

Emergence of Randomness in Temporally Aggregated Financial Tick Sequences

Silvia Onofri^{1*}, Andrey Shternshis^{2†}, Stefano Marmi^{1‡}

1. Scuola Normale Superiore, Piazza dei Cavalieri 7, Pisa, Italy, 56126
2. Uppsala University, Regementsvägen 10, Uppsala, Sweden, 75105

Abstract

Markets efficiency implies that the stock returns are intrinsically unpredictable, a property that makes markets comparable to random number generators. We present a novel methodology to investigate ultra-high frequency financial data and to evaluate the extent to which tick by tick returns resemble random sequences. We extend the analysis of ultra high-frequency stock market data by applying comprehensive sets of randomness tests, beyond the usual reliance on serial correlation or entropy measures. Our purpose is to extensively analyze the randomness of these data using statistical tests from standard batteries that evaluate different aspects of randomness.

We illustrate the effect of time aggregation in transforming highly correlated high-frequency trade data to random streams. More specifically, we use many of the tests in the NIST Statistical Test Suite and in the TestU01 battery (in particular the Rabbit and Alphabit sub-batteries), to prove that the degree of randomness of financial tick data increases together with the increase of the aggregation level in transaction time. Additionally, the comprehensive nature of our tests also uncovers novel patterns, such as non-monotonic behaviors in predictability for certain assets. This study demonstrates a model-free approach for both assessing randomness in financial time series and generating pseudo-random sequences from them, with potential relevance in several applications.

1 Introduction

A central concept in finance, the Efficient Market Hypothesis (EMH) [1], posits that asset prices fully and fairly reflect all available information. In its weak form, the EMH implies that future price movements are unpredictable and follow a random walk or martingale

*Corresponding author: silvia.onofri@sns.it

†andrey.shternshis@it.uu.se

‡stefano.marmi@sns.it

process: conditional on past information, the expected value of future price changes is zero. In this sense, an efficient market behaves like a random number generator (RNG). For instance, [2] demonstrated that random numbers can be generated from Bitcoin price data, while [3] proposed stock prices as a source of randomness for cryptographic applications.

However, empirical evidence shows that the degree of randomness in financial markets is not constant. Market efficiency and thus the randomness of price changes can vary across time periods [4, 5] and across different markets classes [6]. Moreover, randomness depends on the temporal resolution of the data. At intraday frequencies, asset prices exhibit well-documented stylized facts [7], i.e., empirical properties of prices such as higher volatility at market openings and closings. Prices sampled at one-minute intervals already deviate from perfect randomness because of such properties [8], while ultra-high-frequency (tick-by-tick) data are even less random, revealing structured patterns that reflect correlated trading activity [9, 10].

We aim to systematically investigate the randomness of financial time series across different levels of temporal aggregation. While tick-level data contain deterministic microstructure patterns, aggregated prices increasingly resemble random sequences over larger time steps. We assess the randomness of financial data across different levels of temporal aggregation using a diverse set of randomness tests. Our approach treats financial price series as potential outputs of a random number generator and applies batteries of RNG tests to measure their degree of randomness. To the best of our knowledge, this is the first study to apply such test batteries to tick-by-tick financial data.

In [3, 11], authors proposed a methodology of developing randomness beacons from financial data for applications relying on public randomness. Financial time series and sequences generated by RNG were compared in [12]. In [13], authors classified and ranked pseudo-RNGs by their quality of outputs. Prices of stocks were tested as outputs of RNG in [14]. The set of tests on efficiency were applied in [15]. Several authors have used various quantitative measures derived from information theory [16, 17] to assess the randomness or complexity of financial time series [18, 19, 20]. Comparing stock market data to the output of well-tested PRNGs can highlight areas where market behavior deviates from expected randomness. Deviations from ideal randomness can thus reveal underlying market mechanisms, including algorithmic trading behaviors that follow deterministic rules [21, 22].

Our study builds on previous work [23], in which we examined the predictability of ultra-high-frequency financial time series using entropy-based randomness tests. That analysis showed that as tick-by-tick data are aggregated over larger numbers of transactions, their predictability tends to decrease. In the present study, we extend this approach in two main directions. First, instead of relying solely on entropy-based measures, we apply a comprehensive set of randomness tests from several methodological categories, such as frequency, pattern, and spectral analyses. This allows for a broader characterization of the stochastic properties of financial data. Second, we make use of the full available dataset to obtain empirical distributions of the test outcomes, providing a new way to represent how randomness evolves across different temporal scales.

This study offers three main contributions. First, we systematize the methodology for assessing the randomness of financial time series by integrating standard randomness test suites with entropy-based measures. In particular, we employ the NIST Statistical Test Suite [24] and TestU01 [25, 26]. As an initial methodological step, we introduce a procedure for selecting subsets of tests that are appropriate for the specific lengths of financial data sequences, ensuring valid results. Second, we provide a broader perspective on the randomness of binary strings derived from tick-level data. In particular, we have shown that both the rate and monotonicity of convergence toward randomness, as the aggregation level increases, vary between stocks and time periods. Third, we demonstrate that as the data are aggregated, the binarized price sequences undergo a whitening effect, progressively becoming indistinguishable from random sequences. Our methodology enables the generation of ℓ sub-sequences at each aggregation level ℓ . That is, at an appropriate scale depending on the specific financial asset, multiple sequences pass standard randomness tests. Consequently, we propose a model-free approach for generating pseudo-random sequences from financial data, which can be leveraged for further cryptographic and other applications.

The paper is structured as follows. We revise methods used in previous literature and battery tests used to quantify the amount of randomness contained in financial data in Section 2. Here we also explain our approach of converting prices to binary strings, together with the description of dataset. We simulate pseudorandom number sequences and apply randomness batteries in Section 3 to test the methods' suitability. In Section 4, we conduct experiments and classify the results on randomness into several categories. Finally, in Section 5 we discuss our findings and future work.

2 Encoding of prices and randomness tests

We begin by describing our methodology for symbolizing price sequences, which converts financial price data into binary strings.

We focus on executed order prices and compare adjacent transactions to obtain binary strings $\bar{b} \in \{0, 1\}^*$. In more detail, given a sequence of prices $\{s_1, s_2, \dots, s_N\}$ for a certain day and ticker, we construct our first binary sequence \bar{b} considering, for each $i = 2, \dots, N$ the ratio $r = \frac{s_{i-1}}{s_i}$. Then, if $r < 1$, we add a 0 to \bar{b} ; if $r > 1$, we add a 1 to \bar{b} ; if $r = 1$, we do not add any value to the sequence \bar{b} . Then, we construct other sequences by aggregating data by an aggregation level $\ell = 1, \dots, 100$. In order to use the full information even at high aggregation levels, we consider ℓ samplings for each aggregation level ℓ to build binary strings \bar{b}_j , $j = 1, \dots, \ell$ in the way explained before. So, we consider the ratio r and update the string \bar{b}_j in such a way:

$$r = \frac{s_{j+i\ell}}{s_{j+(i-1)\ell}} \begin{cases} < 1 & \text{then } \bar{b}_j \rightarrow \bar{b}_j 0 \\ = 1 & \text{then } \bar{b}_j \rightarrow \bar{b}_j \\ > 1 & \text{then } \bar{b}_j \rightarrow \bar{b}_j 1 \end{cases}$$

Then, from each initial sequence of prices $\{s_1, s_2, \dots, s_N\}$ we get $\sum_{\ell=1}^{100} \ell = 5050$ different binary sequences for each analyzed day. We emphasize that this methodology does not require any future information on the data, since, at the time of transaction i , we only use one past transaction to build the ratios. This means that it can indeed be interpreted as an online approach.

We use randomness tests on these binary sequences to investigate the amount of predictability in prices. Batteries of randomness tests are mainly thought for the testing of cryptographic strings, which are usually very long. Since the strings obtained from each analyzed day are not long enough, we run the tests on strings that contain the concatenation of data for an whole month. That is, for each aggregation level $\ell = 1, \dots, 100$ and for each sample $j = 1, \dots, \ell$, we run the test on a concatenated string obtained as $\bar{b}^m = \bar{b}^{d_1} || \dots || \bar{b}^{d_n}$, where the $\{\bar{b}^{d_i}\}_{i=1,\dots,n}$ are the strings obtained for each analyzed day of the month m .

2.1 Randomness tests

There are different perspectives from which randomness could be determined. The variety of tests to identify random sequences investigates, for example, frequencies of strings, patterns and entropy.

In fact, a randomness test is a statistical test that, given a binary string as input, outputs the decision between accepting or rejecting the hypothesis $H_0=\text{the string is random}$ (or the hypothesis $H_a=\text{the string is non-random}$). To test a certain property, each test uses a reference distribution and tests whether the one produced by the tested string is similar to the one produced by a random string.

For each battery, we fix a level of significance α , that represents the threshold for probability of a false negative error, i.e., $\alpha = \mathbb{P}(\text{accept } H_a | H_0 \text{ is true})$. The final decision between accepting or rejecting the null hypothesis is made by comparing α with a p -value obtained from the computed statistic. The p -value represents the probability that a perfect RNG would have produced a sequence less random than the one tested. So, once we have computed the p -values, we have two possible outcomes. If the test is one-tailed, then:

- if p -value $\geq \alpha$, then accept H_0 .
- if p -value $< \alpha$, then reject H_0 .

If the test is two-tailed, then the choice is made by:

- if $\frac{\alpha}{2} \leq p\text{-value} \leq 1 - \frac{\alpha}{2}$, then accept H_0 .
- if $p\text{-value} < \frac{\alpha}{2}$ or $p\text{-value} > 1 - \frac{\alpha}{2}$, then reject H_0 .

In the following, the value of α is always assumed to be fixed to 0.01.

Thus, the common output of all tests are p -values, allowing us to conclude if the hypothesis on randomness can be rejected. Section 2.1.1 demonstrates how entropy-based statistics are constructed. Other tests are taken from two batteries: NIST Statistical Test Suite and TestU01. In particular, we use NIST STS version 2.1.2 and TestU01 version 1.2.3 (we focus on the Alphabit and Rabbit sub-batteries). Through these batteries, we distinguish five main different points of view on randomness analysis:

- Frequency tests: they evaluate in several ways whether the proportion of zeros and ones is coherent with the one of a uniform distribution.
- Pattern tests: they detect specific local structures, repeated patterns, and correlations between bits.
- Entropy and complexity tests: they assess how difficult a sequence is to compress or predict, using measures such as entropy estimation or linear complexity.
- Spectral tests: they apply discrete Fourier transforms to detect periodic structures or unexpected frequency spikes that would not be present in truly random data.
- Random Walks tests: they analyze the cumulative behavior of sequences interpreted as random walks, checking for imbalances, excursions from the origin, and path regularities.

A classification of the tests we use is given in Table 1. A description of these tests and how we use them is given in the following subsections.

2.1.1 Entropy-based tests

We recall two entropy-based tests introduced in [23]. Let $X = \{x_1, x_2, \dots, x_N\}$ be a realization of a stationary random process with symbols from a binary alphabet: $x_i \in \{0, 1\}$. If the sequence is fully random, then all strings of length $k < N$ have equal probabilities to be met in the sequence X . Then, the hypothesis of full randomness can be tested by evaluating empirical probabilities of appearing all possible blocks of k symbols. Using the probabilities of a finite set of strings, the authors in [23] employ Shannon entropy [16], a measure of uncertainty, to assess the randomness of the sequence. Below, we revise the step-by-step algorithm for estimating randomness via the Shannon entropy. In the following, we refer to this test as ShannonEntropy test.

Category	NIST STS	Rabbit	Alphabit
Frequency tests	Frequency (Mono-bit), Frequency Test within a Block	MultinomialBitsOverlapping, HammingWeight	MultinomialBitsOverlapping (x4)
Pattern tests	Runs Test, Longest Run of Ones in a Block, Non-overlapping Template Matching, Overlapping Template Matching, Serial Test	ClosePairsBitMatch(x2), LongestHeadRun, PeriodsInStrings, HammingCorrelation(x3), HammingIndependence(x3), AutoCorrelation(x2), Runs Test	HammingCorrelation, HammingIndependence (x2)
Entropy and complexity tests	Binary Matrix Rank Test, Maurer's Universal Statistical Test, Linear Complexity Test, Approximate Entropy Test	AppearanceSpacings, LinearComp, LempelZiv, MatrixRank(x3)	—
Spectral tests	Discrete Fourier Transform	Fourier1, Fourier3	—
Random walks tests	Cumulative Sums Test, Random Excursions Test, Random Excursions Variant Test	RandomWalk1(x3), RandomExcursions, RandomExcursionsVariant	RandomWalk1(x2)

Table 1: Taxonomy of tests from NIST STS, Alphabit and Rabbit batteries.

- Choose the block length, k , as suggested in [27]. We choose $k = \lceil 0.5 \log_2 N \rceil$: this is a trade-off between ensuring that there are enough blocks to construct a test statistic and ensuring that the blocks are long enough to capture potential dependencies.
- Divide the sequence into $N_b = \lfloor \frac{N}{k} \rfloor$ non-overlapping blocks, where $\lfloor \cdot \rfloor$ denotes the floor function (rounding down):

$$\hat{x}_t = \{x_{(t-1)k+1}, x_{(t-1)k+2}, \dots, x_{tk}\}, \quad t \in [1, N_b].$$

- Calculate empirical frequencies \hat{f}_j of all blocks of length k using the indicator function I :

$$\hat{f}_j = \sum_{t=1}^{n_b} I(\hat{x}_t = a_j), \quad a_j \in A^k, j \in [1, 2^k]$$

- Estimate the Shannon entropy, which is defined as the averaged measure of uncertainty about a symbol appearing in a sequence:

$$\hat{H} = - \sum_j \frac{\hat{f}_j}{N_b} \ln \frac{\hat{f}_j}{N_b},$$

where $\ln()$ is the natural logarithm with the convention $0 \ln 0 = 0$.

- Test if the estimation is close to the possible maximum of entropy. More precisely, we test whether the difference between the maximum possible entropy, $k \ln 2$, and the estimate follows a χ^2 -distribution with $2^k - 1$ degrees of freedom [28]:

$$Y_1 = 2N_b(k \ln 2 - \hat{H})$$

$$H_0 : Y_1 \sim \chi^2(2^k - 1)$$

In many applications, such as financial ones, the requirement for equiprobable symbols (0 and 1) can be relaxed, and the hypothesis of randomness is defined only in terms of independence. The second entropy-based test that we recall from [23] aims to test the hypothesis $H_0 = \text{The occurrence of a new symbol in the sequence } X \text{ is independent of the sequence's preceding symbols}$, even if they may have different probability of appearance. Moreover, the method considers all overlapping blocks, thereby enriching the dataset for computing the test statistic. The test statistic in this case is the relative entropy between the empirical probabilities and those expected under H_0 . We keep the same value for k and follow the procedure below:

- Define $N_o = N - k + 1$ overlapping blocks

$$\bar{x}_t = \{x_t, x_{t+1}, \dots, x_{t+k-2}\}, \quad t \in [1, N_o]$$

- Calculate empirical frequencies of blocks of length k

$$f_{ij} = \sum_{t=1}^{N_o} I(\bar{x}_t = a_i) I(x_{t+k-1} = a_j), \quad a_i \in A^{k-1}, a_j \in A$$

- Evaluate the test statistic and assess whether it follows a χ^2 -distribution:

$$Y_2 = 2 \sum_{ij} f_{ij} \ln \frac{N_o f_{ij}}{f_{\cdot j} f_{i \cdot}}, \quad f_{\cdot j} = \sum_i f_{ij}, f_{i \cdot} = \sum_j f_{ij}$$

$$H_0 : Y_2 \sim \chi^2(2^{k-1} - 1)$$

The proof for the asymptotic distribution can be found in [23]. In the following, we refer to this test as KL test.

2.1.2 NIST Statistical Test Suite

The NIST Statistical Test Suite (STS) [24] is a set of tests designed to evaluate the randomness of binary sequences and to ensure their suitability for cryptographic applications. Developed by the National Institute of Standards and Technology (NIST), it was first published in 2000 and then revised in 2010. In 2022 NIST announced [29] that they are working on a new revision.

As explained at the beginning of Section 2, we run NIST STS tests¹ on binary strings obtained from months of financial data. Nonetheless, some tests require a length for the string to be tested that is impossible to reach by using financial data with this methodology. The battery also requires the user to select parameters to run the tests. Each test is run on a certain number of subsequences, so we have to decide how to split our string. Some tests also require other parameters (such as block length). In the documentation [24], NIST provides details on how to select such parameters. First, we choose 9 of the 15 tests to run on our sequences, since the length of the strings is not sufficient for the others. Then, we focus on the following tests: Frequency, Block Frequency, Cumulative Sums, Runs, Longest Run of Ones, Discrete Fourier Transform, Non-overlapping Template Matching, Approximate Entropy and Serial Test. We divide the tests into two groups: we run Discrete Fourier Transform and Non-Overlapping Template Matching on substrings of 1000 bits,

¹We apply slight modifications to the file `src/assess.c`: we change the type of the `expCount` variable from an integer to a float.

while the others are run on substrings of 128 bits. Other parameters are chosen according to the suggestions of the documentation and are summarized in Table 4 in Appendix A.

We run a sanity check explained in Section 3 on the tests, that leads to a more fine-grained selection of the tests. Then, we run the tests on each monthly string obtained from financial data and collect the p -values obtained from the output files.

2.1.3 TestU01:

TestU01 [25, 26] offers a wide variety of tests, organized in sub-batteries. We focus in particular on Alphabit and Rabbit sub-batteries.

In particular, Alphabit is composed of 9 different tests. It contains four Multinomial Bits Overlapping tests, which detect correlations between successive bits in blocks of several lengths; two types of Hamming tests (Independence and Correlation tests) that detect correlations between the successive bits of overlapping blocks, and two Random Walks tests.

Rabbit is composed of a wide variety of tests: they are 26 and analyze randomness from very different perspectives. We find again the Multinomial Bits Overlapping test, the Hamming Correlation and Independence tests and the Random Walk tests from the Alphabit battery, while other tests include Close Pairs Bit Match, Appearance Spacings, Linear Compression, Lempel Ziv, Fourier transforms, Longest Head Run, Periods in strings, Hamming Weight, Autocorrelation, Run, Matrix Rank. Some of them correspond to NIST tests if run with certain choices of parameters.

We choose to run the tests with standard parameters. A detail of the tests and the parameters can be found in Table 5 in App. A. Binary strings should be at least 500 bits long to be tested correctly.

We run all the tests on the strings obtained for an whole month, as explained at the beginning of Section 2. In some cases tests fail due to insufficient length of the strings at high aggregation levels. Since the standard battery output just prints the p -values of the failed tests, we slightly modify it to collect all the p -values of all the tests run².

2.2 Dataset

We now introduce the description of the financial dataset used for our randomness tests. We analyze data from limit order books obtained from LOBSTER (www.lobsterdata.com). In a limit order book, traders submit buy and sell orders that either execute immediately by matching with an existing opposite order, or remain active until they are either filled or canceled.

Our study covers 80 trading days between 01-08-2022 and 21-11-2022. Each trading day spans from 9:30 to 16:00, that is 390 minutes in total. Table 2 lists the tickers we selected,

²In particular, we modify the prints of the functions *WritepVal* and *WriteReport* in *testu01/bbattery.c* to print all the p -values.

Table 2: Assets and characteristics of prices

Asset	Ticker	Mean price	Standard deviation of price	Daily trading volume	Daily number of transactions	Average time between transactions
Apple Inc.	AAPL	153.47	0.93	12,184,032	136,136	0.165
Microsoft Corporation	MSFT	251.78	1.37	4,529,093	84,342	0.269
Tesla Inc.	TSLA	388.02	3.81	8,686,354	178,704	0.127
Intel Corporation	INTC	30.15	0.20	7,055,642	38,255	0.595
Eli Lilly and Company	LLY	327.33	1.73	370,050	11,404	2.086
Snap Inc.	SNAP	10.67	0.14	4,967,779	18,521	1.358
Ford Motor Company	F	13.93	0.10	4,468,175	12,954	1.815
Carnival Corporation & plc	CCL	9.24	0.12	5,874,376	15,372	1.518
SPDR S&P 500 ETF	SPY	390.52	1.56	9,136,137	95,181	0.246

Mean price, its standard deviation, trading volume, number of transactions, and average time between transactions are calculated for each day and then are averaged over 80 days. Trading volume is summed up for each day. Average time is given in seconds.

along with their key attributes. The sample includes stocks from a variety of sectors, differing in average price, volatility, transaction counts, inter-trade durations (mostly under one second), and daily volumes (0.3 to 12 million shares). We also include the SPY ETF, which tracks the S&P 500 Index. Transaction timestamps are recorded with nanosecond precision.

For each asset, we downloaded the corresponding message file from LOBSTER. These files contain six fields: order ID, time, price, size (volume), event type, and side (buy or sell). Orders that conceal their size are classified as hidden orders [30]. Our analysis focuses specifically on two types of events: executions of visible limit orders and executions of hidden limit orders.

3 Sanity check with Random Number Generators

Since some tests depend on the choice of parameters, or in some cases the unexpected short length of the strings could lead to inaccurate outputs, we run a procedure to select which

tests we can consider as valid on real data and which not. To do this, we use three RNGs using three different sources of randomness, and we call this procedure *sanity check*.

We run the sanity check on all the tests that we use, in order to be sure that they work properly with string lengths that we get from financial data. In particular, the three generators are: *Quantis QRNG USB* [31] by ID Quantique, that is a quantum RNG; *Linux /dev/urandom*, an operating system-level pseudorandom number generator (PRNG) seeded from environmental noise; and random strings produced by Möbius function. The first is a quantum RNG whose randomness is enhanced by quantum mechanical principles. The second uses environmental noise entropy and is one of the most commonly used generators. The third is believed to be random due to the connections between the Riemann hypothesis and the Möbius function. More details about these generators can be found in App. B.

We generate strings of $N = 50\,000, 100\,000, 500\,000, 1\,000\,000$ bits and handle them similarly as the financial series. This means that, for every generated string $\{s_1, s_2, \dots, s_N\}$, for each level of aggregation $\ell = 1, \dots, 100$ and for each sample $j = 1, \dots, \ell$, we consider the ratios $r = \frac{s_{j+i\ell}}{s_{j+(i-1)\ell}}$ and build the binary string, as before, by adding a 0 or 1 to \bar{b}_j if, respectively, $r < 1$ or $r > 1$. So, again, from each sequence generated from our RNG, we build 5050 different binary strings. We apply all the selected tests from NIST STS, all the test from Alphabit and Rabbit batteries and the two entropy-based tests explained in Section 2.1.1 on the obtained sequences and use the same parameters of the financial series, as explained in Table 4 in App. A.

We fix a threshold of 2% for excluding tests from the analysis. For any test, if the amount of strings that fail the test on all the three RNGs is higher than the threshold, we do not run the test on financial data of the corresponding length, since this is a clue that the test is not working properly with our choice of parameters. Tests are run with the same level of significance of financial data, i.e. $\alpha = 0.01$.

The results of the sanity check applied to NIST STS and TestU01 batteries are shown in Table 3. For ShannonEntropy and KL tests, every string-length passes every test.

The length of the strings generated from financial data can be approximated as follows: $N = 50\,000$ for the tickers CCL, F and SNAP; $N = 100\,000$ for the tickers INTC and LLY; $N = 500\,000$ for the tickers AAPL, MSFT; $N = 1\,000\,000$ for the tickers SPY and TSLA. We then consider tests on these tickers to be valid only if they passed the sanity check on the corresponding string length.

4 Randomness tests applied to tick data

After the selection made by the sanity check, we run randomness tests from NIST STS, Alphabit, Rabbit and the two entropy-based tests on the strings obtained from UHF data. We show some of the results in the following, while the total collection can be found in

NIST STS						RABBIT BATTERY					
N.	Test	50K	100K	500K	1M	N.	Test	50K	100K	500K	1M
1	Frequency	✓				1	MultinomialBitsOverlapping				
2	BlockFrequency	✓	✓			2	ClosePairsBitMatch, $t = 2$	✓	✓	✓	✓
3	CumulativeSums	✓	✓			3	ClosePairsBitMatch, $t = 4$		✓		
4	Runs	✓	✓	✓		4	AppearanceSpacings			✓	✓
5	LongestRun	✓	✓	✓		5	LinearComp	✓	✓	✓	✓
6	Approx. Entropy	✓	✓	✓	✓	6	LempelZiv	✓	✓	✓	✓
7	Serial	✓	✓	✓	✓	7	Fourier1	✓	✓	✓	✓
8	FFT	✓	✓			8	Fourier3	✓	✓	✓	✓
9	NonOverlappingTemplate	✓	✓			9	LongestHeadRun	✓	✓	✓	✓
ALPHABIT BATTERY						10	PeriodsInStrings	✓	✓	✓	✓
1	MultinomialBitsOverlapping, $L = 2$	✓	✓	✓	✓	11	HammingWeight, $L = 32$	✓	✓	✓	✓
2	MultinomialBitsOverlapping, $L = 4$	✓	✓	✓	✓	12	HammingCorrelation, $L = 32$	✓	✓	✓	✓
3	MultinomialBitsOverlapping, $L = 8$	✓	✓	✓	✓	13	HammingCorrelation, $L = 64$	✓	✓	✓	✓
4	MultinomialBitsOverlapping, $L = 16$	✓	✓	✓	✓	14	HammingCorrelation, $L = 128$	✓	✓	✓	✓
5	HammingIndependence, $L = 16$	✓	✓	✓	✓	15	HammingIndependence, $L = 16$	✓	✓	✓	✓
6	HammingIndependence, $L = 32$	✓	✓	✓	✓	16	HammingIndependence, $L = 32$	✓	✓	✓	✓
7	HammingCorrelation, $L = 32$	✓	✓	✓	✓	17	HammingIndependence, $L = 64$	✓	✓	✓	✓
8	RandomWalk1, $L = 64$	✓	✓			18	AutoCorrelation, $d = 1$	✓	✓	✓	✓
9	RandomWalk1, $L = 320$	✓	✓	✓		19	AutoCorrelation, $d = 2$	✓	✓	✓	✓
						20	Run		✓	✓	✓
						21	MatrixRank, 32×32	✓	✓	✓	✓
						22	MatrixRank, 320×320	✓	✓	✓	✓
						23	MatrixRank, 1024×1024	✓	✓	✓	✓
						24	RandomWalk1, $L = 128$	✓			
						25	RandomWalk1, $L = 1024$	✓	✓	✓	✓
						26	RandomWalk1, $L = 10016$	✓	✓	✓	✓

Table 3: Results of sanity check applied to NIST STS, Alphabit, and Rabbit tests on selected string sizes.

our GitHub repository³. Results of the tests are shown as boxplots: we plot aggregation levels on the x-axis, and for every aggregation level ℓ we show the boxplot of $-\log_{10}(p_j)$, $j = 1, \dots, \ell$, where p_j is the p -value of the sample j . The threshold of 2 then represents $\alpha = 0.01$: if $-\log_{10}(p_j)$ is under the threshold, this means that the test has been passed, then the corresponding string can be considered as random; otherwise, the test has rejected the randomness hypothesis.

Our main results can be summarized as follows:

1. Our study aligns with the findings in [23], reaffirming that, generally, randomness tends to increase as the aggregation level grows up.
2. However, novel methods and tests sometimes reveal exceptions; for instance, there are cases where even at high aggregation levels the result of tests reveals predictability.
3. Some of the tests display a novel pattern in predictability-aggregation plots: an example is Fourier3 test, from Rabbit battery. In this case, the maximum of predictability

³<https://github.com/Silvia895/Emergence-of-Randomness-in-Temporally-Aggregated-Financial-Tick-Sequences>

for some stocks occurs at aggregation level higher than 1; after that peak, randomness starts to increase.

4. In another particular case, we find an unexpected behavior. This is the result of the HammingCorrelation tests (both in Alphabit and Rabbit batteries) applied on strings generated from INTC data: in the month of August, we see predictability increasing as the level of aggregation increases. A deeper analysis is run on this example to understand the nature of such behavior.

Every case in this list is analyzed in the following paragraphs.

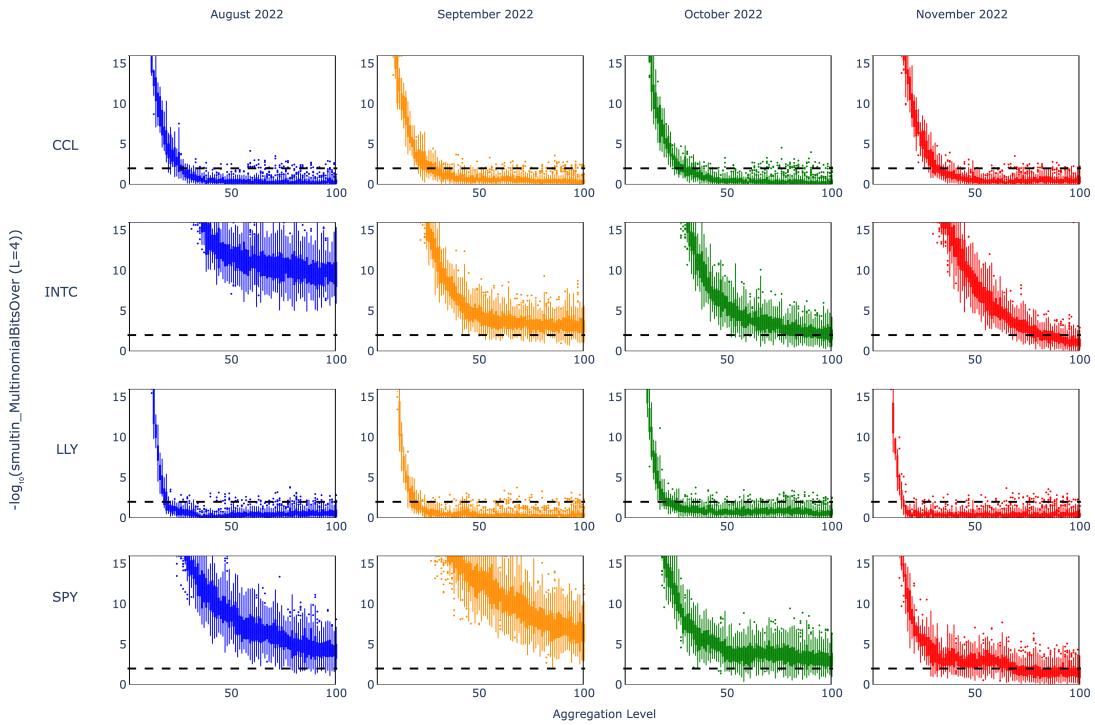


Figure 1: Results of smultin MultinomialBitsOver ($L = 4$) test from Alphabit applied to CCL, INTC, LLY and SPY data.

Cases 1 and 2: In the following, we provide examples of tests that analyse randomness from different points of view: Figure 1 shows the results of smultin_MultinomialBitsOver

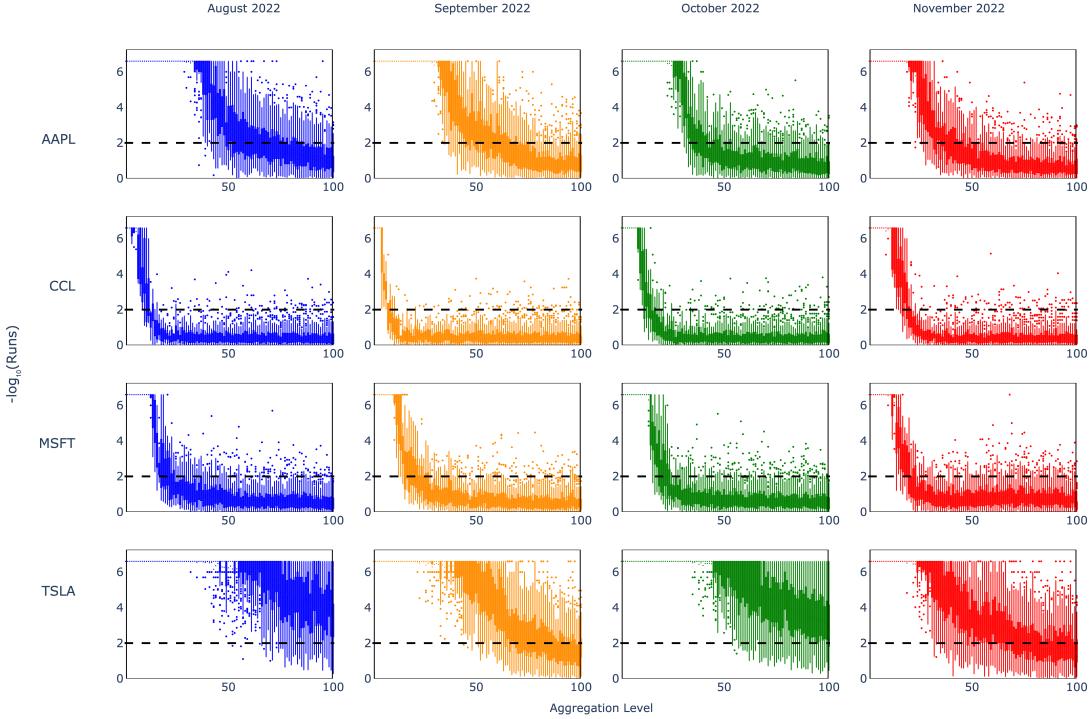


Figure 2: Results of Runs test from NIST STS applied to AAPL, CCL, MSFT and TSLA data.

($L = 4$), a frequency test; Figure 2 represents the results of the pattern test Runs; Figures 3 and 4 give the results of two entropy tests (ApproximateEntropy and KL); Figure 5 shows the results of the spectral test Fourier3, which will be also analysed for Case 3 in the following; finally, Figure 6 gives the results of CumulativeSums test, that is a random walk test. Every figure shows the test applied to four stocks. The general behavior is consistent with expectations: aggregation induces randomness, so the predictability of the strings decreases while increasing the level of the aggregation.

However, we observe in some cases the phenomenon of slow increasing (or not increasing at all) of randomness at high aggregation levels for certain stocks in particular. For example, many tests for AAPL and TSLA highlight predictability also at high aggregation levels (see e.g. Figures 2, 3 and 4).

We attribute this persistent predictability to their high trading activity: AAPL and TSLA have an average of respectively 6 and 8 trades per second (whereas CCL and LLY

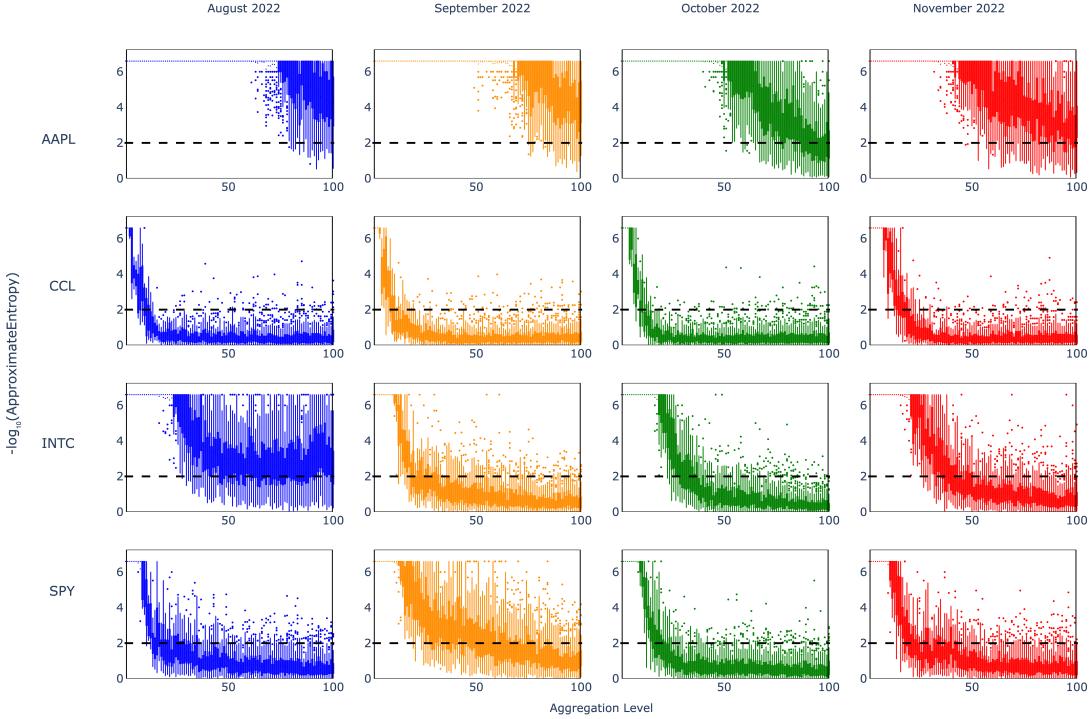


Figure 3: Results of ApproximateEntropy test from NIST STS applied to AAPL, CCL, INTC and SPY data.

on average have 0.6 and 0.5 trades per second). While our results demonstrate that the general trend is towards increasing randomness (decreasing predictability) with aggregation for almost all the stocks in our datasets, sometimes it happens that aggregating until level 100 is not enough to see randomness emerge from the strings.

Also, our methodology shows the usefulness of such a multifaceted approach: different tests analyse different properties of the strings. This means that even if one string appears random with respect to one property, it may not appear random with respect to another. Examples can be seen in Figures 1 and 3. Consider the month of August 2022 for SPY: although the ApproximateEntropy test suggests that randomness is achieved at level 20, aggregation up to level 100 is insufficient to guarantee randomness with respect to the MultinomialBitsOverlapping test with $L = 4$.

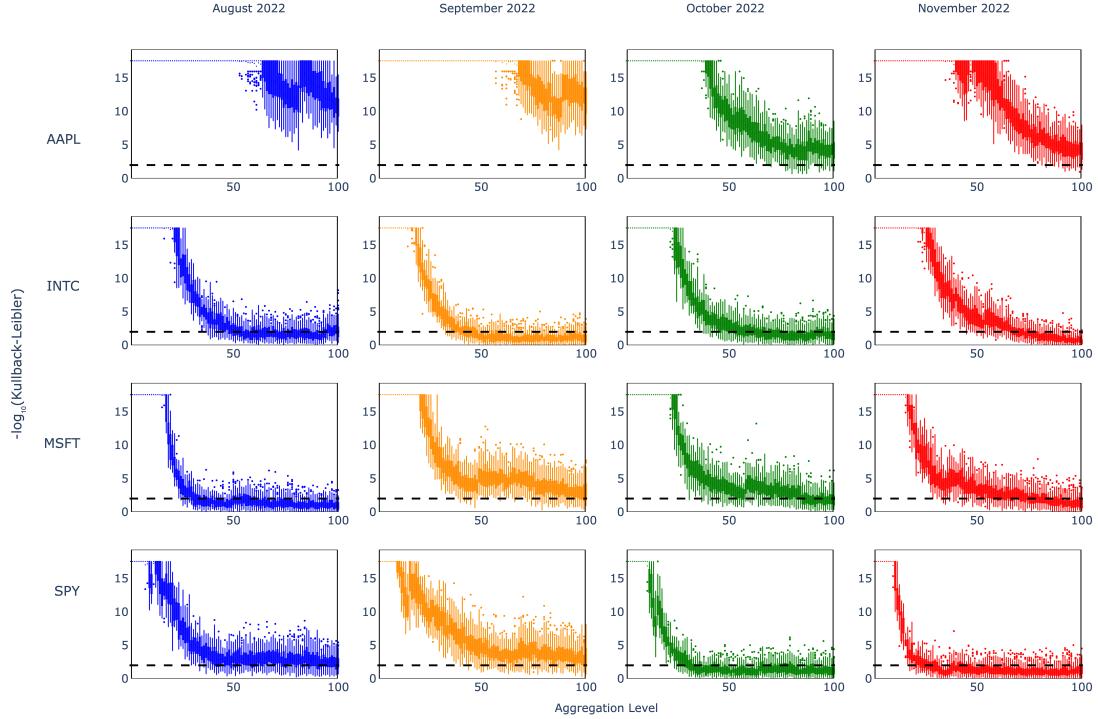


Figure 4: Results of KL test applied to AAPL, INTC, MSFT and SPY data.

Case 3: The Fourier3 spectral test on CCL, F, INTC and SNAP stocks, whose results are shown in Figure 7, shows a rare behavior compared to other tests. Here, the maximum of predictability is reached at a level higher than one, and then decreases. This means that the level of randomness has a non-monotone behavior, that was not detected by previous works.

Such non-monotonic behavior highlighting deterministic patterns at some levels of aggregations can be explained by algorithmic trading. Since the major trading activities are done algorithmically in stock markets, it is possible that the impact of algorithms becomes visible. Actions by a specific algorithm are launched with some periodicity that may coincide with peaks of predictability in the resulting plots. Alternatively, this phenomenon can be attributed to microstructure in high frequency trading. For instance, splitting large orders into smaller pieces happens with some periodicity and not with any transaction.

We note that Fourier3 is not the only test to detect such behavior. Other examples on some of the stocks are given by the AutoCorrelation ($d = 1$) test in Rabbit battery,

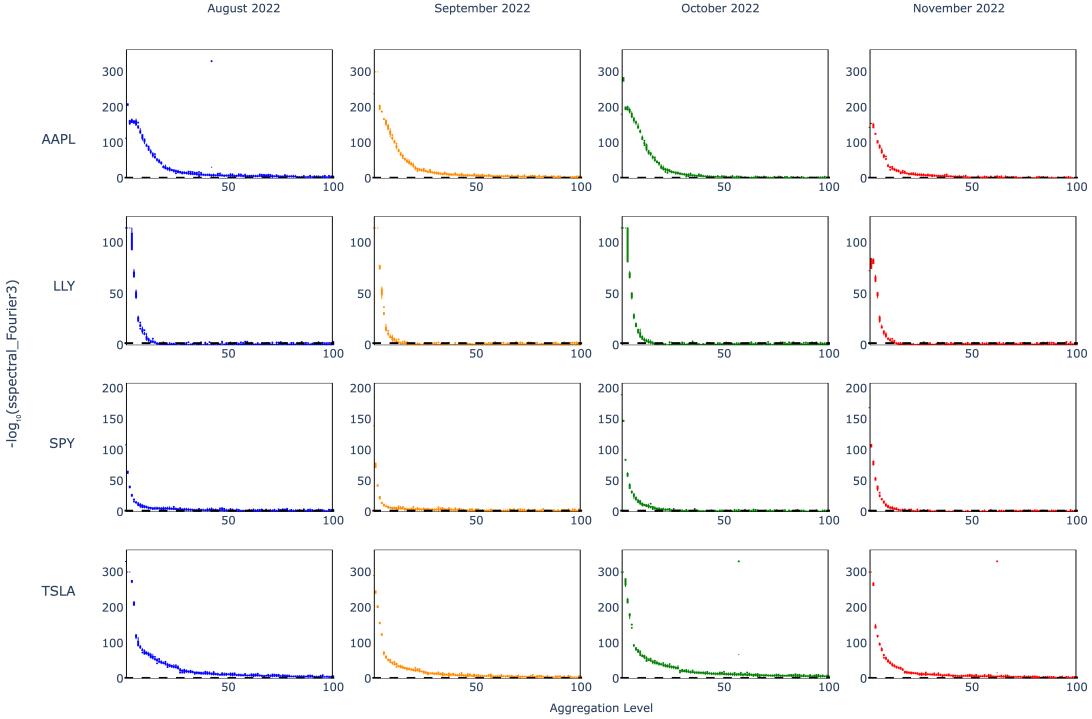


Figure 5: Results of Fourier3 test from Rabbit applied to AAPL, LLY, SPY and TSLA data.

BlockFrequency and CumulativeSums tests in NIST STS, ShannonEntropy and KL tests. Results for these tests are available in our GitHub repository.

Case 4: The HammingCorrelation tests applied on INTC data give results that look different from all the other tests and stocks: we analyze more in detail what could be the reason of such a particular behavior.

HammingCorrelation is a pattern test that analyzes blocks of the sequences and compares their Hamming weights. Then, it is strongly connected with the frequencies of the string and to its distance from having exactly 50% of zeros and 50% of ones. We analyze in more details frequencies of INTC data, in particular in August, where the particular behavior seems to be particularly far from the other tests and stocks. We apply an arithmetic mean test on our data to measure how far they are from having a perfect balance between zeros and ones. We define

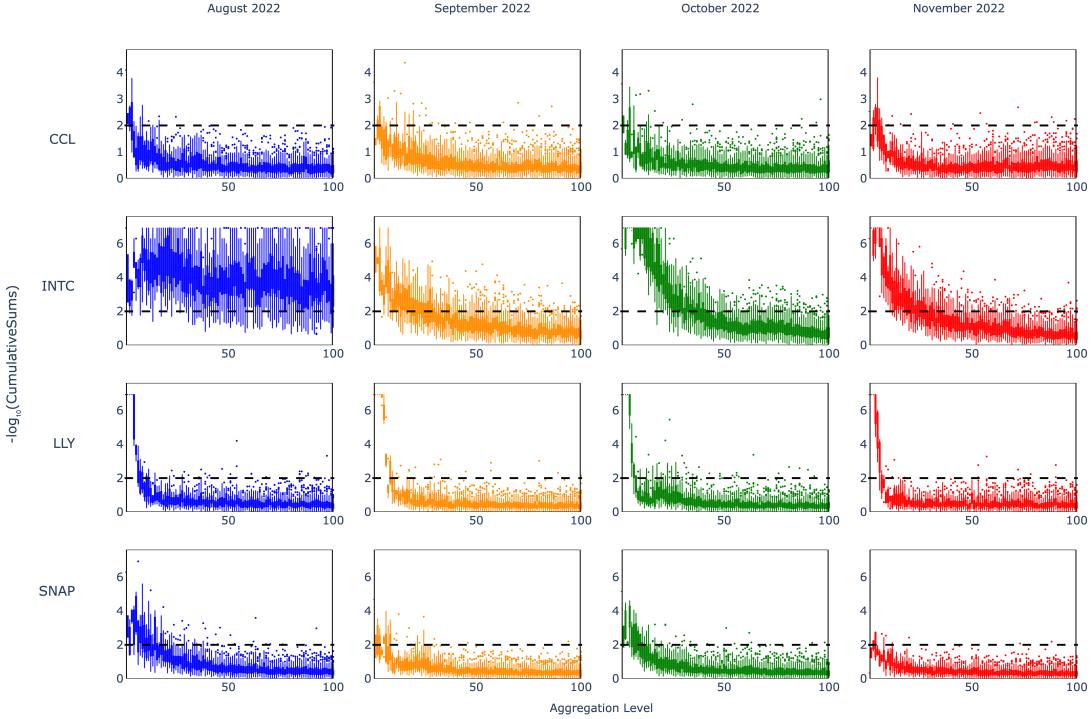


Figure 6: Results of CumulativeSums test from NIST STS applied to CCL, INTC, LLY, and SNAP data.

$$Z = \frac{2}{\sqrt{N}} \cdot \left| c - \frac{N}{2} \right|,$$

where c is the number of zeros in the sequence whose length is N . In Figure 8b, the results of the arithmetic mean test are shown for INTC in the month of August 2022. Frequencies appear really biased, far from equiprobability, and this leads to an obvious rejection of the randomness hypothesis by tests like HammingCorrelation. However, the assumption of equiprobability is not mandatory for some randomness tests. For instance, we assessed the independence of successive symbols using the KL test, checking whether the sequence consists of Bernoulli(p) variables, where p may differ from $\frac{1}{2}$.

In fact, if we try to balance these frequencies and we generate new strings starting from the same data, we get this result to be partly corrected.

To obtain binary strings that contain 50% zeros and 50% ones, we apply a similar

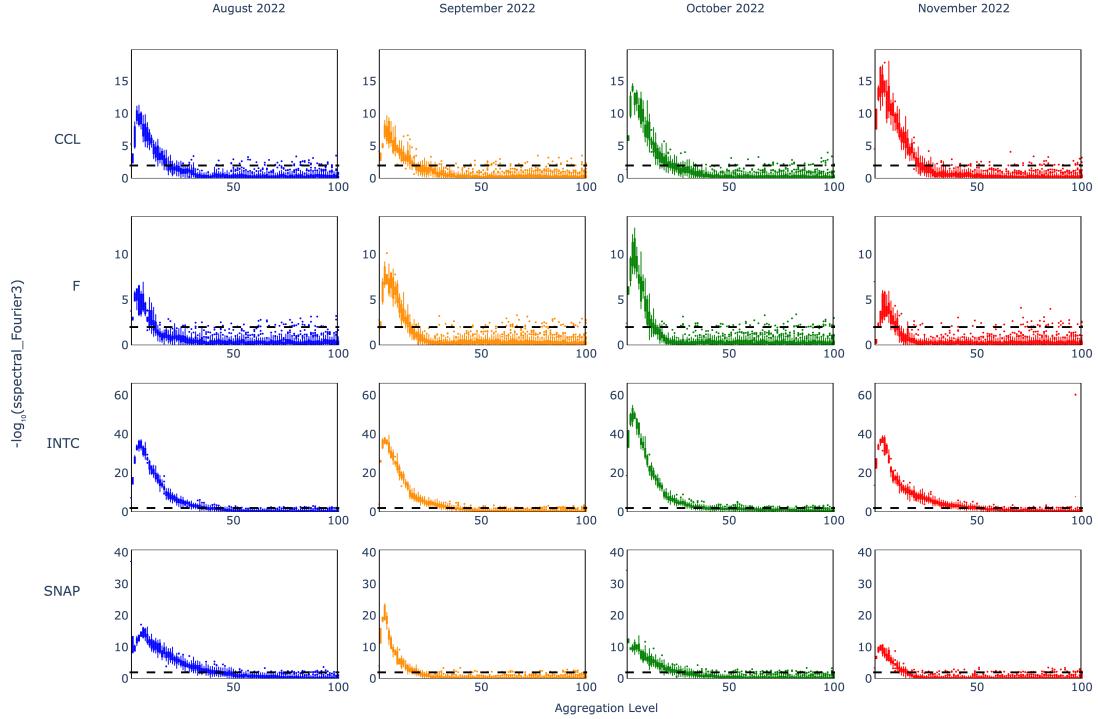


Figure 7: Results of sspectral_Fourier3 test from the Rabbit battery applied to CCL, F, INTC and SNAP.

methodology to the one used before, but this time we need to use the full information on the prices of the day in order to compute the median of them. This means that this is no more an online process, but we have to wait the end of the trading day to compute the binary strings.

Starting from the sequence of prices $\{s_1, s_2, \dots, s_N\}$, we consider the sequence $\{s_m\}$ where $m = j + i\ell$ for each aggregation level $\ell = 1, \dots, 100$ and for each sample $j = 1, \dots, \ell$. Then, we compute the median M of all the ratios $r = \frac{s_{m+1}}{s_m}$ and build the string \bar{b}_j in such a way:

$$r = \frac{s_{m+1}}{s_m} \begin{cases} < M & \text{then } \bar{b}_j \rightarrow \bar{b}_j 0 \\ > M & \text{then } \bar{b}_j \rightarrow \bar{b}_j 1 \end{cases}$$

We highlight that in such a way we get a daily string that is perfectly balanced by construction, but computing the median of the prices of the whole day requires the com-

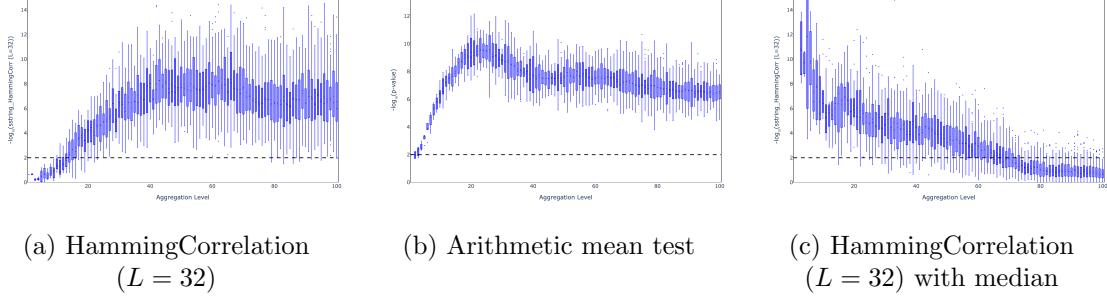


Figure 8: The boxplots show the results of three tests applied to INTC data of August 2022: (a) HammingCorrelation ($L = 32$) from Rabbit battery; (b) Arithmetic mean test; (c) HammingCorrelation ($L = 32$) test after balancing the frequencies.

putation to be done at the end of day. Then, we concatenate daily strings as before to obtain the monthly ones. For each aggregation level $\ell = 1, \dots, 100$ and for each sample $j = 1, \dots, \ell$, we run the test on a concatenated string obtained as $\bar{b}^m = \bar{b}^{d_1} || \dots || \bar{b}^{d_n}$, where the $\{\bar{b}^{d_i}\}_{i=1, \dots, n}$ are the strings obtained for each analyzed day of the month m .

Results of this approach applied to INTC data of August 2022 are shown in Figure 8c. The slow decrease of randomness indeed shows the existence of some patterns inside the sequence, particularly at low aggregation levels; anyway, when we reach level 70, the sequence appears to be random, while if not balancing the frequencies we get a non-decreasing behavior. We emphasize how other tests (for example KL test shown in Figure 9) consider the same string random starting from aggregation level nearly 50, and how the methodology (either balancing the frequencies or not) does not affect the result.

Application: The methodology we introduce in this paper gives us the possibility to let non-random binary strings obtained by UHF data be whitened out by the process of aggregation for the most of the stocks. Our randomizing process can be seen as an *online-process*: it does not require future information on tick data, so we do not have to collect data before computing the strings, but we can update the strings as new data becomes available. So, strings can be computed on-the-fly. Moreover, it does not assume any hypothesis, any model and does not require to run any simulation. Since our monthly strings have lengths of up to 1 000 000 of bits, depending on the stock - i.e., up to 10 000 bits at aggregation level 100 -, and since one month is composed of nearly 20 trading days, we claim that our methodology can produce up to $\frac{10\,000}{20} = 500$ entropy bits per day. If we select stocks not involved in cases 2 and 4, and we use the binary strings that we obtain at level of aggregation 100 as a source of randomness, we have a verifiable source certified to be random by all our tests, that is, by the most common and standard batteries of randomness tests used to certify RNGs. We thus provide an alternative approach for extracting random number

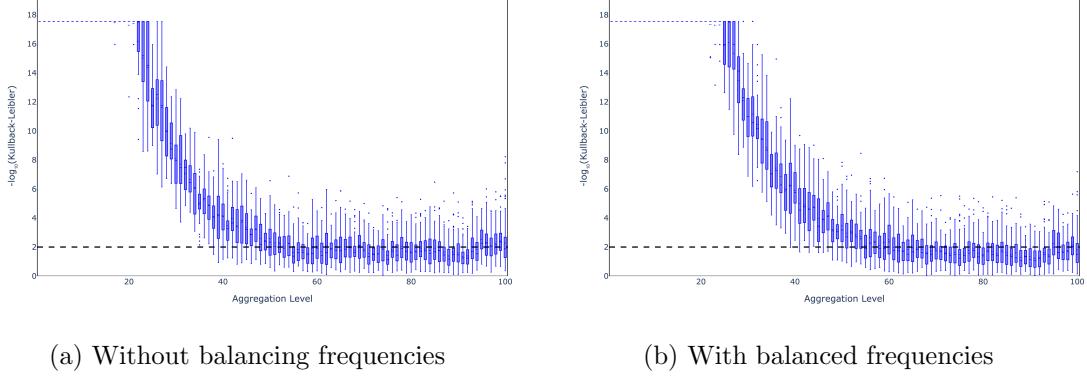


Figure 9: The boxplots show the results of KL test applied to INTC data from August 2022. (a) Results of test on the original string; (b) Results on test after balancing the frequencies of the string.

sequences from financial data [2, 3, 11]. The proposed methodology of whitening financial time series through aggregation can further be examined for its potential applications, for instance, in the implementation of randomness beacons [3]. Our source of randomness from aggregated prices can then be classified and ranked among PRNGs using the framework of [13].

5 Conclusion and future work

In this paper, we introduce a new methodology to manage strings obtained by symbolizing UHF data in order to let the randomness emerge from them.

Systemic analyses are in line with results of the paper on entropy-based approaches [23]. However, with the extension of the methodology and inclusion of standard batteries of tests, we have investigated the data from new perspectives. In particular, some tests identify random numbers sequences well, but determine predictability in stocks such as AAPL and TSLA even after price aggregation. Previous analyses have not discover the slow decay for AAPL stock.

Methods cited in case 3 have displayed a novel pattern in predictability-aggregation plots. The maximum of predictability occurs at aggregation level higher than 1 and then goes down. One of the directions for future work is to investigate such patterns and interpret them. For instance, such a pattern can be caused by algorithmic trading happening after a certain number of transactions. The method shown in case 4 is an example for tests that require equal marginal probabilities of symbols. The future research direction is to construct a battery of tests which do not demand such an assumption. One method with

this property was developed in [23], however, there is a need in developing more methods to make the analysis more systematic. For instance, frequency tests based on χ -squared distribution [32, 33] may be modified for this task. Finally, the plan for future work is to extend the data in applications to larger basket of assets. This task will require a visualization technique allowing to summarize the randomness of each stock at different aggregation levels.

References

- [1] Eugene F. Fama. Efficient capital markets: A review of theory and empirical work. *J. Finance*, 25(2):383–417, 1970. doi: 10.2307/2325486.
- [2] Ayumu Chiba and Shuichi Ichikawa. Random number generation based on cryptocurrency prices and linear feedback shift register. In *2024 Twelfth International Symposium on Computing and Networking (CANDAR)*, pages 142–148. IEEE, 2024.
- [3] Jeremy Clark and Urs Hengartner. On the use of financial data as a random beacon. In *Proceedings of the 2010 International Conference on Electronic Voting Technology/Workshop on Trustworthy Elections*, EVT/WOTE’10, page 1–8, USA, 2010. USENIX Association.
- [4] Andrew W Lo. The adaptive markets hypothesis. *The Journal of Portfolio Management*, 30(5):15–29, 2004. doi: 10.3905/jpm.2004.442611.
- [5] Walid Mensi, Chaker Aloui, Manel Hamdi, and Khuong Nguyen. Crude oil market efficiency: An empirical investigation via the shannon entropy. *Économie internationale*, 129:119–137, 2012. doi: 10.3917/eco.129.0119.
- [6] Wiston Adrián Risso. The informational efficiency: the emerging markets versus the developed markets. *Applied Economics Letters*, 16(5):485–487, 2009. doi: 10.1080/17446540802216219.
- [7] R. Cont. Empirical properties of asset returns: stylized facts and statistical issues. *Quantitative Finance*, 1(2):223–236, 2001. doi: 10.1080/713665670.
- [8] L.M. Calcagnile, F. Corsi, and S. Marmi. Entropy and efficiency of the etf market. *Comput. Econ.*, 55:143–184, 2020. doi: 10.1007/s10614-019-09885-z.
- [9] Fabrizio Lillo and J. Doyne Farmer. The long memory of the efficient market. *Stud. Nonlinear Dyn. Econom.*, 8(3), 2004. doi: 10.2202/1558-3708.1226.
- [10] Jean-Philippe Bouchaud, J. Doyne Farmer, and Fabrizio Lillo. Chapter 2 - how markets slowly digest changes in supply and demand. In Thorsten Hens and Klaus Reiner Schenk-Hoppé, editors, *Handbook of Financial Markets: Dynamics and*

Evolution, Handbooks in Finance, pages 57–160. North-Holland, San Diego, 2009. doi: 10.1016/B978-012374258-2.50006-3.

- [11] Daji Landis and Joseph Bonneau. Randomness beacons from financial data in the presence of an active attacker. *Cryptology ePrint Archive*, 2025.
- [12] Ben Moews. On random number generators and practical market efficiency. *Journal of the Operational Research Society*, 75(5):907–920, 2024.
- [13] Jeaneth Machicao, Quynh Quang Ngo, Vladimir Molchanov, Lars Linsen, and Odemir Bruno. A visual analysis method of randomness for classifying and ranking pseudo-random number generators. *Information Sciences*, 558:1–20, 2021.
- [14] John R Doyle and Catherine H Chen. Patterns in stock market movements tested as random number generators. *European Journal of Operational Research*, 227(1): 122–132, 2013.
- [15] Wei Zhang, Pengfei Wang, Xiao Li, and Dehua Shen. The inefficiency of cryptocurrency and its cross-correlation with dow jones industrial average. *Physica A: Statistical Mechanics and its Applications*, 510:658–670, 2018.
- [16] C. E. Shannon. A mathematical theory of communication. *Bell Labs Tech. J.*, 27(3): 379–423, 1948. doi: 10.1002/j.1538-7305.1948.tb01338.x.
- [17] S. Kullback and R. A. Leibler. On information and sufficiency. *Ann. Math. Stat.*, 22: 79–86, 1951. doi: 10.1214/aoms/1177729694.
- [18] A. Dionisio, R. Menezes, and D Mendes. An econophysics approach to analyse uncertainty in financial markets: an application to the portuguese stock market. *Eur. Phys. J. B*, 50:161–164, 2006. doi: 10.1140/epjb/e2006-00113-2.
- [19] Jose Alvarez-Ramirez and Eduardo Rodriguez. A singular value decomposition entropy approach for testing stock market efficiency. *Phys. A: Stat. Mech. Appl.*, 583:126337, 2021. doi: 10.1016/j.physa.2021.126337.
- [20] Anna Carbone and Linda Ponta. Relative cluster entropy for power-law correlated sequences. *SciPost Physics*, 13(3):076, 2022.
- [21] David A Hsieh. Chaos and nonlinear dynamics: application to financial markets. *The journal of finance*, 46(5):1839–1877, 1991.
- [22] Jose A Scheinkman and Blake LeBaron. Nonlinear dynamics and stock returns. *Journal of business*, pages 311–337, 1989.

- [23] Andrey Shternshis and Stefano Marmi. Price predictability at ultra-high frequency: Entropy-based randomness test. *Communications in Nonlinear Science and Numerical Simulation*, 141:108469, 2025.
- [24] L. E. Bassham, A. L. Rukhin, J. Soto, J. R. Nechvatal, M. E. Smid, E. B. Barker, S. D. Leigh, M. Levenson, M. Vangel, and et al. Banks, D. L. A statistical test suite for random and pseudorandom number generators for cryptographic applications. *NIST Special Publication 800-22 Revision 1a*, 2010. doi: 10.6028/nist.sp.800-22r1a.
- [25] Pierre L'Ecuyer and Richard Simard. Testu01: A C library for Empirical Testing of Random Number Generators. *ACM Transactions on Mathematical Software*, 33, 4 (22), 2007.
- [26] Pierre L'Ecuyer and Richard Simard. Testu01: A software library in ANSI C for empirical testing of random number generators, user's guide. *DIRO*, 2013. URL <http://simul.iro.umontreal.ca/testu01/>.
- [27] Paul C. Shields. *The Ergodic Theory of Discrete Sample Paths*. American Mathematical Society, 1996.
- [28] A. M. Zubkov. Limit distributions for a statistical estimate of the entropy. *Theory Probab. Appl.*, 18(3):611–618, 1974. doi: 10.1137/1118080.
- [29] NIST. Announcement of Proposal to Revise Special Publication 800-22 Revision 1a, January 2022. URL <https://www.nist.gov/news-events/news/2022/01/announcement-proposal-revise-special-publication-800-22-revision-1a>.
- [30] Martin D. Gould, Mason A. Porter, Stacy Williams, Mark McDonald, Daniel J. Fenn, and Sam D. Howison. Limit order books. *Quant. Finance*, 13(11):1709–1742, 2013. doi: 10.1080/14697688.2013.803148.
- [31] IDQuantique. Quantis QRNG USB. URL <https://www.idquantique.com/random-number-generation/products/quantis-random-number-generator/>.
- [32] Pierre L'Ecuyer and Richard Simard. Beware of linear congruential generators with multipliers of the form $a=\pm 2^k \pm 2^l$. *ACM Transactions on Mathematical Software (TOMS)*, 25(3):367–374, 1999.
- [33] Andrew Rukhin, Juan Soto, James Nechvatal, Miles Smid, and Elaine Barker. A statistical test suite for random and pseudorandom number generators for cryptographic applications. Technical report, 2001.
- [34] Aurel Wintner. Random factorizations and Riemann's hypothesis. *Duke Mathematical Journal*, 11(2):267 – 275, 1944. doi: 10.1215/S0012-7094-44-01122-1. URL <https://doi.org/10.1215/S0012-7094-44-01122-1>.

- [35] Giuseppe Mussardo and André LeClair. Randomness of möbius coefficients and brownian motion: Growth of the mertens function and the riemann hypothesis. *Journal of Statistical Mechanics: Theory and Experiment*, 2021(11):11310, Nov 2021. doi: 10.1088/1742-5468/ac22fb.

A Details of randomness tests

Tables 4 and 5 provide details about the tests contained in, respectively, NIST STS and TestU01 Alphabit and Rabbit. In the first case we specify the suggested parameters and our choices, while in the second case we just run tests with the standard parameters.

Randomness tests contained in NIST STS, Alphabit and Rabbit batteries can be divided in five main categories: frequency tests, pattern tests, entropy and complexity tests, spectral tests and random walks tests.

- Frequency tests: in this category we can find Frequency (Monobit) Test and Frequency Test within a Block from the NIST STS, MultinomialBitsOverlapping and HammingWeight from Rabbit and the four MultinomialBitsOverlapping of the Alphabit battery.
- Pattern tests: in this category we have Runs Test, Tests for the Longest-Run-of-Ones in a Block, Non-overlapping Template Matching Test and Overlapping Template Matching Test from the NIST STS; in the battery Rabbit there are the two ClosePairsBitMatch tests, LongestHeadRun, PeriodsInStrings, the three Hamming-Correlation and the three HammingIndependence tests, two AutoCorrelation and the Run tests, while Alphabit analyzes patterns by, again, HammingIndependence and HammingCorrelation tests.
- Entropy and complexity tests: they comprehend Binary Matrix Rank Test, Maurer’s “Universal Statistical” Test, Linear Complexity Test, Serial Test and Approximate Entropy Test from NIST STS and AppearanceSpacings, LinearComp, LempelZiv and the three MatrixRank in Rabbit battery. Moreover, the ShannonEntropy and KL tests described in Section 2.1.1 belong to this category.
- Spectral tests: they comprehend the Discrete Fourier Transform (Spectral) Test of NIST STS and the Fourier1 and Fourier3 tests from the Rabbit battery.
- Random Walks tests: they comprehend the Cumulative Sums Test, the Random Excursions Test and the Random Excursions Variant Test from the NIST STS, the three Random Walk tests in Rabbit and the two Random Walk tests in Alphabit.

B Random Number Generators used for sanity check

As explained in Section 3, we use three RNGs to run a sanity check on the tests of the batteries. Our aim is to validate the efficiency of the tests with given string lengths and parameters: since the batteries are mainly thought for cryptographic purposes, and since we have mandatory choices for some parameters due to the fact that we want to apply tests on real data, we want to certify if tests with such parameters work in a proper way. That is, if a test does not recognize more than the 2% of the strings produced by the chosen RNGs as random, then we assume that the test does not work correctly with such parameters and just avoid to run it on financial data.

We use three RNGs with three different sources of randomness: quantum physics, environmental noise and pure mathematics. For the first case, we use *Quantis QRNG USB* by *ID Quantique* [31]: this is a quantum random number generator, that exploits natural randomness coming from the principles of quantum mechanics to generate random numbers. In particular, this generator uses individual photons sent onto a semi-transparent mirror; each photon has a 50% chance of being reflected or transmitted. These two outcomes are detected and encoded as the binary bits 0 or 1.

The second generator that we use is */dev/urandom*, the random number generator built into the Linux kernel. The Linux */dev/urandom* accumulates environmental noise from device drivers and other system activities into a central entropy pool, while maintaining an internal estimate of the available entropy. Random values are generated by extracting data from this pool, and these are used to reseed periodically a PRNG.

The generator that uses randomness derived from pure mathematics exploits the Möbius function. The Möbius function is defined as

$$\mu(n) = \begin{cases} 1 & \text{if } n = 1 \\ (-1)^k & \text{if } n \text{ is product of } k \text{ distinct primes} \\ 0 & \text{if } n \text{ is divisible by a square } > 1 \end{cases}$$

We manage to obtain binary strings from this function by just considering $\mu(n)$ such that $\mu(n) = \pm 1$, discarding the zeros.

Considerable attention has been given in literature to the analysis of the randomness of the Möbius function.

The Riemann Hypothesis can be reformulated as $|\sum_{n \leq x} \mu(n)| = O(x^{\frac{1}{2} + \epsilon}) \forall \epsilon > 0$. This is equivalent to say that $\mu(n)$ consists of a random walk. So, if the Riemann Hypothesis holds, then it is true that $\mu(n)$ can be considered a source of randomness.

Moreover, in [34], Wintner defined the Rademacher random multiplicative function β :

$$\beta(n) = \begin{cases} -1, +1 \text{ with prob. } \frac{1}{2} & \text{if } n \text{ is prime} \\ \prod_{p|n} \beta(p) & \text{if } n \text{ is product of distinct primes} \\ 0 & \text{if } n \text{ has a repeated prime factor} \end{cases} .$$

For β it has been proven that $|\sum_{n \leq x} \beta(n)| = O(x^{\frac{1}{2} + \epsilon})$, $\epsilon > 0$.

Furthermore, in [35], Mussardo and LeClair test the randomness of Möbius function by tests similar to the ones we use.

Test	Suggested values	Chosen values
Frequency (Monobit) Test	$t \geq 100$	$t = 128$
Frequency Test within a Block	$t \geq 100, t \geq MJ, M \geq 20, M < 0.01\ell, J < 100$	$t = 128, M = 20$
Runs Test	$t \geq 100$	$t = 128$
Tests for the Longest-Run-of-Ones in a Block	$t \geq 128$ and if $128 \leq t \leq 6272$ $M = 8$	$t = 128, M = 8$
Binary Matrix Rank Test	$t \geq 38912$	Not executed
Discrete Fourier Transform (Spectral) Test (FFT)	$t \geq 1000$	$t = 1000$
Non-overlapping Template Matching Test	$m \in \{9, 10\}, J = 8, M > 0.01\ell$	$t = 1000, m = 9$
Overlapping Template Matching Test	$t \geq 10^6$	Not executed
Maurer's "Universal Statistical" Test	$t \geq 387840$	Not executed
Linear Complexity Test	$t \geq 10^6$	Not executed
Serial Test	$m < \lfloor \log_2 t \rfloor - 2$	$t = 128, m = 2$
Approximate Entropy Test	$m < \lfloor \log_2 t \rfloor - 5$	$t = 128, m = 5$
Cumulative Sums (Cusums) Test	$t \geq 100$	$t = 128$
Random Excursions Test	$t \geq 10^6$	Not executed
Random Excursions Variant Test	$t \geq 10^6$	Not executed

Table 4: This table shows tests contained in NIST STS. We summarize here the suggestions contained in the documentation for t =length of substrings in bits, M =block length in bits, $J = \lfloor \frac{t}{M} \rfloor$ number of blocks, m =template length in bits.

N. Test	Alphabit	Rabbit
1	MultinomialBitsOverlapping, $L = 2$	MultinomialBitsOverlapping
2	MultinomialBitsOverlapping, $L = 4$	ClosePairsBitMatch, $t = 2$
3	MultinomialBitsOverlapping, $L = 8$	ClosePairsBitMatch, $t = 4$
4	MultinomialBitsOverlapping, $L = 16$	AppearanceSpacings
5	HammingIndependence, $L = 16$ bits	LinearComp
6	HammingIndependence, $L = 32$ bits	LempelZiv
7	HammingCorrelation, $L = 32$ bits	Fourier1
8	RandomWalk1, $L = 64$	Fourier3
9	RandomWalk1, $L = 320$	LongestHeadRun
10		PeriodsInStrings
11		HammingWeight, $L = 32$ bits.
12		HammingCorrelation, $L = 32$ bits
13		HammingCorrelation, $L = 64$ bits
14		HammingCorrelation, $L = 128$ bits
15		HammingIndependence, $L = 16$ bits
16		HammingIndependence, $L = 32$ bits
17		HammingIndependence, $L = 64$ bits
18		AutoCorrelation, $d = 1$
19		AutoCorrelation, $d = 2$
20		Run
21		MatrixRank, 32×32 matrices
22		MatrixRank, 320×320 matrices
23		MatrixRank, 1024×1024 matrices
24		RandomWalk1, $L = 128$
25		RandomWalk1, $L = 1024$
26		RandomWalk1, $L = 10016$

Table 5: Detail of the tests contained in Alphabit and Rabbit sub-batteries of TestU01. L is the block length (or the random walk length), t is the dimension, d represents the lag.