Efficient Estimation of Bid-Ask Spreads from Open, High, Low, and Close Prices^{*}

1

2

3

5

David Ardia, Emanuele Guidotti, and Tim A. Kroencke[†]

June 22, 2024

Abstract

Popular bid-ask spread estimators are downward biased when trading is infrequent. 6 Moreover, they consider only a subset of open, high, low, and close prices and ne-7 glect potentially useful information to improve the spread estimate. By accounting 8 for discretely observed prices, this paper derives asymptotically unbiased estimators 9 of the effective bid-ask spread. Moreover, we combine them optimally to minimize 10 the estimation variance and obtain an efficient estimator. Through theoretical anal-11 yses, numerical simulations, and empirical evaluations, we show that our efficient 12 estimator dominates other estimators from transaction prices, yields novel insights 13 for measuring bid-ask spreads, and has broad applicability in empirical finance. 14

^{*}We are grateful to the editor and two anonymous referees for their valuable comments. We also thank Jean-Philippe Bouchaud, Victor DeMiguel, Laurent Frésard, Björn Hagströmer, Angelo Ranaldo, Sebastian Stöckl, Dennis Umlandt, and conference and seminar participants at the Annual Meeting of the European Finance Association (EFA, 2022), Annual Meeting of the Financial Management Association (FMA, 2022), International Conference on Quantitative Finance and Financial Econometrics (QFFE, 2022), New Zealand Finance Meeting (NZFM, 2021), Annual Financial Market Liquidity Conference (AFML, 2021), Annual Meeting of the German Finance Association (DGF, 2021), International Risk Management Conference (IRMC, 2021), Goethe University Frankfurt, HEC Montréal, and INSEAD business school for helpful suggestions. We acknowledge financial support from the Institute for Data Valorization (IVADO). This paper is based on the first chapter of Emanuele Guidotti's PhD thesis at the University of Neuchâtel.

[†]David Ardia is at the Department of Decision Sciences and GERAD, HEC Montréal. Emanuele Guidotti is at the Institute of Finance, USI Lugano. Tim A. Kroencke is at the FHNW School of Business, Basel, and at Remaco Asset Management, Basel. E-mail addresses: david.ardia@hec.ca (D. Ardia), emanuele.guidotti@usi.ch (E. Guidotti), tim.kroencke@fhnw.ch (T. A. Kroencke).

Conflict-of-interest disclosure statement

16	DAVID ARDIA
17	I have nothing to disclose.
18	EMANUELE GUIDOTTI
19	I have nothing to disclose.
20	TIM A. KROENCKE
21	Member of the Academic Advisory Council of Vontobel Bank, Zurich.

THE EFFECTIVE BID-ASK SPREAD measures the distance of observed transaction prices 22 from the unobserved fundamental price, and it is a predominant measure of transaction 23 costs in financial markets. The literature on measuring bid-ask spreads has proceeded 24 along two complementary paths that focus on either high-frequency or low-frequency data. 25 The high-frequency literature relies on trade and quote data to obtain an explicit proxy 26 of the fundamental price and calculate the distance of transaction prices from it (Holden 27 and Jacobsen, 2014; Stoikov, 2018; Hagströmer, 2021). The low-frequency literature 28 introduces assumptions about the fundamental price to derive estimators from transaction 29 prices only, without requiring any information about quotes (Roll, 1984; Hasbrouck, 2009; 30 Corwin and Schultz, 2012; Abdi and Ranaldo, 2017). 31

While measures from trades and quotes are typically more accurate, low-frequency 32 estimates are more readily available and are becoming increasingly popular due to the 33 difficulties and costs of obtaining quote data for international markets, historical data 34 samples, and asset classes other than stocks.¹ However, the estimators developed thus 35 far rely on the assumption that prices are observed continuously. In contrast, the number 36 of trades within any time interval is finite in real markets, and prices unfold in discrete 37 time. We show that this assumption causes a downward bias when the number of trades 38 per observation period is small. Moreover, these estimators consider only a subset of open, 39 high, low, and close prices and thus neglect potentially useful information to improve the 40 spread estimate and reduce the estimation variance. Jahan-Parvar and Zikes (2023) show 41

¹For instance, recent use cases of low-frequency estimators include: stock return predictability and asset pricing anomalies (McLean and Pontiff, 2016; Hou, Xue, and Zhang, 2018; Chen, Eaton, and Paye, 2018; Birru, 2018; Hua et al., 2019; Jacobs and Müller, 2020; Patton and Weller, 2020; Amihud and Noh, 2020; Chaieb, Errunza, and Langlois, 2020); municipal and corporate bonds (Schwert, 2017; Bongaerts, de Jong, and Driessen, 2017; Cai et al., 2019; Kaviani et al., 2020; Bali, Subrahmanyam, and Wen, 2021; Ding, Xiong, and Zhang, 2022); bond funds (Goldstein, Jiang, and Ng, 2017; Choi et al., 2020); currency markets (Michaelides, Milidonis, and Nishiotis, 2019; Ranaldo and de Magistris, 2022); OTC derivatives (Loon and Zhong, 2016); interest rates (Ranaldo, Schaffner, and Vasios, 2021); monetary policy (Grosse-Rueschkamp, Steffen, and Streitz, 2019); institutional trading costs (Eaton, Irvine, and Liu, 2021); investor behavior (Li, Subrahmanyam, and Yang, 2018); information and dark pools (Brogaard and Pan, 2021); and machine learning (Easley et al., 2020). See Table I.1 in the Internet Appendix for a survey.

that a larger estimation variance causes a larger upward bias when the spread is small compared to volatility due to the methods employed to guarantee non-negativity of the spread estimates in small samples. In summary, current estimators understate bid-ask spreads when expected to be the largest and overstate bid-ask spreads when expected to be the smallest.

In this paper, we develop an asymptotically unbiased estimator with minimum vari-47 ance by accounting for discretely observed prices and optimally considering the complete 48 information set of open, high, low, and close prices. First, we derive multiple bid-ask 49 spread estimators from several combinations of prices. Our methodology yields estima-50 tors with an analytical term that depends on the probability that opening or closing prices 51 coincide with the highest or lowest prices. Such probability would be zero if prices were 52 observed continuously, and it can be regarded as an analytical correction term account-53 ing for discretely observed prices. To give a sense of the importance of correcting by this 54 term, Figure 1 displays the probability that daily opening or closing prices coincide with 55 the highest or lowest prices for U.S. common stocks from 1926 to 2021. The probability 56 ranges between 25% for large stocks and 75% for small stocks in the last century and 57 decreased in the last two decades. Thus, while correcting by this term is less significant 58 for more recent periods, it becomes essential when analyzing historical samples, and it 59 is increasingly important for smaller stocks. Moreover, this term also depends on the 60 sampling frequency used for estimation. Indeed, if intraday—instead of daily—prices are 61 used, then the number of trades observed per time interval decreases, and the probability 62 that opening or closing prices coincide with the highest or lowest prices increases. Thus, 63 this term corrects a bias that varies in the time series and the cross-section and depends 64 on the sampling frequency of open, high, low, and close prices. 65

66

[Insert Figure 1 about here.]

⁶⁷ Next, we combine our estimators to construct an efficient estimator. All estimators

are asymptotically unbiased, so their efficient combination is obtained by minimizing the 68 estimation variance. We proceed as follows. First, we identify two estimators that achieve 69 minimum variance when the spread is small compared to volatility. Second, we identify 70 two other estimators that exhibit the opposite behavior and achieve minimum variance 71 when the spread is large compared to volatility. Third, we show that these estimators can 72 be written as moment conditions and apply the generalized method of moments (Hansen, 73 1982) to construct our efficient estimator that achieves minimum variance across small 74 and large spreads. By minimizing the estimation variance, our efficient estimation also 75 minimizes the upward bias that arises in small samples due to the methods employed to 76 guarantee non-negativity of the spread estimates (Jahan-Parvar and Zikes, 2023). 77

We compare our efficient estimator with the seminal Roll (1984) estimator and with those by Corwin and Schultz (2012) and Abdi and Ranaldo (2017) as they have been shown to deliver more accurate estimates than previous approaches.

In our simulation experiments, we study the bias and variance of the estimators. In 81 agreement with our theoretical analysis, we find that other estimators understate the 82 spread in simulations that use few trades per period, and the estimate shrinks to zero 83 as the number of trades declines. Instead, our estimator remains unbiased even for 84 simulations where we expect, on average, only a single trade per period. For simulations 85 that use many trades per period, we find that all estimators are asymptotically unbiased, 86 and they correctly estimate the spread used in the simulation. In this case, the best 87 estimator has the lowest variance because it delivers unbiased estimates with higher 88 precision. We find that the estimator by Corwin and Schultz (2012) has a lower variance 89 than Abdi and Ranaldo (2017) for small spreads, while it has a higher variance for large 90 spreads. Our efficient estimator provides the most precise estimates with a variance lower 91 than the other approaches across low and large spreads. In summary, our estimator 92 dominates other approaches by yielding unbiased estimates when other estimators are 93

⁹⁴ biased and achieving minimum variance when all estimators are unbiased.

Our empirical analysis uses the Center for Research in Security Prices (CRSP) U.S. 95 stock database to compute bid-ask spread estimates from daily prices. We compare 96 them with the effective spread computed by matching high-frequency trades with quotes 97 via the NYSE Trades and Quotes (TAQ) database in the sample period 1993–2021. 98 The simulation-based results carry over to the empirical data. Our efficient estimator 99 dominates all other estimators, and it is more correlated and considerably closer to the 100 high-frequency benchmark in each sub-period, in each market venue, for small and large 101 stocks, both in time series and cross-sectional studies, for each sample size and evaluation 102 metric. 103

We illustrate the broad applicability of our estimator in low- and high-frequency both 104 within and outside the U.S. stock market. First, we revisit historical spread estimates 105 from daily prices in the U.S. stock market since 1926. For small stocks, our estimator 106 closely overlaps with the high-frequency benchmark. In contrast, other estimators under-107 state the spread, and their bias increases for older sample periods, mirroring that these 108 estimators are more biased when trading becomes less frequent. Indeed, their bias re-109 duces for larger stocks, which are presumably traded more frequently. For all stocks, we 110 find that the end-of-day quoted spread is higher than our effective spread estimates by a 111 factor of two. Thus, our estimator reproduces previous findings that the quoted spread 112 overstates the effective spread finally paid by traders by up to 100% (Huang and Stoll, 113 1994; Petersen and Fialkowski, 1994; Bessembinder and Kaufman, 1997; Bacidorea, Ross, 114 and Sofianosa, 2003), due to dealers offering a better price than the quotes, also known 115 as trading inside the spread (Lee, 1993). In summary, our estimator makes available the 116 most realistic effective spread estimates for the U.S. stock market from 1926 to the advent 117 of high-frequency data. 118

119

Second, we show that our estimator can exploit intraday prices to improve the spread

estimate significantly and that this approach is more effective than increasing the estima-120 tion sample with more daily data. Using minute prices—instead of daily—increases the 121 correlation of the estimates with the benchmark from 56.17% to 88.79% in the challeng-122 ing sample from October 2003 to December 2021, where the spread is small compared to 123 volatility. The fraction of non-positive estimates reduces from 34.15% to 0.02%, and the 124 upward bias induced by resetting negative estimates to zero essentially vanishes (Jahan-125 Parvar and Zikes, 2023). These results show that our estimator can be applied at any 126 frequency, and, in this sense, it reconciles the high-frequency and low-frequency litera-127 ture. Moreover, by relying on transaction prices only, our estimator is insensitive to the 128 quality of quote data, which causes issues in measuring effective spreads by matching 129 trades with quotes in fast and competitive markets (Holden and Jacobsen, 2014). 130

Third, we apply the estimator outside the stock market and analyze low- and highfrequency estimates for cryptocurrencies. We find that other estimators are dominated by their downward bias in high frequency and produce a tenfold difference between estimates that use daily or intraday prices. Instead, our estimator produces estimates from daily prices that closely overlap with those from hourly and minute prices. We conclude that our efficient estimator can potentially reduce a significant source of non-standard errors (Menkveld et al., 2023) in the measurement of transaction costs.

This paper is structured as follows. Section 1 reviews high- and low-frequency estimators of the effective bid-ask spread. Section 2 introduces our methodology and develops our estimators. Sections 3 and 4 present our simulation and empirical results, respectively. Section 5 illustrates the advantages and wide applicability of our efficient estimator. Finally, Section 6 concludes. To guarantee reproducibility, we make available software for the R statistical environment (R Core Team, 2020) that implements all the results in this paper. To facilitate adoption, we make our efficient estimator available in various ¹⁴⁵ programming languages. We also release open data containing all our spread estimates.²

¹⁴⁶ 1 The Effective Bid-Ask Spread

¹⁴⁷ For a given trade, the relative effective bid-ask spread S is defined as:

$$S = \frac{2D(P - \tilde{P})}{\tilde{P}},\tag{1}$$

where P is the observed transaction price, \tilde{P} is the unobserved fundamental price, and 148 D is a direction of trade indicator taking the value +1 for buyer-initiated trades, and -1149 for seller-initiated trades. As the fundamental price \tilde{P} is unobserved, different ways of 150 estimating the spread exist, which depend on different proxies for \tilde{P} . Here, we review 151 popular measures of the effective bid-ask spread that arise from different proxies. We 152 classify these measures into two groups. First, we discuss measures that require trade 153 and quote data and are typically used in the high-frequency literature. Then, we discuss 154 measures that only require transaction prices and are typically used in the low-frequency 155 literature. 156

¹⁵⁷ 1.1 High-Frequency Measures of Effective Spreads

One way to measure effective spreads is obtaining a proxy of the fundamental price from trade and quote data to plug in Equation (1). This class of estimators measures the distance of transaction prices from the given proxy. Popular proxies are the quoted midpoint, the weighted midpoint, and the microprice.

²The code implementing the estimator is available at https://github.com/eguidotti/bidask. The code to reproduce the paper and the data containing spread estimates will be available upon publication.

162 1.1.1 Quoted Midpoint

¹⁶³ A simple proxy of the fundamental price is the average of the bid and ask prices. The ¹⁶⁴ quoted midpoint P_M is defined as:

$$P_M = \frac{P_A + P_B}{2},\tag{2}$$

where P_A and P_B are the ask and bid prices, respectively. Using $\tilde{P} = P_M$ in Equation (1) we obtain the so-called midpoint effective spread (Hagströmer, 2021). This midpointbased measure is required in U.S. regulations (SEC current Rule 605, Rule 11ac1-5 before 2007) and is often referred to as the effective spread. Here, we use the more precise terminology of Hagströmer (2021) to highlight that effective spreads are not observable and depend on the choice of the fundamental price. The midpoint effective spread is one possible measure of effective spreads.

172 1.1.2 Weighted Midpoint

Hagströmer (2021) challenges using the quoted midpoint as a proxy of the fundamental
price and shows that it leads to overstating effective spreads in markets with discrete
prices and elastic liquidity demand. To overcome this problem, he proposes to use the
weighted midpoint:

$$P_W = \frac{P_A Q_B + P_B Q_A}{Q_A + Q_B},\tag{3}$$

where Q_A and Q_B are the depths quoted at the ask and bid prices, respectively.

178 1.1.3 Microprice

Stoikov (2018) criticizes the midpoint and weighted midpoint as proxies of the fundamental price for generating autocorrelated returns and proposes an alternative proxy—the
microprice—that is a martingale by construction. We refer the reader to Stoikov (2018)

for the construction of the microprice and to Hagströmer (2021) for a comparison of effective spreads obtained with the midpoint, weighted midpoint, and microprice.

¹⁸⁴ 1.2 Low-Frequency Measures of Effective Spreads

Another way to measure effective spreads is by introducing desirable assumptions about the data-generating process to develop an estimator that does not require an explicit proxy of the fundamental price. This class of estimators measures the distance of transaction prices from a fundamental price implicitly defined by the model's assumptions.

Several contributions (Roll, 1984; Hasbrouck, 2009; Corwin and Schultz, 2012; Abdi and Ranaldo, 2017) have proposed to derive an estimator of the effective spread by writing Equation (1) in logarithmic prices $p = \log(P)$ such that:

$$p = \tilde{p} + Z \,, \tag{4}$$

where Z = S/2D is the bid-ask bounce and the basic assumptions are that:

¹⁹³ Assumption 1 Fundamental returns are not serially correlated.

194 Assumption 2 Fundamental returns are uncorrelated with bid-ask bounces.

195 Assumption 3 Bid-ask bounces are uncorrelated and have zero mean.

Assumptions 1–3 are the representative set of assumptions underlying previous contributions. However, they are more general than those required by each of them. For instance, the Roll (1984) model further assumes that buys and sells are equally likely. The Bayesian approach by Hasbrouck (2009) requires that fundamental returns are i.i.d. with normal distribution. Corwin and Schultz (2012) rely on the idea that high prices are buyer-initiated and low prices are seller-initiated and they model the fundamental price with a geometric Brownian motion with zero-mean returns, which is also used by

Abdi and Ranaldo (2017). They further assume that spread and volatility are constant, 203 ruling out stylized facts such as heteroscedasticity and jumps. To mitigate these restric-204 tions, they advocate in favor of measuring the spread over two-day rolling periods and 205 averaging these estimates. However, Jahan-Parvar and Zikes (2023) show this approach 206 produces inconsistent estimators. Finally, one important limitation of all the previous 207 contributions is that they do not account for the discrete nature of trades. Specifically, 208 they require (explicitly or implicitly) the restrictive assumption that there is always at 209 least one trade between two time instants such that prices are observed continuously. 210

Overall, the class of estimators based on Assumptions 1–3 aims at measuring the distance of transaction prices from a fundamental price with serially uncorrelated returns that are not correlated with bid-ask bounces. Such a class of estimators is the central focus of this paper, and we review the most popular approaches below. Other works alter the definition of the fundamental price by adding a dependence between the fundamental returns and the bid-ask bounces to accord an informational role to the trade directions, and they are outside the scope of this paper (see *e.g.*, Chen, Linton, and Yi, 2017).

218 1.2.1 Close Prices

The seminal work by Roll (1984) computes the serial covariance of observed returns to estimate the effective spread from closing prices. He shows that:

$$S^{2} = -4\mathbb{C}\operatorname{ov}[\Delta c_{t}, \Delta c_{t-1}], \qquad (5)$$

where S^2 is the mean squared spread in the estimation sample and $\Delta c_t = c_t - c_{t-1}$ where c_t is the closing log-price of period t. The main limitation of this approach is that it has a large estimation variance, and the squared spread turns out to be negative in 50% of the cases using a yearly sample of daily closing prices (Roll, 1984). To improve the estimation accuracy, Hasbrouck (2009) proposes a Gibbs estimation of the Roll model. However, the method requires an iterative procedure, is computationally expensive, and needs many observations to converge.³

228 1.2.2 High and Low Prices

Corwin and Schultz (2012) propose an alternative estimator from high and low prices 229 with smaller variance than the Roll (1984) estimator. Their methodology is based on 230 the idea that high (low) prices are almost always buy (sell) trades. Hence, the high-low 231 ratio incorporates both the volatility of the fundamental price and the bid-ask spread. 232 As volatility increases with the return interval, while the spread does not, it is possible 233 to derive a spread estimator from the high-low ratios over different time intervals. To 234 link the high-low ratios with volatility, they assume that the fundamental price follows a 235 geometric Brownian motion and use the equations by Parkinson (1980) and Garman and 236 Klass (1980). However, these equations hold only if the price is observed continuously 237 and are biased in practice as the number of trades within any time interval is finite. 238

239 1.2.3 Close, High, and Low Prices

Abdi and Ranaldo (2017) propose an estimator that jointly uses closing and high-low prices to achieve smaller variance than the Roll (1984) estimator and smaller bias than the Corwin and Schultz (2012) estimator. They show that:

$$S^{2} = 4\mathbb{E}[(c_{t-1} - \eta_{t-1})(c_{t-1} - \eta_{t})], \qquad (6)$$

³From Hasbrouck's website (https://pages.stern.nyu.edu/~jhasbrou/Research/GibbsCurrent/ gibbsCurrentIndex.html): "I often receive inquiries regarding Gibbs estimates formed at higher frequencies (e.g., monthly or weekly). I don't provide these estimates due to concerns about their reliability. The 2009 paper describes some of the issues that arise. Briefly, the prior distributions used here are diffuse (to ensure that the posteriors are data-dominated). The priors are generally, however, biased. As the sample size drops, the posteriors start resembling the priors, and the bias problem becomes more acute. The only way out of this is to put more structure on the priors. This is not impractical, but it is application-specific."

where $\eta_t = (h_t + l_t)/2$ is the average of the high and low log-prices. However, their methodology also requires that the fundamental price follows a geometric Brownian motion with continuously observed prices. As a consequence, the estimator is still biased. Moreover, it does not exploit the full information set of open, high, low, and close prices to further improve the spread estimate.

$_{^{248}}$ 2 Methodology

This paper relaxes the assumption that prices are observed continuously—and several other assumptions that were required by previous contributions—by deriving bid-ask spread estimators using Equation (4) under only Assumptions 1–3. By accounting for the discrete nature of trades, we drastically reduce the estimation bias. By exploiting the full information set of open, high, low, and close prices, we minimize the estimation variance.

255 We start by introducing the indicator variable:

$$\tau_t = \begin{cases} 0 & \text{if} \quad h_t = l_t = c_{t-1} \\ 1 & \text{otherwise} \end{cases}$$
(7)

that equals 0 if the highest price matches the lowest price and the previous close, and it equals 1 otherwise. The value $\tau_t = 0$ indicates that either i) all trades in period t are executed at the previous closing price, which is increasingly likely when the number of trades per period is smaller, or ii) there is no trading and the open, high, low, and close prices of period t are filled with the previous close. The value $\tau_t = 1$ is the complementary case and ensures that prices are not forward-filled.

²⁶² We now derive an estimator from close-to-open and open-to-mid (de-meaned) returns

²⁶³ by considering their serial covariance:

$$\mathbb{C}\operatorname{ov}[\overline{\eta_t - o_t}, o_t - c_{t-1}] = \mathbb{E}[(\overline{\eta_t - o_t})(o_t - c_{t-1})], \qquad (8)$$

where $\eta_t = (h_t + l_t)/2$ is the average of the high and low log-prices, o_t is the opening log-price, c_{t-1} is the closing log-price of the previous time interval, and the de-meaned returns are defined as follows:

$$\overline{r_t} = r_t - \tau_t \frac{\mathbb{E}[r_t]}{\mathbb{E}[\tau_t]} \,. \tag{9}$$

²⁶⁷ In Appendix A.1, we prove that the covariance in Equation (8) is equal to:

$$\mathbb{C}ov[\eta_t - o_t, o_t - c_{t-1} \mid \tau_t = 1]\mathbb{P}[\tau_t = 1].$$
(10)

Next, we replace observed prices with fundamental prices and bid-ask bounces as given in Equation (4). As fundamental returns are not autocorrelated (Assumption 1), and they are also uncorrelated with bid-ask bounces (Assumption 2), Equation (10) is equal to:

$$\mathbb{C}\mathrm{ov}[Z_{\eta_t} - Z_{o_t}, Z_{o_t} - Z_{c_{t-1}} \mid \tau_t = 1] \mathbb{P}[\tau_t = 1], \qquad (11)$$

where Z_{o_t} is the bid-ask bounce at the open, $Z_{c_{t-1}}$ is the bid-ask bounce at the previous close, and $Z_{\eta_t} = (Z_{h_t} + Z_{l_t})/2$. Conditional on $\tau_t = 1$, prices are not forward-filled, and thus bid-ask bounces at time t are uncorrelated with bid-ask bounces at time t - 1 by assumption. Moreover, they have zero mean (Assumption 3). Thus, Equation (11) is equal to:

$$\mathbb{E}[Z_{\eta_t} Z_{o_t} - Z_{o_t}^2 \mid \tau_t = 1] \mathbb{P}[\tau_t = 1].$$
(12)

²⁷⁶ We now need to compute the expectation in Equation (12). To this end, we recall that

 $_{277}$ $Z_{o_t}=S_{o_t}/2D_{o_t}$ and thus $Z_{o_t}^2=S_{o_t}^2/4.$ Hence, we have:

$$\mathbb{E}[Z_{o_t}^2 \mid \tau_t = 1] = \mathbb{E}[S_{o_t}^2]/4.$$
(13)

²⁷⁸ and the remaining term is calculated in Appendix A.2:

$$\mathbb{E}[Z_{\eta_t} Z_{o_t} \mid \tau_t = 1] = \frac{\mathbb{E}[S_{o_t}^2]}{4} \frac{\mathbb{P}[o_t = h_t \mid \tau_t = 1] + \mathbb{P}[o_t = l_t \mid \tau_t = 1]}{2}.$$
 (14)

Finally, we substitute Equations (13)–(14) into Equation (12) and solve for the spread. Following the calculations in Appendix A.3, we obtain that the mean squared spread is:

$$S_o^2 = \mathbb{E}[S_{o_t}^2] = \frac{-8\mathbb{E}[(\overline{\eta_t - o_t})(o_t - c_{t-1})]}{\mathbb{P}[o_t \neq h_t, \tau_t = 1] + \mathbb{P}[o_t \neq l_t, \tau_t = 1]}.$$
(15)

281 2.1 Efficient Estimation of Effective Spreads

So far, we have derived an estimator from close-to-open and open-to-mid returns. However, the same methodology can be used to derive estimators from other combinations of prices. This section identifies four estimators that achieve minimum variance under different conditions. Then, we optimally combine the four estimators to minimize the estimation variance under any condition and obtain an efficient estimator.

For illustration, we consider the case where high and low prices always differ from open or close prices, and returns have zero mean such that $\overline{r_t} = r_t$. From Equation (15) we obtain that the spread is proportional to $S_o^2 = \mathbb{E}[(\eta_t - o_t)(o_t - c_{t-1})]$ and thus the estimation variance is proportional to:

$$\mathbb{V}\mathrm{ar}[\hat{S}_{o}^{2}] = \mathbb{V}\mathrm{ar}[(\eta_{t} - o_{t})(o_{t} - c_{t-1})].$$
(16)

²⁹¹ Equation (16) shows that the estimation variance depends on the volatility of observed

returns. Thus, it depends on the volatility of the fundamental price and the size of the
bid-ask spread. We now consider two complementary cases where the spread is either
small or large compared to the volatility of the fundamental price.

In the first case, $S \to 0$ and observed prices p coincide with fundamental prices \tilde{p} . In this case, the estimation variance is proportional to:

$$\operatorname{Var}[\hat{S}_o^2] = \operatorname{Var}[(\tilde{\eta}_t - \tilde{o}_t)(\tilde{o}_t - \tilde{c}_{t-1})] = \operatorname{Var}[\tilde{\eta}_t - \tilde{o}_t] \operatorname{Var}[\tilde{o}_t - \tilde{c}_{t-1}].$$
(17)

Equation (17) shows that the estimation variance decreases with the sampling frequency 297 because the volatility of the fundamental price reduces at higher frequencies and makes 298 the estimation variance smaller. In other words, we obtain that the bid-ask spread should 299 be estimated with the highest frequency data possible and that estimators considering 300 higher time lags are dominated by estimators considering the smallest possible lag. The 301 estimator in Equation (15) is optimal because it considers subsequent close-to-open and 302 open-to-mid returns. An equivalent estimator is obtained by considering subsequent mid-303 to-close and close-to-open returns, as we expect open-to-mid returns to be distributed 304 similarly to mid-to-close returns. All other estimators have larger variance and are dom-305 inated by these two because they require higher time lags. 306

In the second case, $S \to \infty$ and observed returns are driven by bid-ask bounces. Moreover, as the spread is large, high prices are buys, and low prices are sells. In this case, $Z_{\eta_t} = (Z_{h_t} + Z_{l_t})/2 = S/4 - S/4 = 0$ and the estimation variance is proportional to:

$$\operatorname{Var}[\hat{S}_{o}^{2}] = \operatorname{Var}[(Z_{\eta_{t}} - Z_{o_{t}})(Z_{o_{t}} - Z_{c_{t-1}})] = \operatorname{Var}[Z_{o_{t}}]^{2} + \operatorname{Var}[Z_{o_{t}}]\operatorname{Var}[Z_{c_{t-1}}].$$
(18)

Equation (18) shows that the estimation variance can be reduced by using the mid price η_{t-1} ($\mathbb{V}ar[Z_{\eta_{t-1}}] = 0$) instead of the closing price c_{t-1} ($\mathbb{V}ar[Z_{c_{t-1}}] \to \infty$). In this case, the estimation variance becomes $\mathbb{V}ar[Z_{o_t}]^2$, which is strictly smaller than that in Equation (18). In other words, it is convenient to consider subsequent mid-to-open and opento-mid returns when the spread is large compared to volatility. An equivalent estimator is obtained by considering subsequent mid-to-close and close-to-mid returns, as we expect i) mid-to-close returns to be distributed similarly to open-to-mid returns and ii) close-to-mid returns to be distributed similarly to mid-to-open returns.

Table 1 summarizes the four estimators derived from the combinations of prices dis-318 cussed above. We call these estimators Discrete Generalized Estimators (DGEs) because 319 they account for the fact that prices unfold in discrete time and generalize previous ap-320 proaches that rely on continuously observed prices. For instance, the estimator by Abdi 321 and Ranaldo (2017) can be regarded as a particular case of our CHL estimator in Table 1. 322 Indeed, if we require prices to be observed continuously, they are never forward-filled and 323 the closing price always differs from the high or low prices. Therefore, $\pi_c = -4$ and 324 for zero-mean returns the CHL estimator becomes $S^2 = -4\mathbb{E}[(\overline{\eta_t - c_{t-1}})(c_{t-1} - \eta_{t-1})] =$ 325 $4\mathbb{E}[(c_{t-1} - \eta_t)(c_{t-1} - \eta_{t-1})]$, which is identical to the estimator in Equation (6). Thus, 326 our CHL estimator can be regarded as a generalization of the Abdi and Ranaldo (2017) 327 estimator that provides an analytical correction term accounting for discretely observed 328 prices. 329

330

[Insert Table 1 about here.]

Next, we combine our DGEs to minimize the estimation variance and obtain the Efficient DGE (EDGE). To this end, we notice that each DGE can be written as a moment condition so that their efficient combination is obtained by applying the Generalized Methods of Moments (GMM) (Hansen, 1982). As discussed above, the OHL estimator in Table 1 is expected to perform similarly to CHL, and OHLC is expected to perform similarly to CHLO. However, OHL and OHLC measure the spread at the open while CHL and CHLO measure the spread at the close. We thus combine OHL with CHL and ³³⁸ OHLC with CHLO to obtain two moment conditions that measure the average spread at ³³⁹ the open and close:

$$\mathbb{E}\left[2S^2 - \pi_o(\overline{\eta_t - o_t})(o_t - \eta_{t-1}) - \pi_c(\overline{\eta_t - c_{t-1}})(c_{t-1} - \eta_{t-1})\right] = 0,$$
(19)

340

$$\mathbb{E}\left[2S^2 - \pi_o(\overline{\eta_t - o_t})(o_t - c_{t-1}) - \pi_c(\overline{o_t - c_{t-1}})(c_{t-1} - \eta_{t-1})\right] = 0, \qquad (20)$$

where we have set $S^2 = (S_o^2 + S_c^2)/2$ for notational convenience. These moment conditions can be written as $\mathbb{E}[S^2 - x_{i,t}] = 0$ where x is opportunely defined. By applying GMM, the efficient estimator is given by:

$$S_{\rm GMM}^2 = \arg\min_{S^2} \sum_{ij} (S^2 - \mu_i) W_{ij} \left(S^2 - \mu_j\right), \qquad (21)$$

where $\mu_i = \mathbb{E}[x_{i,t}]$ and $W = \Omega^{-1}$ is the inverse of the covariance matrix $\Omega = \mathbb{V} ar[S^2 - x_{i,t}]$, which simplifies to $\Omega = \mathbb{V} ar[x_{i,t}]$ as the variance is translation invariant. Therefore, we have a particular case of GMM where the optimal weighting matrix does not depend on the minimizing variable, and the problem reduces to the minimization of a quadratic form. By differentiating Equation (21), setting the derivative equal to zero, and solving for S^2 , we obtain:

$$S_{\rm GMM}^2 = \sum_i w_i \mu_i \quad \text{with} \quad w_i = \frac{\sum_j W_{ij}}{\sum_{i,j} W_{ij}}.$$
 (22)

Finally, applying GMM in Equation (22) with the two moment conditions above and a diagonal covariance matrix Ω gives our *Efficient Discrete Generalized Estimator* (EDGE):

$$S_{\text{EDGE}}^2 = w_1 \mathbb{E}[x_{1,t}] + w_2 \mathbb{E}[x_{2,t}], \qquad (23)$$

352

$$x_{1,t} = \frac{\pi_o}{2} (\overline{\eta_t - o_t}) (o_t - \eta_{t-1}) + \frac{\pi_c}{2} (\overline{\eta_t - c_{t-1}}) (c_{t-1} - \eta_{t-1}),$$

$$x_{2,t} = \frac{\pi_o}{2} (\overline{\eta_t - o_t}) (o_t - c_{t-1}) + \frac{\pi_c}{2} (\overline{o_t - c_{t-1}}) (c_{t-1} - \eta_{t-1}),$$
(24)

$$w_1 = \frac{\mathbb{V}\mathrm{ar}[x_{2,t}]}{\mathbb{V}\mathrm{ar}[x_{1,t}] + \mathbb{V}\mathrm{ar}[x_{2,t}]}, \quad w_2 = \frac{\mathbb{V}\mathrm{ar}[x_{1,t}]}{\mathbb{V}\mathrm{ar}[x_{1,t}] + \mathbb{V}\mathrm{ar}[x_{2,t}]}.$$
 (25)

³⁵³ For estimation, the usual sample counterparts replace the expectations and variances,³⁵⁴ respectively.

355 2.2 Negative Estimates

Our estimators and those of Roll (1984), Corwin and Schultz (2012), and Abdi and Ranaldo (2017) are formal estimators for the mean squared spread S^2 . However, the estimate \hat{S}^2 may become negative in small samples due to statistical fluctuations. This is an issue because a negative squared spread is not mathematically nor economically meaningful.

To guarantee the non-negativity of spread estimates, it is common to reset negative values to zero by applying the transformation:

$$\hat{S} = \sqrt{\max\left\{0, \hat{S}^2\right\}}.$$
(26)

Although this approach maintains non-negativity, it can lead to a substantial number 363 of zero estimates, which can be problematic for certain applications like portfolio sort-364 ing. In an effort to mitigate this drawback, earlier studies have explored calculating the 365 squared spread across rolling time intervals, resetting negative estimates to zero within 366 these intervals, and subsequently computing the average across the entire estimation pe-367 riod. However, Jahan-Parvar and Zikes (2023) have shown that this strategy introduces 368 a strong upward bias that does not decline as the sample size increases, making the esti-369 mates inconsistent. They document that volatility is the primary driver of the bias and 370 that inconsistent measures fail to replicate some well-known results in empirical finance. 371 Following their recommendations, we apply the transformation in Equation (26) to the 372 final estimate. This ensures that the bias declines as the sample size increases and the 373

374 estimate \hat{S} is consistent.

Another way to produce consistent estimates while avoiding zero values is to take the square root of the modulus of the final estimate \hat{S}^2 . This proposition is motivated by the positive correlation between negative estimates and minus the spread that we have found empirically (see Internet Appendix I.2). However, for the sake of comparability with prior studies, we reset negative estimates to zero within this paper, leaving the exploration of alternative approaches for future research.

381 3 Simulation Results

In this section, we perform Monte Carlo simulations to study the accuracy of EDGE and its building blocks. We compare the results with the seminal Roll (1984) estimator and with the estimators proposed more recently by Corwin and Schultz (2012) and Abdi and Ranaldo (2017). Throughout the paper, we refer to these estimators with ROLL, CS, and AR, respectively. The CS estimator is adjusted for overnight returns as described in Corwin and Schultz (2012).

388 3.1 Setup

For ease of comparison, we use the simulation setup of Corwin and Schultz (2012), also 389 used in Abdi and Ranaldo (2017). Specifically, we simulate 10,000 months of data where 390 each month consists of 21 trading days and each day consists of 390 minutes. For each 391 minute of the day, the fundamental price \tilde{P}_m is simulated as $\tilde{P}_m = \tilde{P}_{m-1}e^{\sigma z}$ with $\tilde{P}_0 = 1$, 392 where σ is the standard deviation per minute and z is a random draw from a standard 393 Gaussian distribution. The daily standard deviation equals 3%, and the standard devia-394 tion per minute equals 3% divided by $\sqrt{390}$. Trade prices are defined as P_m multiplied 395 by one minus (plus) half the assumed bid-ask spread, and we use a 50% chance for bid 396

(ask) prices. Prices are assumed to be observed with a given probability. Daily high and
low prices equal the highest and lowest prices observed during the day. Open and close
prices equal the first and the last prices observed in the day. If no trade is observed for a
given day, then the previous day's closing price is used as the open, high, low, and close
prices for that day.

402 3.2 Results

We start by studying the bias of the various estimators. To this end, we simulate 10,000 months of daily prices and estimate the spread using the whole time series. These simulations use a constant spread of 1%, and the probability of observing a trade ranges from 1/390 to 1, such that the expected number of daily trades ranges from 1 to 390.

Figure 2 shows how the spread estimate varies in function of the trading frequency. 407 We find that all estimators are unbiased and correctly estimate a spread equal to 1%408 when we use 390 trades per day. However, their behavior is substantially different when 409 the trading frequency declines. Indeed, CS estimates a spread of 0.75% in the simulation 410 using 100 trades per day. Moreover, its downward bias increases rapidly as the trading 411 frequency declines further, and it returns an estimate of zero in the simulations that use 412 less than ten trades per day. AR is less sensitive to the trading frequency, but it is still 413 significantly biased in the simulations that use only a few trades per day. Instead, EDGE 414 produces unbiased estimates regardless of the number of trades, suggesting it works well 415 in practice even for assets that trade infrequently. These results demonstrate how CS 416 and AR strongly rely on the assumption that assets are traded continuously and produce 417 downward biased estimates when that assumption is not satisfied. Our more general 418 methodology provides an analytical correction term that accounts for infrequent trading 419 and produces asymptotically unbiased estimates. 420

421

[Insert Figure 2 about here.]

Next, we study the variance of the estimators by computing the standard deviation of monthly spread estimates across 10,000 simulations, where each month consists of 21 trading days. These simulations use 390 trades per day to ensure that all estimators are unbiased. In this setting, an estimator with lower variance is strictly preferable to one with higher variance because it produces unbiased estimates with higher precision.

Figure 3 reports the standard deviation of spread estimates in simulations that use a constant spread ranging from 0.50% to 8.00%. CS is preferable to AR for smaller spreads, while AR is for larger spreads. EDGE is always the best estimator, providing the most precise estimates with minimum variance uniformly across low and large spreads.

To shed light on the performance of EDGE, we also report the behavior of its building 431 blocks. In agreement with the discussion in Section 2.1, Figure 3 shows that OHL is 432 equivalent to CHL, and OHLC is equivalent to CHLO. Moreover, the variance of OHLC 433 and CHLO decreases for smaller spreads. On the contrary, the variance of OHL and CHL 434 decreases for larger spreads. EDGE exploits the opposite behaviors of these estimators 435 to produce estimates with minimum variance uniformly across low and large spreads. 436 Indeed, Equation (25) shows that EDGE puts more weight on the OHLC and CHLO 437 estimators for smaller spreads, while it puts more weight on the OHL and CHL estimators 438 for larger spreads. The result is an estimator that achieves minimum variance across small 439 and large spreads. For an additional comparison, we also report the results for the GMM 440 estimator in Equation (22) where we set the weighting matrix equal to the identity matrix. 441 This estimator has roughly the same variance of EDGE for spreads between 2.00% and 442 5.00%, but its variance is worse for smaller and larger spreads. We conclude that the 443 weighting matrix used for EDGE is effective in minimizing the estimation variance. 444

445

[Insert Figure 3 about here.]

Finally, Table 2 reports the mean and standard deviation of monthly spread estimates
from daily prices across 10,000 simulations. Panel A uses 390 trades per day to simulate

frequent trading. Here, estimators other than ROLL produce mean spreads close to the 448 actual values used in the simulation and are essentially unbiased. ROLL is affected by 449 an upward bias for small spreads that arises from truncating negative estimates and is 450 exacerbated by the large estimation variance. EDGE outperforms all other estimators in 451 these simulations by producing unbiased estimates with the lowest variance across small 452 and large spreads. Panel B introduces infrequent trading in the simulations. We find that 453 EDGE outperforms its building blocks by producing estimates with lower variance and 454 other estimators by producing estimates with lower bias. AR seems to perform similarly 455 to EDGE for simulations that use a spread of 0.50%, but this is due to the downward bias 456 for infrequent trading being counterbalanced by the upward bias induced by truncating 457 negative estimates. Although CS produces estimates with low variance, these estimates 458 are strongly biased. For instance, CS estimates a spread of 0.04% where the actual spread 459 used in the simulation is 1.00%. 460

461

[Insert Table 2 about here.]

In summary, our estimator yields unbiased estimates when other estimators are biased
 and achieves minimum variance when all estimators are unbiased.

464 4 Empirical Results

In this section, we investigate the performance of the estimators on empirical data. To evaluate the performance, we first need to define the ground truth, that is, the spread that serves as the benchmark for the evaluation. Following the literature, we use the effective spread obtained by matching high-frequency trade and quote data to evaluate the performance of the various estimators that only require commonly available daily price data.

471 4.1 Data

To compute bid-ask spread estimates (*i.e.*, EDGE, AR, CS, ROLL), we obtain daily prices 472 from the CRSP US Stock Database in the period 1926–2021 for all NYSE, AMEX, and 473 NASDAQ stocks with CRSP share codes of 10 or 11 (*i.e.*, U.S. common shares). To ensure 474 that all the estimates are obtained from transaction prices, we keep only observations for 475 which the open, high, low, and close prices are directly available. CRSP reports quotes 476 derived from bid and ask prices if transaction prices are unavailable, and a dash in front of 477 the price marks these values. We consider these non-transaction-based prices as missing 478 values. Then, we drop the days where the high, low, or close price is missing. We also 479 drop days where the open or close prices are outside the high-low range or where the low 480 price is higher than the high price. 481

We match CRSP and TAQ daily data using CUSIP identifiers and tickers. First, we 482 reconstruct the time series of CUSIPs for each KYPERMNO in CRSP. Similarly, we recon-483 struct the time series of TICKERs for each KYPERMNO in CRSP. Then, we compute the 484 time series of CUSIPs for each SYMBOL in TAQ using the Monthy TAQ Master files for 485 1993–2009 and the Daily TAQ Master files for 2010–2021. Finally, we merge the daily 486 datasets by matching observations with the same date, with the same CUSIP, and where 487 the TAQ's SYMBOL is equal to the TICKER in CRSP. Our identification strategy allows us 488 to match 99% of the stocks in CRSP. 489

For each stock-month, we estimate the spread from daily prices with EDGE, AR, CS, and ROLL and drop the estimate for all the estimators if it is missing for any of them. For instance, EDGE cannot be computed if open prices are missing, and ROLL cannot be computed if a stock-month contains only two daily observations. In such cases, we drop the corresponding estimate for all estimators. We use no explicit cutoff for the number of observations in a given stock-month. The cutoff is implicitly determined by the requirements of the most stringent estimator. Ultimately, the covariance requires at

least two returns to be computed, meaning we need at least three daily observations in a 497 stock-month. In our CRSP-TAQ merged sample, the frequency of missing estimates for 498 each of the estimators is 1.24% for EDGE, 0.48% for CHL, 1.02% for OHL, 1.17% for 499 CHLO, 1.02% for OHLC, 0.03% for AR, 0.03% for CS, and 0.14% for ROLL. Moreover, 500 when CHL is missing and CS is not, the CS estimate is zero in 90% of the cases. When 501 CHL is missing and AR is not, the AR estimate is zero in 100% of the cases. These are 502 mostly cases when the stock always trades at the same price so that the denominator of 503 our estimators is zero and the estimate is undefined. In such cases, a missing estimate 504 should be preferable to an implicit imputation of zero produced by the other estimators. 505 We rely on the TAQ database from May 1993 to December 2021 to compute the 506 benchmark effective spread. Daily spreads are obtained via the Wharton Research Data 507 Services (WRDS) Intraday Indicators using Monthly TAQ from 1993 to 2003 and Daily 508 TAQ from 2004 onward, according to the methodology described in Holden and Jacobsen 509 (2014). For each month, we winsorize the daily spreads at 99.5% (one-sided) and compute 510 the root mean squared spread for each stock. We refer to this measure as HJ. 511

To ensure that our results are robust to the choice of the benchmark, we also com-512 pute spreads using the weighted midpoint as described in Hagströmer (2021). First, we 513 replicate the daily spread measures from the WRDS Intraday Indicators using the Daily 514 TAQ database in the period 2004–2021 and we recompute our monthly HJ benchmark. 515 The benchmark achieves 99.5% correlation with the one obtained using the estimates 516 pre-computed by WRDS. Next, we replace the midpoint with the weighted midpoint to 517 generate the effective spreads described in Hagströmer (2021). The correlation between 518 the monthly benchmarks using the midpoint and weight-midpoint effective spreads is 519 99.1%. We have evaluated the estimators using both benchmarks, and all the results 520 are fully consistent. Throughout the paper, we use the midpoint benchmark as it is pre-521 computed by WRDS also for the Monthly TAQ database in the period 1993–2003, where 522

the National Best Bid and Offer (NBBO) is not directly available and matching trades with quotes poses several challenges (Holden and Jacobsen, 2014).

525 4.2 Results

⁵²⁶ Our CRSP-TAQ merged sample consists of about 1.6 million stock-month spread esti-⁵²⁷ mates for each estimator in the sample period from May 1993 to December 2021. In ⁵²⁸ Table 3, we report summary statistics and several evaluation metrics for the estimates. ⁵²⁹ EDGE achieves the highest correlation with the HJ benchmark, the lowest mean absolute ⁵³⁰ percentage error (MAPE), root mean squared error (RMSE), and the smallest fraction ⁵³¹ of zero estimates.⁴

The remainder of this section is dedicated to a deeper comparison across the estimators in a cross-sectional, time-series, and panel-data setting.

535 4.3 Cross-Sectional Correlation

Looking at cross-sectional correlations on a month-by-month basis allows us to evaluate 536 the estimators' ability to capture the cross-sectional distribution of spreads in differ-537 ent time periods. Given the effective spread benchmark $S_{i,t}$ for stock i at time t and 538 the corresponding estimate $\hat{S}_{i,t}$, we compute the cross-sectional correlation at time t as 539 $\rho_t = \mathbb{C}\mathrm{or}_i[S_{i,t}, \hat{S}_{i,t}]$. The month-by-month cross-sectional correlations for the various esti-540 mators are displayed in Figure 4. The correlation between EDGE and the effective spread 541 benchmark is consistently higher than the correlations achieved by any other estimator 542 throughout the whole period considered in the analysis. 543

544

[[]Insert Figure 4 about here.]

⁴The MAPE and RMSE are computed on log-spreads as described in Internet Appendix I.4.

545 4.4 Time-Series Correlation

Looking at time-series correlations on a stock-by-stock basis allows us to evaluate the 546 ability of the estimators to capture the time-series distribution of spreads for different 547 kinds of stocks. To this end, we split all stocks in deciles based on their market capital-548 ization. The size deciles are sorted by increasing market capitalization of each stock as 549 its last listing date in CRSP. Then, given the effective spread benchmark $S_{i,t}$ for stock i 550 at time t and the corresponding estimate $\hat{S}_{i,t}$, we compute the time-series correlation for 551 decile d as $\rho_d = \mathbb{C}\operatorname{or}_{t,i\in d}[S_{i,t}, \hat{S}_{i,t}]$. The time-series correlations for each decile obtained 552 with the various estimators are displayed in Figure 5. The correlation between EDGE 553 and the effective spread benchmark is consistently higher than the correlations achieved 554 by any other estimator for all types of stocks. 555

556

[Insert Figure 5 about here.]

557 4.5 Panel-Data Correlation

Next, we analyze the performances across five dimensions: market venues, time periods, 558 market capitalization, spread size, and trading frequency. When analyzing market venues, 559 the groups correspond to NYSE, AMEX, and NASDAQ. For the time periods, we use 560 those defined in Corwin and Schultz (2012) and Abdi and Ranaldo (2017). In addition, 561 we extend the sample and include the more recent sub-period 2016–2021. For size groups, 562 we sort stocks in quintiles based on their market capitalization at their last listing date 563 in CRSP. Spread quintiles are sorted on the average effective spread throughout the life 564 of the stock. For the trading frequency, we split stocks based on their average number of 565 daily trades during the whole sample period. Then, given the effective spread benchmark 566 $S_{i,t}$ for stock *i* at time *t* and the corresponding estimate $\hat{S}_{i,t}$, we compute the correlation 567 for group g as $\rho_g = \mathbb{C}\mathrm{or}_{(i,t)\in g}[S_{i,t}, \hat{S}_{i,t}].$ 568

The results are summarized in Table 4 for market venues (Panel A), time periods 569 (Panel B), market capitalization (Panel C), spread size (Panel D), and trading frequency 570 (Panel E). One clear result emerges: EDGE outperforms all the alternative estimators 571 in each market venue, sub-period, market capitalization, spread size, and for each trad-572 ing frequency by consistently achieving the highest correlation with the effective spread 573 benchmark. To shed light on the performance of EDGE, we also report the behavior of 574 its building blocks. In particular, it is natural to compare CHL with AR as they use the 575 same information set of high, low, and close prices. As the main difference between the 576 two estimators is that CHL accounts for infrequent trading, the outperformance of CHL 577 compared to AR demonstrates the importance of relaxing the assumption that prices are 578 observed continuously to ultimately improve empirical results. Table 4 further shows that 579 any single building block OHL, CHL, OHLC, CHLO outperforms AR, CS, and ROLL. 580 Finally, EDGE optimally combines its building blocks to provide an estimator that is 581 superior to any of them taken individually. 582

583

[Insert Table 4 about here.]

In the Internet Appendix, we provide representative illustrations on individual stocks 584 to investigate further the estimators' performance (Section I.3). We also compare the 585 estimators using additional evaluation metrics such as Spearman's (rank) correlation 586 (Table I.4), MAPE (Table I.5) and RMSE (Table I.6), and the fraction of zero estimates 587 (Table I.7). Overall, EDGE achieves the highest rank correlation with the benchmark, 588 the lowest MAPE and RMSE, and generates the lowest fraction of non-positive estimates. 589 We also find that it achieves the best results when estimating first differences instead of 590 spread levels (Table I.8) and when increasing the estimation window from one month to 591 one year (Table I.9). It is also interesting to note that CS achieves a slightly lower MAPE 592 and RMSE compared to EDGE in the following cases: a) NYSE stocks, b) recent periods, 593 c) large stocks, d) small spreads, and e) frequent trading. Taken together, these are all 594

cases where the bid-ask spread is expected to be small and where the downward bias of the CS estimator may improve the estimate. Indeed, if the spread is expected to be small, then an estimator biased towards zero may yield better results. This observation suggests that, generally, a Bayesian approach may further improve the estimate when a good prior is available for specific applications. We leave such possibility for future research as we focus on an estimator of general applicability here.

5 Applications

To demonstrate the broad applicability of EDGE, we provide three representative examples. The first revisits historical spread estimates from daily prices in the U.S. stock market since 1926. The second studies spread estimates obtained from intraday prices for U.S. stocks. Finally, the third applies the estimator outside the U.S. stock market and compares low- and high-frequency estimates for cryptocurrencies.

⁶⁰⁷ 5.1 Low-Frequency Estimates for the U.S. Stock Market

⁶⁰⁸ Using CRSP data since 1926, we construct, for each month, three portfolios based on ⁶⁰⁹ size according to the following procedure. First, we sort the stocks based on their market ⁶¹⁰ capitalization at the end of each month. Then, we select small-caps, mid-caps, and large-⁶¹¹ caps using the 50th and 80th percentiles as breakpoints. Finally, we compute monthly ⁶¹² spread estimates for individual stocks and construct the average spread for each of the ⁶¹³ three portfolios in each month between 1926–1992 (CRSP sample) and 1993–2021 (CRSP-⁶¹⁴ TAQ merged sample).⁵ The results are reported in Figure 6.

[Insert Figure 6 about here.]

⁶¹⁵

⁵When EDGE cannot be computed, we use the CHL estimator in Table 1 that does not need open prices. Open prices are missing in CRSP from July 1962 through June 1992.

Panel A displays the cross-sectional mean of spread estimates for small stocks. Ac-616 cording to EDGE, the spread was high in the 1930s, spiked in 1933 with peaks between 617 10%-15%, decreased until the 1960s and increased again with a first peak of about 5% 618 in 1963, a second peak of 7.5% in 1975, and a third peak of 10% in the early 1990s. In 619 line with the idea that liquidity evaporates in times of crisis (Nagel, 2012), these years 620 coincide with periods of financial downturn and economic recession, such as the Great 621 Depression between 1929–41, the U.S. Banking Crisis of 1933, the Kennedy Slide of 1962, 622 the 1973–1975 recession following the oil crisis, and the early 1990s recession in the United 623 States. Following the electronization of financial markets in the 2000s, the spread de-624 creased significantly until the global financial crisis, when it spiked again in 2009 with a 625 peak close to 5%. The spread has continued to reduce in the last decade, reaching the 626 lowest level ever as of December 2021. In the CRSP-TAQ merged sample after 1993, 627 the HJ benchmark closely follows this trend and overlaps with EDGE. Instead, CS and 628 AR tend to underestimate the spread, particularly for older periods, mirroring that these 629 estimators are biased when trading becomes increasingly infrequent. In the historical 630 sample before 1993, we find that the gap between EDGE and the alternative estimators 631 widens. EDGE is often larger than AR by a factor of two, and the difference is even 632 more pronounced compared to CS. Given our benchmark result from the recent sample, 633 we conjecture that the alternative estimators considerably underestimate the effective 634 spread in the early sample. 635

Panels B and C report the results for medium and large stocks, respectively. As expected, we find that larger stocks tend to have lower spreads than smaller stocks. Indeed, EDGE estimates an average spread that is typically below 2.5% for medium stocks and below 1% for large stocks. The gap with AR and CS decreases for larger stocks, mirroring that their bias reduces for stocks presumably traded more frequently. Figure 6 also reports end-of-day quoted spreads derived from CRSP. These spreads

are significantly higher than the effective spread benchmark in the sample period between 642 1993 and the early 2000s. The historical sample also supports this finding before 1993, 643 where quoted spreads are often higher than EDGE by a factor of two. We thus confirm 644 earlier studies that the quoted spread overstates the effective spread finally paid by traders 645 by up to 100% (Huang and Stoll, 1994; Petersen and Fialkowski, 1994; Bessembinder and 646 Kaufman, 1997; Bacidorea, Ross, and Sofianosa, 2003), due to dealers offering a better 647 price than the quotes, also known as trading inside the spread (Lee, 1993). We also find 648 that quoted and effective spreads closely overlap in the last two decades, suggesting that 649 this phenomenon has reduced over time and quoted spreads have become a better proxy 650 of effective spreads following the electronization of financial markets. 651

Finally, we notice that estimating spreads from daily prices leads to an upward bias 652 that becomes increasingly evident in more recent periods and for larger stocks. For 653 instance, EDGE estimates an average spread of 0.42% in December 2021 for large-caps, 654 while the HJ benchmark is 0.06%. This bias arises in small samples due to the practice of 655 resetting negative estimates to zero, which leads, on average, to overstating the spread, 656 especially when the spread is small compared to volatility (Jahan-Parvar and Zikes, 2023). 657 A way to mitigate this small-sample bias is to extend the estimation window with more 658 daily observations. For instance, the average EDGE estimate using a yearly sample of 659 daily prices is 0.22% in 2021 for large-caps, and yearly estimates generally achieve a 660 higher correlation with the yearly benchmark (see Table I.9 in the Internet Appendix) 661 compared to monthly estimates with the monthly benchmark (Table 4). Another way to 662 improve the estimation is using intraday prices whenever available, as discussed in the 663 next section. 664

⁶⁶⁵ 5.2 High-Frequency Estimates for the U.S. Stock Market

While the variance component of an asset return is proportional to the return interval, the spread component is not. Hence, we can rely on high-frequency prices to reduce the asset variance without altering the spread component and achieve a better signal-to-noise ratio to improve the spread estimate.

For instance, let N be the sample size and consider estimates derived from a monthly 670 sample of daily data (N = 21), a yearly sample of daily data (N = 252), or a monthly 671 sample of minute data ($N = 21 \times 390$). According to Equation (17), the estimation 672 variance is roughly proportional to σ_1^4/N where σ_1 is the volatility per period, and the 673 standard error is proportional to σ_1^2/\sqrt{N} . For daily prices, $\sigma_1 = \sigma/\sqrt{252}$ where σ is 674 the volatility per year. For minute prices, $\sigma_1 = \sigma/\sqrt{252 \times 390}$. Thus, estimates derived 675 from a yearly sample of daily data have a standard error that is $\sqrt{252/21} = 3.5$ times 676 smaller than that obtained from a monthly sample of daily data. Estimates derived 677 from a monthly sample of minute data have a standard error that is $390^{3/2} = 7,702$ 678 times smaller. To put this in perspective, the enhancement factor of the sample using 679 minute prices would be achieved by a sample of 1,245,699,037 daily prices, equivalent 680 to approximately 5 million years of trading. From this analysis, we conclude that using 681 intraday prices offers a more effective way to improve the spread estimate than increasing 682 the sample size with more daily data. 683

To illustrate how EDGE can substantially improve the estimation of bid-ask spreads using intraday prices, we proceed as follows. First, we aggregate trades into open, high, low, and close prices every minute using the Daily TAQ database from October 2003 to December 2021. Then, we estimate the spread with EDGE from the minute data for each stock-month. Finally, we compare the estimates derived this way with those derived from daily prices and the HJ benchmark.

⁶⁹⁰ Figure 7 reports the results for large-cap stocks. These stocks are featured by tiny

spreads that are difficult to estimate in small samples due to their small signal-to-noise ratio, which causes a large fraction of non-positive estimates and generates an upward bias due to the practice of resetting negative estimates to zero (Jahan-Parvar and Zikes, 2023). Indeed, we find that the EDGE estimates from daily prices are negative in 41% of stock-months, and they are higher than the HJ benchmark by 0.35% (35bps) on average. Instead, estimates derived from minute prices are negative in only 0.05% of stock-months, and their upward bias shrinks to zero (1bps).

698

[Insert Figure 7 about here.]

Next, we analyze the estimates for all stocks. This sample consists of 711,161 stockmonth spread estimates derived from both minute and daily prices. We find that using minute prices reduces the fraction of negative estimates from 34.15% to 0.02% and significantly improves all evaluation metrics. The Pearson's (Spearman's) correlation with the HJ benchmark increases from 56.17% (43.47%) to 88.79% (97.31%). The MAPE (RMSE) reduces from 23.68 (1.80) to 5.17 (0.41).

Finally, we estimate spreads from minute prices using the Monthly TAQ database 705 from May 1993 to July 2014. The Monthly TAQ data are identical to the Daily TAQ 706 data except for two main differences. First, Monthly TAQ only reports raw quotes, while 707 Daily TAQ includes an NBBO file that reports the highest bid price and lowest ask 708 price among all available exchanges at each timestamp. Second, Monthly TAQ data are 709 timestamped to the second while Daily TAQ data are timestamped to the millisecond. 710 While these differences cause several problems in measuring effective spreads by matching 711 trades with quotes (Holden and Jacobsen, 2014), they do not affect EDGE. Indeed, the 712 correlation between the EDGE estimates obtained with Monthly TAQ and Daily TAQ 713 in the overlapping period between October 2003 and July 2014 is 99.8%. This suggests 714 that, by relying on transaction prices only, EDGE is more robust than measuring effective 715 spreads by matching trades with quotes, and it is less sensitive to the quality of the data. 716

In summary, low-frequency estimates can be substantially improved using intraday prices. This is particularly relevant for cases where high-frequency prices are available, but quotes are not, or they cannot be reliably matched with trades. Examples include, but are not limited to, over-the-counter markets, dark pools, and crypto exchanges.

721 5.3 Estimates for Other Markets

⁷²² Our estimator represents a general way to estimate effective spreads, and it is designed ⁷²³ to be applied to a variety of markets. To illustrate its applicability outside the U.S. stock ⁷²⁴ market, we analyze estimates for cryptocurrency pairs listed in Binance.

⁷²⁵ Binance is a leading crypto exchange listing hundreds of cryptocurrencies that can ⁷²⁶ be exchanged for one another via trading pairs. Each trading pair (*e.g.*, ETH/BTC) ⁷²⁷ reports the price of the base currency (*e.g.*, ETH) in units of the quote currency (*e.g.*, ⁷²⁸ BTC). Like other crypto exchanges, Binance provides historical and real-time daily and ⁷²⁹ intraday prices for free, while trade and quote data are subject to subscription fees, and ⁷³⁰ their historical coverage is more limited. As trade and quote data are unavailable to us, ⁷³¹ we cannot compute bid-ask spreads obtained by matching trades with quotes.

To estimate effective spreads from freely available data, we download historical open, high, low, and close prices for all cryptocurrency pairs at the minute, hourly, and daily frequency. We then compute monthly estimates with EDGE, AR, and CS for each pair and each frequency and drop the estimate for all estimators if missing for any of them. Our sample consists of 2,163 crypto pairs and 53,865 pair-month spread estimates for each frequency and estimator in the sample period from July 2017 to December 2021.

We expect AR and CS to overstate the spread when using daily prices due to the upward bias induced by resetting negative estimates to zero. When using intraday prices, we expect them to understate the spread because the number of trades per period reduces at higher frequencies, and their downward bias shrinks the estimate to zero. Instead, we expect EDGE to mitigate these two concerns because its lower variance reduces the
upward bias, and the estimator is unaffected by the downward bias due to infrequent
trading.

Figure 8 reports the time evolution of the average spread across all trading pairs for 745 each estimator. As expected, AR and CS produce different estimates depending on the 746 sampling frequency. Estimates derived from daily prices are significantly higher than 747 those derived from hourly prices, which, in turn, are higher than those derived from 748 minute prices. Depending on the frequency used, the average spread in the whole sample 749 period ranges anywhere between 0.18% (0.02%) and 1.85% (1.45%) according to AR 750 (CS). This tenfold difference makes it impossible to estimate the spread reliably because 751 it is unclear which sampling frequency should be preferred in principle. Instead, EDGE 752 produces estimates less sensitive to the sampling frequency, and estimates from daily 753 prices closely overlap with those from hourly and minute prices. The average spread 754 in the whole sample period remains in the narrow range between 0.68% and 0.70%, 755 depending on whether minute, hourly, or daily prices are used. In 2021, we find that the 756 average spread for crypto pairs is between 0.35%-0.45%. 757

[Insert Figure 8 about here.]

In summary, EDGE is less sensitive to the sampling frequency than other estimators and can potentially reduce a large source of non-standard errors (Menkveld et al., 2023) in the measurement of transaction costs.

762 6 Conclusion

⁷⁶³ Historically, the development of bid-ask spread estimators has evolved along two com⁷⁶⁴ plementary paths that consider either high-frequency or low-frequency data. The former
⁷⁶⁵ exploits trade and quote data to obtain an explicit proxy of the fundamental price and

measure the distance of transaction prices from it. The latter introduces assumptions 766 about the fundamental price to derive an estimator from transaction prices only. While 767 estimates derived from trades and quotes are typically more accurate, low-frequency es-768 timates are more readily available and are becoming increasingly popular. However, low-769 frequency estimators assume that prices are observed continuously. Here, we document 770 that these approaches lead to understating effective spreads, especially for infrequently 771 traded assets that should presumably be associated with high transaction costs. We then 772 develop a novel methodology relaxing the assumption that prices are observed continu-773 ously and derive generalized estimators that correct this downward bias analytically. We 774 show that different estimators are preferable depending on whether the spread is large 775 or small compared to volatility, and we combine them efficiently to produce an unbiased 776 estimator with minimum variance. Through theoretical analyses, numerical simulations, 777 and empirical evaluations, we find that our efficient estimator dominates each generalized 778 estimator taken individually and other estimators from transaction prices. 779

Our efficient estimator has broad applicability for several reasons. First, it is de-780 rived under more general assumptions than other approaches and extends the domain 781 of applicability to various assets and time periods. Second, the estimator is unaffected 782 by the downward bias due to infrequent trading and makes it possible to estimate effec-783 tive spreads for assets traded infrequently, for historical periods, or using high-frequency 784 prices when quotes are unreliable or unavailable. Third, the estimator minimizes the 785 estimation variance and thus also minimizes the upward bias that arises from resetting 786 negative estimates to zero in small samples (Jahan-Parvar and Zikes, 2023). 787

Our results show that other estimators significantly understate effective spreads in the 20th century, while end-of-day quoted spreads overstate effective spreads by up to 100%. Thus, this work makes available the most realistic effective spread estimates for the U.S. stock market from 1926 to the advent of high-frequency data. We further show

that our estimator can substantially improve estimates from daily prices using intraday 792 prices, while other estimators are dominated by their downward bias because trading 793 becomes sparse in high frequency. To demonstrate the generalizability of these results 794 outside the U.S. stock market, we estimate bid-ask spreads for cryptocurrencies. Our 795 efficient estimator produces consistent estimates regardless of whether daily or intraday 796 prices are used, while other estimators produce a tenfold difference between daily and 797 intraday estimates. We conclude that our estimator may reduce a significant source of 798 non-standard errors in applied research (Menkveld et al., 2023). 799

Finally, we provide guidance for future research aimed at estimating transaction costs. 800 First, we have shown that the assumption that prices are observed continuously has far-801 reaching implications and causes biases that generally vary in the cross-section and time 802 series, and they also depend on the sampling frequency of open, high, low, and close 803 prices. Future work should explicitly account for discretely observed prices to avoid this 804 source of bias. Second, our estimator can be applied at any frequency, and, in this sense, 805 it reconciles the high-frequency and low-frequency literature. For this reason, we argue 806 that a better classification is distinguishing between methods that require trade and 807 quote data and those that require transaction prices only. Third, we have constructed 808 an efficient estimator in the class of covariance-based estimators from open, high, low, 809 and close prices. To design more efficient estimators, future work could either consider 810 approaches that are not based on the serial covariance of returns or exploit information 811 other than prices, such as, for instance, the trading volume or a suitable Bayesian prior. 812

References

- Abdi, F., and A. Ranaldo. 2017. A simple estimation of bid-ask spreads from daily close,
 high, and low prices. *Review of Financial Studies* 30:4437–80.
- Amihud, Y., and J. Noh. 2020. Illiquidity and stock returns ii: Cross-section and timeseries effects. *Review of Financial Studies* 34:2101–23.
- ⁸¹⁸ Bacidorea, J., K. Ross, and G. Sofianosa. 2003. Quantifying market order execution quality at the New York stock exchange. *Journal of Financial Markets* 6:281–307.
- Bali, T. G., A. Subrahmanyam, and Q. Wen. 2021. Long-term reversals in the corporate bond market. *Journal of Financial Economics* 139:656–77.
- Bessembinder, H., and H. M. Kaufman. 1997. A comparison of trade execution costs
 for NYSE and NASDAQ-listed stocks. *Journal of Financial and Quantitative Analysis*32:287–310.
- Birru, J. 2018. Day of the week and the cross-section of returns. *Journal of Financial Economics* 130:182–214.
- Bongaerts, D., F. de Jong, and J. Driessen. 2017. An asset pricing approach to liquidity
 effects in corporate bond markets. *Review of Financial Studies* 30:1229–69.
- Brogaard, J., and J. Pan. 2021. Dark pool trading and information acquisition. *Review* of Financial Studies 35:2625–66.
- Cai, F., S. Han, D. Li, and Y. Li. 2019. Institutional herding and its price impact:
 Evidence from the corporate bond market. *Journal of Financial Economics* 131:139–67.
- ⁸³⁴ Chaieb, I., V. Errunza, and H. Langlois. 2020. How is liquidity priced in global markets?
 ⁸³⁵ Review of Financial Studies 34:4216–68.
- ⁸³⁶ Chen, X., O. Linton, and Y. Yi. 2017. Semiparametric identification of the bid–ask spread
 ⁸³⁷ in extended Roll models. *Journal of Econometrics* 200:312–25.
- Chen, Y., G. W. Eaton, and B. S. Paye. 2018. Micro(structure) before macro? the
 predictive power of aggregate illiquidity for stock returns and economic activity. *Journal*of Financial Economics 130:48–73.
- ⁸⁴¹ Choi, J., S. Hoseinzade, S. S. Shin, and H. Tehranian. 2020. Corporate bond mutual
 ⁸⁴² funds and asset fire sales. *Journal of Financial Economics* 138:432–57.
- ⁸⁴³ Corwin, S. A., and P. Schultz. 2012. A simple way to estimate bid-ask spreads from daily
 ⁸⁴⁴ high and low prices. *Journal of Finance* 67:719–60.

- ⁸⁴⁵ Ding, Y., W. Xiong, and J. Zhang. 2022. Issuance overpricing of China's corporate debt
 ⁸⁴⁶ securities. Journal of Financial Economics 144:328–46.
- Easley, D., M. Lopez de Prado, M. O Hara, and Z. Zhang. 2020. Microstructure in the machine age. *Review of Financial Studies* 34:3316–63.
- Eaton, G. W., P. J. Irvine, and T. Liu. 2021. Measuring institutional trading costs and the
 implications for finance research: The case of tick size reductions. *Journal of Financial Economics* 139:832–51.
- Garman, M. B., and M. J. Klass. 1980. On the estimation of security price volatilities from historical data. *Journal of Business* 53:67–78.
- Goldstein, I., H. Jiang, and D. T. Ng. 2017. Investor flows and fragility in corporate bond
 funds. Journal of Financial Economics 126:592–613.
- Grosse-Rueschkamp, B., S. Steffen, and D. Streitz. 2019. A capital structure channel of
 monetary policy. *Journal of Financial Economics* 133:357–78.
- Hagströmer, B. 2021. Bias in the effective bid-ask spread. Journal of Financial Economics
 142:314–37.
- Hansen, L. P. 1982. Large sample properties of generalized method of moments estimators. *Econometrica* 50:1029–54.
- Hasbrouck, J. 2009. Trading costs and returns for us equities: Estimating effective costs
 from daily data. *Journal of Finance* 64:1445–77.
- Holden, C. W., and S. Jacobsen. 2014. Liquidity measurement problems in fast, competitive markets: Expensive and cheap solutions. *Journal of Finance* 69:1747–85.
- Hou, K., C. Xue, and L. Zhang. 2018. Replicating anomalies. *Review of Financial Studies*33:2019–133.
- Hua, J., L. Peng, R. A. Schwartz, and N. S. Alan. 2019. Resiliency and stock returns.
 Review of Financial Studies 33:747–82.
- Huang, R. D., and H. R. Stoll. 1994. Market microstructure and stock return predictions.
 Review of Financial Studies 7:179–213.
- Jacobs, H., and S. Müller. 2020. Anomalies across the globe: Once public, no longer existent? *Journal of Financial Economics* 135:213–30.
- Jahan-Parvar, M. R., and F. Zikes. 2023. When do low-frequency measures really measure effective spreads? Evidence from equity and foreign exchange markets. *Review of Financial Studies* 36:4190–232.

- Kaviani, M. S., L. Kryzanowski, H. Maleki, and P. Savor. 2020. Policy uncertainty and
 corporate credit spreads. *Journal of Financial Economics* 138:838–65.
- Lee, C. M. C. 1993. Market integration and price execution for NYSE-listed securities.
 Journal of Finance 48:1009–38.
- Li, X., A. Subrahmanyam, and X. Yang. 2018. Can financial innovation succeed by catering to behavioral preferences? Evidence from a callable options market. *Journal* of Financial Economics 128:38–65.
- Loon, Y. C., and Z. K. Zhong. 2016. Does Dodd-Frank affect OTC transaction costs and
 liquidity? Evidence from real-time CDS trade reports. *Journal of Financial Economics*119:645–72.
- McLean, R. D., and J. Pontiff. 2016. Does academic research destroy stock return predictability? *Journal of Finance* 71:5–32.
- Menkveld, A. J., A. Dreber, F. Holzmeister, J. Huber, M. Johannesson, M. Kirchler,
 M. Razen, U. Weitzel, D. Abad, M. M. Abudy, et al. 2023. Non-standard errors. *Journal of Finance (forthcoming)*.
- Michaelides, A., A. Milidonis, and G. P. Nishiotis. 2019. Private information in currency
 markets. *Journal of Financial Economics* 131:643–65.
- ⁸⁹⁴ Nagel, S. 2012. Evaporating liquidity. *Review of Financial Studies* 25:2005–39.
- Parkinson, M. 1980. The extreme value method for estimating the variance of the rate
 of return. Journal of Business 53:61–5.
- Patton, A. J., and B. M. Weller. 2020. What you see is not what you get: The costs of
 trading market anomalies. *Journal of Financial Economics* 137:515–49.
- Petersen, M. A., and D. Fialkowski. 1994. Posted versus effective spreads. good prices or
 bad quotes? *Journal of Financial Economics* 35:269–92.
- ⁹⁰¹ R Core Team. 2020. R: A language and environment for statistical computing. R Foun ⁹⁰² dation for Statistical Computing, Vienna, Austria.
- Ranaldo, A., and P. S. de Magistris. 2022. Liquidity in the global currency market.
 Journal of Financial Economics 146:859–83.
- Ranaldo, A., P. Schaffner, and M. Vasios. 2021. Regulatory effects on short-term interest
 rates. Journal of Financial Economics 141:750–70.
- Roll, R. 1984. A simple implicit measure of the effective bid-ask spread in an efficient
 market. Journal of Finance 39:1127–39.

- Schwert, M. 2017. Municipal bond liquidity and default risk. *Journal of Finance* 72:1683– 722.
- 911 Stoikov, S. 2018. The micro-price: A high-frequency estimator of future prices. Quanti-
- ⁹¹² *tative Finance* 18:1959–66.

Estimators										
OHL	$S_o^2 = \pi_o \mathbb{E}[(\overline{\eta_t - o_t})(o_t - \eta_{t-1})]$	$S_c^2 = \pi_c \mathbb{E}[(\overline{\eta_t - c_{t-1}})(c_{t-1} - \eta_{t-1})]$	CHL							
OHLC	$S_o^2 = \pi_o \mathbb{E}[(\overline{\eta_t - o_t})(o_t - c_{t-1})]$	$S_c^2 = \pi_c \mathbb{E}[(\overline{o_t - c_{t-1}})(c_{t-1} - \eta_{t-1})]$	CHLO							
Coefficients										
π_o	$\overline{\tau_o \mid \pi_o = -8 / \left(\mathbb{P}[o_t \neq h_t, \tau_t = 1] + \mathbb{P}[o_t \neq l_t, \tau_t = 1] \right)}$									
π_c	$\pi_c = -8 / \left(\mathbb{P}[c_{t-1} \neq h_{t-1}, \tau_t = 1] \right)$	$+\mathbb{P}[c_{t-1}\neq l_{t-1},\tau_t=1]\big)$								
	Р	rices								
0,h,l,c η	o,h,l,c Open, High, Low, Close log-prices. η Mid-prices computed as $\eta = (l+h)/2$.									

Table 1: Discrete Generalized Estimators

This table reports bid-ask spread estimators obtained from several combinations of open, high, low, and close prices as described in Section 2. The OHL and OHLC estimators measure the spread at the open. The CHL and CHLO estimators measure the spread at the close. The indicator variable τ_t is defined in Equation (7) and the de-meaned returns \bar{r}_t are defined in Equation (9).

		EDGE	OHLC	CHLO	OHL	CHL	AR	CS	ROLL
Panel A: Fr	equent 7	Frading							
S = 0.50%	Mean	0.44	0.46	0.46	0.79	0.79	0.70	0.60	1.44
	(sd)	(0.33)	(0.40)	(0.39)	(0.79)	(0.79)	(0.77)	(0.49)	(1.43)
S=1.00%	Mean	0.90	0.88	0.88	1.03	1.03	0.95	1.03	1.59
	(sd)	(0.42)	(0.55)	(0.55)	(0.86)	(0.86)	(0.85)	(0.58)	(1.49)
S = 3.00%	Mean	2.88	2.87	2.88	2.92	2.93	2.92	2.93	2.95
	(sd)	(0.41)	(0.69)	(0.69)	(0.73)	(0.72)	(0.70)	(0.61)	(1.83)
S=5.00%	Mean	4.87	4.86	4.87	4.92	4.93	4.97	4.90	4.90
	(sd)	(0.42)	(0.81)	(0.81)	(0.62)	(0.62)	(0.58)	(0.61)	(2.14)
S = 8.00%	Mean	7.84	7.78	7.79	7.88	7.89	7.99	7.86	7.93
	(sd)	(0.45)	(1.11)	(1.10)	(0.64)	(0.64)	(0.54)	(0.62)	(2.63)
Panel B: In	frequent	Trading							
S = 0.50%	Mean	0.71	0.77	0.79	0.89	0.91	0.65	0.02	1.44
	(sd)	(0.75)	(0.87)	(0.88)	(0.96)	(0.97)	(0.73)	(0.07)	(1.42)
S = 1.00%	Mean	0.95	0.99	0.99	1.11	1.10	0.81	0.04	1.56
	(sd)	(0.83)	(0.97)	(0.96)	(1.03)	(1.04)	(0.80)	(0.10)	(1.47)
S = 3.00%	Mean	2.89	2.76	2.76	2.86	2.86	2.26	0.35	2.89
	(sd)	(0.83)	(1.23)	(1.23)	(1.20)	(1.19)	(0.92)	(0.36)	(1.82)
S=5.00%	Mean	5.02	4.89	4.92	5.01	5.04	4.04	1.17	4.83
	(sd)	(0.81)	(1.32)	(1.33)	(1.13)	(1.13)	(0.85)	(0.62)	(2.12)
S = 8.00%	Mean	8.19	8.10	8.06	8.23	8.20	6.59	2.66	7.71
	(sd)	(0.96)	(1.59)	(1.62)	(1.24)	(1.26)	(0.94)	(0.96)	(2.65)

Table 2: Monthly Estimates from Simulated Daily Prices

The table reports means and standard deviations (in %) of monthly spread estimates across 10,000 simulations, where each month consists of 21 trading days and each day consists of 390 minutes. For each minute of the day, the fundamental value \tilde{P}_m is simulated as $\tilde{P}_m = \tilde{P}_{m-1}e^{\sigma z}$ with $\tilde{P}_0 = 1$, where σ is the standard deviation per minute and z is a random draw from a standard Gaussian distribution. The daily standard deviation equals 3%, and the standard deviation per minute equals 3% divided by $\sqrt{390}$. Trade prices are defined as \tilde{P}_m multiplied by one minus (plus) half the bid-ask spread S, and we assume a 50% chance that a bid (ask) is observed. Panel A reports the results where the probability of observing a trade is 100%. In Panel B, that probability equals 1%. Daily high and low prices equal the highest and lowest prices observed during the day. Open and close prices equal the first and the last prices observed in the day. If no trade is observed for a given day, then the previous day's closing price is used as the open, high, low, and close prices for that day. Negative spread estimates are set to zero.

	${ m N}$ (1)	Mean (%)	$\begin{array}{c} \operatorname{Med} \\ (\%) \end{array}$	$\begin{array}{c} \mathrm{Sd} \\ (\%) \end{array}$	$\begin{array}{c} \operatorname{Cor}_P\\(\%) \end{array}$	$\begin{array}{c} \operatorname{Cor}_{S} \\ (\%) \end{array}$	MAPE (%)	$\begin{array}{c} \text{RMSE} \\ (1) \end{array}$	$\begin{array}{c} \text{FNPE} \\ (\%) \end{array}$
EDGE	1,637,621	2.11	1.00	3.37	78.86	66.68	16.21	1.23	25.63
OHLC	$1,\!637,\!621$	2.22	1.05	3.57	69.87	57.17	18.83	1.37	29.74
CHLO	$1,\!637,\!621$	2.01	0.85	3.56	74.26	58.77	17.08	1.27	30.97
OHL	$1,\!637,\!621$	2.35	1.21	3.65	69.95	54.64	20.47	1.49	29.97
CHL	$1,\!637,\!621$	2.16	1.03	3.67	73.83	55.44	18.93	1.41	31.30
AR	$1,\!637,\!621$	1.70	0.95	2.50	68.13	53.55	19.90	1.41	31.87
CS	$1,\!637,\!621$	0.66	0.28	1.10	45.55	33.77	35.90	2.61	29.18
ROLL	$1,\!637,\!621$	2.47	1.39	4.09	55.22	41.38	24.53	1.80	32.60
HJ	1,637,621	1.89	0.75	2.73	_	—	_	_	—

Table 3: Summary Statistics

The table reports summary statistics of stock-month spread estimates from daily prices in the sample period 1993–2021 (CRSP-TAQ merged sample). Negative spread estimates are set to zero, and we drop the stock-month estimate for all the estimators if it is missing for any of them. The table reports the number of stock-months (N), the mean (Mean), median (Med), and standard deviation (Sd) of the estimates, their Pearson's (Cor_P) and Spearman's (Cor_S) correlation with the HJ benchmark, the mean absolute percentage error (MAPE) and the root mean squared error (RMSE) computed on the log-spreads (see Internet Appendix I.4), and the Fraction of Non-Positive Estimates (FNPE). The highest correlations, the lowest errors, and the lowest fraction of non-positive estimates are in bold.

	EDGE	OHLC	CHLO	OHL	CHL	AR	CS	ROLL
Panel A: Analysis by	/ Market]	Exchange	,					
NYSE	64.94	53.22	61.84	52.15	58.43	46.79	45.87	29.59
AMEX	68.99	57.77	67.84	59.26	68.23	61.05	38.32	48.15
NASDAQ	78.16	68.62	73.24	68.87	72.99	67.03	41.09	54.81
Panel B: Analysis by	7 Time Pe	riod						
1993-1996	82.93	75.94	76.83	77.68	78.57	75.31	46.94	70.23
1997 - 2000	78.47	68.24	73.40	68.95	73.41	69.20	45.06	60.19
2001 - 2002	73.04	60.47	70.32	61.44	69.59	67.31	40.70	59.00
2003 - 2007	67.65	57.56	63.45	57.26	61.97	57.34	33.87	38.10
2008-2011	69.89	62.16	64.00	61.49	62.26	59.17	33.99	43.70
2012 - 2015	60.78	51.41	55.93	52.09	55.64	53.14	37.29	29.24
2016 - 2021	53.98	46.72	43.97	46.48	43.11	40.78	39.34	22.11
Panel C: Analysis by	v Market (Capitaliza	ation					
Size quintile 1	74.35	63.60	69.90	64.86	70.44	65.00	37.16	56.08
Size quintile 2	71.29	60.36	66.68	60.39	66.33	56.08	30.20	44.31
Size quintile 3	75.13	65.09	70.32	63.12	67.28	57.22	38.41	40.11
Size quintile 4	72.55	62.93	68.07	59.63	63.44	53.02	44.60	32.90
Size quintile 5	66.65	57.77	61.32	54.24	56.17	47.31	47.24	30.31
Panel D: Analysis by	v Spread S	Size						
Spread quintile 1	17.84	15.62	16.56	15.43	15.21	14.18	12.64	9.60
Spread quintile 2	45.66	39.59	41.73	34.67	34.15	30.35	32.79	15.06
Spread quintile 3	61.98	52.28	57.80	48.88	52.46	44.72	40.82	24.64
Spread quintile 4	67.76	55.55	64.44	55.32	63.21	55.22	37.74	38.98
Spread quintile 5	71.38	60.78	66.08	62.57	67.24	61.83	33.15	55.06
Panel E: Analysis by	^r Trading	Frequenc	У					
Numtrd quintile 1	74.77	65.88	69.12	67.68	70.37	67.81	40.02	65.10
Numtrd quintile 2	79.15	69.32	74.45	69.55	74.17	69.58	51.59	52.00
Numtrd quintile 3	75.41	65.77	70.94	63.92	67.92	60.98	50.42	40.53
Numtrd quintile 4	67.17	58.53	62.78	56.00	58.71	52.54	48.98	32.41
Numtrd quintile 5	55.48	48.05	50.02	45.23	45.78	39.36	43.91	22.29

Table 4: Pearson's Correlation with the HJ Benchmark

The table reports Pearson's correlations (in %) with the HJ benchmark for stock-month spread estimates from daily prices in the sample period 1993–2021 (CRSP-TAQ merged sample). The highest correlation per group is in bold. Negative spread estimates are set to zero, and we drop the stock-month estimate for all the estimators if it is missing for any of them. The size quintiles are sorted by increasing market capitalization at the last observed period for each individual stock. The spread quintiles are sorted by increasing average HJ spreads during the whole sample period. The trade quintiles are sorted by increasing average number of daily trades during the whole sample period.



Figure 1: Probability that Open/Close Prices are High/Low Prices

The probability is computed for each stock-month as the average across: i) the fraction of days such that the open matches the high, ii) the fraction of days such that the open matches the low, iii) the fraction of days such that the close matches the high, iv) the fraction of days such that the close matches the low. Stocks are sorted into small-caps, mid-caps, and large-caps based on their market capitalization at the end of each month and using the 50th and 80th percentiles as breakpoints. The figure reports the average probability across stocks for each month and size group from 1926 to 2021. Open prices are missing from July 1962 through June 1992.





The figure reports spreads estimated from 10,000 months of simulated data where each month consists of 21 trading days and each day consists of 390 minutes. For each minute of the day, the fundamental value \tilde{P}_m is simulated as $\tilde{P}_m = \tilde{P}_{m-1}e^{\sigma z}$ with $\tilde{P}_0 = 1$, where σ is the standard deviation per minute and z is a random draw from a standard Gaussian distribution. The daily standard deviation equals 3%, and the standard deviation per minute equals 3% divided by $\sqrt{390}$. Trade prices are defined as \tilde{P}_m multiplied by one minus (plus) half the bid-ask spread, where the spread equals 1.00%, and we assume a 50% chance that a bid (ask) is observed. The probability of observing a trade ranges from 1/390 to 1, and the average number of trades per day is reported on the x-axis. Daily high and low prices equal the highest and lowest prices observed for a given day, then the previous day's closing price is used as the open, high, low, and close prices for that day. Negative spread estimates are set to zero.





The figure reports the standard deviation of monthly spread estimates across 10,000 simulations, where each month consists of 21 trading days and each day consists of 390 minutes. For each minute of the day, the fundamental value \tilde{P}_m is simulated as $\tilde{P}_m = \tilde{P}_{m-1}e^{\sigma z}$ with $\tilde{P}_0 = 1$, where σ is the standard deviation per minute and z is a random draw from a standard Gaussian distribution. The daily standard deviation equals 3%, and the standard deviation per minute equals 3% divided by $\sqrt{390}$. Trade prices are defined as \tilde{P}_m multiplied by one minus (plus) half the bid-ask spread, where the spread is reported on the x-axis, and we assume a 50% chance that a bid (ask) is observed. Daily high and low prices equal the highest and lowest prices observed during the day. Open and close prices equal the first and the last prices observed in the day. Negative spread estimates are set to zero. The OHL, CHL, OHLC, and CHLO estimators are defined in Table 1, and GMM is their GMM-combination in Equation (22) where the weighting matrix is the identity matrix.



Figure 4: Cross-Sectional Correlation with the HJ Benchmark

The figure shows cross-sectional Pearson's correlations between stock-month spread estimates from daily prices and the HJ benchmark for each month in the sample period 1993–2021 (CRSP-TAQ merged sample). Negative spread estimates are set to zero, and we drop the stock-month estimate for all the estimators if it is missing for any of them.





The figure shows time-series Pearson's correlations between stock-month spread estimates from daily prices and the HJ benchmark for size deciles in the sample period 1993–2021 (CRSP-TAQ merged sample). Size deciles are sorted by increasing market capitalization at the last observed period for each individual stock. Negative spread estimates are set to zero, and we drop the stock-month estimate for all the estimators if it is missing for any of them.





The figure reports the average spread across stocks for each month and size group from 1926 to 2021. Stocks are sorted into small-caps, mid-caps, and large-caps based on their market capitalization at the end of each month and using the 50th and 80th percentiles as breakpoints. Spreads are estimated for each stock-month using daily prices. Negative spread estimates are set to zero, and we drop the stock-month estimate for all the estimators if it is missing for any of them. EDGE is replaced with CHL when open prices are missing in CRSP. End-of-day quoted spreads (QS) are missing from July 1962 to October 1982. The HJ benchmark obtained from TAQ data is available since May 1993.



Figure 7: High-Frequency Estimates for U.S. Stocks

The figure reports the average spread across large-cap stocks for each month from October 2003 to December 2021. Spreads are estimated for each stock-month using daily (EDGE) or minute (EDGE/HF) prices. Negative spread estimates are set to zero, and we drop the stock-month estimate for all the estimators if it is missing for any of them. HJ is the benchmark spread obtained from TAQ data.



Figure 8: Low- and High-Frequency Estimates for Cryptocurrencies

The figure reports the average spread across trading pairs listed in Binance for each month from July 2017 to December 2021. Spreads are estimated for each pair-month using daily, hourly, or minute prices. Negative spread estimates are set to zero, and we drop the pair-month estimate for all the estimators if it is missing for any of them.

913 A Appendix

914 A.1 Proof of Equation 10

The de-meaned returns defined in Equation (9) have mean zero conditional on τ_t , for any return r_t computed in the time interval between the end of period t-1 and the end of period t. Indeed, r_t is identically zero conditional on $\tau_t = 0$ because $h_t = l_t = c_{t-1}$ and thus $\mathbb{E}[\bar{r}_t \mid \tau_t = 0] = 0$. Moreover, $\mathbb{E}[\bar{r}_t \mid \tau_t = 1] = \mathbb{E}[r_t \mid \tau_t = 1] - \mathbb{E}[r_t]/\mathbb{E}[\tau_t] = 0$ because $\mathbb{E}[r_t \mid \tau_t = 0] = 0$. In summary, it holds that $\mathbb{E}[\bar{r}_t \mid \tau_t] = 0$ and using the law of total covariance we have:

$$\begin{aligned} \mathbb{C}\operatorname{ov}[\overline{r_t}, r_s] &= \mathbb{E}[\mathbb{C}\operatorname{ov}[\overline{r_t}, r_s \mid \tau_t]] + \mathbb{C}\operatorname{ov}[\mathbb{E}[\overline{r_t} \mid \tau_t], \mathbb{E}[r_s \mid \tau_t]] \\ &= \mathbb{E}[\mathbb{C}\operatorname{ov}[\overline{r_t}, r_s \mid \tau_t]] \\ &= \mathbb{C}\operatorname{ov}[\overline{r_t}, r_s \mid \tau_t = 1]\mathbb{P}[\tau_t = 1] + \mathbb{C}\operatorname{ov}[\overline{r_t}, r_s \mid \tau_t = 0]\mathbb{P}[\tau_t = 0] \\ &= \mathbb{C}\operatorname{ov}[r_t, r_s \mid \tau_t = 1]\mathbb{P}[\tau_t = 1]. \end{aligned}$$
(A.1)

The last equality follows from the fact that $\overline{r_t} = r_t = 0$ conditional on $\tau_t = 0$, while $\overline{r_t} = r_t + const$. conditional on $\tau_t = 1$ and the constant is irrelevant for the calculation of the covariance.

924 A.2 Proof of Equation 14

⁹²⁵ We need to compute:

$$\mathbb{E}[Z_{\eta_t} Z_{o_t} \mid \tau_t = 1] = \frac{\mathbb{E}[Z_{h_t} Z_{o_t} \mid \tau_t = 1] + \mathbb{E}[Z_{l_t} Z_{o_t} \mid \tau_t = 1]}{2}.$$
 (A.2)

We start by considering high prices, and we condition on whether or not the opening price o_t is equal to the highest price h_t :

$$\mathbb{E}[Z_{h_t} Z_{o_t} \mid \tau_t = 1] = \mathbb{E}[Z_{h_t} Z_{o_t} \mid o_t = h_t, \tau_t = 1] \mathbb{P}[o_t = h_t \mid \tau_t = 1] + \mathbb{E}[Z_{h_t} Z_{o_t} \mid o_t \neq h_t, \tau_t = 1] \mathbb{P}[o_t \neq h_t \mid \tau_t = 1].$$
(A.3)

⁹²⁸ If $o_t = h_t$, then the opening price is selected as the highest price in the period, and the

⁹²⁹ bid-ask bounces $Z_{h_t} = Z_{o_t} = S_{o_t}/2D_{o_t}$ coincide. Thus, we have:

$$\mathbb{E}[Z_{h_t} Z_{o_t} \mid o_t = h_t, \tau_t = 1] = \mathbb{E}[S_{o_t}^2]/4.$$
(A.4)

If $o_t \neq h_t$, then Z_{h_t} and Z_{o_t} are uncorrelated by Assumption 3 and:

$$\mathbb{E}[Z_{h_t} Z_{o_t} \mid o_t \neq h_t, \tau_t = 1] = \mathbb{E}[Z_{h_t} \mid o_t \neq h_t, \tau_t = 1] \mathbb{E}[Z_{o_t} \mid o_t \neq h_t, \tau_t = 1] = 0, \quad (A.5)$$

because $\mathbb{E}[Z_{o_t} \mid o_t \neq h_t, \tau_t = 1] = \mathbb{E}[Z_{o_t}] = 0$ if we consider that the bid-ask bounce at the open is independent from whether the opening price is the highest price in the period. Substituting Equations (A.4)–(A.5) into Equation (A.3) gives:

$$\mathbb{E}[Z_{h_t} Z_{o_t} \mid \tau_t = 1] = \mathbb{E}[S_{o_t}^2] \mathbb{P}[o_t = h_t \mid \tau_t = 1]/4.$$
(A.6)

The same equation holds for low prices by replacing Z_{h_t} with Z_{l_t} and h_t with l_t . Substituting Equation (A.6) for high and low prices into Equation (A.2) yields Equation (14).

936 A.3 Proof of Equation 15

Substituting Equations (13)-(14) in Equation (12) yields:

$$\mathbb{E}[(\overline{\eta_t - o_t})(o_t - c_{t-1})] = \frac{\mathbb{E}[S_{o_t}^2]}{4} \left(\frac{\mathbb{P}[o_t = h_t \mid \tau_t = 1] + \mathbb{P}[o_t = l_t \mid \tau_t = 1]}{2} - 1 \right) \mathbb{P}[\tau_t = 1]$$

$$= \frac{\mathbb{E}[S_{o_t}^2]}{4} \left(-\frac{\mathbb{P}[o_t \neq h_t \mid \tau_t = 1] + \mathbb{P}[o_t \neq l_t \mid \tau_t = 1]}{2} \right) \mathbb{P}[\tau_t = 1]$$

$$= \frac{\mathbb{E}[S_{o_t}^2]}{4} \left(-\frac{\mathbb{P}[o_t \neq h_t, \tau_t = 1] + \mathbb{P}[o_t \neq l_t, \tau_t = 1]}{2} \right).$$
(A.7)

Solving Equation (A.7) for $\mathbb{E}[S_{o_t}^2]$ gives Equation (15).

939	Internet Appendix
940	Efficient Estimation of Bid-Ask Spreads from Open,
941	High, Low, and Close Prices
942	David Ardia, Emanuele Guidotti, Tim A. Kroencke

Reference	Journal	Area	Estimator
McLean and Pontiff (2016)	JF	AP, stocks	CS
Loon and Zhong (2016)	JFE	MM, OTC derivatives	Roll
Bongaerts, de Jong, and Driessen (2017)	RFS	AP, corporate bonds	Roll
Schwert (2017)	JF	AP, muni bonds	Roll
Goldstein, Jiang, and $Ng (2017)$	JFE	AP, bond funds	Roll
Hou, Xue, and Zhang (2018)	RFS	AP, stocks	CS
Chen, Eaton, and Paye (2018)	JFE	AP, time-series predictability	Roll, CS
Li, Subrahmanyam, and Yang (2018)	JFE	AP, investor behavior	CS
Birru (2018)	JFE	AP, stocks	CS
Grosse-Rueschkamp, Steffen, and Streitz (2019)	JFE	CF, capital structure	CS
Michaelides, Milidonis, and Nishiotis (2019)	JFE	MM, currency markets	CS
Cai et al. (2019)	JFE	AP, corporate bonds	Roll
Hua et al. (2019)	RFS	AP, stocks	Roll, CS
Easley et al. (2020)	JFE	MM, information	CS
Jacobs and Müller (2020)	JFE	AP, stocks	CS
Patton and Weller (2020)	JFE	AP, stocks	CS
Kaviani et al. (2020)	JFE	AP, corporate bonds	CS
Amihud and Noh (2020)	m RFS	AP, stocks	AR
Choi et al. (2020)	JFE	AP, bond funds	Roll
Chaieb, Errunza, and Langlois (2020)	RFS	AP, stocks	CS, AR
Bali, Subrahmanyam, and Wen (2021)	JFE	AP, corporate bonds	Roll
Eaton, Irvine, and Liu (2021)	JFE	MM, trading costs	Roll, CS
Ranaldo, Schaffner, and Vasios (2021)	JFE	MM, interest rates	Roll
Brogaard and Pan (2021)	m RFS	MM, information	CS
Ranaldo and de Magistris (2022)	JFE	MM, currencies	Roll, CS
Ding, Xiong, and Zhang (2022)	JFE	AP, corporate bonds	CS

Table I.1: Bid-Ask Spread Estimators Used in Top-Three Finance Publications: 2016–2022

Total number of papers (count multiple usages):

26(31)

To be listed, a paper must have been published in the Journal of Financial Economics, Journal of Finance, or Review of Financial Studies Corwin and Schultz (2012), AR: Abdi and Ranaldo (2017). Articles are classified in the areas of interest: market microstructure (MM), in 2016–2022 and must reference and apply one of the bid-ask spread estimators as indicated in the last column; Roll: Roll (1984), CS: asset pricing (AP), and corporate finance (CF).

⁹⁴³ I.1 Literature Review

944 I.2 Signed Estimates

⁹⁴⁵ We define the signed spread estimate as follows:

$$\hat{S} = \operatorname{sign}(\hat{S}^2) \times \sqrt{|\hat{S}^2|}.$$
(I.1)

In Table I.2, we report the number of negative spread estimates and their correlation
with minus the HJ benchmark.

	EDGE	OHLC	CHLO	OHL	CHL	AR	CS	ROLL
Ν	419,617	486,690	507,074	490,702	512,542	521,928	465,952	$533,\!543$
Cor_P	54.23%	57.43%	59.41%	53.84%	55.56%	36.01%	21.47%	38.57%
Cor_S	45.95%	48.29%	49.73%	49.24%	50.61%	41.02%	21.78%	40.41%

 Table I.2: Summary Statistics of Negative Estimates

The table reports summary statistics of negative stock-month spread estimates from daily prices in the sample period 1993–2021 (CRSP-TAQ merged sample). We drop the stock-month estimate for all the estimators if it is missing for any of them. The table reports the number of negative estimates (N), and their Pearson's (Cor_P) and Spearman's (Cor_S) correlation with minus the HJ benchmark.

948 I.3 Individual Stocks

This section reports two illustrative cases for individual stocks. Specifically, Table I.3 949 reports signed bid-ask spread estimates for the monthly samples of open, high, low, and 950 close daily prices displayed in Figure I.1. Panel A is the case of a stock traded infrequently 951 in May 2001, where the high and low prices often match the open and close prices. In this 952 situation, EDGE, OHLC, CHLO, OHL, and CHL are close to the benchmark spread of 953 2.53%. AR gives a downward biased estimate of 1.14%, and the CS estimate is essentially 954 zero. ROLL estimates a negative spread, likely due to its large estimation variance. 955 Panel B gives an example in December 2021 where the stock trades frequently, the high 956 and low prices differ from the open and close prices, but the spread is small (0.56%). In 957 this case, some building blocks of EDGE give positive estimates (OHLC, OHL, CHL) and 958 others give a negative estimate (CHLO). By weighting them optimally, EDGE estimates 959 a spread of 0.58%, while other methods produce upward biased estimates of 1.87% (AR), 960

Panel	HJ	EDGE	OHLC	CHLO	OHL	CHL	AR	CS	ROLL
A	2.53	2.50	2.52	2.47	2.53	2.48	1.14	-0.00	-0.58
В	0.56	0.58	1.13	-0.99	2.46	1.93	1.87	2.35	2.87

Table I.3: Spread Estimates for Individual Stocks

The table reports monthly spread estimates from daily prices for the examples in Figure I.1. The estimates are signed as defined by Equation (I.1).



Figure I.1: Daily Open, High, Low, and Close Prices for Individual Stocks Examples of daily open, high, low, and close prices. The highest and lowest extremes of each bar are the high and low prices. The left and right segments are the open and close prices, respectively. Bars are green if the close is higher than the open, and red otherwise. Stocks are identified with their CRSP PERMNO. Panel A reports stock 75272 in May 2001. Panel B reports stock 21543 in December 2021.

⁹⁶² I.4 Additional Evaluation Metrics

The empirical distribution of benchmark spreads is highly skewed, as displayed in Figure I.2. For this reason, we evaluate the MAPE and RMSE on the log-spreads, which are more symmetrically distributed. As the logarithm is defined for positive values, we ignore non-positive estimates in these calculations.

MAPE =
$$\frac{1}{N} \sum_{i=1}^{N} \left| \frac{\log(S_i) - \log(\hat{S}_i)}{\log(S_i)} \right|$$
, RMSE = $\sqrt{\frac{1}{N} \sum_{i=1}^{N} (\log(S_i) - \log(\hat{S}_i))^2}$. (I.2)



Figure I.2: Distribution of the HJ Benchmark The histograms show the empirical distribution of the HJ benchmark. The left panel reports the distribution of the spreads. The right panel reports the distribution of the natural logarithm of the spreads.

We report several evaluation metrics for monthly and yearly spread estimates from daily prices. The following tables report metrics for monthly estimates: Table I.4 reports Spearman's correlations; Tables I.5 and I.6 report mean absolute percentage errors and root mean squared errors; Table I.7 reports the fractions of non-positive estimates. Table I.8 reports Pearson's correlations for first-differences of monthly estimates. Table I.9 reports Pearson's correlations for yearly estimates.

	EDGE	OHLC	CHLO	OHL	CHL	AR	CS	ROLL	
Panel A: Analysis by	v Market 1	Exchange	,						
NYSE	22.98	17.85	21.99	17.47	19.20	17.37	2.69	10.62	
AMEX	49.52	37.02	47.01	38.11	46.21	43.39	17.72	37.99	
NASDAQ	72.29	60.68	64.21	58.34	60.89	57.64	30.80	45.92	
Panel B: Analysis by Time Period									
1993–1996	83.91	74.62	75.72	75.67	76.90	74.39	53.19	66.76	
1997 - 2000	72.18	59.47	64.86	57.64	61.96	58.87	40.30	48.53	
2001 - 2002	63.53	47.01	60.06	45.36	55.81	53.32	28.03	42.23	
2003 - 2007	56.71	49.83	42.68	47.34	41.38	41.20	23.52	29.37	
2008-2011	58.82	49.68	43.98	51.79	46.13	47.51	31.53	33.63	
2012 - 2015	54.14	46.60	40.58	47.87	43.71	44.96	40.68	31.10	
2016 - 2021	49.35	43.41	35.91	42.95	36.48	37.96	42.80	25.06	
Panel C: Analysis by	v Market	Capitaliza	ation						
Size quintile 1	71.50	56.30	65.57	57.39	65.39	60.85	28.17	53.76	
Size quintile 2	64.83	51.84	57.97	50.52	55.71	50.15	19.71	42.05	
Size quintile 3	60.82	52.18	52.64	46.92	46.03	42.84	23.87	29.00	
Size quintile 4	43.85	38.26	38.56	33.48	32.10	30.76	21.68	18.82	
Size quintile 5	26.88	24.15	23.85	22.18	20.83	20.17	14.53	14.00	
Panel D: Analysis by	v Spread S	Size							
Spread quintile 1	13.68	11.53	12.37	12.46	12.00	11.40	7.22	8.83	
Spread quintile 2	31.60	27.75	26.98	23.00	20.76	19.56	13.95	10.09	
Spread quintile 3	50.83	43.85	44.23	37.35	36.23	33.68	23.49	19.15	
Spread quintile 4	61.52	49.05	55.05	47.35	52.10	47.51	25.32	35.88	
Spread quintile 5	69.49	53.41	62.53	55.90	64.05	58.29	21.81	56.12	
Panel E: Analysis by	Trading	Frequenc	У						
Numtrd quintile 1	72.68	58.35	64.90	60.83	66.75	64.32	33.88	63.10	
Numtrd quintile 2	72.25	61.54	63.17	59.19	60.14	59.00	43.17	45.44	
Numtrd quintile 3	61.22	54.60	51.83	48.60	44.82	44.27	36.75	29.87	
Numtrd quintile 4	45.85	42.07	37.16	37.95	33.01	32.89	32.74	22.75	
Numtrd quintile 5	27.49	25.90	22.83	23.64	20.30	19.80	19.78	15.75	

Table I.4: Spearman's Correlation with the HJ Benchmark

The table reports Spearman's correlations (in %) with the HJ benchmark for stock-month spread estimates from daily prices in the sample period 1993–2021 (CRSP-TAQ merged sample). The highest correlation per group is in bold. Negative spread estimates are set to zero, and we drop the stock-month estimate for all the estimators if it is missing for any of them. The size quintiles are sorted by increasing market capitalization at the last observed period for each stock. The spread quintiles are sorted by increasing average HJ spreads during the whole sample period. The trade quintiles are sorted by increasing average number of daily trades during the whole sample period.

	EDGE	OHLC	CHLO	OHL	CHL	AR	CS	ROLL		
Panel A: Analysis by	v Market	Exchang	e							
NYSE	21.12	23.96	22.06	26.77	25.55	25.19	20.72	32.92		
AMEX	13.93	15.80	13.64	16.59	14.37	17.00	50.66	17.56		
NASDAQ	14.77	17.41	15.62	18.64	16.98	18.32	39.91	21.94		
Panel B: Analysis by Time Period										
1993–1996	9.29	11.14	10.55	11.39	10.71	13.83	58.27	14.15		
1997 - 2000	11.32	13.76	12.77	14.64	13.77	15.05	48.19	18.24		
2001 - 2002	14.49	16.99	15.46	18.27	17.05	17.79	43.41	22.29		
2003 - 2007	19.09	22.13	19.15	24.10	21.50	21.81	25.22	27.15		
2008-2011	23.23	25.93	24.36	28.21	26.90	27.11	28.91	33.10		
2012 - 2015	20.38	23.35	21.29	25.42	23.40	23.53	22.32	30.04		
2016-2021	21.65	24.88	22.83	27.31	25.64	25.52	20.72	33.26		
Panel C: Analysis by	v Market	Capitaliz	zation							
Size quintile 1	13.37	15.82	14.38	16.42	14.92	16.96	61.84	18.54		
Size quintile 2	12.93	15.31	13.65	16.09	14.44	16.88	43.95	18.27		
Size quintile 3	14.10	16.63	14.74	18.44	16.83	17.48	28.17	22.26		
Size quintile 4	18.76	21.63	19.51	24.39	22.77	22.60	21.19	30.01		
Size quintile 5	24.43	27.35	25.41	29.84	28.48	28.12	20.21	36.14		
Panel D: Analysis by	v Spread	Size								
Spread quintile 1	26.09	29.24	27.09	31.93	30.47	30.13	19.15	38.17		
Spread quintile 2	19.56	22.42	20.23	25.44	23.84	23.44	18.79	31.33		
Spread quintile 3	15.09	17.86	15.38	19.93	17.81	17.86	21.93	24.05		
Spread quintile 4	11.96	14.45	12.54	15.16	13.23	14.54	34.55	17.20		
Spread quintile 5	12.28	14.26	13.65	14.33	13.58	17.24	78.46	16.06		
Panel E: Analysis by	⁷ Trading	Frequen	cy							
Numtrd quintile 1	11.48	12.88	12.66	12.78	12.46	16.98	82.35	14.37		
Numtrd quintile 2	11.25	13.72	12.10	14.63	12.98	13.43	34.14	17.13		
Numtrd quintile 3	15.21	18.13	15.59	20.48	18.34	18.30	22.29	24.70		
Numtrd quintile 4	20.50	23.66	21.17	26.53	24.61	24.31	19.45	32.25		
Numtrd quintile 5	26.51	29.65	27.63	32.17	30.73	30.38	20.57	38.57		

Table I.5: Mean Absolute Percentage Error with the HJ Benchmark

The table reports Mean Absolute Percentage Errors (in %) with the HJ benchmark as defined in Equation (I.2), for stock-month spread estimates from daily prices in the sample period 1993–2021 (CRSP-TAQ merged sample). The lowest error per group is in bold. We drop the stock-month estimate for all the estimators if it is missing for any of them. The size quintiles are sorted by increasing market capitalization at the last observed period for each stock. The spread quintiles are sorted by increasing average HJ spreads during the whole sample period. The trade quintiles are sorted by increasing average number of daily trades during the whole sample period.

	EDGE	OHLC	CHLO	OHL	CHL	AR	CS	ROLL	
Panel A: Analysis by N	Market Ex	change							
NYSE	1.76	1.95	1.82	2.11	2.02	2.00	1.60	2.52	
AMEX	0.79	0.88	0.75	0.93	0.79	0.88	3.17	0.98	
NASDAQ	1.03	1.16	1.05	1.26	1.17	1.18	2.83	1.49	
Panel B: Analysis by Time Period									
1993-1996	0.52	0.55	0.52	0.56	0.53	0.64	3.91	0.70	
1997 - 2000	0.63	0.73	0.68	0.79	0.75	0.79	3.19	1.01	
2001 - 2002	0.87	1.00	0.91	1.10	1.02	1.04	3.04	1.35	
2003 - 2007	1.39	1.57	1.43	1.70	1.59	1.58	1.74	1.99	
2008 - 2011	1.69	1.85	1.78	1.98	1.94	1.91	1.88	2.37	
2012 - 2015	1.61	1.80	1.67	1.92	1.81	1.79	1.55	2.28	
2016-2021	1.70	1.90	1.80	2.05	1.97	1.95	1.54	2.48	
Panel C: Analysis by M	Market Ca	pitalizati	on						
Size quintile 1	0.72	0.81	0.72	0.85	0.76	0.82	4.03	0.96	
Size quintile 2	0.78	0.89	0.80	0.95	0.87	0.93	2.80	1.10	
Size quintile 3	0.98	1.12	1.01	1.24	1.16	1.16	1.92	1.51	
Size quintile 4	1.45	1.63	1.49	1.79	1.70	1.68	1.56	2.17	
Size quintile 5	2.04	2.23	2.10	2.39	2.30	2.28	1.64	2.83	
Panel D: Analysis by S	Spread Size	e							
Spread quintile 1	2.11	2.32	2.18	2.48	2.38	2.36	1.59	2.91	
Spread quintile 2	1.51	1.68	1.55	1.85	1.76	1.74	1.39	2.23	
Spread quintile 3	1.04	1.19	1.06	1.33	1.23	1.22	1.43	1.60	
Spread quintile 4	0.71	0.82	0.72	0.88	0.78	0.82	2.06	1.01	
Spread quintile 5	0.58	0.62	0.58	0.62	0.57	0.71	4.63	0.68	
Panel E: Analysis by T	Trading Fr	equency							
Numtrd quintile 1	0.58	0.58	0.57	0.57	0.55	0.73	4.72	0.62	
Numtrd quintile 2	0.62	0.73	0.65	0.79	0.72	0.73	2.14	0.94	
Numtrd quintile 3	1.01	1.16	1.02	1.31	1.21	1.20	1.56	1.59	
Numtrd quintile 4	1.51	1.69	1.55	1.86	1.77	1.75	1.42	2.25	
Numtrd quintile 5	2.14	2.34	2.22	2.50	2.41	2.39	1.66	2.94	

Table I.6: Root Mean Squared Error with the HJ Benchmark

The table reports Root Mean Squared Errors with the HJ benchmark as defined in Equation (I.2), for stock-month spread estimates from daily prices in the sample period 1993–2021 (CRSP-TAQ merged sample). The lowest error per group is in bold. We drop the stock-month estimate for all the estimators if it is missing for any of them. The size quintiles are sorted by increasing market capitalization at the last observed period for each stock. The spread quintiles are sorted by increasing average HJ spreads during the whole sample period. The trade quintiles are sorted by increasing average number of daily trades during the whole sample period.

	EDGE	OHLC	CHLO	OHL	CHL	AR	CS	ROLL	
Panel A: Analysis by Market Exchange									
NYSE	40.90	44.50	42.34	42.75	42.19	44.10	42.01	40.06	
AMEX	25.84	30.89	32.43	30.69	31.83	31.29	40.30	32.75	
NASDAQ	18.34	22.58	25.40	23.81	26.06	26.13	21.76	29.04	
Panel B: Analysis by Time Period									
1993–1996	15.50	19.02	18.76	21.34	21.99	22.26	26.12	26.31	
1997 - 2000	22.47	27.31	25.51	29.82	29.51	29.99	34.76	31.85	
2001 - 2002	24.55	31.40	27.03	32.85	30.49	30.92	35.42	32.26	
2003 - 2007	27.07	30.58	35.40	29.86	33.47	33.67	28.59	35.13	
2008-2011	26.38	32.84	32.91	30.42	31.30	31.81	26.79	33.10	
2012 - 2015	31.47	33.89	39.30	32.89	36.36	37.41	25.16	35.89	
2016 - 2021	34.70	37.15	40.84	35.54	37.98	39.18	27.33	35.02	
Panel C: Analysis by Market Capitalization									
Size quintile 1	15.64	21.09	21.48	22.21	22.74	22.21	24.72	25.95	
Size quintile 2	17.42	22.34	24.38	23.00	24.72	24.33	25.52	27.61	
Size quintile 3	23.13	26.99	29.70	27.48	29.97	30.37	26.56	32.96	
Size quintile 4	33.18	36.35	37.55	35.94	37.31	38.75	31.95	37.07	
Size quintile 5	38.77	41.93	41.76	41.23	41.77	43.71	37.16	39.42	
Panel D: Analysis by Spread Size									
Spread quintile 1	41.85	44.73	44.18	43.29	43.50	45.65	38.06	40.27	
Spread quintile 2	34.63	37.99	39.02	37.57	38.92	40.46	33.63	38.69	
Spread quintile 3	24.94	28.44	31.89	29.56	32.62	33.24	27.54	35.55	
Spread quintile 4	15.97	20.17	23.32	21.88	24.74	24.29	23.92	29.14	
Spread quintile 5	10.75	17.33	16.44	17.54	16.72	15.70	22.74	19.34	
Panel E: Analysis by Trading Frequency									
Numtrd quintile 1	13.74	21.14	17.64	20.67	17.98	17.42	29.40	19.92	
Numtrd quintile 2	16.24	20.48	22.94	22.51	25.13	24.82	25.12	30.33	
Numtrd quintile 3	24.40	28.08	32.00	28.29	31.76	32.07	25.37	34.58	
Numtrd quintile 4	33.61	36.45	39.20	36.08	38.36	39.64	29.58	37.96	
Numtrd quintile 5	40.15	42.55	43.10	42.30	43.28	45.42	36.43	40.23	

Table I.7: Fraction of Non-Positive Estimates

The table reports fractions of non-positive stock-month spread estimates (in %) from daily prices in the sample period 1993–2021 (CRSP-TAQ merged sample). The lowest fraction per group is in bold. We drop the stock-month estimate for all the estimators if it is missing for any of them. The size quintiles are sorted by increasing market capitalization at the last observed period for each stock. The spread quintiles are sorted by increasing average HJ spreads during the whole sample period. The trade quintiles are sorted by increasing average number of daily trades during the whole sample period.

	EDGE	OHLC	CHLO	OHL	CHL	AR	CS	ROLL	
Panel A: Analysis by Market Exchange									
NYSE	15.44	11.68	13.67	11.40	12.48	11.74	13.10	7.78	
AMEX	21.99	16.33	19.13	17.75	20.05	20.41	19.40	9.49	
NASDAQ	30.82	23.73	25.19	24.43	25.81	25.04	18.75	16.39	
Panel B: Analysis by Time Period									
1993-1996	32.83	25.15	25.94	26.68	27.59	27.78	19.49	22.22	
1997 - 2000	32.29	24.63	26.64	25.05	27.06	27.28	21.16	20.41	
2001 - 2002	31.51	23.49	28.62	23.98	27.88	30.43	24.34	22.16	
2003 - 2007	23.70	18.42	18.38	18.77	18.65	17.04	13.26	2.99	
2008 - 2011	25.30	21.08	19.11	21.49	19.24	19.80	13.12	11.77	
2012 - 2015	13.48	11.73	9.31	12.61	10.35	10.43	10.54	2.04	
2016 - 2021	13.79	11.85	8.95	13.08	10.16	9.38	11.51	2.95	
Panel C: Analysis by Market Capitalization									
Size quintile 1	31.34	23.96	26.04	24.86	26.72	27.54	21.17	18.77	
Size quintile 2	27.51	21.34	22.01	22.18	22.80	21.10	16.19	13.33	
Size quintile 3	27.16	20.97	21.99	21.28	21.84	19.53	15.35	11.20	
Size quintile 4	22.96	17.77	19.76	17.15	18.99	16.24	14.30	7.98	
Size quintile 5	16.42	13.44	13.86	12.76	13.05	11.04	9.64	8.39	
Panel D: Analysis by Spread Size									
Spread quintile 1	4.22	4.87	3.47	4.77	3.55	4.24	3.67	3.94	
Spread quintile 2	12.58	10.76	10.96	10.36	10.01	10.11	9.76	6.32	
Spread quintile 3	21.42	17.00	17.80	16.71	17.32	15.88	16.18	7.65	
Spread quintile 4	28.97	22.01	24.09	22.33	23.97	22.98	19.40	12.13	
Spread quintile 5	31.43	24.24	25.58	25.43	26.60	27.77	20.20	19.39	
Panel E: Analysis by Trading Frequency									
Numtrd quintile 1	29.94	23.17	24.67	24.33	25.64	27.07	16.56	22.59	
Numtrd quintile 2	34.18	25.77	27.45	26.66	28.09	28.09	24.89	14.46	
Numtrd quintile 3	27.56	21.46	22.43	21.30	21.76	20.36	20.74	8.91	
Numtrd quintile 4	20.15	15.97	16.65	15.79	16.29	14.59	15.93	6.74	
Numtrd quintile 5	12.39	10.96	10.77	10.22	10.07	8.93	9.93	4.83	

Table I.8: Pearson's Correlation of First-Differences with the HJ Benchmark

The table reports Pearson's correlations (in %) with the HJ benchmark for first-differences of stock-month spread estimates from daily prices in the sample period 1993–2021 (CRSP-TAQ merged sample). The highest correlation per group is in bold. Negative spread estimates are set to zero, and we drop the stock-month estimate for all the estimators if it is missing for any of them. The size quintiles are sorted by increasing market capitalization at the last observed period for each stock. The spread quintiles are sorted by increasing average HJ spreads during the whole sample period. The trade quintiles are sorted by increasing average number of daily trades during the whole sample period.

	EDGE	OHLC	CHLO	OHL	CHL	AR	CS	ROLL	
Panel A: Analysis by Market Exchange									
NYSE	74.17	66.98	70.93	66.76	69.06	56.35	58.89	35.72	
AMEX	76.85	67.56	75.57	69.50	76.40	64.59	45.33	42.23	
NASDAQ	86.54	79.71	83.60	80.20	83.09	78.43	49.89	62.94	
Panel B: Analysis by Time Period									
1993-1996	88.27	83.89	85.02	85.30	86.09	83.02	52.46	80.14	
1997 - 2000	87.52	79.94	84.25	81.09	84.06	81.77	54.39	74.52	
2001 - 2002	85.73	73.21	84.72	75.27	84.14	82.02	52.91	74.98	
2003 - 2007	76.21	67.34	71.30	68.64	70.32	64.25	47.23	34.06	
2008-2011	80.29	76.57	74.68	76.05	73.03	68.32	48.62	43.58	
2012 - 2015	70.78	63.56	63.91	64.94	64.32	63.62	51.92	32.02	
2016 - 2021	63.73	58.34	54.63	58.54	53.62	50.81	52.09	27.82	
Panel C: Analysis by Market Capitalization									
Size quintile 1	82.80	73.76	79.97	75.38	80.19	74.87	42.42	62.40	
Size quintile 2	80.64	72.48	77.33	73.38	77.50	64.95	36.99	49.91	
Size quintile 3	82.49	76.15	79.01	74.79	76.45	64.40	46.87	43.24	
Size quintile 4	82.55	75.85	81.54	72.88	76.99	65.84	54.12	31.89	
Size quintile 5	78.31	71.00	76.70	68.34	71.59	61.95	58.96	39.53	
Panel D: Analysis by Spread Size									
Spread quintile 1	23.91	22.14	21.58	19.70	17.65	16.01	17.99	6.95	
Spread quintile 2	52.04	50.50	46.69	43.11	37.40	32.46	39.33	13.30	
Spread quintile 3	67.97	62.29	62.80	58.62	57.75	46.96	49.02	21.11	
Spread quintile 4	77.37	66.28	76.75	66.68	76.04	66.42	45.87	42.26	
Spread quintile 5	80.43	71.18	76.82	73.31	77.50	71.45	36.99	62.02	
Panel E: Analysis by Trading Frequency									
Numtrd quintile 1	83.59	76.27	79.77	78.11	80.38	78.24	47.76	77.30	
Numtrd quintile 2	88.24	81.52	85.93	82.04	85.95	81.45	63.30	55.62	
Numtrd quintile 3	83.86	77.56	81.10	75.93	78.43	69.37	61.20	40.14	
Numtrd quintile 4	74.16	70.80	67.81	68.32	64.25	57.04	62.84	33.50	
Numtrd quintile 5	68.26	62.99	62.02	60.06	56.85	48.64	55.44	26.63	

Table I.9: Pearson's Correlation of Yearly Estimates with the HJ Benchmark

The table reports Pearson's correlations (in %) with the yearly HJ benchmark for stock-year spread estimates from daily prices in the sample period 1993–2021 (CRSP-TAQ merged sample). The yearly HJ benchmark is the root mean squared monthly HJ benchmark within the year. The highest correlation per group is in bold. Negative spread estimates are set to zero, and we drop the stock-year estimate for all the estimators if it is missing for any of them. The size quintiles are sorted by increasing market capitalization at the last observed period for each stock. The spread quintiles are sorted by increasing average HJ spreads during the whole sample period. The trade quintiles are sorted by increasing average number of daily trades during the whole sample period.