}\gamma(0)$ . Consider  $\iota = 0.8$ . Then,  $(1 - \iota)\gamma(0) = 0.2\gamma(0)$  and  $\frac{1+\iota}{1-\iota}\gamma(0) = 9\gamma(0)$  while  $(1 + \iota)\gamma(0) = 1.8\gamma(0)$  and  $\frac{1-\iota}{1+\iota}\gamma(0) = \frac{1}{9}\gamma(0)$ . Therefore, the difference in the asymptotic bias is smaller when the noise is negatively autocorrelated; consequently, the IV estimates by  $\widetilde{IV}_n^{(1)}$  and  $\widetilde{IV}_n^{(2)}$  are close, see also Tables 1 and 2 in our simulation study.

$1 - \pi$ ),<sup>18</sup> so that the deviation of transaction prices from fundamentals becomes an AR(1) process. Fitting the autocorrelation function recovered by BCRV in Figure 5 to that of an AR(1) model, we obtain an estimate of the AR(1) coefficient equal to  $\hat{\iota} = 0.73$  and the probability of order continuation is  $1 - \hat{\pi} = (1 + \hat{\iota})/2 = 0.87$ . That is, the estimated probability that a buy (or sell) order follows another buy (or sell) order is 0.87. In view of the extensive empirical results in Huang and Stoll (1997) (see Table 5 therein), this is a reasonable estimate.

One possible interpretation of the positively autocorrelated order flows is that a large *order* is often executed as a series of smaller *trades* to reduce the price impact, or conducted against multiple *trades* from stale limit orders. However, such positive autocorrelation contradicts the prediction of inventory models, in which market makers induce negatively autocorrelated order flows to equilibrate inventories; see Ho and Stoll (1981). Consequently, according to inventory models the probability of order reversal would be  $\pi > 0.5$ . One remedy, suggested by Huang and Stoll (1997), is to collapse multiple *trades* at the same price into one *order*, which is exactly the tick time sampling scheme considered in Section 7.5.2. Exploiting the estimates by BCRV presented in Figure 10, we obtain an estimate of the probability of order reversal equal to  $\hat{\pi} = 0.84$ , which is very close to the average probability 0.87 in Huang and Stoll (1997). We emphasize that we recover these probabilities without any prior knowledge or estimates of the order flows.

The dependence structure of microstructure noise affects the estimation of the IV. Popular de-noise methods that assume i.i.d. noise work reasonably well with relatively sparse regular time samples or tick time samples. However, this discards a substantial amount of the original transaction data.<sup>19</sup> Instead, we can directly estimate the IV from the original data using our multi-step estimators that explicitly take the potential dependence in noise into account.

In our empirical study, we have also illustrated that bias corrections play an essential role in recovering the statistical properties of noise and in estimating the IV. Our multi-step estimators are specifically designed to conduct such bias corrections.

## 8 Conclusion

In high-frequency financial data the efficient price is contaminated by microstructure noise, which is usually assumed to be independently and identically distributed. This simple distributional assumption is challenged by both microeconomic financial models and various empirical facts. In this paper, we deviate from the i.i.d. assumption by allowing noise to be dependent in a general setting. We then

<sup>18</sup>It is the probability that a buy (sell) order follows another sell (buy) order.

<sup>19</sup>To obtain the Citigroup tick time sample and the 1-second regular time sample, we delete roughly 70% and 90% of the original transaction data, respectively.

develop econometric tools to recover the dynamic properties of microstructure noise and design improved approaches for the estimation of the integrated volatility.

This paper makes four contributions. First, it develops nonparametric estimators of the second moments of microstructure noise in a general setting. Second, it provides robust estimators of the integrated volatility, without assuming serially independent noise. Third, it reveals the importance of both asymptotic and finite sample bias analysis and develops simple and readily implementable multi-step estimators. Empirically, it characterizes the dependence structures of noise at several time scales and provides intuitive economic interpretations; it also investigates the impact of the dynamic properties of microstructure noise on integrated volatility estimation.

This paper thus introduces a robust and accurate method to effectively separate the two components of high-frequency financial data — the efficient price and microstructure noise. The robustness lies in its flexibility to accommodate rich dependence structures of microstructure noise motivated by various economic models and trading practices, whereas the accuracy is achieved by the finite sample refinement. As a result, we discover dynamic properties of microstructure noise consistent with microstructure theory and obtain accurate volatility estimators.

## Acknowledgements

We are very grateful to the Editor, the Associate Editor and two anonymous referees for their helpful comments and suggestions that have significantly improved this paper. We are also grateful to Federico Bandi, Peter Boswijk, Peter Reinhard Hansen, Siem Jan Koopman, Oliver Linton, and Xiye Yang for their comments and discussions on earlier versions of this paper. This research was funded in part by the Netherlands Organization for Scientific Research under grant NWO VIDI 2009 (Laeven).

## References

- AÏT-SAHALIA, Y. AND J. JACOD (2014): *High-frequency Financial Econometrics*, Princeton University Press.
- AÏT-SAHALIA, Y., P. A. MYKLAND, AND L. ZHANG (2005): “How often to sample a continuous-time process in the presence of market microstructure noise,” *Review of Financial Studies*, 18, 351–416.
- (2011): “Ultra high frequency volatility estimation with dependent microstructure noise,” *Journal of Econometrics*, 160, 160–175.

- BANDI, F. M. AND J. R. RUSSELL (2006): “Separating microstructure noise from volatility,” *Journal of Financial Economics*, 79, 655–692.
- (2008): “Microstructure noise, realized variance, and optimal sampling,” *Review of Economic Studies*, 75, 339–369.
- BARNDORFF-NIELSEN, O. E., P. R. HANSEN, A. LUNDE, AND N. SHEPHARD (2008): “Designing realized kernels to measure the ex post variation of equity prices in the presence of noise,” *Econometrica*, 76, 1481–1536.
- BLACK, F. (1986): “Noise,” *Journal of Finance*, 41, 529–543.
- BRADLEY, R. C. (2007): *Introduction to Strong Mixing Conditions*, Kendrick Press.
- CHAKER, S. (2017): “On high frequency estimation of the frictionless price: The use of observed liquidity variables,” *Journal of Econometrics*, 201, 127–143.
- CHOI, J. Y., D. SALANDRO, AND K. SHASTRI (1988): “On the estimation of bid-ask spreads: Theory and evidence,” *Journal of Financial and Quantitative Analysis*, 23, 219–230.
- CHRISTENSEN, K., R. C. OOMEN, AND M. PODOLSKIJ (2014): “Fact or friction: Jumps at ultra high frequency,” *Journal of Financial Economics*, 114, 576–599.
- CLINET, S. AND Y. POTIRON (2017): “Estimation for high-frequency data under parametric market microstructure noise,” Tech. rep.
- (2018): “Efficient asymptotic variance reduction when estimating volatility in high frequency data,” *Journal of Econometrics*, 206, 103–142.
- (2019): “Testing if the market microstructure noise is fully explained by the informational content of some variables from the limit order book,” *Journal of Econometrics*.
- DA, R. AND D. XIU (2019): “When Moving-Average Models Meet High-Frequency Data: Uniform Inference on Volatility,” Tech. rep.
- DUFFIE, D. (2010): *Dynamic Asset Pricing Theory*, Princeton University Press.
- GARMAN, M. B. (1976): “Market microstructure,” *Journal of Financial Economics*, 3, 257–275.
- GLOSTEN, L. R. AND P. R. MILGROM (1985): “Bid, ask and transaction prices in a specialist market with heterogeneously informed traders,” *Journal of Financial Economics*, 14, 71–100.

- GLOTER, A. AND J. JACOD (2001a): “Diffusions with measurement errors. I. Local asymptotic normality,” *ESAIM: Probability and Statistics*, 5, 225–242.
- (2001b): “Diffusions with measurement errors. II. Optimal estimators,” *ESAIM: Probability and Statistics*, 5, 243–260.
- GRIFFIN, J. E. AND R. C. OOMEN (2008): “Sampling returns for realized variance calculations: tick time or transaction time?” *Econometric Reviews*, 27, 230–253.
- GROSS-KLUSSMANN, A. AND N. HAUTSCH (2013): “Predicting bid–ask spreads using long-memory autoregressive conditional Poisson models,” *Journal of Forecasting*, 32, 724–742.
- HANSEN, P. R., J. LARGE, AND A. LUNDE (2008): “Moving average-based estimators of integrated variance,” *Econometric Reviews*, 27, 79–111.
- HANSEN, P. R. AND A. LUNDE (2006): “Realized variance and market microstructure noise,” *Journal of Business & Economic Statistics*, 24, 127–161.
- HARRIS, L. (1990): “Estimation of stock price variances and serial covariances from discrete observations,” *Journal of Financial and Quantitative Analysis*, 25, 291–306.
- HASBROUCK, J. (2007): *Empirical Market Microstructure: The Institutions, Economics, and Econometrics of Securities Trading*, Oxford University Press.
- HASBROUCK, J. AND T. S. HO (1987): “Order arrival, quote behavior, and the return-generating process,” *Journal of Finance*, 42, 1035–1048.
- HAUTSCH, N. AND M. PODOLSKIJ (2013): “Preaveraging-based estimation of quadratic variation in the presence of noise and jumps: Theory, implementation, and empirical evidence,” *Journal of Business & Economic Statistics*, 31, 165–183.
- HENDERSHOTT, T., C. M. JONES, AND A. J. MENKVELD (2013): “Implementation shortfall with transitory price effects,” in *High-Frequency Trading: New Realities for Trades, Markets and Regulators*, ed. by D. Easley, M. L. de Prado, and M. O’Hara, London: Risk Books.
- HO, T. AND H. R. STOLL (1981): “Optimal dealer pricing under transactions and return uncertainty,” *Journal of Financial Economics*, 9, 47–73.
- HUANG, R. D. AND H. R. STOLL (1997): “The components of the bid–ask spread: A general approach,” *Review of Financial Studies*, 10, 995–1034.

- JACOD, J., Y. LI, P. A. MYKLAND, M. PODOLSKIJ, AND M. VETTER (2009): “Microstructure noise in the continuous case: The pre-averaging approach,” *Stochastic Processes and Their Applications*, 119, 2249–2276.
- JACOD, J., Y. LI, AND X. ZHENG (2017): “Statistical properties of microstructure noise,” *Econometrica*, 85, 1133–1174.
- (2019): “Estimating the integrated volatility with tick observations,” *Journal of Econometrics*, 208, 80–100.
- JACOD, J. AND P. A. MYKLAND (2015): “Microstructure noise in the continuous case: Approximate efficiency of the adaptive pre-averaging method,” *Stochastic Processes and their Applications*, 125, 2910–2936.
- JACOD, J., M. PODOLSKIJ, AND M. VETTER (2010): “Limit theorems for moving averages of discretized processes plus noise,” *Annals of Statistics*, 38, 1478–1545.
- JACOD, J. AND A. N. SHIRYAEV (2003): *Limit Theorems for Stochastic Processes*, vol. 288, Springer-Verlag Berlin.
- KALNINA, I. (2011): “Subsampling high frequency data,” *Journal of Econometrics*, 161, 262–283.
- LEE, S. S. AND P. A. MYKLAND (2012): “Jumps in equilibrium prices and market microstructure noise,” *Journal of Econometrics*, 168, 396–406.
- LI, Y., S. XIE, AND X. ZHENG (2016): “Efficient estimation of integrated volatility incorporating trading information,” *Journal of Econometrics*, 195, 33–50.
- LI, Z. M., R. J. A. LAEVEN, AND M. H. VELLEKOOP (2019): “Supplementary material to “Dependent microstructure noise and integrated volatility estimation from high-frequency data,”” Tech. rep.
- MANCINI, C. (2001): “Disentangling the jumps of the diffusion in a geometric jumping Brownian motion,” *Giornale dell’Istituto Italiano degli Attuari*, 64, 19–47.
- MOKKADEM, A. (1988): “Mixing properties of ARMA processes,” *Stochastic Processes and Their Applications*, 29, 309–315.
- MYKLAND, P. A. AND L. ZHANG (2012): “The econometrics of high frequency data,” in *Statistical Methods for Stochastic Differential Equations*, ed. by M. Kessler, A. Lindner, and M. Sørensen, New York: Chapman and Hall/CRC Press, chap. 2, 109–190.

- PODOLSKIJ, M. AND M. VETTER (2009a): “Bipower-type estimation in a noisy diffusion setting,” *Stochastic Processes and Their Applications*, 119, 2803–2831.
- (2009b): “Estimation of volatility functionals in the simultaneous presence of microstructure noise and jumps,” *Bernoulli*, 15, 634–658.
- REISS, M. (2011): “Asymptotic equivalence for inference on the volatility from noisy observations,” *Annals of Statistics*, 39, 772–802.
- ROLL, R. (1984): “A simple implicit measure of the effective bid-ask spread in an efficient market,” *Journal of Finance*, 39, 1127–1139.
- STOLL, H. R. (1989): “Inferring the components of the bid-ask spread: theory and empirical tests,” *Journal of Finance*, 44, 115–134.
- TSAY, R. S. (2005): *Analysis of Financial Time Series*, vol. 543, John Wiley & Sons.
- XIU, D. (2010): “Quasi-maximum likelihood estimation of volatility with high frequency data,” *Journal of Econometrics*, 159, 235–250.
- ZHANG, L. (2006): “Efficient estimation of stochastic volatility using noisy observations: A multi-scale approach,” *Bernoulli*, 12, 1019–1043.
- ZHANG, L., P. A. MYKLAND, AND Y. AÏT-SAHALIA (2005): “A tale of two time scales: Determining integrated volatility with noisy high-frequency data,” *Journal of the American Statistical Association*, 100, 1394–1411.
- ZHOU, B. (1996): “High-frequency data and volatility in foreign-exchange rates,” *Journal of Business & Economic Statistics*, 14, 45–52.

## Tables and Figures

$\iota$	-0.7	-0.3	0	0.3	0.7
$\widehat{IV}_n$	-22.37 (14.15)	-22.36 (14.17)	-22.36 (14.18)	-22.40 (14.20)	-22.87 (14.27)
$\widehat{IV}_n^{(1)}$	-1.71 (4.19)	-0.97 (4.21)	-0.23 (4.24)	0.85 (4.29)	4.33 (4.47)
$\widehat{IV}_n^{(2)}$	-0.94 (5.58)	-0.55 (5.61)	-0.19 (5.65)	0.31 (5.72)	1.57 (5.98)
$\widehat{IV}_n^{(3)}$	-0.55 (6.32)	-0.35 (6.35)	-0.17 (6.40)	0.04 (6.48)	0.19 (6.79)
$\widetilde{IV}_n$	-22.66 (13.94)	-22.66 (13.96)	-22.68 (13.96)	-22.74 (13.97)	-23.30 (13.99)
$\widetilde{IV}_n^{(1)}$	-2.00 (3.07)	-1.27 (3.08)	-0.55 (3.09)	0.51 (3.10)	3.90 (3.18)
$\widetilde{IV}_n^{(2)}$	-1.37 (3.60)	-1.01 (3.60)	-0.67 (3.60)	-0.20 (3.61)	0.93 (3.69)
$\widetilde{IV}_n^{(3)}$	-1.06 (3.89)	-0.88 (3.89)	-0.73 (3.90)	-0.56 (3.91)	-0.55 (3.99)
$\widetilde{IV}_n^{\text{JLZ}}$	-11.74 (7.63)	-11.65 (7.63)	-11.65 (7.64)	-11.65 (7.65)	-11.19 (7.68)
$\widehat{IV}_n^{\text{QMLE}}$	0.83 (10.13)	-0.19 (3.16)	-0.18 (3.43)	0.04 (3.52)	1.08 (4.35)

Table 1: Estimation of the IV. We take  $\Delta = 1$  sec and the number of observations is  $n = 23,400$ . We report the estimation results of three groups of IV estimators: our pre-averaging estimator and its multi-step versions based on non-overlapping intervals  $\widehat{IV}_n, \widehat{IV}_n^{(1)}, \widehat{IV}_n^{(2)}$  and  $\widehat{IV}_n^{(3)}$ ; our pre-averaging estimator and its multi-step versions based on overlapping intervals  $\widetilde{IV}_n, \widetilde{IV}_n^{(1)}, \widetilde{IV}_n^{(2)}$  and  $\widetilde{IV}_n^{(3)}$ ; the estimator  $\widetilde{IV}_n^{\text{JLZ}}$  based on the pre-averaging method proposed in [Jacod et al. \(2019\)](#) and the estimator  $\widehat{IV}_n^{\text{QMLE}}$  based on the QMLE method in [Da and Xiu \(2019\)](#). The numbers represent the centered mean estimates based on 1,000 simulations with standard deviations between parentheses. All numbers in the table have been multiplied by  $10^5$ . The tuning parameters for the first eight estimators are  $j_n = 20$ ,  $\ell_n = 10$  and  $\theta = 0.4$ , and we use the triangular kernel. For the estimator in [Jacod et al. \(2019\)](#) we used the choices suggested in that paper:  $\bar{h}_n = 0.5/\sqrt{\Delta_n}$ ,  $k_n = (\Delta_n)^{-1/5}$  and  $k'_n = (\Delta_n)^{-1/8}$ . In [Da and Xiu \(2019\)](#) the parameter  $q$  of the fitted  $\text{MA}(q)$  model was found by optimization over  $q \in \{8, 9, 10\}$  only for each sample in order to save time, since test runs indicated that the optimal order was usually around  $q = 9$ .



$\iota$	-0.7	-0.3	0	0.3	0.7
$\widehat{\text{IV}}_n$	-4.49 (3.87)	-4.49 (3.88)	-4.49 (3.90)	-4.50 (3.93)	-4.62 (4.05)
$\widehat{\text{IV}}_n^{(1)}$	-1.47 (2.83)	-0.72 (2.85)	0.02 (2.87)	1.13 (2.91)	4.98 (3.06)
$\widehat{\text{IV}}_n^{(2)}$	-0.11 (3.07)	-0.03 (3.09)	0.04 (3.11)	0.13 (3.15)	0.40 (3.32)
$\widehat{\text{IV}}_n^{(3)}$	0.02 (3.09)	0.04 (3.12)	0.04 (3.14)	0.03 (3.18)	-0.06 (3.35)
$\widetilde{\text{IV}}_n$	-4.83 (3.48)	-4.83 (3.48)	-4.83 (3.49)	-4.85 (3.50)	-5.00 (3.55)
$\widetilde{\text{IV}}_n^{(1)}$	-1.80 (2.13)	-1.06 (2.13)	-0.32 (2.14)	0.78 (2.15)	4.60 (2.20)
$\widetilde{\text{IV}}_n^{(2)}$	-0.48 (2.28)	-0.41 (2.29)	-0.34 (2.29)	-0.25 (2.31)	-0.02 (2.37)
$\widetilde{\text{IV}}_n^{(3)}$	-0.35 (2.30)	-0.34 (2.30)	-0.34 (2.31)	-0.36 (2.32)	-0.48 (2.39)
$\widetilde{\text{IV}}_n^{\text{JLZ}}$	-3.79 (3.11)	-3.68 (3.12)	-3.68 (3.12)	-3.68 (3.13)	-3.09 (3.17)
$\widehat{\text{IV}}_n^{\text{QMLE}}$	0.50 (3.61)	-0.69 (2.64)	-0.76 (3.16)	-0.80 (3.28)	0.28 (4.74)

Table 2: Estimation of the IV. We take  $\Delta = 0.2$  sec and the number of observations is  $n = 117,000$ . We report the estimation results of three groups of IV estimators: our pre-averaging estimator and its multi-step versions based on non-overlapping intervals  $\widehat{\text{IV}}_n, \widehat{\text{IV}}_n^{(1)}, \widehat{\text{IV}}_n^{(2)}$  and  $\widehat{\text{IV}}_n^{(3)}$ ; our pre-averaging estimator and its multi-step versions based on overlapping intervals  $\widetilde{\text{IV}}_n, \widetilde{\text{IV}}_n^{(1)}, \widetilde{\text{IV}}_n^{(2)}$  and  $\widetilde{\text{IV}}_n^{(3)}$ ; the estimator  $\widetilde{\text{IV}}_n^{\text{JLZ}}$  based on the pre-averaging method proposed in [Jacod et al. \(2019\)](#) and the estimator  $\widehat{\text{IV}}_n^{\text{QMLE}}$  based on the QMLE method in [Da and Xiu \(2019\)](#). The numbers represent the centered mean estimates based on 1,000 simulations with standard deviations between parentheses. All numbers in the table have been multiplied by  $10^5$ . The tuning parameters for the first eight estimators are  $j_n = 20$ ,  $\ell_n = 10$  and  $\theta = 0.4$ , and we use the triangular kernel. For the estimator in [Jacod et al. \(2019\)](#) we used the choices suggested in that paper:  $\bar{h}_n = 0.5/\sqrt{\Delta_n}$ ,  $k_n = (\Delta_n)^{-1/5}$  and  $k'_n = (\Delta_n)^{-1/8}$ . In [Da and Xiu \(2019\)](#) the parameter  $q$  of the fitted  $\text{MA}(q)$  model was found by optimization over  $q \in \{8, 9, 10\}$  only for each sample in order to save time, since test runs indicated that the optimal order was usually around  $q = 9$ .

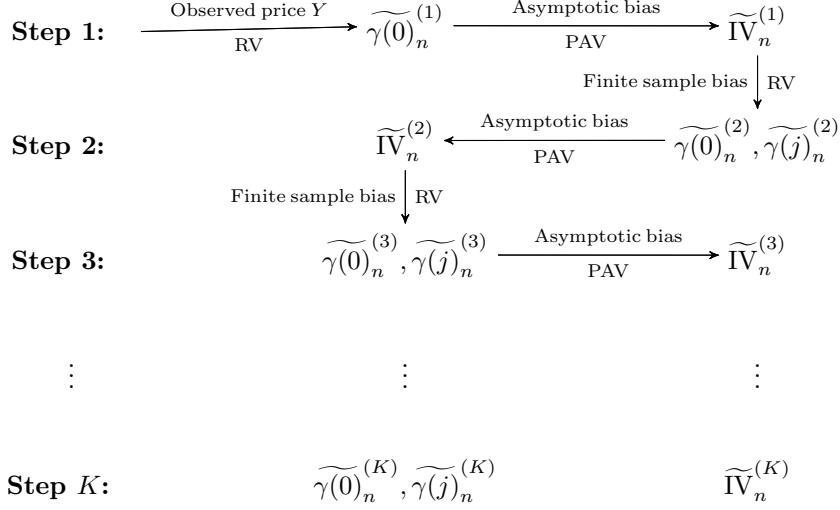


Figure 1: Illustration of the construction of the multi-step estimators. In the first step, we use realized volatility (RV) to obtain an estimator of the variance of (possibly misspecified) i.i.d. noise,  $\widetilde{\gamma(0)}_n^{(1)}$ . Next, this estimator is used to correct the asymptotic bias of the pre-averaging estimator (PAV) to derive the first-step estimator of the IV,  $\widetilde{IV}_n^{(1)}$ . In the second step, we use  $\widetilde{IV}_n^{(1)}$  to obtain finite sample bias corrected estimators of the variance and covariances of noise,  $\widetilde{\gamma(0)}_n^{(2)}$  and  $\widetilde{\gamma(j)}_n^{(2)}$ , which are then used to remove the asymptotic bias in PAV, leading to the second-step IV estimator,  $\widetilde{IV}_n^{(2)}$ . Iterating this procedure will lead to  $K$ -step estimators  $\widetilde{\gamma(0)}_n^{(K)}, \widetilde{\gamma(j)}_n^{(K)}, \widetilde{IV}_n^{(K)}$ .

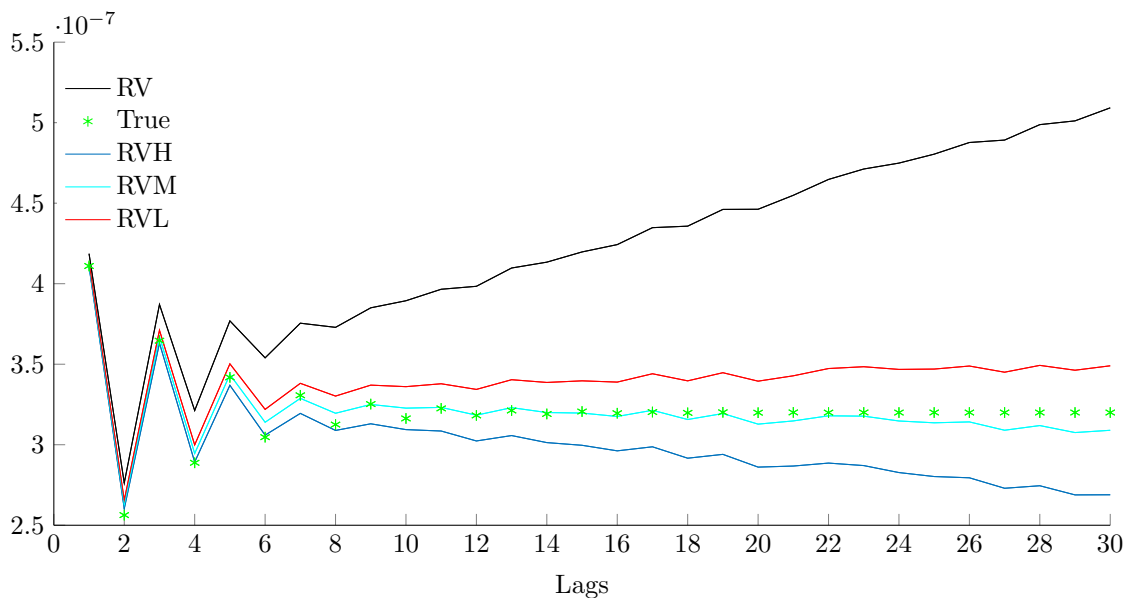


Figure 2: Realized volatility estimators against the number of lags  $j$ , based on a single simulated sample, without and with finite sample bias correction, cf. (6) and (10). Here, RV:  $\widehat{\langle Y, Y \rangle}(j)_n$ ; RVL:  $\widehat{\langle Y, Y \rangle}(j)_n - \frac{0.8jIV}{2(n-j+1)}$ ; RVM:  $\widehat{\langle Y, Y \rangle}(j)_n - \frac{jIV}{2(n-j+1)}$ ; and RVH:  $\widehat{\langle Y, Y \rangle}(j)_n - \frac{1.2jIV}{2(n-j+1)}$ . We take  $\Delta = 1$  sec, the number of observations is 23,400, and  $\iota = -0.7$ . The designation “True” corresponds to the stochastic limit  $\gamma(0) - \gamma(j)$ .

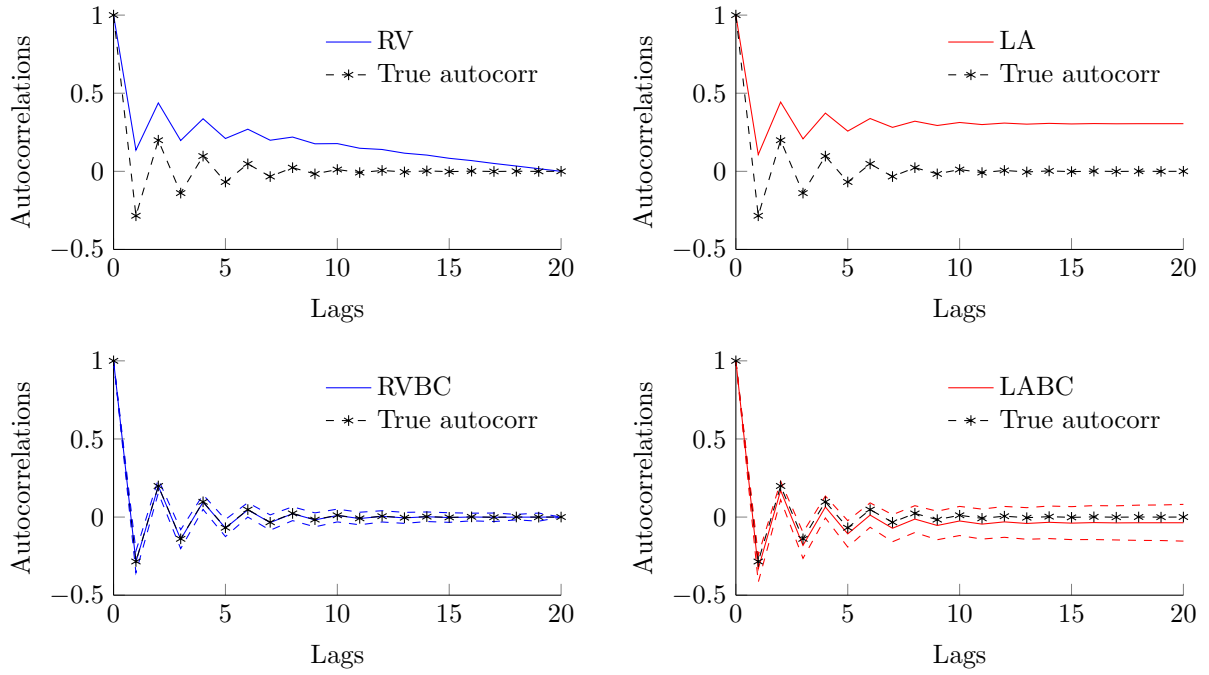


Figure 3: Realized volatility (RV) and local averaging (LA) estimators of the autocorrelations of noise against the number of lags  $j$ , averaged over 1,000 simulated samples. Top panel: RV and LA estimators without finite sample bias corrections. Bottom panel: RV and LA estimators with finite sample bias corrections (RVBC, LABC). The dashed lines are the 95% simulated confidence intervals. We take  $\Delta = 1$  sec, the number of observations is 23,400, and  $\iota = -0.7$ . The tuning parameters of the RV and LA estimators are  $j_n = 20$  and  $K_n = 6$ , respectively.

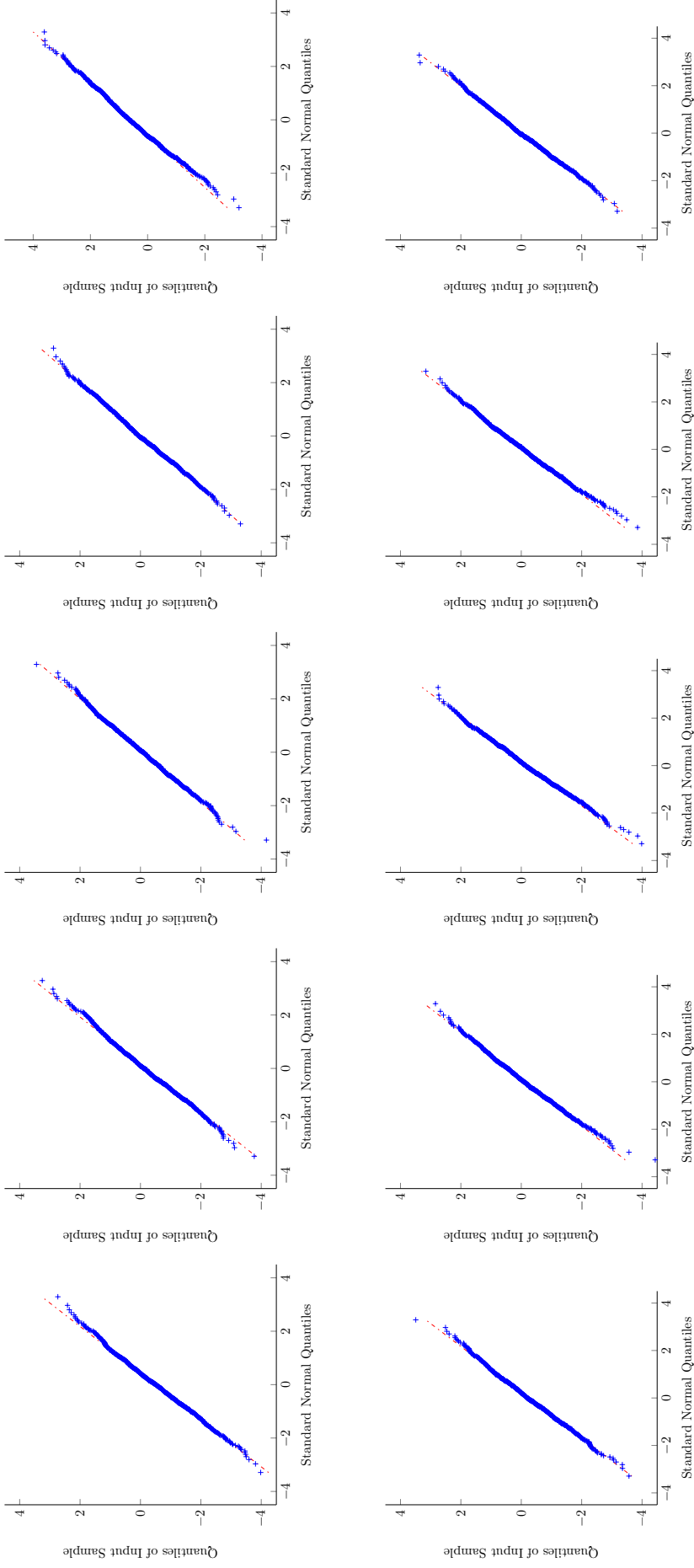


Figure 4: Standard normal QQ-plots of the second-step IV estimators. Top panel:  $\Delta_n^{-\frac{1}{4}} \left( \widehat{\text{IV}}_n^{(2)} - \text{IV} \right) / \sqrt{\widehat{\Sigma}_{\text{IV}_n}^{(2)}}$ . Bottom panel:  $\Delta_n^{-\frac{1}{4}} \left( \widetilde{\text{IV}}_n^{(2)} - \text{IV} \right) / \sqrt{\widetilde{\Sigma}_{\text{IV}_n}^{(2)}}$ . The AR(1) coefficient  $\iota$  of the noise process, from left to right, is  $-0.7, -0.3, 0, 0.3$ , and  $0.7$ . The number of simulations is 1,000, the data frequency is  $\Delta = 0.1$  sec, and the number of observations is 234,000. The tuning parameter of the RV estimator is  $j_n = 20$ , and  $\ell_n = 10$ . The tuning parameter  $\theta$  equals  $1/3$ .

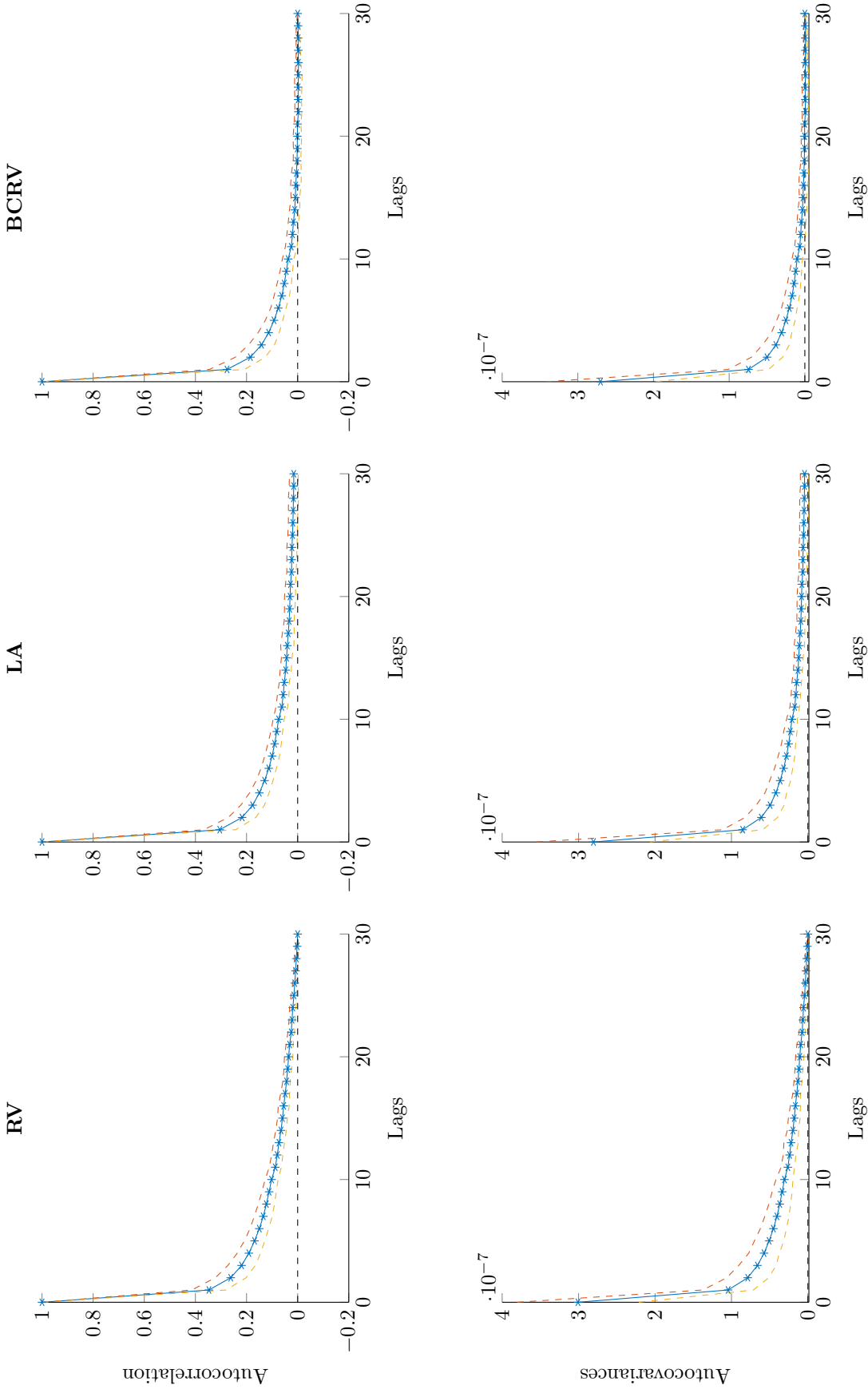


Figure 5: From the left to the right, we display the realized volatility (RV), local averaging (LA), and the bias corrected realized volatility (BCRV) estimators of the autocorrelations (top panel) and autocovariances (bottom panel) of noise against the number of lags  $j$  based on transaction data for Citigroup. Sample period: January, 2011. On average there are 10.5 observations per second in the sample. The three estimators are applied to and then averaged over each of the 20 trading days. The stars indicate the means of the 20 estimates. The dashed lines are 2 standard deviations away from the mean. The tuning parameter of the RV estimator is  $j_n = 30$ .

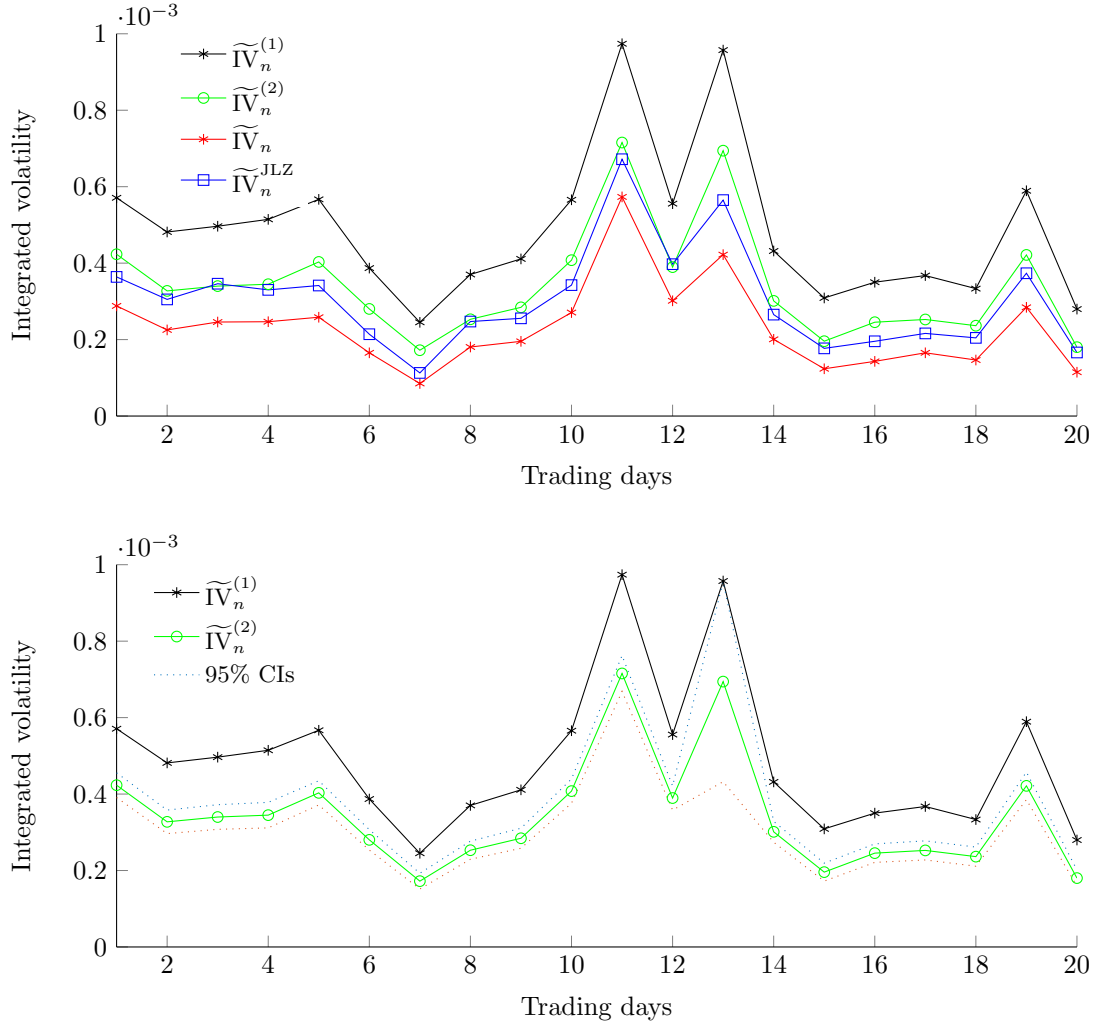


Figure 6: Estimation of the IV based on transaction data for Citigroup. Sample period: January, 2011, consisting of 20 trading days. On average there are 10.5 observations per second in the sample. The estimators  $\widetilde{IV}_n^{(1)}$ ,  $\widetilde{IV}_n^{(2)}$ , and  $\widetilde{IV}_n$  are given by (36), (41), and (25). The  $\widetilde{IV}_n^{JLZ}$  estimator is proposed in Jacod et al. (2019). In the bottom panel, the asymptotic confidence intervals (CIs) are based on the limit distribution in Theorem 5.1. The tuning parameter of the RV estimator is  $j_n = 30$ , and  $\ell_n = 10$ .  $\theta$  is selected according to (28).

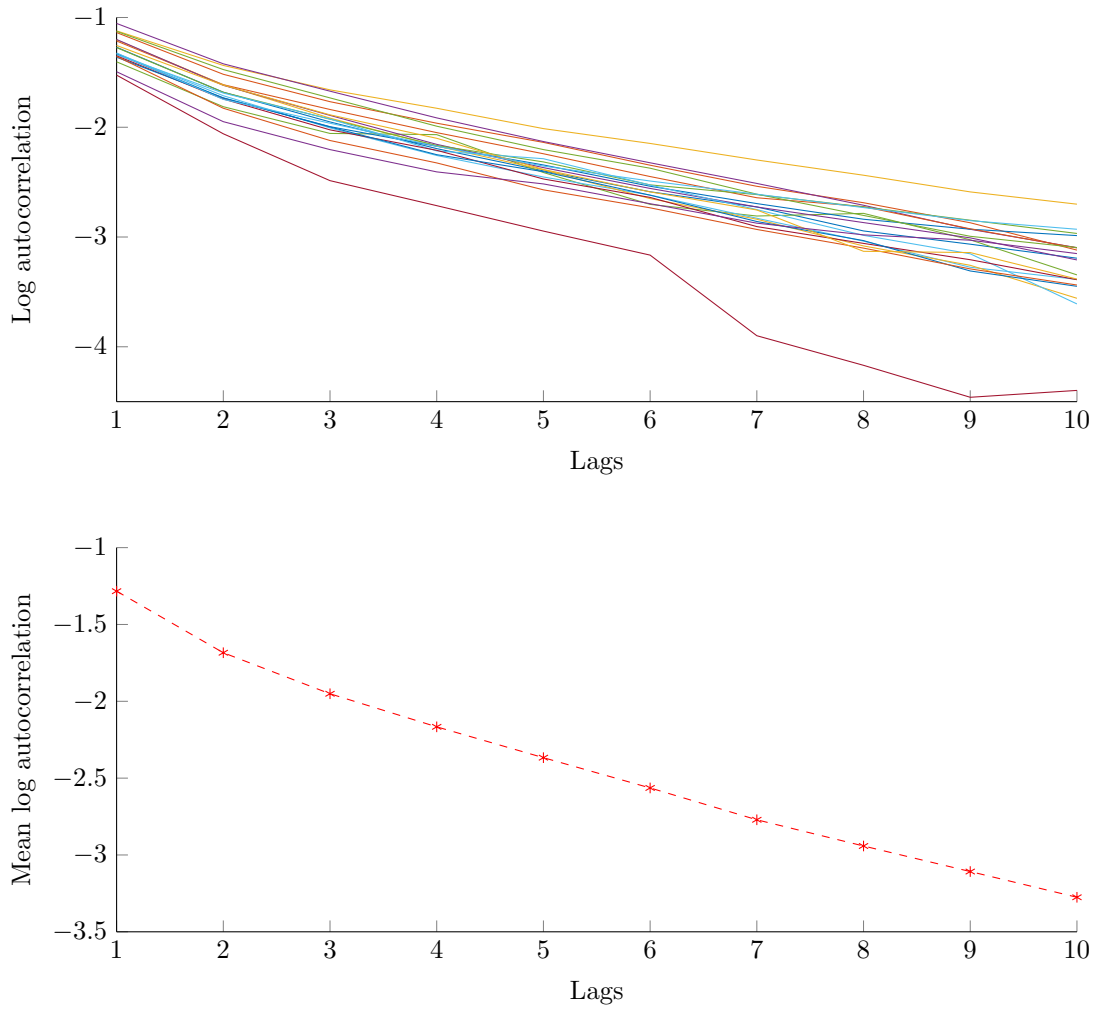


Figure 7: Top panel: Logarithmic autocorrelations of noise against the number of lags  $j$  estimated by BCRV for each trading day based on transaction data for Citigroup. Bottom panel: Means of the logarithmic autocorrelations of noise and a linear regression line. Sample period: January, 2011, consisting of 20 trading days. On average there are 10.5 observations per second in the sample. The tuning parameter of the RV estimator is  $j_n = 30$ .



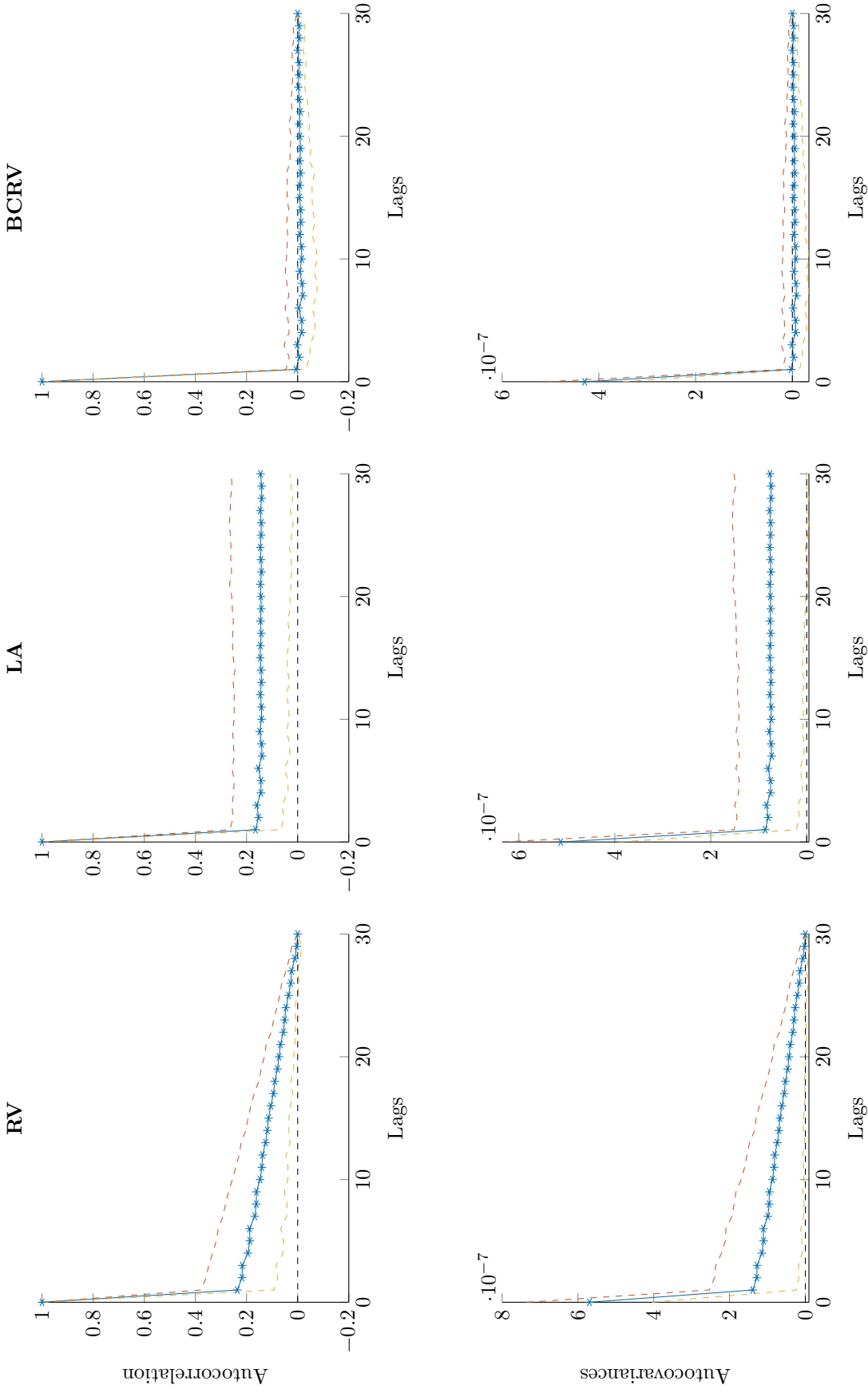


Figure 8: From the left to the right, we display the realized volatility (RV), local averaging (LA), and the bias corrected realized volatility (BCRV) estimators of the autocorrelations (top panel) and autocovariances (bottom panel) of noise against the number of lags  $j$  based on a subsample of the transaction data for Citigroup. Sample period: January, 2011. The subsample is recorded on a 1-sec time scale. The three estimators are applied to and next averaged over each of the 20 trading days. The stars indicate the means of the 20 estimates. The dashed lines are 2 standard deviations away from the mean. The tuning parameter of the RV estimator is  $j_n = 30$ .

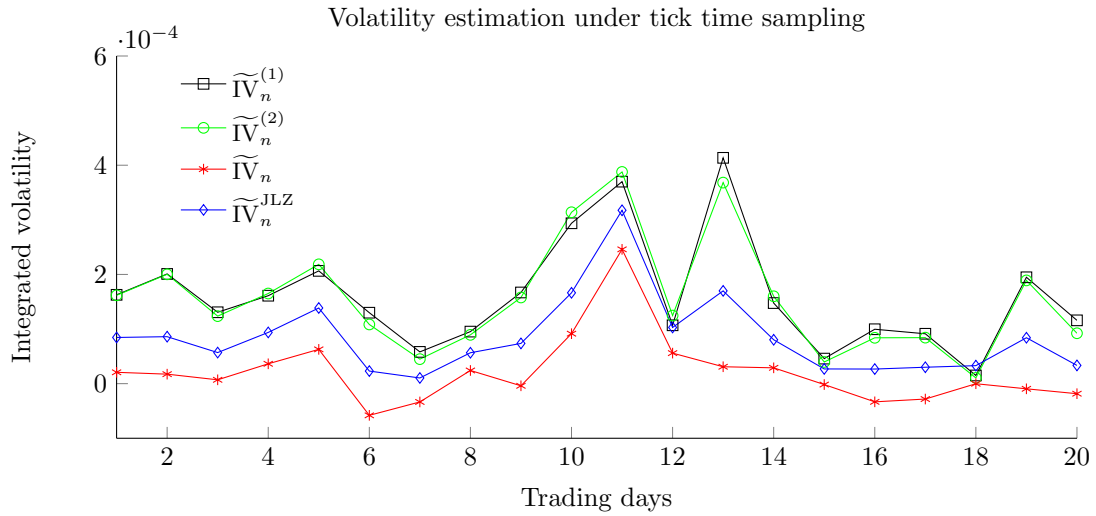
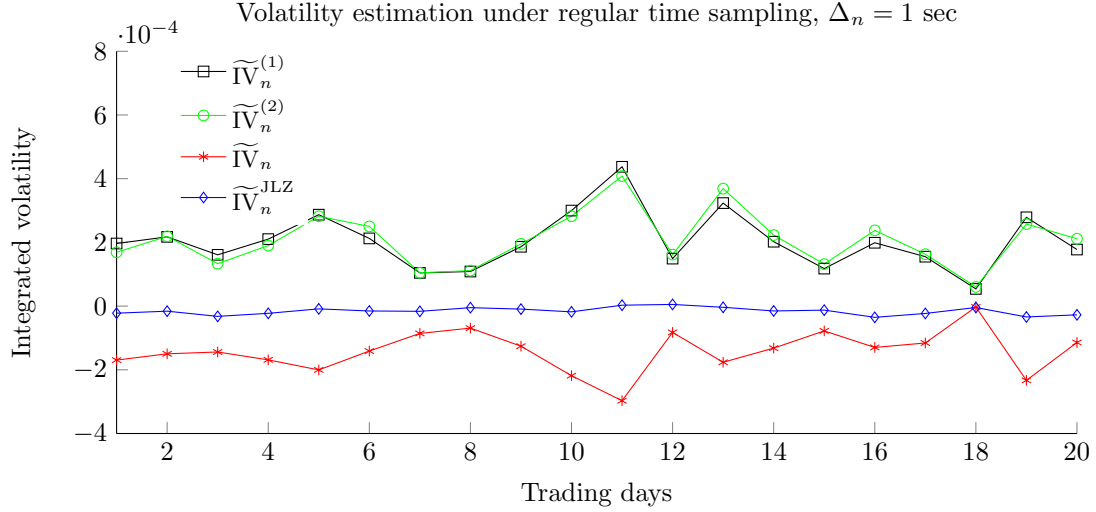


Figure 9: Estimation of the IV based on subsamples of the transaction data for Citigroup. Sample period: January, 2011, consisting of 20 trading days. In the top panel, the estimation is performed on a subsample that is recorded on a 1-sec time scale. In the bottom panel, the estimation is performed on a subsample that is recorded at tick time; on average there are 3.2 observations per second in the sample. The estimators  $\widetilde{IV}_n^{(1)}$ ,  $\widetilde{IV}_n^{(2)}$ , and  $\widetilde{IV}_n$  are given by (36), (41), and (25). The  $\widetilde{IV}_n^{\text{JLZ}}$  estimator is proposed in Jacod et al. (2019). The tuning parameter of the RV estimator is  $j_n = 30$ , and  $\ell_n = 4$  for the 1-sec sample and  $\ell_n = 6$  for the tick time sample.  $\theta$  is selected according to (28).

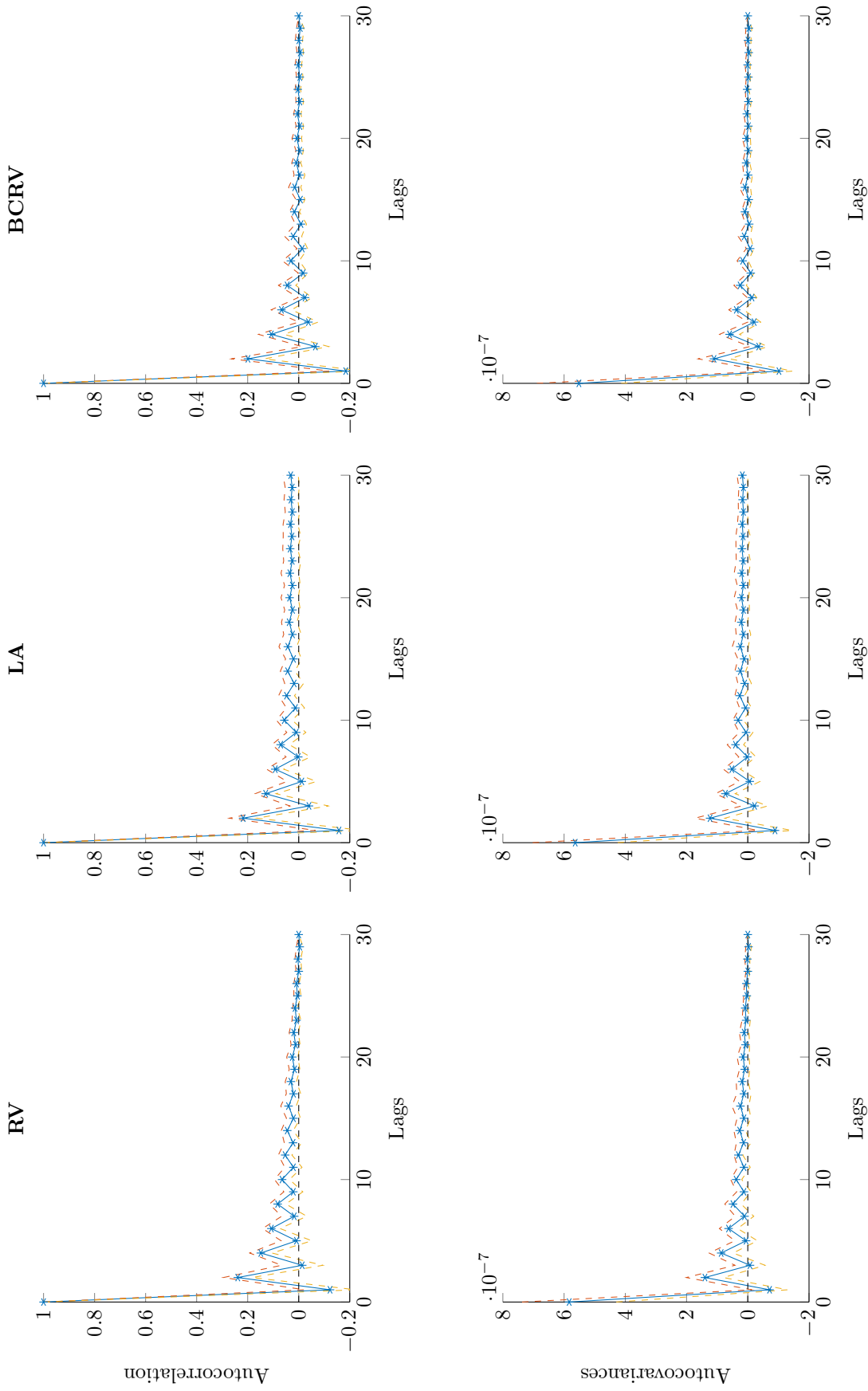


Figure 10: From the left to the right, we display the realized volatility (RV), local averaging (LA), and the bias corrected realized volatility (BCRV) estimators of the autocorrelations (top panel) and autocovariances (bottom panel) of noise against the number of lags  $j$  based on a subsample of the transaction data for Citigroup. Sample period: January, 2011. The subsample is recorded at tick time. On average there are 3.2 observations per second in the sample. The three estimators are applied to and then averaged over each of the 20 trading days. The stars indicate the means of the 20 estimates. The dashed lines are 2 standard deviations away from the mean. The tuning parameter of the RV estimator is  $j_n = 30$ .