

# Sprintz: Time Series Compression for the Internet of Things

DAVIS BLALOCK, Massachusetts Institute of Technology

SAMUEL MADDEN, Massachusetts Institute of Technology

JOHN GUTTAG, Massachusetts Institute of Technology

Thanks to the rapid proliferation of connected devices, sensor-generated time series constitute a large and growing portion of the world's data. Often, this data is collected from distributed, resource-constrained devices and centralized at one or more servers. A key challenge in this setup is reducing the size of the transmitted data without sacrificing its quality. Lower quality reduces the data's utility, but smaller size enables both reduced network and storage costs at the servers and reduced power consumption in sensing devices. A natural solution is to compress the data at the sensing devices. Unfortunately, existing compression algorithms either violate the memory and latency constraints common for these devices or, as we show experimentally, perform poorly on sensor-generated time series.

We introduce a time series compression algorithm that achieves state-of-the-art compression ratios while requiring less than 1KB of memory and adding virtually no latency. This method is suitable not only for low-power devices collecting data, but also for servers storing and querying data; in the latter context, it can decompress at over 3GB/s in a single thread, even faster than many algorithms with much lower compression ratios. A key component of our method is a high-speed forecasting algorithm that can be trained online and significantly outperforms alternatives such as delta coding.

Extensive experiments on datasets from many domains show that these results hold not only for sensor data but also across a wide array of other time series.

CCS Concepts: • **Information systems** → **Data compression**; • **Computer systems organization** → *Embedded systems*; • **Mathematics of computing** → Time series analysis;

Additional Key Words and Phrases: Data Compression, Time Series, Embedded Systems

## ACM Reference Format:

Davis Blalock, Samuel Madden, and John Guttag. 2018. Sprintz: Time Series Compression for the Internet of Things. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.* 2, 3, Article 93 (September 2018), 23 pages. <https://doi.org/10.1145/3264903>

## 1 INTRODUCTION

Thanks to the proliferation of smartphones, wearables, autonomous vehicles, and other connected devices, it is becoming common to collect large quantities of sensor-generated time series. Once this data is centralized in servers, many tools exist to analyze and create value from it [10, 22, 59, 66, 68, 71, 75, 76]. However, centralizing it can be challenging because of power constraints on the devices collecting the data. In particular, transmitting data wirelessly is extremely power-intensive—on a representative set of chips [38, 39], transmitting data over Bluetooth Low Energy (BLE) costs tens of *milliwatts*, while computing at full power costs only tens of *microwatts*, and sitting idle costs close to 1 microwatt.

Authors' addresses: Davis Blalock, Massachusetts Institute of Technology, [dblalock@mit.edu](mailto:dblalock@mit.edu); Samuel Madden, Massachusetts Institute of Technology, [madden@csail.mit.edu](mailto:madden@csail.mit.edu); John Guttag, Massachusetts Institute of Technology, [guttag@mit.edu](mailto:guttag@mit.edu).

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, or post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

© 2018 Copyright held by the owner/author(s). Publication rights licensed to ACM.

2474-9567/2018/9-ART93 \$15.00

<https://doi.org/10.1145/3264903>

Proc. ACM Interact. Mob. Wearable Ubiquitous Technol., Vol. 2, No. 3, Article 93. Publication date: September 2018.

One strategy for reducing this power consumption is to extract information locally and only transmit summaries [4, 14, 74]. This can be effective in some cases, but requires both a predefined use case for which a summary is sufficient and an appropriate method of constructing this summary. Devising such a method can be a significant endeavor; e.g., Verma et al. [74] conducted a research project to identify features in their data that would allow for accurate classification of brain activity.

A complementary and more general approach is to compress the data before transmitting it [4, 9, 12, 37, 72]. This allows arbitrary subsequent analysis and does not require elaborate summary construction algorithms. Unfortunately, existing compression methods either 1) are only applicable for specific types of data, such as timestamps [6, 47, 55], audio [1, 17, 52, 61] or EEG [49, 74] recordings; or 2) use algorithms that are ill-suited to sensor-generated time series.

More specifically, existing methods (e.g., [19, 20, 25, 27, 30, 41, 45, 48, 51]) violate one or more of the following design requirements:

- (1) **Small block size.** On devices with only a few kilobytes of memory, it is not possible to buffer large amounts of data before compressing it. Moreover, even with more memory, buffering can add unacceptable latency; for example, a smartwatch transmitting nine axes of 8-bit motion data at 20Hz to a smartphone would need to wait  $10000/(9 \times 1 \times 20) = 56$  seconds to fill even a 10KB buffer. This precludes using this data for gesture recognition and would add tremendous user interface latency for step counting, activity recognition, or most other purposes.
- (2) **High decompression speed.** While the device collecting the data may not need to decompress it, it is desirable to have an algorithm that could also function well in a central database. This eliminates the need to transcode the data at the server and simplifies the application. In a database, time series workloads are not only read-heavy [6, 10, 14], but often necessitate materializing data (or downsampled versions thereof) for visualization, clustering, computing correlations, or other operations [14]. At the same time, writing is often append-only [14, 55]. As a result, decompression speed is paramount, while compression speed need only be fast enough to keep up with the rate of data ingestion.
- (3) **Lossless.** Given that time series are almost always noisy and often oversampled, it might not seem necessary to compress them losslessly. However, noise and oversampling 1) tend to vary across applications, and 2) can be addressed in an application-specific way as a preprocessing step. Consequently, instead of assuming that some level of downsampling or some particular smoothing will be appropriate for all data, it is better for the compression algorithm to preserve what it is given and leave preprocessing up to the application developer.

The primary contribution of this work is SPRINTZ, a compression algorithm for time series that offers state-of-the-art compression ratios and speed while also satisfying all of the above requirements. It requires <1KB of memory, can use blocks of data as small as eight samples, and can decompress at up to 3GB/s in a single thread. SPRINTZ's effectiveness stems from exploiting 1) temporal correlations in each variable's value and variance, and 2) the potential for parallelization across different variables, realized through the use of vector instructions. The main limitation of SPRINTZ is that it operates directly only on integer time series. However, as we discuss in Section 5.8, straightforward preprocessing allows it to be applied to most floating point time series as well.

A key component of SPRINTZ's operation is a novel, vectorized forecasting algorithm for integers. This algorithm can simultaneously train online and generate predictions at close to the speed of memcpy, while significantly improving compression ratios compared to delta coding.

A second contribution is an empirical comparison of a range of algorithms currently used to compress time series, evaluated across a wide array of public datasets. We also make available code to easily reproduce these experiments, including the plots and statistical tests in the paper. To the best of our knowledge, this constitutes the largest public benchmark for time series compression.

The remainder of this paper is structured as follows. In Section 2, we introduce relevant definitions, background, and details regarding the problem we consider. In Section 3, we survey related work and what distinguishes SPRINTZ. In Sections 4 and 5, we describe SPRINTZ and evaluate it across a number of publicly available datasets. We also discuss when SPRINTZ is advantageous relative to other approaches.

## 2 DEFINITIONS AND BACKGROUND

Before elaborating upon how our method works, we introduce necessary definitions and provide relevant information regarding the problem being solved.

### 2.1 Definitions

**DEFINITION 2.1. *Sample.*** A sample is a vector  $\mathbf{x} \in \mathbb{R}^D$ .  $D$  is the sample's **dimensionality**. Each element of the sample is an integer represented using a number of bits  $w$ , the **bitwidth**. The bitwidth  $w$  is shared by all elements.

**DEFINITION 2.2. *Time Series.*** A time series  $\mathbf{X}$  of length  $T$  is a sequence of  $T$  samples,  $\mathbf{x}_1, \dots, \mathbf{x}_T$ . All samples  $\mathbf{x}_t$  share the same bitwidth  $w$  and dimensionality  $D$ . If  $D = 1$ ,  $\mathbf{X}$  is called **univariate**; otherwise it is **multivariate**.

**DEFINITION 2.3. *Rows, Columns.*** When represented in memory, we assume that each sample of a time series is one row and each dimension is one column. Because data arrives as samples and memory constraints may limit how many samples can be buffered, we assume that the data is stored in row-major order—i.e., such that each sample is stored contiguously.

### 2.2 Hardware Constraints

Many connected devices are powered by batteries or harvested energy [34]. This results in strict power budgets and, in order to satisfy them, omission of certain functionality. In particular, many devices lack hardware support for floating point operations, SIMD (vector) instructions, and integer division. Moreover, they often have no more than a few kilobytes of memory, clocks in the tens of MHz at most, and 8-, 16-, or 32-bit processors instead of 64-bit [3, 38, 39].

In contrast, we assume that the hardware used to decompress the data does not share these limitations. It is likely a modern x86 server with SIMD instructions, gigabytes of RAM, and a multi-GHz clock. However, because the amount of data it must store and query can be large, compression ratio and decompression speed are still important.

### 2.3 Data Characteristics

From a compression perspective, time series have four attributes uncommon in other data.

- (1) **Lack of exact repeats.** In text or structured records, there are many sequences of bytes—often corresponding to words or phrases—that will exactly repeat many times. This makes dictionary-based methods a natural fit. In time series, however, the presence of noise makes exact repeats less common [11, 57].
- (2) **Multiple variables.** Real-world time series often consist of multiple variables that are collected and accessed together. For example, the Inertial Measurement Unit (IMU) in modern smartphones collects three-dimensional acceleration, gyroscope, and magnetometer data, for a total of nine variables sampled at each time step. These variables are also likely to be read together, since each on its own is insufficient to characterize the phone's motion.
- (3) **Low bitwidth.** Any data collected by a sensor will be digitized into an integer by an Analog-to-Digital Converter (ADC). Nearly all ADCs have a precision of 32 bits or fewer [2], and typically 16 or fewer of these bits are useful. For example, even lossless audio codecs store only 16 bits per sample [17, 61]. Even data

that is not collected from a sensor can often be stored using six or fewer bits without loss of performance for many tasks [36, 48, 57].

- (4) **Temporal correlation.** Since the real world usually evolves slowly relative to the sampling rate, successive samples of a time series tend to have similar values. However, when multiple variables are present and samples are stored contiguously, this correlation is often present only with a lag—e.g., with nine IMU variables, every ninth value is similar. Such lag correlations violate the assumptions of most compressors, which treat adjacent bytes as the most likely to be related.

Much of the reason SPRINTZ outperforms existing methods is that it exploits or accounts for all of these characteristics, while existing methods do not.

### 3 RELATED WORK

SPRINTZ draws upon ideas from time series compression, time series forecasting, integer compression, general-purpose compression, and high-performance computing. From a technical perspective, SPRINTZ is unusual or unique in its abilities to:

- (1) Bit pack with extremely small block sizes
- (2) Bit pack low-bitwidth integers effectively
- (3) Efficiently exploit correlation between successive samples in *multivariate* time series
- (4) Naturally integrate both run-length encoding and bit packing
- (5) Exploit vectorized hardware through forecaster, learning algorithm, and bit packing method co-design

From an application perspective, SPRINTZ is distinct in that it enables higher-ratio lossless compression with far less memory and latency than competing methods.

#### 3.1 Compression of Time Series

Most work on compressing time series has focused on lossy techniques. The most common approach is to approximate the data as a sequence of low-order polynomials [27, 40–42, 45, 72]. An alternative, commonly seen in the data mining literature, is to discretize the time series using Symbolic Aggregate Approximation (SAX) [48] or its variations [15, 65]. These approaches are designed to preserve enough information about the time series to support indexing or specific data mining algorithms (e.g. [43, 56, 65]), rather than to compress the time series *per se*. As a result, they are extremely lossy; a hundred-sample time series might be compressed into one or two bytes, depending on the exact discretization parameters.

For audio time series specifically, there are a large number of lossy codecs [1, 52, 61, 73], as well as a small number of lossless [8, 17] codecs. In principle, some of these could be applied to non-audio time series. However, modern codecs make such strong assumptions about the possible numbers of channels, sampling rates, bit depths, or other characteristics that it is infeasible to use them on non-audio time series.

Many fewer algorithms exist for lossless time series compression. For floating-point time series, the only algorithm of which we are aware is that of the Gorilla database [55]. This method XORs each value with the previous value to obtain a diff, and then bit packs the diffs. In contrast to our approach, it assumes that time series are univariate and have 64-bit floating-point elements.

For lossless compression of integer time series (including timestamps), existing approaches include directly applying general-purpose compressors [14, 35, 63, 68, 75], (double) delta encoding and then applying an integer compressor [10, 55], or predictive coding and byte-packing [44]. These approaches can work well, but tend to offer both less compression and less speed than SPRINTZ.

### 3.2 Compression of Integers

The fastest methods of compressing integers are generally based on bit packing—i.e., using at most  $b$  bits to represent values in  $\{0, 2^b - 1\}$ , and storing these bits contiguously [47, 67, 78]. Since  $b$  is determined by the largest value that must be encoded, naively applying this method yields limited compression. To improve it, one can encode fixed-size blocks of data at a time, so that  $b$  can be set based on the largest values in a block instead of the whole dataset [47, 62, 78]. A further improvement is to ignore the largest few values when setting  $b$  and store their omitted bits separately [47, 78].

SPRINTZ bit packing differs significantly from existing methods in two ways. First, it compresses much smaller blocks of samples. This reduces its throughput as compared to, e.g., [47], but significantly improves compression ratios (c.f. Section 5). This is because large values only increase  $b$  for a few samples instead of for many. Second, SPRINTZ is designed for 8 and 16-bit integers, rather than 32 or 64-bit integers. Existing methods are often inapplicable to lower-bitwidth data (unless converted to higher-bitwidth data) thanks to strong assumptions about bitwidth and data layout.

A common [17, 61] alternative to bit packing is Golomb coding [31], or its special case Rice coding [60]. The idea is to assume that the values follow a geometric distribution, often with a rate constant fit to the data.

Both bit packing and Golomb coding are bit-based methods in that they do not guarantee that encoded values will be aligned on byte boundaries. When this is undesirable, one can employ byte-based methods such as 4-Wise Null Suppression [62], LEB128 [21], or Varint-G8IU [69]. These methods reduce the number of bytes used to store each sample by encoding in a few bits how many bytes are necessary to represent its value, and then encoding only that many bytes. Some, such as Simple8B [7] and SIMD-GroupSimple [77], allow fractional bytes to be stored while preserving byte alignment for groups of samples.

### 3.3 General-Purpose Compression

A reasonable alternative to using a time series compressor would be to apply a general-purpose compression algorithm, possibly after delta coding or other preprocessing. Thanks largely to the development of Asymmetric Numeral Systems (ANS) [26] for entropy coding, general purpose compressors have advanced greatly in recent years. In particular, Zstd [20], Brotli [5], LZ4 [19] and others have attained speed-compression tradeoffs significantly better than traditional methods such as GZIP [30], LZ0 [54], etc. However, these methods have much higher memory requirements than SPRINTZ and, empirically, often do not compress as well and/or decompress as quickly.

### 3.4 Predictive Filtering

For numeric data such as time series, there are four types of predictive coding commonly in use: predictive filtering [13], delta coding [47, 67], double-delta coding [10, 55], and XOR-based encoding [55]. In predictive filtering, each prediction is a linear combination of a fixed number of recent samples. This can be understood as an autoregressive model or the application of a Finite Impulse Response (FIR) filter. When the filter is learned from the data, this is termed “adaptive filtering.” Many audio compressors use some form of adaptive filtering [1, 17, 61].

Delta coding is a special case of predictive filtering where the prediction is always the previous value. Double-delta coding, also called delta-delta coding or delta-of-deltas coding, consists of applying delta coding twice in succession. XOR-based encoding is similar to delta coding, but replaces subtraction of the previous value with the XOR operation. This modification is often desirable for floating-point data [55].

Our forecasting method can be understood as a special case of adaptive filtering. While adaptive filtering is a well-studied mathematical problem in the signal processing literature, we are unaware of a practical algorithm that attains speed within an order of magnitude of our own. I.e., our method’s primary novelty is as a vectorized

*algorithm* for fitting and predicting multivariate time series, rather than as a mathematical *model* of multivariate time series. That said, it does incorporate different modeling assumptions than other compression algorithms for time series in that it reduces the model to one parameter and omits a bias term.

## 4 METHOD

To describe how SPRINTZ works, we first provide an overview of the algorithm, then discuss each of its component in detail.

### 4.1 Overview

SPRINTZ is a bit packing-based predictive coder. It consists of four components:

- (1) **Forecasting.** SPRINTZ employs a forecaster to predict each sample based on previous samples. It encodes the difference between the next sample and the predicted sample, which is typically closer to zero than the next sample itself.
- (2) **Bit packing.** SPRINTZ then bit packs the errors as a “payload” and prepends a header with sufficient information to invert the bit packing.
- (3) **Run-length encoding.** If a block of errors is all zeros, SPRINTZ waits for a block in which some error is nonzero and then writes out the number of all-zero blocks instead of the (otherwise empty) payload.
- (4) **Entropy coding.** SPRINTZ Huffman codes the headers and payloads.

These components are run on blocks of eight samples (motivated in Section 4.3), and can be modified to yield different compression-speed tradeoffs. Concretely, one can 1) skip entropy coding for greater speed and 2) choose between delta coding and our online learning method as forecasting algorithms. The latter is slightly slower but often improves compression.

We chose these steps since they allow for high speed and exploit the characteristics of time series. Forecasting leverages the high correlation of successive samples to reduce the entropy of the data. Run-length encoding allows for extreme compression in the (common) scenario that there is no change in the data—e.g., a user’s smartphone may be stationary for many hours while the user is asleep. Our method of bit packing exploits temporal correlation in the variability of the data by using the same bitwidth for points that are within the same block. Huffman coding is not specific to time series but has low memory requirements and improves compression ratios.

An overview of how SPRINTZ compresses one block of samples is shown in Algorithm 1. In lines 2-5, SPRINTZ predicts each sample based on the previous sample and any state stored by the forecasting algorithm. For the first sample in a block, the previous sample is the last element of the previous block, or zeros for the initial block. In lines 6-8, SPRINTZ determines the number of bits required to store the largest error in each column and then bit packs the values in that column using that many bits. (Recall that each column is one variable of the time series). If all columns require 0 bits, SPRINTZ continues reading in blocks until some error requires >0 bits (lines 11-13). At this point, it writes out a header of all 0s and then the number of all-zero blocks. Finally, it writes out the number of bits required by each column in the latest block as a header, and the bit packed data as a payload. Both header and payload are compressed with Huffman coding.

SPRINTZ begins decompression (Algorithm 2) by decoding the Huffman-coded bitstream into a header and a payload. Once decoded, these two components are easy to separate since the header is always first and of fixed size. If the header is all 0s, the payload indicates the length of a run of zero errors. In this case, SPRINTZ runs the predictor until the corresponding number of samples have been predicted. Since the errors are zero, the forecaster’s predictions are the true sample values. In the nonzero case, SPRINTZ unpacks the payload using the number of bits specified for each column by the header.



**Algorithm 1** encodeBlock( $\{x_1, \dots, x_B\}$ , forecaster)

---

```

1: Let buff be a temporary buffer
2: for  $i \leftarrow 1, \dots, B$  do                                // For each sample
3:    $\hat{x}_i \leftarrow \text{forecaster.predict}(x_{i-1})$ 
4:    $\text{err}_i \leftarrow x_i - \hat{x}_i$ 
5:    $\text{forecaster.train}(x_{i-1}, x_i, \text{err}_i)$ 
6: for  $j \leftarrow 1, \dots, D$  do                                // For each column
7:    $\text{nbits}_j \leftarrow \max_i \{\text{requiredNumBits}(\text{err}_{ij})\}$ 
8:    $\text{packed}_j \leftarrow \text{bitPack}(\{\text{err}_{1j}, \dots, \text{err}_{Bj}\}, \text{nbits}_j)$ 
9: // Run-length encode if all errors are zero
10: if  $\text{nbits}_j == 0, 1 \leq j \leq D$  then
11:   repeat                                                    // Scan until end of run
12:     Read in another block and run lines 2-8
13:   until  $\exists_j [\text{nbits}_j \neq 0]$ 
14:   Write  $D$  0s as headers into buff
15:   Write number of all-zero blocks as payload into buff
16:   Output  $\text{huffmanCode}(\text{buff})$ 
17: Write  $\text{nbits}_j, j = 1, \dots, D$  as headers into buff
18: Write  $\text{packed}_j, j = 1, \dots, D$  as payload into buff
19: Output  $\text{huffmanCode}(\text{buff})$ 

```

---

**Algorithm 2** decodeBlock(bytes,  $B, D$ , forecaster)

---

```

1:  $\text{nbits}, \text{payload} \leftarrow \text{huffmanDecode}(\text{bytes}, B, D)$ 
2: if  $\text{nbits}_j == 0 \forall j$  then                                // Run-length encoded
3:    $\text{numblocks} \leftarrow \text{readRunLength}()$ 
4:   for  $i \leftarrow 1, \dots, (B \cdot \text{numblocks})$  do
5:      $x_i \leftarrow \text{forecaster.predict}(x_{i-1})$ 
6:     Output  $x_i$ 
7: else                                                    // Not run-length encoded
8:   for  $i \leftarrow 1, \dots, B$  do
9:      $\hat{x}_i \leftarrow \text{forecaster.predict}(x_{i-1})$ 
10:     $\text{err}_i \leftarrow \text{unpackErrorVector}(i, \text{nbits}, \text{payload})$ 
11:     $x_i \leftarrow \text{err}_i + \hat{x}_i$ 
12:    Output  $x_i$ 
13:     $\text{forecaster.train}(x_{i-1}, x_i, \text{err}_i)$ 

```

---

## 4.2 Forecasting

SPRINTZ forecasting can use either delta coding or FIRE (Fast Integer REgression), a novel online forecasting algorithm we introduce.

**4.2.1 Delta Coding.** Forecasting with delta coding consists of predicting each sample  $x_i$  to be equal to the previous sample  $x_{i-1}$ , where  $x_0 \triangleq 0$ . This method is stateless given  $x_{i-1}$  and is extremely fast. It is particularly fast when combined with run-length encoding, since it yields a run of zero errors if and only if the data is constant.

This means that decompression of runs requires only copying a fixed vector, with no additional forecasting or training. Moreover, when answering queries, one can sometimes avoid decompression entirely—e.g., one can compute the max of all samples in the run by computing the max of only the first value.

**4.2.2 FIRE.** Forecasting with FIRE is slightly more expensive than delta coding but often yields better compression. The basic idea of FIRE is to model each value as a linear combination of a fixed number of previous values and learn the coefficients of this combination. Specifically, we learn an autoregressive model of the form:

$$x_i = ax_{i-1} + bx_{i-2} + \varepsilon_i \quad (1)$$

where  $x_i$  denotes the value of some variable at time step  $i$  and  $\varepsilon_i$  is a noise term.

Different values of  $a$  and  $b$  are suitable for different data characteristics. If  $a = 2$ ,  $b = -1$ , we obtain double-delta coding, which extrapolates linearly from the previous two points and works well when the time series is smooth. If  $a = 1$ ,  $b = 0$ , we recover delta coding, which models the data as a random walk. If  $a = \frac{1}{2}$ ,  $b = \frac{1}{2}$ , we predict each value to be the average of the previous two values, which is optimal if the  $x_i$  are i.i.d. Gaussians. In other words, these cases are appropriate for successively noisier data.

The reason FIRE is effective is that it learns online what the best coefficients are for each variable. To make prediction and learning as efficient as possible, FIRE restricts the coefficients to lie within a useful subspace. Specifically, we exploit the observation that all of the above cases can be written as:

$$x_i = x_{i-1} + \alpha x_{i-1} - \alpha x_{i-2} + \varepsilon_i \quad (2)$$

for  $\alpha \in [-\frac{1}{2}, 1]$ . Letting  $\delta_i \triangleq x_i - x_{i-1}$  and subtracting  $x_{i-1}$  from both sides, this is equivalent to

$$\delta_i = \alpha \delta_{i-1} + \varepsilon_i \quad (3)$$

This means that we can capture all of the above cases by predicting the next delta as a rescaled version of the previous delta. This requires only a single addition and multiplication, and reduces the learning problem to that of finding a suitable value for a single parameter.

To train and predict using this model, we use the functions shown in Algorithm 3. First, to initialize a FIRE forecaster, one must specify three values: the number of columns  $D$ , the learning rate  $\eta$ , and the bitwidth  $w$  of the integers stored in the columns. Internally, the forecaster also maintains an accumulator for each column (line 4) and the difference (delta) between the two most recently seen samples (line 5). The accumulator is a scaled version of the current  $\alpha$  value with a bitwidth of  $2w$ . It enables fast updates of  $\alpha$  with greater numerical precision than would be possible if modifying  $\alpha$  directly. The accumulators and deltas are both initialized to zeros.

To predict, the forecaster first derives the coefficient  $\alpha$  for each column based on the accumulator. By right shifting the accumulator  $\log_2(\eta)$  bits, the forecaster obtains a learning rate of  $2^{-\log_2(\eta)} = \eta$ . It then estimates the next deltas as the elementwise product (denoted  $\odot$ ) of these coefficients and the previous deltas. It predicts the next sample to be the previous sample plus these estimated deltas.

Because all values involved are integers, the multiplication is done using twice the bitwidth  $w$  of the data type—e.g., using 16 bits for 8 bit data. The product is then right shifted by an amount equal to the bit width. This has the effect of performing a fixed-point multiplication with step size equal to  $2^{-w}$ .

The forecaster trains by performing a gradient update on the L1 loss between the true and predicted samples. I.e., given the loss:

$$\mathcal{L}(x_i, \hat{x}_i) = |x_i - \hat{x}_i| = |x_i - (x_{i-1} + \frac{\alpha}{2^w} \cdot \delta_{i-1})| \quad (4)$$

$$= |\delta_i - \frac{\alpha}{2^w} \cdot \delta_{i-1}| \quad (5)$$



**Algorithm 3** FIRE\_Forecaster Class

---

```

1: function INIT( $D, \eta, w$ )
2:   self.learnShift  $\leftarrow \lg(\eta)$ 
3:   self.bitWidth  $\leftarrow w$  // 8-bit or 16-bit
4:   self.accumulators  $\leftarrow \text{zeros}(D)$ 
5:   self.deltas  $\leftarrow \text{zeros}(D)$ 
6: function PREDICT( $\mathbf{x}_{i-1}$ )
7:   alphas  $\leftarrow \text{self.accumulators} \gg \text{self.learnShift}$ 
8:    $\hat{\delta} \leftarrow (\text{alphas} \odot \text{self.deltas}) \gg \text{self.bitWidth}$ 
9:   return  $\mathbf{x}_{i-1} + \hat{\delta}$ 
10: function TRAIN( $\mathbf{x}_{i-1}, \mathbf{x}_i, \text{err}_i$ )
11:   gradients  $\leftarrow -\text{sign}(\text{err}_i) \odot \text{self.deltas}$ 
12:   self.accumulators  $\leftarrow \text{self.accumulators} - \text{gradients}$ 
13:   self.deltas  $\leftarrow \mathbf{x}_i - \mathbf{x}_{i-1}$ 

```

---

for one column's value  $x_i = x_{ij}$  for some  $j$  and coefficient  $\alpha$ , the gradient is:

$$\frac{\partial}{\partial \alpha} |\delta_i - \frac{\alpha}{2^w} \cdot \delta_{i-1}| = \begin{cases} -2^{-w} \delta_{i-1} & \mathbf{x}_i > \hat{\mathbf{x}}_i \\ 2^{-w} \delta_{i-1} & \mathbf{x}_i \leq \hat{\mathbf{x}}_i \end{cases} \quad (6)$$

$$= -\text{sign}(\varepsilon) \cdot 2^{-w} \delta_{i-1} \quad (7)$$

$$\propto -\text{sign}(\varepsilon) \cdot \delta_{i-1} \quad (8)$$

where we define  $\varepsilon \triangleq \mathbf{x}_i - \hat{\mathbf{x}}_i$  and ignore the  $2^{-w}$  as a constant that can be absorbed into the learning rate. In all experiments reported here, we set the learning rate to  $\frac{1}{2}$ . This value is unlikely to be ideal for any particular dataset, but preliminary experiments showed that it consistently worked reasonably well.

In practice, FIRE differs from the above pseudocode in three ways. First, instead of computing the coefficient for each sample, we compute it once at the start of each block. Second, instead of performing a gradient update after each sample, we average the gradients of all samples in each block and then perform one update. Finally, we only compute a gradient for every other sample, since this has little or no effect on the accuracy and slightly improves speed.

### 4.3 Bit Packing

An illustration of SPRINTZ's bit packing is given in Figure 1. The prediction errors from delta coding or FIRE are zigzag encoded [32] and then the minimum number of bits required is computed for each column. Zigzag encoding is an invertible transform that interleaves positive and negative integers such that each integer is represented by twice its absolute value, or twice its absolute value minus one for negative integers. This makes all values nonnegative and maps integers farther from zero to larger numbers.

Given the zigzag encoded errors, the number of bits  $w'$  required in each column can be computed as the bitwidth minus the fewest leading zeros in any of that column's errors. E.g., in Figure 1a, the first column's largest encoded value is 16, represented as 00010000, which has three leading zeros. This means that we require  $w' = 8 - 3 = 5$  bits to store the values in this column. One can find this value by ORing all the values in a column together and then using a built-in function such as GCC's `__builtin_clz` to compute the number of leading zeros in a single assembly instruction (c.f. [47]). This optimization motivates our use of zigzag encoding to make all values nonnegative.

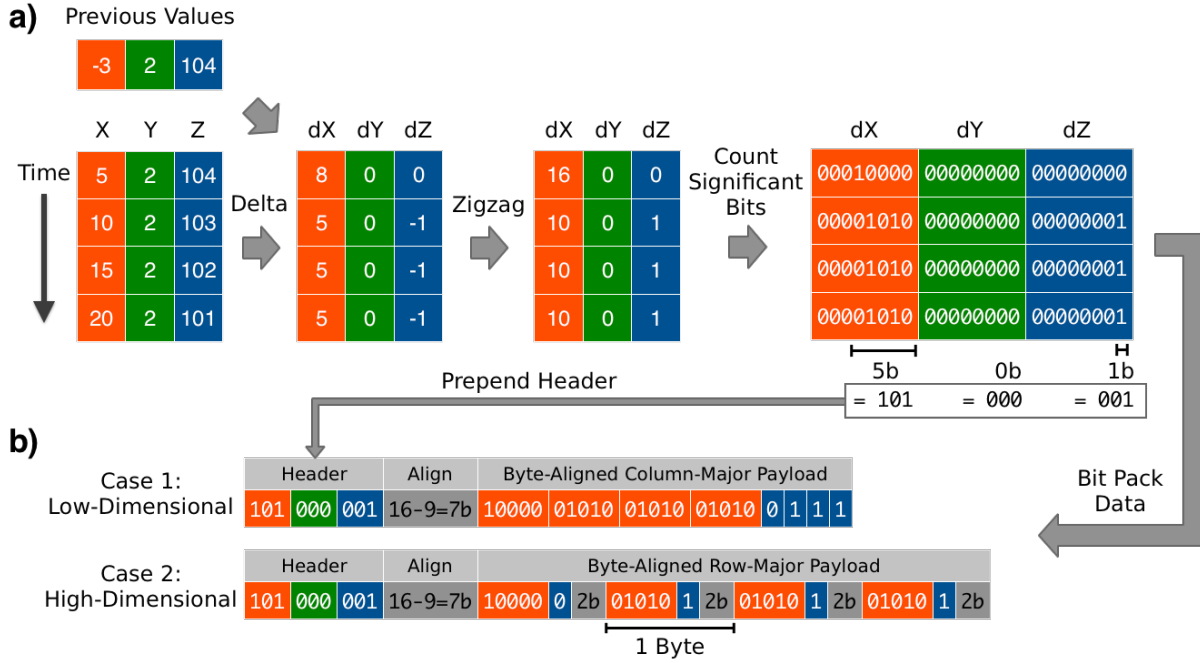


Fig. 1. Overview of SPRINTZ using a delta coding predictor. *a)* Delta coding of each column, followed by zigzag encoding of resulting errors. The maximum number of significant (nonzero) bits is computed for each column. *b)* These numbers of bits are stored in a header, and the original data is stored as a (byte-aligned) payload, with leading zeros removed. When there are few columns, each column's data is stored contiguously. When there are many columns, each row is stored contiguously, possibly with padding to ensure alignment on a byte boundary.

Once the number of bits  $w'$  required for each column is known, the zigzag-encoded errors can be bit packed. First, SPRINTZ writes out a header consisting of  $D$  unsigned integers, one for each column, storing the bitwidths. Each integer is stored in  $\log_2(w)$  bits, where  $w$  is the bitwidth of the data. Since there are  $w + 1$  possible values of  $w'$  (including 0), width  $w - 1$  is treated as a width of  $w$  by both the encoder and decoder. E.g., 8-bit data that could only be compressed to 7 bits is both stored and decoded with a bitwidth of 8.

After writing the headers, SPRINTZ takes the appropriate number of low bits from each element and packs them into the payload. When there are few columns, all the bits for a given column are stored contiguously (i.e., column-major order). When there are many columns, the bits for each *sample* are stored contiguously (i.e., row-major order). In the latter case, up to seven bits of padding are added at the end of each row so that all rows begin on a byte boundary. This means that the data for each column begins at a fixed bit offset within each row, facilitating vectorization of the decompressor. The threshold for choosing between the two formats is a sample width of 32 bits.

The reason for this threshold is as follows. Because the block begins in row-major order and we seek to reconstruct it the same way, the row-major bit packing case is more natural. For small numbers of columns, however, the row padding can significantly reduce the compression ratio. Indeed, for univariate 8-bit data, it makes compression ratios greater than 1 impossible. This gives rise to the column-major case; using a block size of eight samples and column-major order, each column's data always falls on a byte boundary without any padding. The downside of this approach is that both encoder and decoder must transpose the block. However, for

up to four 8-bit columns or two 16-bit columns, this can be done quickly using SIMD shuffling instructions.<sup>1</sup> This gives rise to the cutoff of 32 bit sample width for choosing between the formats.

As a minor bit packing optimization, one can store the headers for two or more blocks contiguously, so that there is one group of headers followed by one group of payloads. This allows many headers to share one set of padding bits between the headers and payload. Grouping headers does not require buffering more than one block of raw input, but it does require buffering the appropriate number of blocks of compressed output. In addition to slightly improving the compression ratio, it also enables more headers to be unpacked with a given number of vector instructions in the decompressor. Microbenchmarks show up to 10% improvement in decompression speed as the number of blocks in a group grows towards eight. However, we use groups of only two in all reported experiments to ensure that our results tend towards pessimism and are applicable under even the most extreme buffer size constraints.

#### 4.4 Entropy Coding

We entropy code the bit packed representation of each block using Huff0, an off-the-shelf Huffman coder [18]. This encoder treats individual bytes as symbols, regardless of the bitwidth of the original data. We use Huffman coding instead of Finite-State Entropy [18] or an arithmetic coding scheme since they are slower, and we never observed a meaningful increase in compression ratio.

The benefit of adding Huffman coding to bit packing stems from bit packing's inability to optimally encode individual bytes. For a given packed bitwidth  $w$ , bit packing models its input as being uniformly distributed over an interval of size  $2^w$ . Appropriately setting  $w$  allows it to exploit the similar variances of nearby values, but does not optimally encode individual values (unless they truly are uniformly distributed within the interval). Huffman coding is complementary in that it fails to capture relationships between nearby bytes but optimally encodes individual bytes.

We Huffman code after bit packing, instead of before, for two reasons. First, doing so is faster. This is because the bit packed block is usually shorter than the original data, so less data is fed to the Huffman coding routines. These routines are slower than the rest of SPRINTZ, so minimizing their input size is beneficial. Second, this approach increases compression. Bit packed bytes benefit from Huffman coding, but Huffman coded bytes do not benefit from bit packing, since they seldom contain large numbers of leading zeros. This absence of leading zeros is unsurprising since Huffman codes are not byte-aligned and use ones and zeros in nearly equal measure.

#### 4.5 Vectorization

Much of SPRINTZ's speed comes from vectorization. For headers, the fixed bitwidths for each field and fixed number of fields allows for packing and unpacking with a mix of vectorized byte shuffles, shifts, and masks. For payloads, delta (de)coding, zigzag (de)coding, and FIRE all operate on each column independently, and so naturally vectorize. Because the packed data for all rows is the same length and aligned to a byte boundary (in the high-dimensional case), the decoder can compute the bit offset of each column's data one time and then use this information repeatedly to unpack each row. In the low-dimensional case, all packed data fits in a single vector register which can be shuffled/masked appropriately for each possible number of columns. This is possible since there are at most four columns in this case. On an x86 machine, bit packing and unpacking can be accelerated with the pext and pdep instructions, respectively.

<sup>1</sup>For recent processors with AVX-512 instructions, one could double these column counts, but we refrain from assuming that these instructions will be available.

## 5 EXPERIMENTAL RESULTS

To assess SPRINTZ's effectiveness, we compared it to a number of state-of-the-art compression algorithms on a large set of publicly available datasets. All of our code and raw results are publicly available on the SPRINTZ website.<sup>2</sup> This website also contains additional experiments, as well as documentation of both our code and experimental setups. All experiments use a single thread on a 2013 Macbook Pro with a 2.6GHz Intel Core i7-4960HQ processor.

All reported timings and throughputs are the best of ten runs. We use the best, rather than average, since this is 1) desirable in the presence of the non-random, purely additive noise characteristic of microbenchmarks, and, 2) consequently, a best practice in microbenchmarking [46]. The best values are nearly always within 10% of the averages.

### 5.1 Datasets

- **UCR [16]** — The UCR Time Series Archive is a repository of 85 univariate time series datasets from various domains, commonly used for benchmarking time series algorithms. Because each dataset consists of many (often short) time series, we concatenate all the time series from each dataset to form a single longer time series. This is to allow dictionary-based methods to share information across time series (instead of compressing each in isolation). To mitigate artificial jumps in value from the end of one time series to the start of another, we linearly interpolate five samples between each pair.
- **PAMAP [58]** — The PAMAP dataset consists of inertial motion and heart rate data from wearable sensors on subjects performing everyday actions. It has 31 variables, most of which are accelerometer and gyroscope readings.
- **MSRC-12 [29]** — The MSRC-12 dataset consists of 80 variables of (x, y, z, depth) positions of human joints captured by a Microsoft Kinect. The subjects performed various gestures one might perform when interacting with a video game.
- **UCI Gas [28]** — This dataset consists of 18 columns of gas concentration readings and ground truth concentrations during a chemical experiment.
- **AMPDs [50]** — The Almanac of Minutely Power Datasets describes electricity, water, and natural gas consumption recorded once per minute for two years at a single home.

For datasets stored as delimited files, we first parsed the data into a contiguous, numeric array and then dumped the bytes as a binary file. Before obtaining any timing results, we first load each dataset into main memory. Because the datasets are provided as floating point values (despite most reflecting analog-to-digital converter output that was originally integer-valued), we quantized them into integers before operating on them. We did so by linearly rescaling them such that the largest and smallest values corresponded to the largest and smallest values representable with the number of bits tested—e.g., 0 and 255 for 8 bits—and then applying the floor function. Note that this is the worst case scenario for our method since it maximizes the number of bits required to represent the data.

For multivariate datasets, we allowed all methods but our own to operate on the data one variable at a time; i.e., instead of interleaving values for every variable, we store all values for each variable contiguously. This corresponds to allowing them an unlimited buffer size in which to store incoming data before compressing it. We allow these ideal conditions in order to ensure that our results for existing methods err towards optimism and to eliminate buffer size as a lurking variable.

### 5.2 Comparison Algorithms

- **SIMD-BP128 [47]** — The fastest known method of compressing integers.
- **FastPFOR [47]** — An algorithm similar to SIMD-BP128, but with better compression ratios.

<sup>2</sup><https://smarturl.it/sprintz>

- **Simple8b** [7] — An integer compression algorithm used by the popular time series database InfluxDB [10].
- **Snappy** [33] — A general-purpose compression algorithm developed by Google and used by InfluxDB, KairosDB [35], OpenTSDB [68], RocksDB [70], the Hadoop Distributed File System [66] and numerous other projects.
- **Zstd** [20] — Facebook’s state-of-the-art general purpose compression algorithm. It is based on LZ77 and entropy codes using a mix of Huffman coding and Finite State Entropy (FSE) [18]. It is available in RocksDB [70].
- **LZ4** [19] — A widely-used general-purpose compression algorithm optimized for speed and based on LZ77. It is used by RocksDB and ChronicleDB [63].
- **Zlib** [25] — A popular implementation of the DEFLATE [24] dictionary coder, which also underlies gzip [30].

For Zlib and Zstd, we use a compression level of 9 unless stated otherwise. This level heavily prioritizes compression ratio at the expense of increased compression time. We use it to improve the results for these methods in experiments in which compression time is not penalized.

We also assess three variations of SPRINTZ, corresponding to different speed/ratio tradeoffs:

- (1) **SprintzFIRE+Huf**. The full algorithm described in Section 4.
- (2) **SprintzFIRE**. Like SprintzFIRE+Huf, but without Huffman coding.
- (3) **SprintzDelta**. Like SprintzFIRE, but with delta coding instead of FIRE as the forecaster.

### 5.3 Compression Ratio

In order to rigorously assess the compression performance of both SPRINTZ and existing algorithms, it is desirable to evaluate each on a large corpus of time series from heterogeneous domains. Consequently, we use the UCR Time Series Archive [16]. This corpus contains dozens of datasets and is almost universally used for evaluating time series classification and clustering algorithms in the data mining community.

The distributions of compression ratios on these datasets for the above algorithms are shown in Figure 2. SPRINTZ exhibits consistently strong performance across almost all datasets. High-speed codecs such as Snappy, LZ4, and the integer codecs (FastPFOR, SIMDBP128, Simple8B) hardly compress most datasets at all.

Perhaps counter-intuitively, 8-bit data tends to yield higher compression ratios than 16-bit data. This is a product of the fact that the number of bits that are “predictable” is roughly constant. I.e., suppose that an algorithm can correctly predict the four most significant bits of a given value; this enables a 2:1 compression ratio in the 8-bit case, but only a  $16:12 = 4:3$  ratio in the 16-bit case. Interestingly, the fact that trailing bits tend to be too noisy to compress also suggests that one could use a lower bitwidth with little loss of information.

To assess SPRINTZ’s performance statistically, we use a Nemenyi test [53] as recommended in [23]. This test compares the mean rank of each algorithm across all datasets, where the highest-ratio algorithm is given rank 1, the second-highest rank 2, and so on. The intuition for why this test is desirable is that it not only accounts for multiple hypothesis testing in making pairwise comparisons, but also prevents a small number of large or highly compressible datasets from dominating the results.

The results of the Nemenyi test are shown in the Critical Difference Diagrams [23] in Figure 3. These diagrams show the mean rank of each algorithm on the x-axis and join methods that are not statistically significantly different with a horizontal line. SPRINTZ on high compression settings is significantly better than any existing algorithm. On lower settings, it is still as effective as the best current methods (Zlib and Zstd).

In addition to this overall comparison, it is important to assess whether FIRE improves performance compared to delta coding. Since this is a single hypothesis with matched pairs, we assess it using a Wilcoxon signed rank test. This yields p-values of .0094 in the 8-bit case and  $4.09e-12$  in the 16-bit case. As a more interpretable measure, FIRE obtains better compression on 51 of 85 datasets using 8 bits and 74 of 85 using 16. These results suggest that FIRE is generally beneficial on 8-bit data but even more beneficial on 16-bit data.

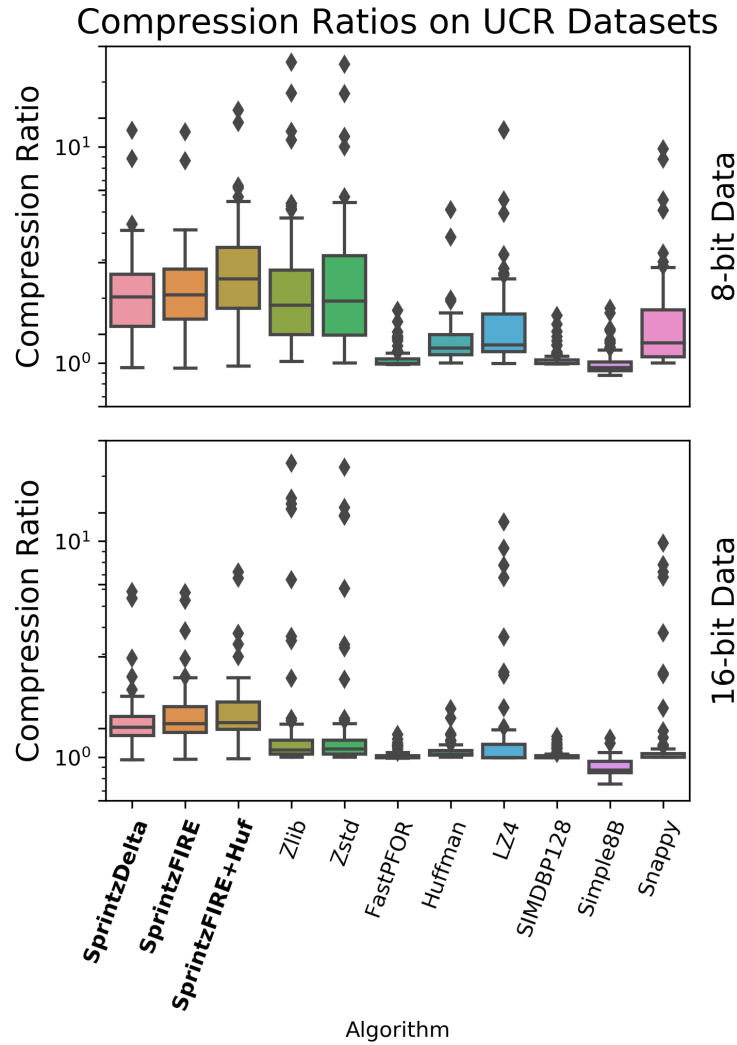


Fig. 2. Boxplots of compression performance of different algorithms on the UCR Time Series Archive. Each boxplot captures the distribution of one algorithm across all 85 datasets.

To understand why 16-bit data benefits more, we examined datasets where FIRE gives differing benefits in the two cases. The difference most commonly occurs when the data is highly compressible with just delta coding. With 8 bits and  $\sim 4\times$  compression, the forecaster's task is effectively to guess whether the next delta is -1, 0, or 1 given a current delta drawn from this same set. The Bayes error rate is high for this problem, and FIRE's attempt to learn adds variance compared to the delta coding approach of always predicting 0. In contrast, with 16 bits, the deltas span many more values and retain continuity that FIRE can exploit.



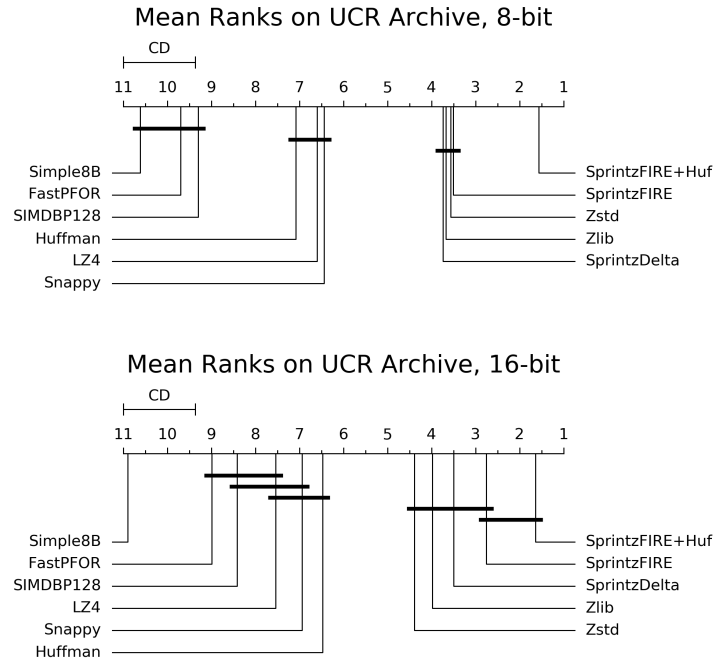


Fig. 3. Compression performance of different algorithms on the UCR Time Series Archive. The x-axis is the mean rank of each method, where rank 1 on a given dataset had the highest ratio. Methods joined with a horizontal black line are not statistically significantly different.

#### 5.4 Decompression Speed

To systematically assess the speed of SPRINTZ, we ran it on time series with varying numbers of columns and varying levels of compressibility. Because real datasets have a fixed and limited number of columns, we ran this experiment on synthetic data. Specifically, we generated a synthetic dataset of 100 million random values uniformly distributed across the full range of those possible for the given bitwidth. This data is incompressible and thus provides a worst-case estimate of SPRINTZ's speed (though in practice, we find that the speed is largely consistent across levels of compressibility).

We compressed the data with SPRINTZ set to treat it as if it had 1 through 80 columns. Numbers that do not evenly divide the data size result in SPRINTZ memcpy-ing the trailing bytes.

While using this synthetic data cannot tell us anything about SPRINTZ's compression ratio, it is suitable for throughput measurement. This is because both SPRINTZ's sequence of instructions executed and memory access patterns are effectively independent of the data distribution—SPRINTZ's core loop has no conditional branches and SPRINTZ's memory accesses are always sequential. Moreover, it exhibits throughputs on real data matching or slightly exceeding the numbers below for the corresponding number of columns (c.f. Figure 7).

As shown in Figure 4, SPRINTZ becomes faster as the number of columns increases and as the number of columns approaches multiples of 32 for 8-bit data or 16 for 16-bit data. These values correspond to the 256-bit width of a SIMD register on the machine used for testing. There is small but consistent overhead associated with using FIRE over delta coding, but both approaches are extremely fast. Without Huffman coding, SPRINTZ decompresses at multiple GB/s once rows exceed ~16B. With Huffman coding, the other components of SPRINTZ are no longer the bottleneck and SPRINTZ consistently decompresses at over 500MB/s. Note that we omit comparison to other

algorithms in this section since their speed varies with compressibility, not number of columns; see Section 5.7 for a direct comparison. Further note that the speed's dependence on number of columns is not an artifact of more columns yielding larger blocks of data. The limiting factor is serial dependence between decoding one sample and predicting the next one; this is accelerated by having wider samples that fill a vector register, but not by having longer blocks.

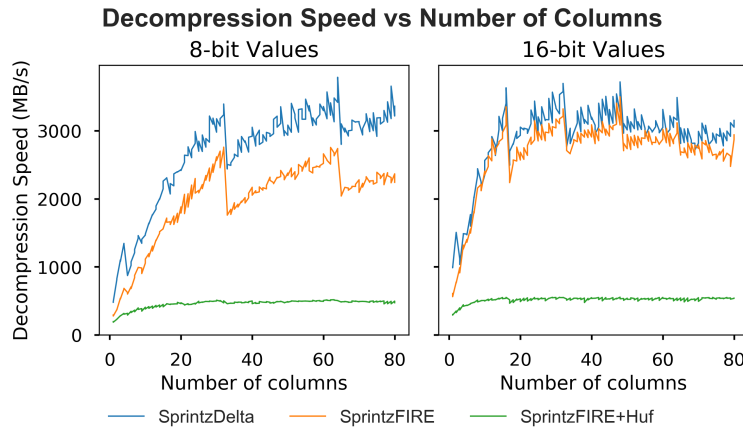


Fig. 4. SPRINTZ becomes faster as the number of columns increases and as the width of each sample approaches multiples of 32B (on a machine with 32B vector registers).

### 5.5 Compression Speed

It is important that SPRINTZ's compression speed be fast enough to keep up with the rate of data ingestion. We measured SPRINTZ's compression speed using the same methodology as decompression speed. As shown in Figure 5, SPRINTZ compresses 8-bit data at over 200MB/s on the highest-ratio setting and 600MB/s on the fastest setting. These numbers are roughly 50% larger on 16-bit data. We refrained from vectorizing this prototype implementation because 1) 200MB/s is already fast enough to run in real time even if every thread were fed data from its own gigabit network connection, and 2) low-power devices often lack vector instructions, so the measured speeds are more indicative of the rate at which these devices could compress (if scaled to the appropriate clock frequency). We again omit comparison to other compressors for the same reason as in the previous section.

The dips after 4 columns in 8-bit data and 2 columns in 16-bit data correspond to the switch from column-major bit packing to rowmajor bit packing.

### 5.6 FIRE Speed

To characterize the speed of the FIRE we repeated the above throughput experiments with both it and two other predictors commonly seen in the literature: delta and double delta coding. As shown in Figure 6, FIRE can encode at up to 5GB/s and decode at up to 6GB/s. This is nearly the same speed as the competing methods and close to the 7.5 GB/s speed of memcpy on the tested machine. Note that “encode” and “decode” here mean converting raw samples to errors and reconstructing samples from sequences of errors, respectively. These operations do not change the data size, but are the subroutines run in the SPRINTZ compressor and decompressor. The reason that there is less discrepancy between delta and FIRE encoding in isolation versus when embedded in SPRINTZ compression (Figure 5) is that, in this experiment, the implementations are vectorized.

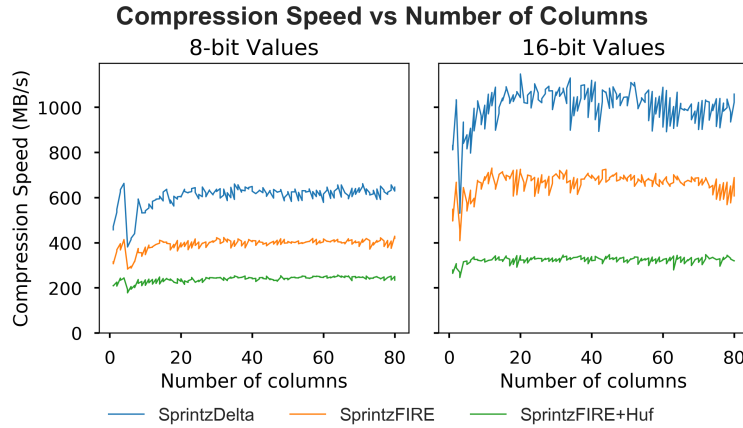


Fig. 5. SPRINTZ compresses at hundreds of MB/s even in the slowest case: its highest-ratio setting with incompressible 8-bit data. On lower settings with 16-bit data, it can exceed 1GB/s.

### 5.7 When to Use Sprintz

The above experiments provide a characterization of SPRINTZ’s speed and a statistically meaningful assessment of its compression ratio in general. However, because one often wants to obtain the best results on a particular type of data, it is helpful to know when SPRINTZ is likely to work well or poorly.

Regarding speed, SPRINTZ is most desirable when there are many variables to be compressed. We have found that the speed is largely insensitive to compression ratio, so the results in Sections 5.4 and 5.5 offer a good estimate of the speed one could expect on similar hardware. The exception to this is if the data contains long runs of constants (or constant slopes if using FIRE). In this case, the decompression speed approaches the speed of memcpy for SprintzDelta or the speed of FIRE for SprintzFIRE and SprintzFIRE+Huf.

Regarding compression ratio, the dominant requirement is that the data must have relatively strong correlations between consecutive values. This occurs when the sampling rate is fast relative to the time scale over which the measured quantity changes—the typical case when one seeks reasonably high-quality measurements. When these correlations are absent, predictive filtering (with only a two-component filter) has little value. Indeed, it can even be counterproductive. Consider the case of data that has an isolated nonzero value every few samples—e.g., the sequence  $\{0, -1, 0, 0\}$ . When delta coded, this yields  $\{0, -1, 1, 0\}$ , which requires an extra bit for SPRINTZ bit packing. In general, SPRINTZ has to pay the cost of abrupt changes twice—once when they happen, and once when they “revert” to the previous level.

Another specific case in which SPRINTZ is undesirable is when the data distribution tends to switch between discrete states. For example, in electricity consumption data, an appliance tends to use little or no electricity when it is off and a relatively constant amount when it is on. Switches between these states are expensive for SPRINTZ, and predictive filtering offers little benefit on sequences of samples that are already almost constant. SPRINTZ can still achieve reasonably good compression in this situation, but dictionary-based compressors will likely perform better. This is because they suffer no penalty from state changes, and runs of constants are their best-case input in terms of both ratio and speed. Their ratio benefits because they can often run-length encode the number of repeated values, and their speed benefits because they can decode runs at memory speed by memcpy-ing the repeated values.

As an illustration of when SPRINTZ is and is not preferable, we ran it and the comparison algorithms on several real-world datasets with differing characteristics. In Figure 7, we use the MSRC-12, PAMAP and UCI Gas datasets.

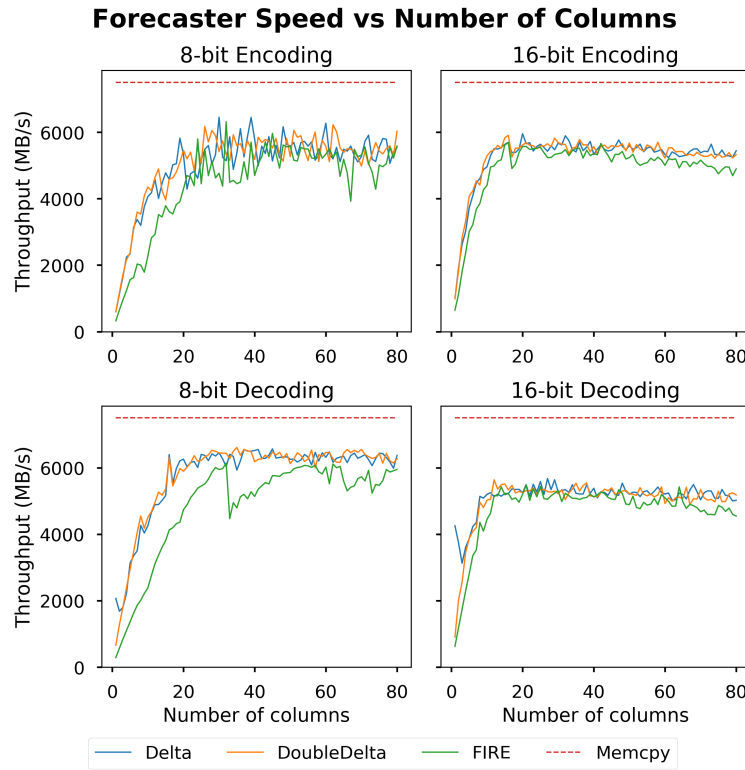


Fig. 6. FIRE is nearly as fast as delta and double delta coding. For a moderate number of columns, it runs at 5-6GB/s on a machine with 7.5GB/s memcpy speed.

These datasets contain time series that change slowly relative to the sampling rate and have 80, 31, and 18 variables, respectively. SPRINTZ achieves an excellent ratio-speed tradeoff on all three datasets, and the highest compression of any method *even on its lowest-compression setting* on the MSRC-12 dataset.

In contrast, SPRINTZ performs poorly on the AMPD Gas and AMPD Water datasets (Figure 8). These datasets chronicle the natural gas and water consumption of a house over a year, and often switch between discrete states and/or have isolated nonzero values. They also have only three and two variables, respectively. SPRINTZ achieves more than 10 $\times$  compression, but dictionary-based methods such as Zstd and LZ4 achieve even greater compression, while also decompressing faster.

### 5.8 Generalizing to Floats

While floating-point values are not the focus of this work, it is possible to apply SPRINTZ to floats by first quantizing the floating-point data. The downside of doing this is that, because floating-point numbers are not uniformly distributed along the real line, such quantization is lossy. To assess the degree of loss, we carried out an experiment to measure the error induced when quantizing real data. Note that this experiment does not assess whether SPRINTZ is the *best* means of compressing floats—it merely suggests that using integer compressors like SPRINTZ as lossy floating-point compressors is reasonable and could be a fruitful avenue for future work.

We assessed the magnitude of typical quantization errors by quantizing the UCR time series datasets. Specifically, we linearly offset and rescaled the time series in each dataset such that the minimum and maximum values in any

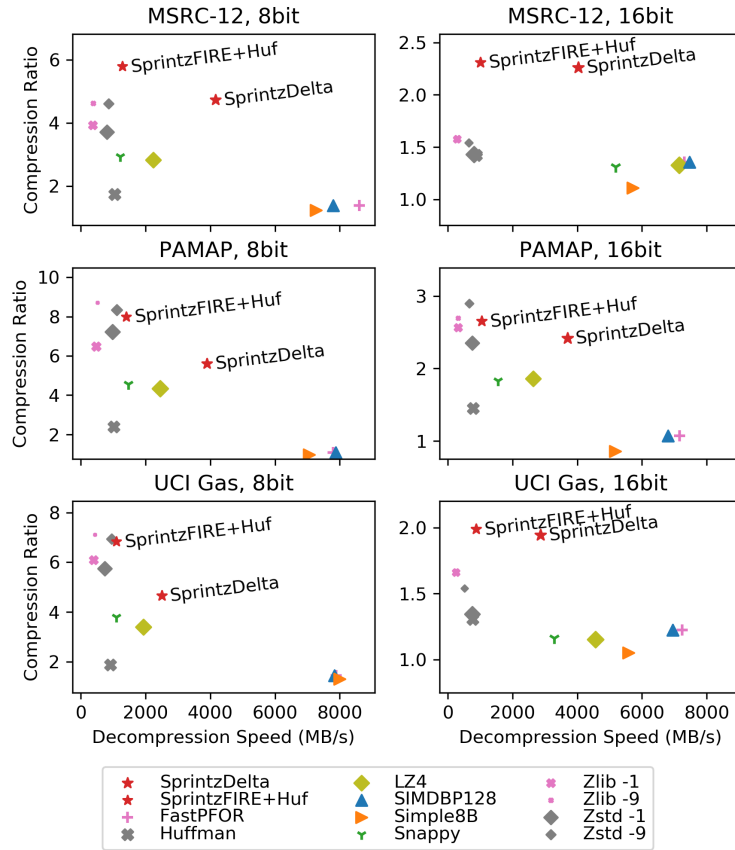
**Decompression Speed vs Compression Ratio, Success Cases**

Fig. 7. SPRINTZ achieves excellent compression ratios and speed on relatively slow-changing time series with many variables.

time series correspond to (0,255) for 8-bit quantization or (0,65535) for 16-bit quantization. We then obtained the quantized data by applying the floor function to this linearly transformed data.

To measure the error this introduced, we then inverted the linear transformation and computed the mean squared error between the original and the “reconstructed” data. The resulting error values for each dataset, normalized by the dataset’s variance, are shown in Figure 9. These normalized values can be understood as signal-to-noise ratio measurements, where the noise is the quantization error. As the figure illustrates, the quantization error is orders of magnitude smaller than the variance for nearly all datasets, and never worse than  $10\times$  smaller, even for 8-bit quantization.

This of course does not indicate that all time series can be safely quantized. Two counterexamples of which we are aware are 1) timestamps where microsecond or nanosecond resolution matters, and 2) GPS coordinates, where small decimal places may correspond to many meters. However, the above results suggest that quantization is a suitable means of applying SPRINTZ to floating-point data in many applications. This is bolstered by previous work showing that quantization even to a mere six bits [57] rarely harms classification accuracy, and quantizing to two bits is enough to support many data mining tasks [43, 48, 64, 65].

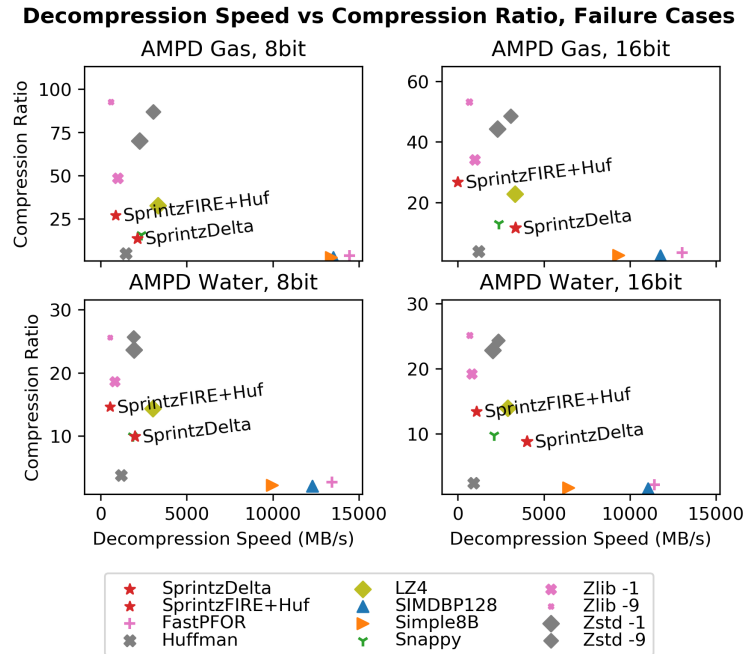


Fig. 8. SPRINTZ is less effective than other methods when the time series has large, abrupt changes and few variables.

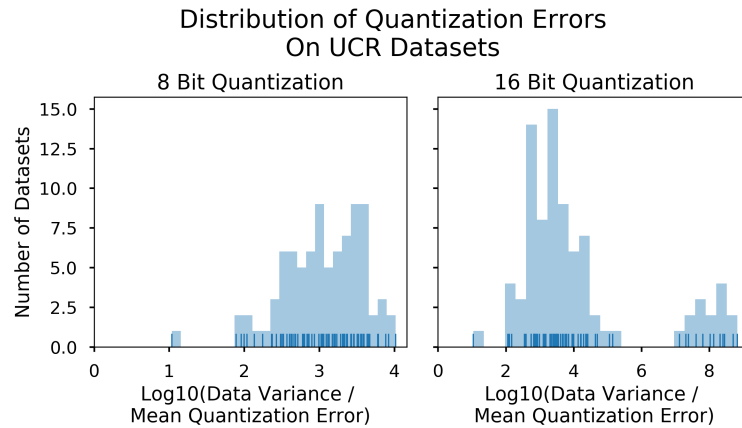


Fig. 9. Quantizing floating-point time series to integers introduces error that is orders of magnitude smaller than the variance of the data. Even with 8 bits, quantization introduces less than 1% error on 82 of 85 datasets.

## 6 CONCLUSION

We introduce SPRINTZ, a compression algorithm for multivariate integer time series that achieves state-of-the-art compression ratios across a large number of publicly available datasets. It also attains speeds of up to 3GB/s in a single thread and predictable performance as a function of the number of variables being compressed. Moreover,



it only needs to buffer eight samples at a time, enabling low latency for continuously arriving data. Finally, SPRINTZ has extremely low memory requirements, making it feasible to run even on resource-constrained devices.

As part of evaluating SPRINTZ, we also conducted what is, to the best of our knowledge, the largest empirical investigation of time series compression that has been reported. To both ensure reproducibility of our work and facilitate future research in this area, we make available all of our experiments as a public benchmark.

In future work, we hope to characterize the relationship between compression and power savings, both for SPRINTZ and for other methods. The savings are upper bounded by the compression ratio in the limit of data transmission consuming all power, but real-world systems have various overheads that cause significant deviation from this idealized model.

## ACKNOWLEDGMENTS

This material is based upon work supported by the National Science Foundation Graduate Research Fellowship under Grant No. 1122374. Any opinion, findings, and conclusions or recommendations expressed in this material are those of the authors(s) and do not necessarily reflect the views of the National Science Foundation.

## REFERENCES

- [1] Information technology – generic coding of moving pictures and associated audio information – part 7: Advanced audio coding (aac), 2006. <https://www.iso.org/standard/43345.html>.
- [2] Digi-key electronics, 2017. <https://www.digikey.com/products/en/integrated-circuits-ics/data-acquisition-analog-to-digital-converters-adc/700?k=adc&k=&pkeyword=adc&pv1989=0>.
- [3] Intel quark microcontrollers, 2017. <https://www.intel.com/content/www/us/en/embedded/products/quark/overview.html>.
- [4] N. Aharony, W. Pan, C. Ip, I. Khayal, and A. Pentland. Social fmri: Investigating and shaping social mechanisms in the real world. *Pervasive and Mobile Computing*, 7(6):643–659, 2011.
- [5] J. Alakuijala and Z. Szabadka. Brotli compressed data format. Technical report, 2016.
- [6] M. P. Andersen and D. E. Culler. Btrdb: Optimizing storage system design for timeseries processing. In *FAST*, pages 39–52, 2016.
- [7] V. N. Anh and A. Moffat. Index compression using 64-bit words. *Software: Practice and Experience*, 40(2):131–147, 2010.
- [8] Apple. Apple lossless audio codec, 2011. <https://github.com/macOSforge/alac>.
- [9] S. Arrabi and J. Lach. Adaptive lossless compression in wireless body sensor networks. In *Proceedings of the Fourth International Conference on Body Area Networks*, page 19. ICST (Institute for Computer Sciences, Social-Informatics and Telecommunications Engineering), 2009.
- [10] S. Beckett. Influxdb, 2017. <https://influxdata.com>.
- [11] D. W. Blalock and J. V. Guttag. Extract: Strong examples from weakly-labeled sensor data. In *Data Mining (ICDM), 2016 IEEE 16th International Conference on*, pages 799–804. IEEE, 2016.
- [12] T. Bose, S. Bandyopadhyay, S. Kumar, A. Bhattacharyya, and A. Pal. Signal characteristics on sensor data compression in iot-an investigation. In *Sensing, Communication, and Networking (SECON), 2016 13th Annual IEEE International Conference on*, pages 1–6. IEEE, 2016.
- [13] T. Boutell. Png (portable network graphics) specification version 1.0. 1997.
- [14] M. Buevich, A. Wright, R. Sargent, and A. Rowe. Respawn: A distributed multi-resolution time-series datastore. In *Real-Time Systems Symposium (RTSS), 2013 IEEE 34th*, pages 288–297. IEEE, 2013.
- [15] A. Camerra, T. Palpanas, J. Shieh, and E. Keogh. isax 2.0: Indexing and mining one billion time series. In *Data Mining (ICDM), 2010 IEEE 10th International Conference on*, pages 58–67. IEEE, 2010.
- [16] Y. Chen, E. Keogh, B. Hu, N. Begum, A. Bagnall, A. Mueen, and G. Batista. The ucr time series classification archive, July 2015. [www.cs.ucr.edu/~eamonn/time\\_series\\_data/](http://www.cs.ucr.edu/~eamonn/time_series_data/).
- [17] J. Coalson. Flac-free lossless audio codec, 2008. <http://flac.sourceforge.net>.
- [18] Y. Collet. Finite state entropy. <https://github.com/Cyan4973/FiniteStateEntropy>.
- [19] Y. Collet. Lz4-extremely fast compression, 2017. <https://github.com/Cyan4973/lz4>.
- [20] Y. Collet. Zstandard - fast real-time compression algorithm, 2017. <https://facebook.github.io/zstd/>.
- [21] D. D. I. F. Committee et al. Dwarf debugging information format, version 4. *Free Standards Group*, 2010.
- [22] J. Dean and S. Ghemawat. Mapreduce: simplified data processing on large clusters. *Communications of the ACM*, 51(1):107–113, 2008.
- [23] J. Demšar. Statistical comparisons of classifiers over multiple data sets. *Journal of Machine learning research*, 7(Jan):1–30, 2006.
- [24] L. P. Deutsch. Deflate compressed data format specification version 1.3. 1996.
- [25] P. Deutsch and J.-L. Gailly. Zlib compressed data format specification version 3.3. Technical report, 1996.

- [26] J. Duda. Asymmetric numeral systems: entropy coding combining speed of huffman coding with compression rate of arithmetic coding. *arXiv preprint arXiv:1311.2540*, 2013.
- [27] F. Eichinger, P. Efron, S. Karnouskos, and K. Böhm. A time-series compression technique and its application to the smart grid. *The VLDB Journal*, 24(2):193–218, 2015.
- [28] J. Fonollosa, S. Sheik, R. Huerta, and S. Marco. Reservoir computing compensates slow response of chemosensor arrays exposed to fast varying gas concentrations in continuous monitoring. *Sensors and Actuators B: Chemical*, 215:618–629, 2015.
- [29] S. Fothergill, H. M. Mentis, P. Kohli, and S. Nowozin. Instructing people for training gestural interactive systems. In J. A. Konstan, E. H. Chi, and K. Höök, editors, *CHI*, pages 1737–1746. ACM, 2012.
- [30] J.-L. Gailly and M. Adler. The gzip home page, 2003. <https://www.gzip.org/>.
- [31] S. Golomb. Run-length encodings. *IEEE transactions on information theory*, 12(3):399–401, 1966.
- [32] Google. Protocol buffers encoding, 2001. <https://developers.google.com/protocol-buffers/docs/encoding#types>.
- [33] S. Gunderson. Snappy: A fast compressor/decompressor, 2015. <https://code.google.com/p/snappy>.
- [34] M. A. Hanson, H. C. Powell Jr, A. T. Barth, K. Ringgenberg, B. H. Calhoun, J. H. Aylor, and J. Lach. Body area sensor networks: Challenges and opportunities. *Computer*, 42(1), 2009.
- [35] B. Hawkins. Kairos db: Fast time series database on cassandra. <https://github.com/kairosdb/kairosdb>.
- [36] B. Hu, T. Rakthanmanon, Y. Hao, S. Evans, S. Lonardi, and E. Keogh. Discovering the intrinsic cardinality and dimensionality of time series using mdl. In *Data Mining (ICDM), 2011 IEEE 11th International Conference on*, pages 1086–1091. IEEE, 2011.
- [37] N. Q. V. Hung, H. Jeung, and K. Aberer. An evaluation of model-based approaches to sensor data compression. *IEEE Transactions on Knowledge and Data Engineering*, 25(11):2434–2447, 2013.
- [38] T. Instruments. 2.4-ghz bluetooth® low energy system-on-chip, 2013. <http://www.ti.com/lit/ds/symlink/cc2540.pdf>.
- [39] T. Instruments. Cc2640 simplelink bluetooth wireless mcu, 2016. <http://www.ti.com/lit/ds/swrs176b/swrs176b.pdf>.
- [40] E. Keogh, K. Chakrabarti, M. Pazzani, and S. Mehrotra. Dimensionality reduction for fast similarity search in large time series databases. *Knowledge and information Systems*, 3(3):263–286, 2001.
- [41] E. Keogh, K. Chakrabarti, M. Pazzani, and S. Mehrotra. Locally adaptive dimensionality reduction for indexing large time series databases. *ACM Sigmod Record*, 30(2):151–162, 2001.
- [42] E. Keogh, S. Chu, D. Hart, and M. Pazzani. An online algorithm for segmenting time series. In *Data Mining, 2001. ICDM 2001, Proceedings IEEE International Conference on*, pages 289–296. IEEE, 2001.
- [43] E. Keogh, J. Lin, and A. Fu. Hot sax: Efficiently finding the most unusual time series subsequence. In *Data mining, fifth IEEE international conference on*, pages 8–pp. Ieee, 2005.
- [44] E. Lazin. Akumuli time-series database. <https://akumuli.org>.
- [45] D. Lemire. A better alternative to piecewise linear time series segmentation. In *Proceedings of the 2007 SIAM International Conference on Data Mining*, pages 545–550. SIAM, 2007.
- [46] D. Lemire. Microbenchmarking calls for idealized conditions, 2018. <https://lemire.me/blog/2018/01/16/microbenchmarking-calls-for-idealized-conditions/>.
- [47] D. Lemire and L. Boytsov. Decoding billions of integers per second through vectorization. *Software: Practice and Experience*, 45(1):1–29, 2015.
- [48] J. Lin, E. Keogh, S. Lonardi, and B. Chiu. A symbolic representation of time series, with implications for streaming algorithms. In *Proceedings of the 8th ACM SIGMOD workshop on Research issues in data mining and knowledge discovery*, pages 2–11. ACM, 2003.
- [49] Y. Liu, M. De Vos, and S. Van Huffel. Compressed sensing of multichannel eeg signals: the simultaneous cosparsity and low-rank optimization. *IEEE Transactions on Biomedical Engineering*, 62(8):2055–2061, 2015.
- [50] S. Makonin, B. Ellert, I. V. Bajic, and F. Popowich. Electricity, water, and natural gas consumption of a residential house in Canada from 2012 to 2014. *Scientific Data*, 3(160037):1–12, 2016.
- [51] S.-G. Miaou and H.-L. Yen. Multichannel eeg compression using multichannel adaptive vector quantization. *IEEE transactions on biomedical engineering*, 48(10):1203–1207, 2001.
- [52] J. Moffitt. Ogg vorbis. *Linux journal*, 2001(81es):9, 2001.
- [53] P. Nemenyi. Distribution-free multiple comparisons. In *Biometrics*, volume 18, page 263. INTERNATIONAL BIOMETRIC SOC 1441 I ST, NW, SUITE 700, WASHINGTON, DC 20005-2210, 1962.
- [54] M. Oberhumer. Lzo-a real-time data compression library. <http://www.oberhumer.com/opensource/lzo/>, 2008.
- [55] T. Pelkonen, S. Franklin, J. Teller, P. Cavallaro, Q. Huang, J. Meza, and K. Veeraraghavan. Gorilla: A fast, scalable, in-memory time series database. *Proceedings of the VLDB Endowment*, 8(12):1816–1827, 2015.
- [56] T. Rakthanmanon and E. Keogh. Fast shapelets: A scalable algorithm for discovering time series shapelets. In *Proceedings of the 2013 SIAM International Conference on Data Mining*, pages 668–676. SIAM, 2013.
- [57] T. Rakthanmanon, E. J. Keogh, S. Lonardi, and S. Evans. Time series epenthesis: Clustering time series streams requires ignoring some data. In *Data Mining (ICDM), 2011 IEEE 11th International Conference on*, pages 547–556. IEEE, 2011.

- [58] A. Reiss and D. Stricker. Towards global aerobic activity monitoring. In *Proceedings of the 4th International Conference on Pervasive Technologies Related to Assistive Environments*, page 12. ACM, 2011.
- [59] S. Rhea, E. Wang, E. Wong, E. Atkins, and N. Storer. Littletable: a time-series database and its uses. In *Proceedings of the 2017 ACM International Conference on Management of Data*, pages 125–138. ACM, 2017.
- [60] R. F. Rice. Some practical universal noiseless coding techniques, part 3, module psl14, k+. 1991.
- [61] T. Robinson. Shorten: Simple lossless and near-lossless waveform compression, 1994.
- [62] B. Schlegel, R. Gemulla, and W. Lehner. Fast integer compression using simd instructions. In *Proceedings of the Sixth International Workshop on Data Management on New Hardware*, pages 34–40. ACM, 2010.
- [63] M. Seidemann and B. Seeger. Chronicledb: A high-performance event store. In *EDBT*, pages 144–155, 2017.
- [64] P. Senin and S. Malinchik. Sax-vsm: Interpretable time series classification using sax and vector space model. In *Data Mining (ICDM), 2013 IEEE 13th International Conference on*, pages 1175–1180. IEEE, 2013.
- [65] J. Shieh and E. Keogh. isax: disk-aware mining and indexing of massive time series datasets. *Data Mining and Knowledge Discovery*, 19(1):24–57, 2009.
- [66] K. Shvachko, H. Kuang, S. Radia, and R. Chansler. The hadoop distributed file system. In *Mass storage systems and technologies (MSST), 2010 IEEE 26th symposium on*, pages 1–10. IEEE, 2010.
- [67] H. Siedelmann, A. Wender, and M. Fuchs. High speed lossless image compression. In *German Conference on Pattern Recognition*, pages 343–355. Springer, 2015.
- [68] B. Sigoure. Opentsdb: The distributed, scalable time series database. *Proc. OSCON*, 11, 2010.
- [69] A. A. Stepanov, A. R. Gangolli, D. E. Rose, R. J. Ernst, and P. S. Oberoi. Simd-based decoding of posting lists. In *Proceedings of the 20th ACM international conference on Information and knowledge management*, pages 317–326. ACM, 2011.
- [70] F. D. E. Team. Rocksdb: A persistent key-value store for fast storage environments. <http://rocksdb.org>.
- [71] A. Thusoo, J. S. Sarma, N. Jain, Z. Shao, P. Chakka, N. Zhang, S. Antony, H. Liu, and R. Murthy. Hive-a petabyte scale data warehouse using hadoop. In *Data Engineering (ICDE), 2010 IEEE 26th International Conference on*, pages 996–1005. IEEE, 2010.
- [72] A. Ukil, S. Bandyopadhyay, and A. Pal. Iot data compression: Sensor-agnostic approach. In *Data Compression Conference (DCC), 2015*, pages 303–312. IEEE, 2015.
- [73] J.-M. Valin, K. Vos, and T. Terriberry. Definition of the opus audio codec. Technical report, 2012.
- [74] N. Verma, A. Shoeb, J. Bohorquez, J. Dawson, J. Guttag, and A. P. Chandrakasan. A micro-power eeg acquisition soc with integrated feature extraction processor for a chronic seizure detection system. *IEEE Journal of Solid-State Circuits*, 45(4):804–816, 2010.
- [75] F. Yang, E. Tschetter, X. Léauté, N. Ray, G. Merlino, and D. Ganguli. Druid: A real-time analytical data store. In *Proceedings of the 2014 ACM SIGMOD international conference on Management of data*, pages 157–168. ACM, 2014.
- [76] M. Zaharia, M. Chowdhury, M. J. Franklin, S. Shenker, and I. Stoica. Spark: Cluster computing with working sets. *HotCloud*, 10(10-10):95, 2010.
- [77] W. X. Zhao, X. Zhang, D. Lemire, D. Shan, J.-Y. Nie, H. Yan, and J.-R. Wen. A general simd-based approach to accelerating compression algorithms. *ACM Transactions on Information Systems (TOIS)*, 33(3):15, 2015.
- [78] M. Zukowski, S. Heman, N. Nes, and P. Boncz. Super-scalar ram-cpu cache compression. In *Data Engineering, 2006. ICDE'06. Proceedings of the 22nd International Conference on*, pages 59–59. IEEE, 2006.

Received May 2018; revised July 2018; accepted September 2018