

# Learned Compression of Nonlinear Time Series With Random Access

Andrea Guerra\*  
University of Pisa  
Pisa, Italy

Giorgio Vinciguerra\*  
University of Pisa  
Pisa, Italy

Antonio Boffa†  
EPFL  
Lausanne, Switzerland

Paolo Ferragina†  
Sant’Anna School for Advanced Studies  
Pisa, Italy

**Abstract**—Time series play a crucial role in many fields, including finance, healthcare, industry, and environmental monitoring. The storage and retrieval of time series can be challenging due to their unstoppable growth. In fact, these applications often sacrifice precious historical data to make room for new data.

General-purpose compressors like Xz and Zstd can mitigate this problem with their good compression ratios, but they lack efficient random access on compressed data, thus preventing real-time analyses. Ad-hoc streaming solutions, instead, typically optimise only for compression and decompression speed, while giving up compression effectiveness and random access functionality. Furthermore, all these methods lack awareness of certain special regularities of time series, whose trends over time can often be described by some linear and nonlinear functions.

To address these issues, we introduce NeaTS, a randomly-accessible compression scheme that approximates the time series with a sequence of nonlinear functions of different kinds and shapes, carefully selected and placed by a partitioning algorithm to minimise the space. The approximation residuals are bounded, which allows storing them in little space and thus recovering the original data losslessly, or simply discarding them to obtain a lossy time series representation with maximum error guarantees.

Our experiments show that NeaTS improves the compression ratio of the state-of-the-art lossy compressors that use linear or nonlinear functions (or both) by up to 14%. Compared to lossless compressors, NeaTS emerges as the only approach to date providing, simultaneously, compression ratios close to or better than the best existing compressors, a much faster decompression speed, and orders of magnitude more efficient random access, thus enabling the storage and real-time analysis of massive and ever-growing amounts of (historical) time series data.

## I. INTRODUCTION

Time series are pervasive across a multitude of fields, including finance, healthcare, industry, and environmental monitoring. These sorted sequences of time-stamped data points represent a wide variety of dynamic phenomena, from market prices to patient vitals and sensor readings, and they have become invaluable for decision-making, trend analysis, and forecasting.

Unsurprisingly, the efficient storage, transmission, and analysis of time series have become more and more challenging as their volume has grown exponentially [1], [2], leading to the development of numerous time series databases [3]–[7].

This work has been accepted for publication in Proceedings of the 41st IEEE International Conference on Data Engineering (ICDE 2025)

\*Equal contribution. Correspondence to: andrea.guerra@phd.unipi.it, giorgio.vinciguerra@unipi.it

†Work done while the author was at the University of Pisa.

Data compression is the key strategy to lower the cost of time series storage and transmission [8]. The easiest way to approach it is to use one of the off-the-shelf general-purpose compressors (such as Brotli [9], Zstd [10], Xz [11], Lz4 [12], Snappy [13], etc.). These tools are capable of achieving commendable compression ratios, but they require a substantial computational overhead, both in terms of CPU and memory usage, which often makes them unsuitable on hardware- and energy-constrained devices such as smartphones, smart wearable, IoT or edge devices.

Motivated by this shortcoming, several new special-purpose compressors have been developed for time series, often reducing the computational overhead at the expense of lower compression ratios. Most notably, several works [14]–[17] have shown how to compress and decompress floating-point time series data much faster than general-purpose compressors, enabling both high ingestion rates and efficient scans. However, not much attention has been given to the design and the evaluation of the random access operation to single values of the time series [18], [19], even in benchmarking studies [20]. This is quite surprising given that the most fundamental queries in time series databases ultimately rely on accessing data within a specific time interval [21], [22], which from a compressed storage perspective boils down to combining a random access operation (to retrieve the first data point) with a scan (to retrieve the subsequent data points within the interval). However, providing efficient random access is challenging, and it often conflicts with achieving good compression ratios.

Furthermore, none of the above compressors can harness a key peculiarity of time series data: its trends over time can often be described by some linear and nonlinear functions [23], [24]. Indeed, although there is a rich literature on approximating and indexing a time series via linear [25]–[27], polynomial and other functions [28]–[31], or via Fourier and wavelet transforms [32], all these methods are lossy [8] and thus inapplicable in cases where we need to reconstruct the original data for accurate analyses. A step in this direction has been made by some learned compressors that are not specifically designed for time series [33]–[35]. But these approaches either exploit linear functions only [33], [34] or use sub-optimal partitioning algorithms and non-error-bounded approximations [33], [35], so they fall short of reaching the best possible compression efficacy.

**Our contribution.** We contribute to the long line of research on time series compression as follows:

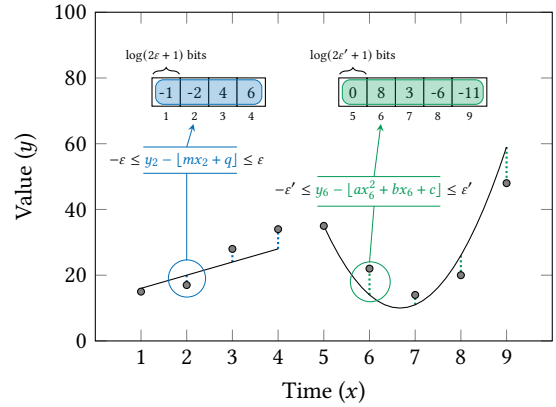
- We show how to compute piecewise approximations using several kinds of *nonlinear* functions (such as quadratic, radical, exponential, logarithmic, and Gaussian) under a given error bound  $\varepsilon$  optimally, i.e. in linear time and with the guarantee that the number of pieces is minimised. This generalises the classic algorithm to compute piecewise *linear* approximations [36].
- We introduce an algorithm to partition a time series into variable-sized fragments, each associated with a *different* nonlinear approximation, so that the space of the output is minimised. This generalises a previous result for increasing linear functions only [34].
- By combining the above two results with proper succinct data structures, we design NeaTS, a new randomly-accessible compression scheme that approximates the time series with a sequence of nonlinear functions of different kinds and shapes. The residuals of the approximation are bounded, which allows storing them in little space and thus recovering the original data losslessly, or discarding them to obtain a lossy time series representation with maximum error guarantees. Figure 1 shows an example of NeaTS.
- We conduct a thorough experimental evaluation on 16 real-world time series, whose size ranges from thousands to hundreds of millions of data points, comparing our NeaTS against 2 lossy compressors, as well as 5 general-purpose and 7 special-purpose lossless compressors, including the recent ALP [17] and LeCo [35]. Our results show that the lossy version of NeaTS improves uniformly the compression ratio of previous error-bounded approximations based on linear or nonlinear functions (or both), with an improvement of up to 14%. Compared to lossless compressors, NeaTS emerges as the only approach to date delivering, simultaneously, compression ratios close to or better than the existing compressors (i.e. the best compression ratio among the special-purpose compressor on 14/16 datasets, and the best overall on 4/16 dataset), a much faster decompression speed, and up to 3 orders of magnitude more efficient random access. No other compressor to date can achieve such a good performance in one of these factors without significantly sacrificing others. We finally show that NeaTS delivers superior performance across range queries of different sizes, thus benefiting the wide variety of queries in time series databases that access data within specific time intervals.

**Outline.** Section II gives the background and definitions. Section III introduces our NeaTS. Section IV presents our experimental results. Section V discusses related work. Section VI concludes the paper and suggests some open problems.

## II. BACKGROUND

We now provide some background information, starting with a definition of the data we compress.

**Definition 1** (Time series). A time series  $T$  is a sequence of  $n$  data points of the form  $(x_k, y_k)$ , where  $x_k \in \mathbb{N}$  is the



**Fig. 1:** NeaTS represents fragments of the time series via linear or nonlinear functions learned from the data. The residuals of the approximation are bounded by a value  $\varepsilon$  so, if lossless compression is needed, we store them in packed arrays (shown on top).

timestamp, and  $y_k \in \mathbb{Z}$  is the value associated with it, ordered increasingly by time, i.e.  $T = [(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)]$  where  $x_1 < x_2 < \dots < x_n$ . A fragment of  $T$  is a subsequence  $T[i, j] = [(x_i, y_i), \dots, (x_j, y_j)]$  for any two indexes  $i, j$  such that  $1 \leq i \leq j \leq n$ .

We require the values to be integers, which is common in practice. In fact, single/double-precision floating-point values can be interpreted as 32/64-bit integers, or better, since the values in real-world time series typically have a fixed number  $x$  of significant digits after the decimal point, we can multiply them by the constant  $10^x$  and turn them into integers [19].

We focus on functional approximations for time series but, unlike known approaches [8], we use them to design not only a lossy compressor but also a lossless one.

To illustrate, consider a time series  $T$  on a Cartesian plane where the horizontal axis represents time and the vertical axis represents values. Any function  $f$  passing through all the data points  $(x_1, y_1), \dots, (x_n, y_n)$  is a lossless encoding of  $T$  because, for a given timestamp  $x_k$ , we can recover the corresponding value as  $y_k = f(x_k)$ , thus requiring us to store only the (parameters of the) function  $f$ . However, the number of parameters of a function passing exactly through all the data points could be so large to result in no compression (consider e.g. a polynomial interpolation of the  $n$  data points, which generally requires storing  $n$  coefficients). Therefore, we allow  $f$  to make some “errors” but in a controlled way, namely, we bound the infinity norm of the errors.

**Definition 2** ( $\varepsilon$ -approximations). Let  $\varepsilon \geq 0$  be an integer. A function  $f$  is said to be an  $\varepsilon$ -approximation of a time series  $T$  (or fragment  $T[i, j]$ ) if we have  $|f(x_k) - y_k| \leq \varepsilon$  for every data point  $(x_k, y_k)$  in  $T$  (resp.  $T[i, j]$ ).

A function  $f$  defined in this way is a lossy representation of  $T$ . To make it lossless, as observed in [34], it is sufficient to also store the “corrections” (i.e. residuals)  $c_k = y_k - \lfloor f(x_k) \rfloor$  in  $\lceil \log(2\varepsilon + 1) \rceil$  bits each, which thus allows recovering  $y_k$  as  $\lfloor f(x_k) \rfloor + c_k$ . Intuitively, the smaller the value of  $\varepsilon$ , the less

space is needed to store the corrections; however, at the same time, more space may be required to store the parameters of  $f$ , due to it being more complex to better fit the data points.

In this scenario, there are two issues to solve. The first is how to compute a function that  $\varepsilon$ -approximates (a fragment of)  $T$ . The second is how to choose  $\varepsilon$  so that the storage of both the function  $f$  and the corrections takes minimal space, possibly using a different  $\varepsilon$ -value for different fragments of  $T$ . In the case  $f$  is a linear function, these two issues have already been solved [34], [36].

We now recall how to compute a linear  $\varepsilon$ -approximation, which will be the starting point of our extension to nonlinear functions. Given an integer  $\varepsilon \geq 0$  and an index  $i$  into a time series  $T$ , O'Rourke's algorithm [36] finds the longest fragment  $T[i, j]$  that admits a linear  $\varepsilon$ -approximation  $f(x) = mx + b$  in optimal  $O(j - i)$  time.<sup>1</sup> The algorithm works by maintaining a set  $P$  of feasible coefficients in the 2D space with  $m$  on the horizontal axis and  $b$  on the vertical axis, and by processing the data points in  $T[i, j]$  left-to-right by shrinking  $|P|$  at each data point. Regarding  $P$ , notice that, according to Definition 2, the linear function  $f$  must satisfy the inequality  $|f(x_k) - y_k| = |mx_k + b - y_k| \leq \varepsilon$  for every  $k = i, \dots, j$ , which can be rewritten as  $b \leq (-x_k)m + y_k + \varepsilon$  and  $b \geq (-x_k)m + y_k - \varepsilon$ . Therefore,  $P$  is a convex polygon in the 2D space defined by the intersection of the  $2(j - i + 1)$  half-planes arising from these two inequalities for  $k = i, \dots, j$ .

When adding the  $(j + 1)$ th data point would cause  $P$  to be empty (i.e. the fragment  $T[i, j]$  cannot be made any longer), the algorithm stops and picks a pair  $(m, b) \in P$  as the coefficients of the linear function  $f$  that  $\varepsilon$ -approximates the fragment  $T[i, j]$ .

Let us call  $P_k$  the polygon at a generic step  $k$ . Crucial for the time-efficiency of the algorithm is the fact that each edge of  $P_k$  has a slope that lies in the range  $[-x_k, 0)$ , which is easy to see since each half-plane defining  $P_k$  has slope of the form  $-x_l$  (due to the above inequalities) and that  $0 < x_l < x_k$  for  $l < k$  (due to Definition 1). This fact allows updating  $P_k$  with the inequalities arising from the  $(k + 1)$ th data point in constant amortised time.

We refer the reader to the seminal paper [36] for more details and anticipate that we will generalise this algorithm to work with several kinds of nonlinear functions.

### III. THE NEATS COMPRESSOR

We now introduce our compression scheme for time series. We call it NeaTS (Nonlinear error-bounded approximation for Time Series) since it exploits the potential of nonlinear functions to compress time series. We design it in three steps.

- 1) We describe how to compute an  $\varepsilon$ -approximation for a fragment of the time series using several kinds of nonlinear functions, such as quadratic, radical, exponential, logarithmic, and Gaussian functions (Section III-A).

<sup>1</sup>More recent papers address this same problem [25]–[27], [37], [38], sometimes proposing algorithms that are equivalent, or sub-optimal in terms of time complexity, or that find shorter fragments compared to the earlier algorithm by O'Rourke [36] we consider here.

- 2) We introduce an algorithm to divide the input time series into variable-sized fragments, each with an associated nonlinear approximation with a different  $\varepsilon$ -value, so that the space of the compressed output is minimised. While our focus is on lossless compression, we also show how to adapt this algorithm to produce a (lossy) piecewise nonlinear  $\varepsilon$ -approximation of a time series in linear time (Section III-B).
- 3) We show how to support random access to individual values of the time series by combining the compressed output with proper succinct data structures (Section III-C).

#### A. Computing a nonlinear $\varepsilon$ -approximation

Linear functions have surely been the most widely used representations for time series due to their simplicity [25]–[27], [37], [38]. Nonetheless, they may not always be the best choice to approximate a time series because of the possible presence of nonlinear patterns in real-world data [23], [24].

Let us be given an integer  $\varepsilon \geq 0$  and an index  $i$  into a time series  $T$ . We now show how to find the longest fragment  $T[i, j]$  that admits an  $\varepsilon$ -approximation for some kinds of nonlinear functions with two parameters  $\theta_1$  and  $\theta_2$ . We do so by generalising the classic algorithm by O'Rourke [36] for linear functions (recalled in Section II).

We start with exponential functions of the form  $f(x) = \theta_2 e^{\theta_1 x}$ . According to Definition 2, we must ensure that  $f$  satisfies the inequality  $|f(x_k) - y_k| = |\theta_2 e^{\theta_1 x_k} - y_k| \leq \varepsilon$  for every data point  $k = i, \dots, j$ , which can be rewritten as

$$\begin{aligned} \ln \theta_2 &\leq (-x_k)\theta_1 + \ln(y_k + \varepsilon) \\ \ln \theta_2 &\geq (-x_k)\theta_1 + \ln(y_k - \varepsilon), \end{aligned}$$

under the assumption  $y_k - \varepsilon > 0$ .<sup>2</sup>

The above inequalities define a pair of half-planes in the Cartesian plane with  $\theta_1$  on the horizontal axis and  $\ln \theta_2$  on the vertical axis. Therefore, their intersection for  $k = i, \dots, j$  originates a convex polygon of feasible parameters for the exponential function  $f$ . Since the slope of each polygon edge is  $-x_k$ , a direct reduction to the algorithm by O'Rourke gives an optimal  $O(j - i)$ -time algorithm to compute an exponential  $\varepsilon$ -approximation of  $T[i, j]$ .

In general, we can show the following.

**Theorem 1.** *Let  $T = [(x_1, y_1), \dots, (x_n, y_n)]$  be a time series, let  $\varepsilon \geq 0$  be an integer, and let  $f$  be a function with two parameters  $\theta_1$  and  $\theta_2$ . If, for any  $k$ , the inequalities  $-\varepsilon \leq f(x_k) - y_k \leq \varepsilon$  can be transformed into inequalities of the form  $\alpha_k \leq t_k m + b \leq \omega_k$ , where:*

- 1)  $\alpha_k$ ,  $t_k$ , and  $\omega_k$  can be computed in constant time from  $\varepsilon$  and  $T[k]$ ;
- 2)  $m$  and  $b$  are derived from  $\theta_1$  and  $\theta_2$ , respectively, via a change of variable, i.e.  $m = \phi(\theta_1)$  and  $b = \psi(\theta_2)$  for some invertible functions  $\phi$  and  $\psi$ ;
- 3)  $t_k$  is a positive increasing function of  $x_k$ .

*Then, for any  $i$ , we can compute the longest fragment  $T[i, j]$  that admits an  $\varepsilon$ -approximation  $f$  in optimal  $O(j - i)$  time.*

<sup>2</sup>If not satisfied, just add  $\varepsilon + 1 - \min_k y_k$  to all  $y_k$ s in the time series.

**TABLE I:** Some examples of two-parameter functions  $f$  that we can use in Theorem 1, together with the terms  $m, b, t_k, \alpha_k$  and  $\omega_k$  defining the transformed inequalities.

$f(x)$	$m$	$b$	$t_k$	$\alpha_k$	$\omega_k$
$\theta_2 e^{\theta_1 x}$	$\theta_1$	$\ln \theta_2$	$x_k$	$\ln(y_k - \varepsilon)$	$\ln(y_k + \varepsilon)$
$\theta_2 x^{\theta_1}$	$\theta_1$	$\ln \theta_2$	$\ln x_k$	$\ln(y_k - \varepsilon)$	$\ln(y_k + \varepsilon)$
$\ln(\theta_2 x^{\theta_1})$	$\theta_1$	$\ln \theta_2$	$\ln x_k$	$y_k - \varepsilon$	$y_k + \varepsilon$
$\theta_1 x + \theta_2$	$\theta_1$	$\theta_2$	$x_k$	$y_k - \varepsilon$	$y_k + \varepsilon$
$\theta_1 \sqrt{x} + \theta_2$	$\theta_1$	$\theta_2$	$\sqrt{x_k}$	$y_k - \varepsilon$	$y_k + \varepsilon$
$\theta_1 x^2 + \theta_2$	$\theta_1$	$\theta_2$	$x_k^2$	$y_k - \varepsilon$	$y_k + \varepsilon$
$\theta_1 x^2 + \theta_2 x$	$\theta_1$	$\theta_2$	$x_k$	$(y_k - \varepsilon)/x_k$	$(y_k + \varepsilon)/x_k$
$\theta_1 x^3 + \theta_2 x$	$\theta_1$	$\theta_2$	$x_k^3$	$(y_k - \varepsilon)/x_k$	$(y_k + \varepsilon)/x_k$
$\theta_1 x^3 + \theta_2 x^2$	$\theta_1$	$\theta_2$	$x_k$	$(y_k - \varepsilon)/x_k^2$	$(y_k + \varepsilon)/x_k^2$

*Proof.* We reduce the computation of the feasible parameters  $\theta_1$  and  $\theta_2$  of the (possibly nonlinear) function  $f$  to the intersection of half-planes in the 2D space with  $m$  on the horizontal axis and  $b$  on the vertical axis, for which we can use the algorithm of O'Rourke [36]. Using the same notation as in [36], let us rewrite the transformed inequalities as  $b \geq (-t_k)m + \alpha_k$  and  $b \leq (-t_k)m + \omega_k$ . These inequalities clearly represent half-planes in that 2D space.

Now, let  $P_k$  be the convex polygon resulting from the intersection of these inequalities where the subscript ranges as  $i, i+1, \dots, k$ . Analogously to [36, Lemma 1], we need to establish the property that the slope of each edge of  $P_k$  belongs to  $[-t_k, 0)$ . The polygon is specified by edges of the form  $b = (-t_l)m + \alpha_l$  and  $b = (-t_l)m + \omega_l$  for  $l = i, \dots, k$ . By assumption (3), the value  $t_l$  is the result of applying a positive increasing function to  $x_l$ . This, combined with the fact that  $0 < x_l < x_k$  (by Definition 1), implies that each edge of  $P_k$  has slope  $-t_l \geq -t_k$  and that  $-t_l < 0$ , thus each slope belongs to  $[-t_k, 0)$ . This property is enough to guarantee the correctness and time complexity of the algorithm that maintains  $P_k$  (cf. proofs of [36, Theorems 1 and 2]). We thus conclude by observing that once the next data point  $T[j+1]$  causes the polygon to be empty, we choose a pair  $(m, b) \in P_j$  and return  $\theta_1 = \phi^{-1}(m)$  and  $\theta_2 = \psi^{-1}(b)$  as the parameters of  $f$ .  $\square$

Taking again as an example exponential functions of the form  $f(x) = \theta_2 e^{\theta_1 x}$ , we first transform the inequalities  $-\varepsilon \leq f(x_k) - y_k \leq \varepsilon$  via simple algebraic manipulations to

$$\underbrace{\ln(y_k - \varepsilon)}_{\alpha_k} \leq \underbrace{x_k}_{t_k} \underbrace{\theta_1}_m + \underbrace{\ln \theta_2}_b \leq \underbrace{\ln(y_k + \varepsilon)}_{\omega_k},$$

and then we apply Theorem 1, which gives the desired exponential  $\varepsilon$ -approximation for a fragment  $T[i, j]$  in optimal  $O(j-i)$  time.<sup>3</sup> Table I shows other examples with linear, exponential, power, logarithmic, and radical functions.<sup>4</sup>

<sup>3</sup>The logarithm and other operations can be computed in constant time with mild assumptions on the model of computation [39].

<sup>4</sup>It is straightforward (and sometimes useful to better approximate the data) to compute a function whose graph is horizontally shifted to the first timestamp  $x_i$  of  $T[i, j]$ : we simply store  $x_i$ , subtract it from the timestamps in  $T[i, j]$ , then we apply Theorem 1 to compute a function  $g$  with domain  $[0, x_j - x_i]$  and output  $f(x) = g(x - x_i)$ .

In some cases, we can use Theorem 1 even for functions  $f$  with three parameters, provided that we add some constraints to reduce the number of *free* parameters to two, since otherwise the set of feasible parameters for  $f$  becomes a polyhedron  $P_k$  in a 3D space (i.e. one dimension for each parameter), which we cannot handle in linear time [30], [31], [40], [41].

Take, for example, quadratic functions of the form  $f(x) = \theta_1 x^2 + \theta_2 x + \theta_3$ . By forcing the function to pass through the first data point  $T[i]$ , i.e. by setting  $f(x_i) = y_i$  and thus fixing  $\theta_3 = y_i - \theta_1 x_i^2 - \theta_2 x_i$  (which we store explicitly), we can transform the inequalities  $-\varepsilon \leq f(x_k) - y_k \leq \varepsilon$  via simple algebraic manipulations to

$$\underbrace{\frac{y_k - y_i - \varepsilon}{x_k - x_i}}_{\alpha_k} \leq \underbrace{(x_k + x_i)}_{t_k} \underbrace{\theta_1}_m + \underbrace{\theta_2}_b \leq \underbrace{\frac{y_k - y_i + \varepsilon}{x_k - x_i}}_{\omega_k},$$

and then we apply Theorem 1. A similar derivation can be done for Gaussian-like functions of the form  $f(x) = e^{\theta_1 x^2 + \theta_2 x + \theta_3}$ .

We conclude this section by observing that a repeated application of Theorem 1 from  $T[1]$  to  $T[n]$  allows partitioning  $T$  into the longest fragments associated with an  $\varepsilon$ -approximation, thus giving the following result.

**Corollary 1.** *Given a time series  $T = [(x_1, y_1), \dots, (x_n, y_n)]$ , a value  $\varepsilon \geq 0$ , and a function  $f$  satisfying the assumptions of Theorem 1, we can compute a piecewise  $\varepsilon$ -approximation of  $T$  with the smallest number of functions of the  $f$ -kind in  $O(n)$  time.*

Corollary 1 directly yields a lossy error-bounded (in terms of infinity norm) representation of  $T$ . As discussed in Section II, this can be made lossless by storing the corrections  $y_k - \lfloor f(x_k) \rfloor$  in  $\lceil \log(2\varepsilon + 1) \rceil$  bits each. In the next section, we describe a more powerful partitioning algorithm to orchestrate different types of nonlinear functions and error bounds.

### B. Partitioning a time series with nonlinear $\varepsilon$ -approximations

Let us be given a set  $\mathcal{F}$  of functions that satisfy the assumptions of Theorem 1, and a set  $\mathcal{E}$  of error bounds. We now turn our attention to the problem of partitioning a time series  $T$  into fragments, each  $\varepsilon$ -approximated (with  $\varepsilon \in \mathcal{E}$ ) by a function from  $\mathcal{F}$ , with the goal of minimising the overall space of the lossless representation of  $T$ , which is given by the storage of the corrections and the parameters of the functions.

At a high level, our approach computes, for each  $f \in \mathcal{F}$  and each  $\varepsilon \in \mathcal{E}$ , the piecewise  $\varepsilon$ -approximation of  $T$  composed of functions of the  $f$ -kind, and then produces the desired partition of  $T$  by stitching together properly-chosen fragments (possibly adjusting their start and end points) taken from the  $|\mathcal{F}| \cdot |\mathcal{E}|$  different piecewise approximations of  $T$ . This generalises a previous result for increasing linear functions only [34].

More in detail, we define a graph  $\mathcal{G}$  with one node for each data point in  $T$ , plus one sink node denoting the end of the time series. Each fragment  $T[i, j-1]$  that is  $\varepsilon$ -approximated by a function  $f \in \mathcal{F}$  produces an edge  $(i, j)$  of  $\mathcal{G}$  whose weight  $w_{f, \varepsilon}(i, j)$  is defined as the bit-size of the compression

of  $T[i, j - 1]$  via  $f$  and the  $j - i$  corrections stored in  $\lceil \log(2\varepsilon + 1) \rceil$ -bits each, i.e.  $w_{f,\varepsilon}(i, j) = (j - i) \lceil \log(2\varepsilon + 1) \rceil + \kappa_f$ , where  $\kappa_f$  is the space in bits taken by the parameters of  $f$  (plus some small metadata, such as the function kind, encoded as an index from  $\{1, \dots, |\mathcal{F}|\}$ ). Moreover, since  $f$  is also an  $\varepsilon$ -approximation of any prefix and suffix of  $T[i, j - 1]$ , other than the edge  $(i, j)$  we add to  $\mathcal{G}$  also the prefix edge  $(i, k)$  and the suffix edge  $(k, j)$ , for all  $k = i, \dots, j - 1$  [34]. It is not difficult to conclude that the shortest path from node 1 to node  $n + 1$  gives the desired partition of  $T$ .

It is well-known that, in the case of a directed acyclic graph (like  $\mathcal{G}$ ), the shortest path can be computed by taking the nodes in order  $1, \dots, n$  and relaxing their outgoing edges, i.e. checking whether these edges can improve the shortest path found so far [42]. Furthermore, generalising what has been done in [34], instead of precomputing all the  $|\mathcal{F}| \cdot |\mathcal{E}|$  different piecewise  $\varepsilon$ -approximations, we only keep track of the  $|\mathcal{F}| \cdot |\mathcal{E}|$  edges of the form  $(i, j)$  that overlap the currently visited node  $k$ , i.e.  $i \leq k < j$ , and split them on-the-fly into prefix and suffix edges of the forms  $(i, k)$  and  $(k, j)$ , respectively.

Algorithm 1 formalises this description. We use  $distance[k]$  to store an upper bound on the cost of the shortest path from node 1 to  $k$ , and  $previous[k]$  to store the previous node and corresponding fragment in the shortest path. We use  $J_{f,\varepsilon}$  to keep track of the start/end positions of the fragment overlapping  $k$  and the parameters of the corresponding function of the  $f$ -kind that  $\varepsilon$ -approximates it. We initialise and update  $J_{f,\varepsilon}$  in Line 10 with a call to `MAKEAPPROXIMATION( $T, k, f, \varepsilon$ )`, which runs the algorithm of Theorem 1 starting from the data point  $T[k]$ . Lines 12–15 and 17–20 relax prefix and suffix edges, respectively, and Lines 21–26 conclude the algorithm by reading and returning the shortest path.

*Complexity analysis:* We now discuss the time complexity of Algorithm 1. For a fixed  $f$  and  $\varepsilon$ , the overall contribution of Line 10 to the time complexity is  $O(n)$ , since it eventually computes via Theorem 1 the piecewise  $\varepsilon$ -approximation of  $T$  composed of a function of the  $f$ -kind. Since there are  $|\mathcal{F}| \cdot |\mathcal{E}|$  possible pairs of  $f$  and  $\varepsilon$ , the overall computation of piecewise approximations takes  $O(|\mathcal{F}| |\mathcal{E}| n)$  time. It is easy to see that the relaxation of all the prefix and suffix edges runs within that same asymptotic time bound, thus the overall time complexity of Algorithm 1 is  $O(|\mathcal{F}| |\mathcal{E}| n)$ .

Concerning  $|\mathcal{F}|$ , we can assume that a real-world time series can be approximated well by a fixed number of function kinds, such as those in Table I, and thus it holds  $|\mathcal{F}| = O(1)$ . Concerning  $|\mathcal{E}|$ , instead, we can pessimistically bound it as follows. Let  $\Delta$  be one plus the difference between the maximum value and the minimum value  $\hat{y}$  in  $T$ . Then, each value  $y_k$  of  $T$  can be stored in  $\lceil \log \Delta \rceil$  bits by just encoding the binary representation of  $y_k - \hat{y}$ . This, in turn, entails that we can restrict our attention to the set  $\mathcal{E} = \{0, 2^1, \dots, 2^{\lceil \log \Delta \rceil}\}$ , since higher values of  $\varepsilon$  would not pay off, i.e. even the most trivial constant function can  $\varepsilon$ -approximate the whole time series. Given that such a set  $\mathcal{E}$  has size  $O(\log \Delta)$ , the time complexity of Algorithm 1 under these conditions is  $O(n \log \Delta)$ .

The average value of  $\log \Delta$  for the diverse dataset we use

---

**Algorithm 1** Partitioning a time series with NeaTS.

---

**In:** Time series  $T[1, n]$ , set  $\mathcal{F}$  of functions, set  $\mathcal{E}$  of error bounds

**Out:** A partitioning of  $T$  into fragments, each associated with an  $\varepsilon$ -approximation  $f$  (with  $\varepsilon \in \mathcal{E}$ ,  $f \in \mathcal{F}$ ), that minimises the size of the NeaTS encoding of  $T$

```

1:  $distance[1, n + 1] \leftarrow [\infty, \dots, \infty]$ 
2:  $previous[1, n + 1] \leftarrow [NULL, \dots, NULL]$ 
3: for all  $(f, \varepsilon) \in \mathcal{F} \times \mathcal{E}$  do  $\triangleright$  Initialise edges and functions
4:    $J_{f,\varepsilon}.start \leftarrow -\infty$ 
5:    $J_{f,\varepsilon}.end \leftarrow -\infty$ 
6:    $J_{f,\varepsilon}.params \leftarrow NULL$ 
7: for  $k \leftarrow 1$  to  $n$  do
8:   for all  $(f, \varepsilon) \in \mathcal{F} \times \mathcal{E}$  do
9:     if  $J_{f,\varepsilon}.end \leq k$  then  $\triangleright$  A new edge overlaps  $k$ 
10:       $J_{f,\varepsilon} \leftarrow \text{MAKEAPPROXIMATION}(T, k, f, \varepsilon)$ 
11:     else
12:        $i \leftarrow J_{f,\varepsilon}.start$   $\triangleright$  Relax prefix edge  $(i, k)$ 
13:       if  $distance[k] > distance[i] + w_{f,\varepsilon}(i, k)$  then
14:          $distance[k] \leftarrow distance[i] + w_{f,\varepsilon}(i, k)$ 
15:          $previous[k] \leftarrow (i, J_{f,\varepsilon})$ 
16:       for all  $(f, \varepsilon) \in \mathcal{F} \times \mathcal{E}$  do
17:          $j \leftarrow J_{f,\varepsilon}.end$   $\triangleright$  Relax suffix edge  $(k, j)$ 
18:         if  $distance[j] > distance[k] + w_{f,\varepsilon}(k, j)$  then
19:            $distance[j] \leftarrow distance[k] + w_{f,\varepsilon}(k, j)$ 
20:            $previous[j] \leftarrow (k, J_{f,\varepsilon})$ 
21:  $result \leftarrow$  an empty dynamic array
22:  $k \leftarrow n + 1$ 
23: while  $k \neq 1$  do  $\triangleright$  Read the shortest path backwards
24:    $result.PUSHFRONT(previous[k])$ 
25:    $k \leftarrow$  the first element of  $previous[k]$ 
26: return  $result$ 

```

---

in Section IV is 28.8, which is a small constant. Moreover, we do not actually need to use all  $\mathcal{F} \times \mathcal{E}$  pairs in Algorithm 1 but rather those surviving a model-selection procedure. For instance, we can initially run Algorithm 1 on a small sample of  $T$  (chosen e.g. according to the seasonality of  $T$ ) and select just the pairs that are used in the result, as these are likely to be effective and enough for the whole time series too. Our experiments show that this model-selection procedure improves the compression speed by an order of magnitude, with little impact on the compression ratio (see Section IV-C1).

*Partitioning for lossy compression:* We can easily modify Algorithm 1 to obtain a lossy representation of the time series  $T$  with a given  $\varepsilon$ -bound on the error, still using functions from a given set  $\mathcal{F}$  and minimising the space, which is given this time by just the storage of the functions' parameters (since we drop the corrections). It is enough to set  $\mathcal{E} = \{\varepsilon\}$  and define the edge weight  $w_f(i, j)$  to be equal to the space in bits taken by the parameters of  $f$ . The resulting algorithm runs in  $O(|\mathcal{F}| n)$ , so in linear time if  $|\mathcal{F}| = O(1)$ .

Our experiments will show that, for a fixed  $\varepsilon$ -bound, this algorithm produces more succinct lossy representations of time

series than known algorithms based on linear or nonlinear functions (namely, the algorithm by O’Rourke [36] and the Adaptive Approximation algorithm [30], [31], respectively).

### C. Designing the NeaTS compressor

We now describe the layout of the compressed time series and how to support the random access operation. As common in the literature [8], [14]–[17], we focus on the storage of the values  $y_1, \dots, y_n$  and assume the timestamps are  $1, \dots, n$ .<sup>5</sup>

Let us assume that the output of Algorithm 1 is a sequence of  $m$  tuples having the form  $\langle f_i, params_i, \varepsilon_i, start_i, end_i \rangle$ , where each tuple indicates a fragment  $T[start_i, end_i]$  of  $T[1, n]$  that is  $\varepsilon_i$ -approximated by a function of kind  $f_i \in \mathcal{F}$  with parameters  $params_i$ . We encode these  $m$  tuples and the values in their corresponding time series fragments via:

- An integer array  $S[1, m]$  storing in  $S[i]$  the starting position of the  $i$ th fragment, i.e.  $S[i] = start_i$ . To obtain the index of the fragment that covers a certain data point  $T[k]$ , we use the  $S.rank(k)$  operation, which returns the number of elements in  $S$  that are smaller than or equal to  $k$ . Since  $S$  is an increasing integer sequence, we compress it via the Elias-Fano encoding [45], [46], which supports accessing an element in  $O(1)$  time and  $S.rank$  in  $O(\min(\log m, \log \frac{n}{m}))$  time [47].
- An integer array  $B[1, m]$  storing in  $B[i]$  the bit size of the corrections of the  $i$ th fragment, i.e.  $B[i] = \lceil \log(2\varepsilon_i + 1) \rceil$ .
- An integer array  $O[1, m + 1]$  storing in  $O[i]$  the cumulative bit size of the corrections in the fragments preceding the  $i$ th one, i.e.  $O[i] = \sum_{j=1}^{i-1} B[j] (end_j - start_j + 1)$ . Notice that  $O[m + 1]$  denotes the overall bit size of the corrections. Similarly to  $S$ , we compress  $O$  via the Elias-Fano encoding.
- A bit string  $C[1, O[m + 1]]$  storing in  $C[O[i], O[i + 1] - 1]$  the correction values  $y_j - \lfloor f_i(x_j) \rfloor$  of the  $i$ th fragment, where  $j \in [start_i, end_i]$ .
- An integer array  $K[1, m]$  storing in  $K[i]$  the function kind for the  $i$ th fragment, i.e.  $K[i] = f_i$ . We regard  $K$  as a string over the alphabet  $\{1, \dots, |\mathcal{F}|\}$  and represent it as a wavelet tree data structure [47], [48]. This allows us to compute the  $K.rank_f(i)$  operation, which returns the number of occurrences of the function kind  $f$  in  $K[1, i]$  in  $O(\log |\mathcal{F}|) = O(1)$  time.
- For each  $f \in \mathcal{F}$ , an array  $P_f$  concatenating the parameters  $params_i$  of the functions of the same kind  $f$ . This way, the parameters  $params_i$  of the  $i$ th fragment can be found in  $P_{f_i}[K.rank_{f_i}(i)]$ .

All the above arrays use cells whose bit size is just enough to contain the largest value stored in them. If  $\mathcal{F}$  contains functions with the same number of parameters (recall from Section III-A that we can use functions with more than two parameters), we can simplify the above encoding by avoiding the use of a wavelet tree for  $K$  and by concatenating all the

<sup>5</sup>The timestamps  $x_1, \dots, x_n$  form an increasing sequence of integers that can be easily mapped to  $1, \dots, n$  via monotone minimal perfect hash functions [43] or compressed rank data structures [34], [44]: the former are very succinct (about 3 bits per integer), the latter take more space but enable range queries over timestamps.

---

### Algorithm 2 Full decompression in NeaTS.

---

**In:** The NeaTS encoding  $\langle S, B, O, C, K, P \rangle$  of  $T$

**Out:** The uncompressed values of  $T$

```

1:  $o \leftarrow 1$  ▷ Bit-offset to the correction
2: for  $i \leftarrow 1$  to  $m$  do ▷ For each fragment
3:    $start \leftarrow S[i]$  ▷ First data point index
4:    $end \leftarrow S[i + 1] - 1$  ▷ Last data point index
5:    $f \leftarrow K[i]$  ▷ Function kind
6:    $params \leftarrow P_f[K.rank_f(i)]$  ▷ Function parameters
7:    $b \leftarrow B[i]$  ▷ Correction bit size
8:   for  $k \leftarrow start$  to  $end$  do
9:      $\tilde{y} \leftarrow \text{compute } \lfloor f(k) \rfloor$  using  $params$ 
10:    output  $\tilde{y} + \text{int}(C[o, o + b - 1])$ 
11:     $o \leftarrow o + b$ 

```

---



---

### Algorithm 3 Random access in NeaTS.

---

**In:** An index  $k$ , the NeaTS encoding  $\langle S, B, O, C, K, P \rangle$  of  $T$

**Out:** The value of  $T[k]$

```

1:  $i \leftarrow S.rank(k)$  ▷ Index of the fragment
2:  $start \leftarrow S[i]$  ▷ First data point index
3:  $f \leftarrow K[i]$  ▷ Function kind
4:  $params \leftarrow P_f[K.rank_f(i)]$  ▷ Function parameters
5:  $b \leftarrow B[i]$  ▷ Correction bit size
6:  $\tilde{y} \leftarrow \text{compute } \lfloor f(k) \rfloor$  using  $params$ 
7:  $o \leftarrow O[i] + (k - start)b$  ▷ Bit-offset to the correction
8: return  $\tilde{y} + \text{int}(C[o, o + b - 1])$ 

```

---

functions’ parameters  $P_f$  into a single array, which is accessed simply through the index of the queried fragment.

Having defined how we represent  $T$  in compressed form via a tuple  $\langle S, B, O, C, K, P \rangle$  of data structures, we are now ready to discuss the decompression and random access operations.

Algorithm 2 shows how to decompress the whole time series. For each fragment, we first decode the associated boundaries and kind of approximation function (Lines 3–7), and then we output all the values  $y_k$  within the fragment’s boundaries by applying the function to index  $k$  and adding the corresponding correction value (Lines 8–11). It is easy to see that the time complexity of Algorithm 2 is  $O(n)$  given that all the involved operations take constant time. Furthermore, since each data point is decompressed independently from the others, the algorithm could be parallelised trivially by decompressing different fragments with different workers, and the computation of the function within a fragment could be implemented via SIMD instructions.

Algorithm 3 shows how to perform the random access operation to the value of  $T[k]$ , for a given index  $k \in \{1, 2, \dots, n\}$ . We start by identifying the index of the fragment where  $T[k]$  falls into (Line 1), then we decode the index of the first data point in that fragment and the (kind and parameters of) function associated with that fragment (Lines 2–5), and finally we apply the function to position  $k$  and add the corresponding correction value (Lines 6–8). The time complexity of Algorithm 3 is dominated by the operation  $S.rank$  at

Line 1, which takes  $O(\min(\log m, \log \frac{n}{m}))$  time. We can easily achieve  $O(1)$  time by representing  $S$  as a bitvector of length  $n$  with a 1 in each position  $start_i$ , and then using the well-known constant-time rank/select operations [49], [50].

#### IV. EXPERIMENTS

##### A. Experimental setting

We run our experiments on a machine with 1.17 TiB of RAM and an Intel Xeon Gold 6140M CPU, running CentOS 7. Our code is in C++23, compiled with GCC 13.2.1, and publicly available at <https://github.com/and-gue/NeaTS>. We refer to our lossy and lossless approaches as NeaTS-L and NeaTS, respectively. We use four types of functions—namely, linear, exponential, quadratic, and radical—which turned out to be sufficient to capture the trends in our real-world datasets well. We use vector instructions in our decompression procedures via the `std::experimental::simd` library, and succinct data structures from the `sdsl` [51] and `sux` [52] libraries.

1) *Datasets*: We use 16 real-world time series datasets out of which 13 were sourced by Chimp [15], and the remaining 3 were obtained from the Geolife project [53] and a study on arrhythmia [54], [55]. Consistent with previous studies [14]–[17], we ignore timestamps as they are either consecutive increasing integers or can be transformed into such with other ad hoc data structures (see Footnote 5). All the datasets report values in textual fixed-precision format, therefore, unless the compressor is designed for doubles, we transform them into 64-bit integers by multiplying each value by a factor  $10^x$ , where  $x$  is the number of fractional digits.

- **IR-bio-temp (IT)** [56] contains about 477M biological temperature observations from an infrared sensor, with 2 fractional digits.
- **Stocks-USA (US)**, **Stocks-UK (UK)**, and **Stocks-DE (GE)** [57] contain about 282M, 59M, and 43M stock exchange prices of USA, UK, and Germany, with 2, 1, and 3 fractional digits, respectively.
- **Electrocardiogram (ECG)** [54], [55] contains about 226M electrocardiogram signals of over 45K patients, with 3 fractional digits.
- **Wind-direction (WD)** [58] contains about 199M wind direction observations, with 2 fractional digits.
- **Air-pressure (AP)** [59] contains about 138M timestamped values of barometric pressure corrected to sea level and surface level, with 5 fractional digits.
- **Geolife-longitude (LON)**, and **Geolife-latitude (LAT)** [53] contain about 25M timestamped longitude and latitude values of 182 users’ GPS trajectories, with 4 fractional digits.
- **Dewpoint-temp (DP)** [60] contains about 5M relative dew point temperature observations, with 3 fractional digits.
- **City-temp (CT)** [61] contains about 3M temperature observations of cities around the world, with 1 fractional digit.
- **PM10-dust (DU)** [62] contains about 334K measurements of PM10 in the atmosphere, with 3 fractional digits.
- **Basel-wind (BW)**, and **Basel-temp (BT)** [63] contain about 130K records of wind speed and temperature data of Basel (Switzerland), with 7 and 9 fractional digits, respectively.

- **Bird-migration (BM)** [64] contains about 18K positions of birds, with 5 fractional digits.
- **Bitcoin-price (BP)** [64] contains about 7K prices of Bitcoin in the dollar exchange rate, with 4 fractional digits.

2) *Competitors*: Regarding the lossy compressors, we compare our NeaTS-L against 2 functional approximation algorithms: the optimal Piecewise Linear Approximation algorithm (PLA) [36], and the Adaptive Approximation algorithm (AA) [30], [31] that combines linear, exponential and quadratic functions. We implemented the PLA and the AA algorithms in C++ since their code is not publicly available.

Regarding the lossless compressors, we compare our NeaTS against 5 widely-used general-purpose compressors—namely, Xz [11], Brotli [9], Zstd [10], Lz4 [12], and Snappy [13]—and 7 state-of-the-art special-purpose compressors—namely, Chimp and Chimp128 [15], TSXor [65], DAC [66], Gorilla [14], LeCo [35], and ALP [17]. We use the Squash library [67] for all the general-purpose compressors. We use the public implementations of TSXor and Gorilla available in the repository of [65], the implementation of DAC available in `sdsl` [51], and the original implementations of LeCo and ALP. We ported Chimp and Chimp128 to C++ since their original implementations are in Java.

Following [15], [16], we apply compressors that do not natively support random access (thus excluding DAC, LeCo, and NeaTS) to blocks of 1000 consecutive values. We then maintain an array that maps each block index to a pointer referencing the starting byte of the block in the compressed output.

##### B. On the lossy compressors

To evaluate the lossy compressors with a meaningful error-bound parameter  $\varepsilon$ , we determined the smallest  $\varepsilon$  such that NeaTS-L achieves better compression than our lossless compressor NeaTS. We express the resulting  $\varepsilon$  as a % of the range of values (i.e. largest minus smallest value) in a dataset, and the compression ratio as the size of the compressed output divided by the size of the original data.

The results in Table II show that NeaTS-L outperforms in compression ratio both the PLA and the AA algorithms on all datasets. On average, NeaTS-L improves the compression ratio of PLA by 7.02% and the one of AA by 11.77%. This demonstrates that, under the same  $\varepsilon$ -bound, the use of nonlinear functions allows achieving better compression compared to linear functions alone (as in the widely-used PLA). In turn, despite employing nonlinear approximations, AA is worse than PLA for nearly all datasets due to the use of a heuristic technique to partition the time series into fragments and of a sub-optimal algorithm for (non)linear  $\varepsilon$ -approximations, two issues that our NeaTS-L solve.

For the approximation accuracy, we report that the Mean Absolute Percentage Error (MAPE)—i.e. the mean of the absolute relative errors between the approximated and the actual values, expressed as a percentage—is 2.47% for AA, 2.85% for NeaTS-L, and 4.37% for PLA (on average over all datasets). Therefore, NeaTS-L has a much better accuracy than PLA and a slightly worse accuracy than AA. This is because



**TABLE II:** Compression ratios of the 3 experimented lossy approaches —i.e. AA, PLA, and NeaTS-L— on the 16 datasets.

Dataset	$\varepsilon$ (%)	Compression ratio (%)			NeaTS-L improv. (%)	
		AA	PLA	NeaTS-L	wrt AA	wrt PLA
IT	1.15E-1	12.11	12.07	<b>11.07</b>	8.57	8.29
US	2.40E-3	7.96	7.41	<b>6.99</b>	12.09	5.65
ECG	5.43E-2	15.03	13.46	<b>12.97</b>	13.71	3.64
WD	6.36E-0	28.09	26.94	<b>24.76</b>	11.88	8.11
AP	3.08E-3	21.90	20.00	<b>19.17</b>	12.49	4.16
UK	9.53E-3	9.82	9.21	<b>8.69</b>	11.50	5.63
GE	9.12E-3	13.95	12.79	<b>12.08</b>	13.35	5.52
LAT	7.00E-6	25.40	23.59	<b>22.09</b>	13.03	6.35
LON	1.40E-5	19.92	18.32	<b>17.26</b>	13.37	5.78
DP	6.32E-2	17.51	16.89	<b>15.87</b>	9.35	6.07
CT	3.88E0	16.19	14.45	<b>13.92</b>	14.03	3.69
DU	6.00E-3	10.04	10.32	<b>9.15</b>	8.93	11.39
BT	4.85E-1	59.62	61.29	<b>53.77</b>	9.81	12.26
BW	3.16E-3	52.19	48.28	<b>45.01</b>	13.75	6.77
BM	1.42E-2	27.13	25.32	<b>23.29</b>	14.15	8.00
BP	3.61E-1	43.05	41.76	<b>38.52</b>	10.54	7.76

AA creates more time series fragments than NeaTS-L, and its functions pass through the first data point of each fragment: two factors that together yield zero errors on many data points.

In terms of compression speed, PLA is the fastest at 123.36 MB/s, followed by AA at 63.11 MB/s, and NeaTS-L at 18.23 MB/s. These results reflect the higher computational effort of NeaTS-L in achieving better compression ratios.

In terms of decompression speed, PLA is the fastest at 2997.00 MB/s, followed by NeaTS at 2561.31 MB/s, and AA at 2420.20 MB/s. This can be attributed to the fact that the linear models in PLA are faster to evaluate, and that NeaTS uses fewer fragments than AA, thus reducing the overhead associated with switching between fragments.

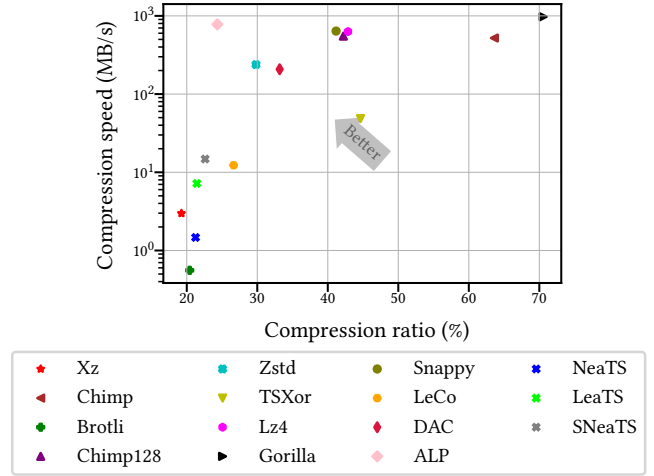
### C. On the lossless compressors

We now compare our NeaTS against the 5 lossless general-purpose compressors (i.e. Xz, Brotli, Zstd, Lz4, and Snappy) and the 7 lossless special-purpose compressors (i.e. Chimp, Chimp128, TSXor, DAC, Gorilla, LeCo, and ALP).

Table III reports the compression ratio, decompression speed, and random access speed of all compressors on each dataset, where the best result in each family of compressors is in bold, and the best result overall is underlined. Moreover, we plot the trade-offs compression ratio vs compression speed, compression ratio vs decompression speed, and compression ratio vs random access speed in Figures 2 to 3 and dig into them in Sections IV-C1 to IV-C3, respectively. Then, in Section IV-C4, we explore the benefits of our approach from a data management perspective by focusing on the important case of range queries. Finally, we provide a summary of our experiments in Section IV-C5.

1) *Compression ratio vs compression speed:* Figure 2 shows the trade-off between compression speed and compression ratio of the lossless compressors.

First, we notice that Xz, followed by Brotli, achieve on average (but not always, see below) the best compression ratio at the cost of a slow compression speed, indeed, they are at the bottom-left of Figure 2. The opposite extreme in this trade-off



**Fig. 2:** The trade-off between compression ratio and speed of the lossless compressors, averaged on the 16 datasets.

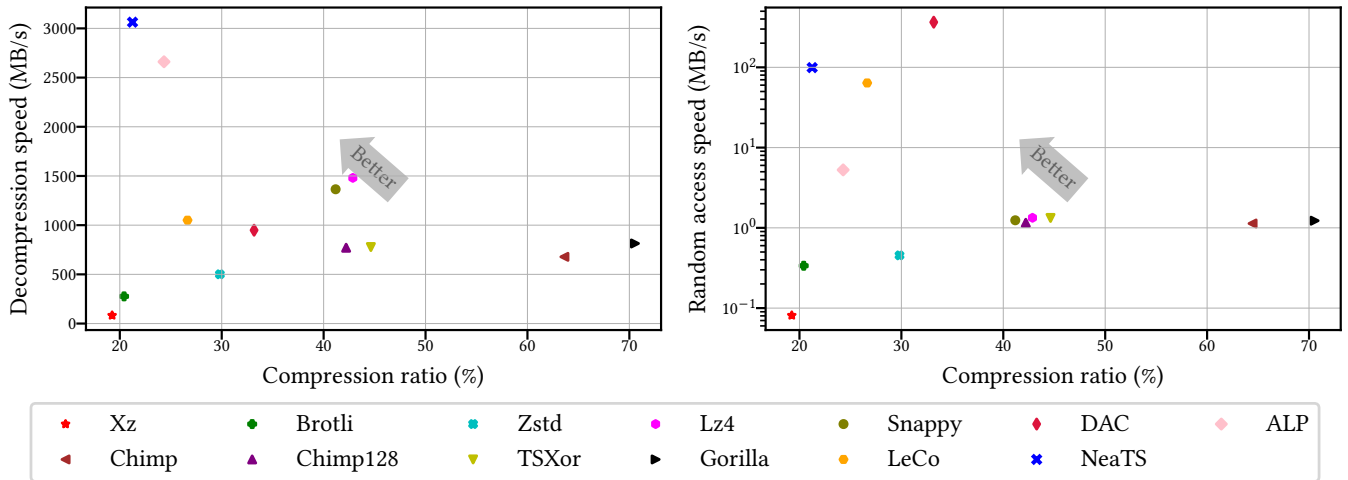
is occupied by the special-purpose compressor Gorilla (at the top-right of Figure 2), which is 3 orders of magnitude faster in compression speed than Brotli but achieves a compression ratio above 70% on average. In between these two extremes, we notice that ALP is on the Pareto front of this trade-off, dominating (in order of increasingly higher compression speeds and worse compression ratios) LeCo, TSXor, DAC, Zstd, Chimp, Chimp128, Lz4 and Snappy, but not our NeaTS.

Indeed, we observe from Table III that NeaTS achieves the best compression ratio among the special-purpose compressors on 14/16 datasets, and the best compression ratio overall on 4/16 datasets. Its compression speed is low but still 164.14% faster than Brotli with just a 4.09% worse compression ratio on average. Moreover, NeaTS always achieves better compression ratios than Lz4 and Snappy (the fastest general-purpose compressors in terms of compression speed) by 52.77% and 50.20% on average, respectively. NeaTS also achieves better compression ratios than Zstd for almost all the datasets (except for AP), with an average improvement of 28.49%.

Compared to the special-purpose compressors, the only 2/16 datasets in which NeaTS does not achieve the best compression ratio are BP and BT, where ALP is slightly better. However, these two datasets are also among the smallest ones, and NeaTS achieves a better compression ratio than ALP by 16.36% on average.

We conclude this section by experimenting with two variants of NeaTS that improve the compression speed at the cost of worse compression. The first, named LeaTS, reduces the set of functions considered by Algorithm 1 to linear functions only. The second, named SNeaTS, reduces the set of functions and error bounds considered by Algorithm 1 to those surviving a model-selection procedure (included in the construction time) that picks the top-5 most-used pairs in the first 10% of the dataset. The results, depicted in Figure 2, show that LeaTS and SNeaTS achieve a compression speed that is  $5.22\times$  and  $12.86\times$  that of NeaTS, and a compression ratio that is 0.89%





**Fig. 3:** The trade-off between compression ratio and decompression speed (left plot), and between compression ratio and random access speed (right plot) of the lossless compressors, averaged on the 16 datasets. Note that the vertical axis of the right plot is logarithmic.

and 8.18% worse than NeaTS, respectively.<sup>6</sup> The latter variant, in particular, is both faster in compression speed than LeCo by 36.26% and better in compression ratio by 12.80% on average.

Despite their better compression ratios, these variants are still not as fast as ALP or Gorilla in compression speed. However, we anticipate from the next subsections that NeaTS also excels in decompression and random access speed, thus making it the most competitive compressor in a query-intensive and space-constrained scenario. Moreover, if compression speed is key for the underlying application, we could imagine using a lightweight compressor like ALP or Gorilla when the time series is first ingested, and running NeaTS later on (or in the background) to provide much more effective compression and efficient query operations in the long run.

2) *Compression ratio vs decompression speed:* In an analytical scenario, the decompression speed is a crucial performance metric. The middle of Table III shows the decompression speed of all the compressors on all the datasets, while Figure 3 shows the trade-off between compression ratio and decompression speed averaged on all the datasets.

First, we notice from Table III that NeaTS achieves the fastest decompression speed on 10/16 datasets thanks to its cache-friendly and vectorised decompression procedure. Compared to ALP, which obtains better performance on the remaining 6/16 datasets, NeaTS is both 16.36% better in compression ratio and 27.08% faster in decompression speed on average. If we consider instead Xz and Brotli (i.e. the closest competitors to NeaTS in terms of compression ratio, as commented above), their decompression speeds are  $44.92\times$  and  $12.27\times$  lower on average than that of NeaTS, respectively.

Finally, NeaTS dominates all the other special-purpose compressors in this compression ratio vs decompression

<sup>6</sup>In the lossless scenario, the space is clearly dominated by the storage of the corrections rather than the function parameters, which is why the improvement in compression ratio of nonlinear functions over linear ones is not as high as 12%, as experienced in the lossy scenario (Table II).

speed trade-off, as Figure 3 clearly shows by placing them at the bottom-right of NeaTS in the decompression plot. For instance, compared to LeCo, NeaTS is 18.23% better in compression ratio and 201.11% faster in decompression speed.

3) *Compression ratio vs random access speed:* Another key performance metric for the efficient analysis of time series is the random access speed. The bottom of Table III shows the average random access speed (for 10M queries) of all the compressors on all the datasets, while Figure 3 shows the trade-off between compression ratio and random access speed averaged on all the datasets.

First, we notice from Table III that DAC, followed by NeaTS, achieves the best random access speed. However, we remark that NeaTS is much more effective than DAC in terms of compression ratio, i.e. 37.25% better on average and up to 67.86% better overall. This is why both NeaTS and DAC occupy a prominent position in the compression ratio vs random access speed trade-off, as Figure 3 clearly shows by placing them close to the top-left edge of the random access plot.

LeCo is the only other compressor supporting random access natively (i.e. without the block-wise approach described in Section IV-A2). NeaTS is both 118.55% faster in random access speed and 18.23% better in compression ratio than LeCo.

The remaining compressors are from 2 to 3 orders of magnitude slower in random access speed than NeaTS. In particular, NeaTS dominates all other special-purpose compressors and the vast majority of general-purpose compressors in the compression ratio vs random access speed trade-off. The only exceptions are Xz and Brotli that, compared to NeaTS, can provide slightly better compression ratios (except for the 4/16 dataset where NeaTS is better) but much slower random access speeds (always), which is why they are at the bottom-left of every plot.

4) *Range queries:* The most fundamental queries in time series databases — such as trend analysis, anomaly detection, correlation analysis, and data aggregation — ultimately rely on

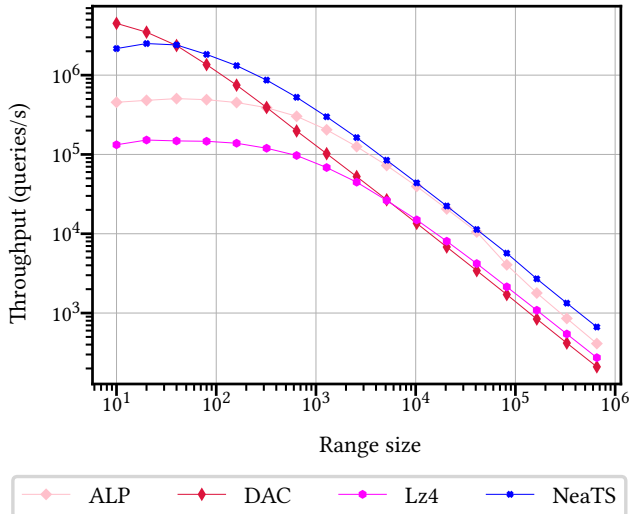


Fig. 4: Range queries throughput across different range sizes.

accessing data within a specific time interval (i.e. a range query) [21], [22], which boils down to a random access operation (to retrieve the first data point) followed by a scan (to retrieve the subsequent data points within the interval). We thus now focus on the best compressors in terms of random access or decompression speed (i.e. ALP, DAC, Lz4, and NeaTS) and evaluate their range query performance at different range sizes, from  $10 \cdot 2^0$  to  $10 \cdot 2^{16}$  data points. Figure 4 shows the throughput in queries per second measured on 10K random range queries and averaged over the 11 largest datasets. For range sizes smaller than 40, DAC is the fastest solution, followed by NeaTS, which remains an order of magnitude more efficient than other compressors. For larger range sizes, NeaTS clearly outperforms all competitors. This demonstrates the ability of NeaTS to provide a full spectrum of efficient data access, from small to large ranges, thus benefiting a wide variety of queries in time series databases.

5) *Summary*: Our results on 16 real-world time series show that NeaTS emerges as the only approach to date providing, simultaneously, compression ratios close to or better than the existing compressors (i.e. the best compression ratio among the special-purpose compressor on 14/16 datasets, and the best overall on 4/16 dataset), a much faster decompression speed, and up to 3 orders of magnitude more efficient random access.

No other compressor can strike such a good trade-off among all these factors together. For example, Xz achieves the best compression ratio on 9/16 datasets, but on average its decompression and random access speeds are  $44.92\times$  and  $1657.10\times$  slower than that of NeaTS, respectively. ALP achieves the fastest compression speed (ignoring Gorilla, whose compression ratio is not very competitive), but on average NeaTS achieves 16.36% better compression ratio, 27.08% faster decompression speed, and at least one order of magnitude faster random access speed than ALP. DAC achieves the fastest random access speed, but on average NeaTS is 37.25% better in

compression ratio and 234.89% faster in decompression speed than DAC. Furthermore, NeaTS outperforms the other compressors for range queries involving 40 or more data points.

This is evidence that NeaTS has the potential to be the compressor of choice for the storage and real-time analysis of massive and ever-growing amounts of time series data.

## V. RELATED WORK

We now review the literature of general- and special-purpose lossless compressors for time series, and of lossy compressors. For the latter, we focus on approaches based on error-bounded functional approximations, which are relevant to our work.

a) *General-purpose lossless compressors*: These compressors are not specifically designed for time series but can be applied to any byte sequence. We discuss below the ones based on the LZ77-parsing [68], which currently offer the best combination of compression ratio and (de)compression speed.

Brotli [9] relies on a modern variant of the LZ77-parsing of the input file that uses a pseudo-optimal entropy encoder based on second-order context modelling. Xz [11] achieves effective compression by using Markov chain modelling and range coding of the LZ77-parsing. Zstd [10] achieves very fast (de)compression speed and good compression ratios via a tabled asymmetric numeral systems encoding. Finally, Lz4 [12] and Snappy [13] trade compression effectiveness with speed by adopting a faster byte-oriented encoding format for the LZ77-parsing. Given this plethora of approaches offering a variety of trade-offs, we tested them all in our experiments.

b) *Special-purpose lossless compressors*: Most recent compressors for time series are often based on encoding the result of bitwise XOR operations between close or adjacent floating-point values. Their compression ratio is strongly influenced by data fluctuations: the more severe the fluctuations, the less effective the compression. On the other hand, these algorithms offer very fast (de)compression speeds.

For instance, Gorilla [14] improves earlier floating-point compressors [69]–[71] by simply computing the XOR between consecutive values of the time series and properly encoding the number of leading zeros and significant bits of the result.

Chimp [15] improves both the compression ratio and speed of Gorilla by using different encoding modes based on the number of trailing and leading zeros of the XOR result. Chimp128 [15] and TSXor [65] use a window of 128 values to choose the best reference value for the XOR computation: Chimp128 uses the value that produces the most trailing zeros, while TSXor selects the value with the most bits in common.

Elf [16] performs an erasing operation on the floats before XORing them, which makes the resulting value more compressible. We do not experiment with Elf because ALP [72] (described next and included in our experiments) was shown to achieve better compression ratios on average, and always faster compression and decompression speeds than Elf.

ALP [72], unlike the above approaches, does not use XOR operations but rather builds on the idea of encoding a double  $x$  via the storage of the significant digits  $d$  and an exponent  $e$ , i.e.  $d = \text{round}(x \cdot 10^e)$ , also known as Pseudodecimal

**TABLE III:** Compression ratio (top), decompression speed (middle), and random access speed (bottom) achieved by the 5 general-purpose and the 8 special-purpose lossless compressors (including our NeaTS) on 16 datasets, sorted by decreasing size. We highlight in bold the best result in each family, and in underline the best result overall.

Dataset	General-purpose compressors					Special-purpose compressors								
	Xz	Brotli	Zstd	Lz4	Snappy	Chimp128	Chimp	TSXor	DAC	Gorilla	LeCo	ALP	NeaTS	
Compression ratio (%)	IT	<b>12.86</b>	14.25	23.46	41.31	36.96	29.43	72.30	30.76	23.83	78.60	13.62	16.86	<b>11.88</b>
	US	9.18	<b>8.70</b>	12.82	27.09	21.51	18.94	54.55	18.89	24.95	57.54	9.16	10.50	<b>8.02</b>
	ECG	<b>12.12</b>	<b>12.12</b>	17.04	26.14	33.75	54.11	43.18	20.03	25.39	45.26	15.58	16.23	<b>12.96</b>
	WD	<b>23.60</b>	27.60	33.78	52.70	54.19	43.38	84.09	46.42	25.75	91.02	24.71	24.90	<b>24.37</b>
	AP	<b>12.35</b>	12.69	17.87	26.50	24.82	30.00	35.76	34.78	41.13	37.67	23.52	25.74	<b>19.27</b>
	UK	9.42	<b>9.06</b>	12.99	26.94	21.41	23.13	46.95	15.85	25.79	53.92	10.83	11.64	<b>9.09</b>
	GE	11.07	<b>11.04</b>	15.27	30.25	23.94	21.08	66.90	21.44	29.01	71.49	13.43	13.88	<b>12.11</b>
	LON	<b>17.03</b>	18.63	32.72	49.71	49.28	58.64	61.70	71.64	47.27	63.09	20.74	26.87	<b>17.53</b>
	LAT	<b>21.51</b>	23.67	40.77	52.12	51.44	58.09	61.44	71.93	47.27	65.02	25.56	26.70	<b>22.22</b>
	DP	<b>16.37</b>	17.02	29.35	48.61	47.54	49.53	77.17	60.91	26.95	83.53	17.83	22.04	<b>16.10</b>
	CT	<b>15.72</b>	16.37	25.33	42.92	37.31	36.09	73.25	30.96	19.14	87.11	17.91	15.27	<b>14.20</b>
	DU	8.21	<b>7.78</b>	11.37	23.00	18.62	21.68	39.74	18.31	11.14	44.49	28.54	13.34	<b>9.46</b>
	BT	<b>45.66</b>	45.69	58.12	67.20	68.64	46.90	84.01	53.88	57.07	92.88	58.15	<b>46.25</b>	54.01
	BW	<b>36.17</b>	41.49	50.24	58.74	58.79	71.27	87.16	82.32	45.91	99.72	56.99	50.01	<b>45.21</b>
	BM	<b>19.67</b>	20.70	29.52	43.58	39.39	40.96	61.96	48.18	37.42	74.67	50.72	30.80	<b>23.44</b>
	BP	<b>36.97</b>	39.85	66.43	69.03	71.22	72.09	67.84	87.86	42.79	82.72	39.03	<b>38.37</b>	39.89
Decompression speed (MB/s)	IT	90.83	304.65	459.91	<b>1405.36</b>	1207.61	725.74	598.36	743.69	999.45	795.22	1082.45	2249.26	<b>2549.04</b>
	US	133.37	396.11	643.89	1609.06	<b>1928.74</b>	1109.06	692.86	1084.44	896.80	839.86	1097.74	2295.05	<b>2982.09</b>
	ECG	94.39	253.72	512.45	1325.59	<b>1473.75</b>	559.59	645.71	773.99	1082.36	790.96	1306.75	2344.27	<b>2897.08</b>
	WD	56.94	227.86	490.46	1443.34	<b>1191.21</b>	732.56	606.28	754.89	864.86	854.76	1035.08	<b>2253.39</b>	1936.42
	AP	106.07	332.69	683.65	<b>1740.56</b>	1599.33	885.36	893.27	885.33	705.69	978.22	1013.16	2116.29	<b>2944.05</b>
	UK	131.95	392.03	634.25	1645.54	<b>1815.86</b>	853.25	670.05	1102.54	829.75	863.34	1087.50	2312.32	<b>3015.90</b>
	GE	115.76	350.47	594.39	<b>1611.32</b>	1761.48	1008.27	656.04	986.84	962.52	829.82	1062.78	2307.16	<b>3243.73</b>
	LAT	53.10	171.93	376.31	<b>1289.07</b>	1033.77	758.29	785.61	492.87	1019.50	622.10	960.91	2144.64	<b>2935.03</b>
	LON	61.55	212.49	375.82	<b>1198.19</b>	963.81	776.44	801.64	493.81	1025.93	625.60	925.42	2113.61	<b>2870.15</b>
	DP	75.14	293.10	429.39	<b>1456.28</b>	1098.91	581.50	626.82	618.88	1056.31	782.95	963.67	<b>2057.49</b>	1953.35
	CT	79.93	298.52	458.10	<b>1467.29</b>	1225.00	760.50	556.01	768.86	1070.76	841.18	969.50	<b>3290.46</b>	3269.65
	DU	147.58	436.39	696.34	1701.39	<b>2031.31</b>	923.52	805.54	1093.90	658.38	948.68	746.33	<b>5410.27</b>	4007.91
	BT	33.43	173.70	414.01	<b>1418.67</b>	1039.45	662.70	590.54	672.81	783.96	805.45	846.57	<b>4241.33</b>	2570.33
	BW	41.74	153.77	378.60	<b>1527.78</b>	1020.14	607.71	611.78	655.01	1189.51	849.33	1501.59	<b>4006.59</b>	3675.40
	BM	69.31	235.46	494.08	<b>1450.65</b>	1283.83	758.16	626.94	700.60	802.93	790.78	1441.36	1677.20	<b>4508.69</b>
	BP	36.85	187.29	366.83	<b>1400.66</b>	1163.14	624.77	682.54	612.99	1231.05	803.20	768.30	1752.34	<b>3649.14</b>
Random access speed (MB/s)	IT	0.09	0.36	0.39	<b>1.21</b>	1.04	1.05	0.92	1.28	<b>137.93</b>	1.09	19.94	5.43	48.19
	US	0.11	0.45	0.52	1.24	<b>1.57</b>	1.55	1.10	1.76	<b>108.11</b>	1.22	22.15	6.36	57.14
	ECG	0.11	0.36	0.49	1.18	<b>1.46</b>	1.16	1.32	1.29	<b>153.85</b>	1.50	30.08	3.78	50.00
	WD	0.05	0.26	0.42	<b>1.23</b>	1.01	1.12	1.02	1.21	<b>135.59</b>	1.28	20.05	5.19	47.06
	AP	0.10	0.39	0.57	<b>1.46</b>	1.30	1.29	1.60	1.50	<b>78.43</b>	1.50	19.02	5.31	53.69
	UK	0.12	0.50	0.56	1.46	<b>1.63</b>	1.20	1.05	2.01	<b>131.15</b>	1.26	25.94	6.26	76.92
	GE	0.12	0.42	0.51	1.38	1.46	1.42	1.02	1.65	<b>186.05</b>	1.18	30.75	4.23	80.00
	LAT	0.06	0.25	0.37	<b>1.22</b>	1.01	1.14	1.24	0.98	<b>210.53</b>	0.98	22.15	5.19	76.19
	LON	0.07	0.32	0.38	<b>1.22</b>	0.99	1.16	1.26	0.97	<b>210.53</b>	0.99	21.07	5.12	80.01
	DP	0.07	0.33	0.37	<b>1.31</b>	0.97	0.82	0.93	0.94	<b>571.43</b>	1.05	53.48	5.56	123.08
	CI	0.07	0.34	0.41	<b>1.28</b>	1.10	1.15	0.90	1.19	<b>666.67</b>	1.16	57.22	6.04	140.35
	DU	0.15	0.58	0.66	1.62	<b>1.95</b>	1.38	1.38	1.92	<b>363.64</b>	1.50	96.39	7.87	142.46
	BT	0.03	0.18	0.39	<b>1.31</b>	0.97	1.05	1.06	1.15	<b>666.67</b>	1.23	145.45	4.49	140.35
	BW	0.04	0.17	0.37	<b>1.46</b>	0.99	0.93	1.08	1.05	<b>800.00</b>	1.23	131.15	4.38	148.15
	BM	0.07	0.27	0.47	<b>1.40</b>	1.25	1.05	1.06	1.18	<b>533.33</b>	1.13	145.45	4.61	160.00
	BP	0.04	0.22	0.38	<b>1.37</b>	1.17	1.06	1.17	1.18	<b>888.89</b>	1.35	181.82	4.55	163.26

Encoding [73]. It finds a single best exponent for a block of 1024 values and bit-packs the resulting significant digits via the frame-of-reference integer code. Values failing to be losslessly encoded as pseudodecimals are stored uncompressed separately. Further optimisations (such as cutting trailing zeros and using vector instructions) are applied to improve the compression ratio and speed.

BUFF [74] compresses a float by eliminating the less significant bits based on a given precision, splitting it into the integer and fractional parts, and then compressing the two parts separately with a fixed-length encoding. We do not experiment with BUFF because its average compression ratio on time series was shown to be worse than that of Chimp (which,

in turn, is always worse than NeaTS in our experiments) and its compression and decompression speeds were shown to be no more than  $6\times$  that of Chimp [20] (which, in turn, are outmatched by those of ALP by one order of magnitude [17]).

Sprintz [75] encodes time series using four components: forecasting, bit packing, run-length encoding, and entropy coding. Sprintz focuses on 8- or 16-bit integers, which is a limitation for our datasets with 64-bit data. Also, it was shown to be worse than BUFF in compression ratio and speed [74].

There are also floating-point compressors targeted to scientific simulation and observational data, such as fpzip [76] and ndzip [77], but we do not experiment with them since their compression ratios were shown to be poor on time series [20].

Concerning the random access operation to time-series values, this is not directly offered by most compressor implementations. Therefore, the typical approach is to compress blocks of the time series separately, and then access a single value by decompressing just the corresponding block. This is often insufficient to guarantee a reasonable speed and use of computational resources (as our experiments confirm), which is why Brisaboa et al. [66] introduced the Directly Addressable Codes (DAC) scheme that enables fast access to individual values. Given this feature, although designed for generic integer sequences, we included DAC in our experiments.

DACTS [19] uses Re-Pair [78] on top of DAC to better capture repeating patterns. This additional compression step is effective when the time series is highly repetitive but slows down the access time, so DACTS proved to be useful on the so-called industrial time series originating from sensors producing long sequences of constant values. This is a restrictive situation, so we did not experiment with DACTS.

Titchy [18] focuses on random access to IoT data and relies on a dictionary-based approach combined with a partitioning of the time series into chunks. Each chunk is represented with a pair  $\langle \textit{base}, \textit{deviation} \rangle$ , where *base* indicates the item of the dictionary to copy, and *deviation* encodes what makes the chunk slightly different from the other. We could not experiment with Titchy because its source code is not available.

Finally, LeCo [35] is a recent proposal (not specific for time series) that, similarly to our NeaTS and earlier work [33], [34], is based on the idea of lossless compression via functional approximations and the storage of residuals. LeCo uses a learned model to choose a kind of function suitable for the values at hand, and then it uses a heuristic partitioning algorithm that greedily splits fragments (each associated with a function learned through regression methods) and merges neighbouring ones if this improves an estimate of the compression ratio. Our NeaTS, instead, compresses data with different function types, learns each function optimally under a given error bound, and employs a rigorous partitioning algorithm to minimise the actual compression ratio. These features, together with the more query-efficient compressed layout, make NeaTS better than LeCo in compression ratio, random access, and decompression speed, as our experiments show.

*c) Functional approximation for lossy compression:* The idea behind lossy functional approximation is to represent a time series as a sequence of functions over time, often under a chosen error metric [8], [79]. In the case of  $L_2$ -norm, the goal is to minimise the *sum of the squared residuals* (i.e. the vertical distances between the true and the approximated values). In the case of the  $L_1$ -norm, the goal is to minimise the *sum of the absolute values* of the residuals. The focus of our paper is instead to bound the  $L_\infty$ -norm of the residuals, thus bounding the *maximum* absolute residual.

As discussed in Section II, the optimal linear-time algorithm to solve this problem using a Piecewise Linear Approximation (PLA) with the minimum number of segments was first proposed by O’Rourke [36] (see also [25]–[27], [37], [38]). Interestingly, ModelarDB [80] has shown how PLAs

can be used in a fully-fledged distributed time series database. ModelarDB could benefit from our results since it computes PLAs via an algorithm [38] that was shown to use more segments than the optimal one that we use as a baseline [27], [36]. Other works [81], [82] proposed to further compress similar segments, which is a post-processing step that we can apply to our techniques too.

A few works [29]–[31] address the problem of partitioning a time series using nonlinear functions. We compared our approach to the most recent one that employs nonlinear  $\varepsilon$ -approximations [30], [31], called Adaptive Approximation (AA). AA heuristically partitions a time series using quadratic, linear, and exponential functions. However, as our experiments show, despite the use of nonlinear functions, AA is almost always worse than PLA, and in turn worse than our approach. In fact, NeaTS-L finds a better partition of the time series into fragments and selects the best approximating function within a larger set of nonlinear ones.

Finally, we mention HIRE [83], which focuses on the distinct problem of constructing a single encoding of a time series that can be decompressed at different  $L_\infty$  error bounds. Since HIRE relies on piecewise constant approximations as building blocks, we believe it could benefit from our more general nonlinear approximations.

## VI. CONCLUSIONS AND FUTURE WORK

We introduced new lossy and lossless compressors that harness the trends and patterns in time series data via a sequence of error-bounded linear and nonlinear approximations of different kinds and shapes. Our approaches experimentally proved to offer new trade-offs in terms of random access speed, decompression speed, and compression ratio compared to existing compressors on 16 diverse time series datasets.

For future work, we suggest further compressing the nonlinear approximation models by exploiting similarities between functions as introduced in [81], [84]. Another interesting research direction is to exploit the information encoded by the functions to efficiently answer aggregate queries on the time series data. Finally, it would be interesting to investigate the impact of our new techniques for computing error-bounded nonlinear approximations in the design of learned data structures [43], [44], [81], [85]–[95].

**Acknowledgments.** This work was supported by the European Union – Horizon 2020 Program under the scheme “INFRAIA-01-2018-2019 – Integrating Activities for Advanced Communities”, Grant Agreement n. 871042, “SoBigData++: European Integrated Infrastructure for Social Mining and Big Data Analytics” (<http://www.sobigdata.eu>), by the NextGenerationEU – National Recovery and Resilience Plan (Piano Nazionale di Ripresa e Resilienza, PNRR) – Project: “SoBigData.it - Strengthening the Italian RI for Social Mining and Big Data Analytics” – Prot. IR0000013 – Avviso n. 3264 del 28/12/2021, by the spoke “FutureHPC & BigData” of the ICSC – Centro Nazionale di Ricerca in High-Performance Computing, Big Data and Quantum Computing funded by European Union – NextGenerationEU – PNRR, by Regione Toscana under POR FSE 2021/27.

## REFERENCES

- [1] S. K. Jensen, T. B. Pedersen, and C. Thomsen, "Time series management systems: A survey," *IEEE Transactions on Knowledge and Data Engineering*, vol. 29, no. 11, pp. 2581–2600, 2017.
- [2] C. C. Aggarwal, Ed., *Managing and Mining Sensor Data*. Springer, 2013.
- [3] C. Wang, J. Qiao, X. Huang, S. Song, H. Hou, T. Jiang, L. Rui, J. Wang, and J. Sun, "Apache IoTDB: A time series database for IoT applications," *Proc. ACM Manag. Data*, vol. 1, no. 2, jun 2023.
- [4] InfluxData Inc., "InfluxDB," 2023. [Online]. Available: <https://www.influxdata.com/>
- [5] The OpenTSDB Authors, "OpenTSDB," 2023. [Online]. Available: <http://opentsdb.net/>
- [6] The Prometheus Authors, "Prometheus," 2023. [Online]. Available: <https://prometheus.io/>
- [7] Timescale Inc., "Timescale," 2023. [Online]. Available: <https://www.timescale.com/>
- [8] G. Chiarot and C. Silvestri, "Time series compression survey," *ACM Comput. Surv.*, aug 2022.
- [9] J. Alakuijala, A. Farruggia, P. Ferragina, E. Kliuchnikov, R. Obyrk, Z. Szabadka, and L. Vandevenne, "Brotli: A general-purpose data compressor," *ACM Transactions on Information Systems*, vol. 37, no. 1, pp. 1–30, 2018.
- [10] Meta Platforms, Inc., "Zstandard - fast real-time compression algorithm," 2016, retrieved March 19, 2023 from <https://github.com/facebook/zstd>.
- [11] The Tukaani Project, "The .xz file format," 2023, retrieved March 19, 2023 from <https://tukaani.org/xz/xz-file-format.txt>.
- [12] Y. Collet *et al.*, "Lz4: Extremely fast compression algorithm," 2013, retrieved March 19, 2023 from <https://github.com/lz4/lz4>.
- [13] Google, "Snappy: A fast compressor/decompressor," 2023, retrieved March 19, 2023 from <https://github.com/google/snappy>.
- [14] T. Pelkonen, S. Franklin, J. Teller, P. Cavallaro, Q. Huang, J. Meza, and K. Veeraraghavan, "Gorilla: A fast, scalable, in-memory time series database," *Proc. VLDB Endowment*, vol. 8, no. 12, pp. 1816–1827, 2015.
- [15] P. Liakos, K. Papakonstantinou, and Y. Kotidis, "Chimp: Efficient lossless floating point compression for time series databases," *PVLDB*, vol. 15, no. 11, pp. 3058–3070, sep 2022.
- [16] R. Li, Z. Li, Y. Wu, C. Chen, and Y. Zheng, "Elf: Erasing-based lossless floating-point compression," *PVLDB*, vol. 16, no. 7, pp. 1763–1776, may 2023.
- [17] A. Afroozeh, L. X. Kuffo, and P. Boncz, "ALP: Adaptive lossless floating-point compression," *Proc. ACM Manag. Data*, vol. 1, no. 4, dec 2023.
- [18] R. Vestergaard, Q. Zhang, M. Á. Sipos, and D. E. Lucani, "Titchy: Online time-series compression with random access for the internet of things," *IEEE Internet Things J.*, vol. 8, no. 24, pp. 17 568–17 583, 2021.
- [19] A. Gómez-Brandón, J. R. Paramá, K. Villalobos, A. Illarramendi, and N. R. Brisaboa, "Lossless compression of industrial time series with direct access," *Computers in Industry*, vol. 132, p. 103503, 2021.
- [20] X. Chen, J. Tian, I. Beaver, C. Freeman, Y. Yan, J. Wang, and D. Tao, "FCBench: Cross-domain benchmarking of lossless compression for floating-point data," *PVLDB*, vol. 17, no. 6, pp. 1418–1431, may 2024.
- [21] A. Khelifati, M. Khayati, A. Dignös, D. E. Difallah, and P. Cudré-Mauroux, "TSM-Bench: Benchmarking time series database systems for monitoring applications," *PVLDB*, vol. 16, no. 11, pp. 3363–3376, 2023.
- [22] Y. Hao, X. Qin, Y. Chen, Y. Li, X. Sun, Y. Tao, X. Zhang, and X. Du, "TS-Benchmark: A benchmark for time series databases," in *Proc. 37th International Conference on Data Engineering (ICDE)*, 2021, pp. 588–599.
- [23] K. F. Turkman, M. G. Scotto, and P. de Zea Bermudez, *Non-Linear Time Series*. Springer Cham, 2014.
- [24] J. Fan and Q. Yao, *Nonlinear Time Series*. Springer New York, NY, 2008.
- [25] E. J. Keogh, S. Chu, D. M. Hart, and M. J. Pazzani, "An online algorithm for segmenting time series," in *Proc. 2001 IEEE International Conference on Data Mining (ICDM)*, 2001, pp. 289–296.
- [26] X. Liu, Z. Lin, and H. Wang, "Novel online methods for time series segmentation," *IEEE Transactions on Knowledge and Data Engineering*, vol. 20, no. 12, pp. 1616–1626, 2008.
- [27] Q. Xie, C. Pang, X. Zhou, X. Zhang, and K. Deng, "Maximum error-bounded piecewise linear representation for online stream approximation," *The VLDB Journal*, vol. 23, no. 6, pp. 915–937, 2014.
- [28] E. Fuchs, T. Gruber, J. Nitschke, and B. Sick, "Online segmentation of time series based on polynomial least-squares approximations," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 32, no. 12, pp. 2232–2245, 2010.
- [29] F. Eichinger, P. Efron, S. Karnouskos, and K. Böhm, "A time-series compression technique and its application to the smart grid," *The VLDB Journal*, vol. 24, no. 2, pp. 193–218, 2015.
- [30] Z. Xu, R. Zhang, K. Ramamohanarao, and U. Parampalli, "An adaptive algorithm for online time series segmentation with error bound guarantee," in *Proc. 15th International Conference on Extending Database Technology (EDBT)*, 2012, pp. 192–203.
- [31] J. Qi, R. Zhang, K. Ramamohanarao, H. Wang, Z. Wen, and D. Wu, "Indexable online time series segmentation with error bound guarantee," *World Wide Web*, vol. 18, no. 2, pp. 359–401, 2015.
- [32] E. J. Keogh, K. Chakrabarti, M. J. Pazzani, and S. Mehrotra, "Dimensionality reduction for fast similarity search in large time series databases," *Knowledge and Information Systems*, vol. 3, no. 3, pp. 263–286, 2001.
- [33] N. Ao, F. Zhang, D. Wu, D. S. Stones, G. Wang, X. Liu, J. Liu, and S. Lin, "Efficient parallel lists intersection and index compression algorithms using graphics processing units," *PVLDB*, vol. 4, no. 8, pp. 470–481, 2011.
- [34] A. Boffa, P. Ferragina, and G. Vinciguerra, "A learned approach to design compressed rank/select data structures," *ACM Transactions on Algorithms*, 2022.
- [35] Y. Liu, X. Zeng, and H. Zhang, "LeCo: lightweight compression via learning serial correlations," *Proc. ACM Manag. Data*, vol. 2, no. 1, mar 2024.
- [36] J. O'Rourke, "An on-line algorithm for fitting straight lines between data ranges," *Communications of the ACM*, vol. 24, no. 9, pp. 574–578, 1981.
- [37] M. Dalai and R. Leonardi, "Approximations of one-dimensional digital signals under the  $l^\infty$  norm," *IEEE Transactions on Signal Processing*, vol. 54, no. 8, pp. 3111–3124, 2006.
- [38] H. Elmeleegy, A. K. Elmagarmid, E. Cecchet, W. G. Aref, and W. Zwaenepoel, "Online piece-wise linear approximation of numerical streams with precision guarantees," *PVLDB*, vol. 2, no. 1, pp. 145–156, aug 2009.
- [39] É. Grandjean and L. Jachiet, "Which arithmetic operations can be performed in constant time in the RAM model with addition?" 2023.
- [40] F. P. Preparata and D. E. Muller, "Finding the intersection of  $n$  half-spaces in time  $o(n \log n)$ ," *Theoretical Computer Science*, vol. 8, pp. 45–55, 1979.
- [41] M. H. Overmars, *The Design of Dynamic Data Structures*, ser. Lecture Notes in Computer Science. Springer, 1983, vol. 156.
- [42] T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein, *Introduction to Algorithms*, 3rd ed. The MIT Press, 2009.
- [43] P. Ferragina, H.-P. Lehmann, P. Sanders, and G. Vinciguerra, "Learned monotone minimal perfect hashing," in *Proc. 31st Annual European Symposium on Algorithms (ESA)*, 2023, pp. 46:1–46:17.
- [44] P. Ferragina, G. Manzini, and G. Vinciguerra, "Compressing and querying integer dictionaries under linearities and repetitions," *IEEE Access*, vol. 10, pp. 118 831–118 848, 2022.
- [45] P. Elias, "Efficient storage and retrieval by content and address of static files," *Journal of the ACM*, vol. 21, no. 2, pp. 246–260, 1974.
- [46] R. M. Fano, *On the number of bits required to implement an associative memory*. Massachusetts Institute of Technology, Project MAC, 1971.
- [47] G. Navarro, *Compact data structures: a practical approach*. Cambridge University Press, 2016.
- [48] R. Grossi, A. Gupta, and J. S. Vitter, "High-order entropy-compressed text indexes," in *Proc. 14th Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, 2003, pp. 841–850.
- [49] G. Jacobson, "Space-efficient static trees and graphs," in *Proc. 30th IEEE Symposium on Foundations of Computer Science (FOCS)*, 1989, pp. 549–554.
- [50] D. R. Clark, "Compact pat trees," Ph.D. dissertation, University of Waterloo, Canada, 1996.
- [51] S. Gog, T. Beller, A. Moffat, and M. Petri, "From theory to practice: Plug and play with succinct data structures," in *Proc. 13th International Symposium on Experimental Algorithms (SEA)*, 2014, pp. 326–337, code available at <https://github.com/simongog/sdsl-lite>.
- [52] S. Vigna, "Broadword implementation of rank/select queries," in *Proc. 7th International Workshop on Experimental Algorithms (WEA)*, 2008, pp. 154–168, code available at <https://github.com/vigna/sux>.

- [53] Y. Zheng, H. Fu, X. Xie, W.-Y. Ma, and Q. Li, "Geolife GPS trajectory dataset – user guide," July 2011. [Online]. Available: <https://www.microsoft.com/en-us/research/publication/geolife-gps-trajectory-dataset-user-guide/>
- [54] J. Zheng, H. Guo, and H. Chu, "A large scale 12-lead electrocardiogram database for arrhythmia study (version 1.0.0)," 2022.
- [55] J. Zheng, H. Chu, D. Struppa, J. Zhang, S. M. Yacoub, H. El-Askary, A. Chang, L. Ehwerhemuepha, I. Abudayyeh, A. Barrett, G. Fu, H. Yao, D. Li, H. Guo, and C. Rakovski, "Optimal multi-stage arrhythmia classification approach," *Scientific Reports*, vol. 10, no. 1, p. 2898, 2020.
- [56] National Ecological Observatory Network (NEON), "Ir biological temperature (dp1.00005.001)," 2022, retrieved March 19, 2023 from <https://data.neonscience.org/data-products/DP1.00005.001/RELEASE-2022>.
- [57] INFORE project, "Financial data set used in infores project," 2023, retrieved March 19, 2023 from <https://zenodo.org/record/3886895>.
- [58] National Ecological Observatory Network (NEON), "2d wind speed and direction (dp1.00001.001)," 2022, retrieved March 19, 2023 from <https://data.neonscience.org/data-products/DP1.00001.001/RELEASE-2022>.
- [59] —, "Barometric pressure (dp1.00004.001)," 2022, retrieved March 19, 2023 from <https://data.neonscience.org/data-products/DP1.00004.001/RELEASE-2022>.
- [60] —, "Relative humidity above water on-buoy (dp1.20271.001)," 2022, retrieved March 19, 2023 from <https://data.neonscience.org/data-products/DP1.20271.001/RELEASE-2022>.
- [61] University of Dayton, "Daily temperature of major cities," 2023, retrieved March 19, 2023 from <https://www.kaggle.com/sudalairajkumar/daily-temperature-of-major-cities>.
- [62] National Ecological Observatory Network (NEON), "Dust and particulate size distribution (dp1.00017.001)," 2022, retrieved March 19, 2023 from <https://data.neonscience.org/data-products/DP1.00017.001/RELEASE-2022>.
- [63] Meteoblue, "Historical weather data download," 2023, retrieved March 19, 2023 from [https://www.meteoblue.com/en/weather/archive/export/basel\\_switzerland](https://www.meteoblue.com/en/weather/archive/export/basel_switzerland).
- [64] InfluxData Inc., "InfluxDB 2.0 sample data," 2023, retrieved March 19, 2023 from <https://github.com/influxdata/influxdb2-sample-data>.
- [65] A. Bruno, F. M. Nardini, G. E. Pibiri, R. Trani, and R. Venturini, "TSXor: A simple time series compression algorithm," in *Proc. 28th International Symposium on String Processing and Information Retrieval (SPIRE)*, 2021, pp. 217–223, code available at <https://github.com/andybb Bruno/TSXor/>.
- [66] N. R. Brisaboa, S. Ladra, and G. Navarro, "DACs: bringing direct access to variable-length codes," *Information Processing & Management*, vol. 49, no. 1, pp. 392–404, 2013.
- [67] E. Nemerson, "Squash - compression abstraction library," 2015, <https://quixdb.github.io/squash>.
- [68] J. Ziv and A. Lempel, "A universal algorithm for sequential data compression," *IEEE Trans. Inf. Theory*, vol. 23, no. 3, pp. 337–343, 1977.
- [69] P. Ratanaworabhan, J. Ke, and M. Burtscher, "Fast lossless compression of scientific floating-point data," in *Proc. 16th Data Compression Conference (DCC)*, 2006, pp. 133–142.
- [70] P. Lindstrom and M. Isenburg, "Fast and efficient compression of floating-point data," *IEEE Transactions on Visualization and Computer Graphics*, vol. 12, no. 5, pp. 1245–1250, 2006.
- [71] M. Burtscher and P. Ratanaworabhan, "High throughput compression of double-precision floating-point data," in *2007 Data Compression Conference (DCC)*, 2007, pp. 293–302.
- [72] A. Afrozeh and P. Boncz, "The FastLanes compression layout: Decoding >100 billion integers per second with scalar code," *PVLDB*, vol. 16, no. 9, pp. 2132–2144, may 2023.
- [73] M. Kuschewski, D. Sauerwein, A. Alhomssi, and V. Leis, "BtrBlocks: efficient columnar compression for data lakes," *Proc. ACM Manag. Data*, vol. 1, no. 2, jun 2023.
- [74] C. Liu, H. Jiang, J. Paparrizos, and A. J. Elmore, "Decomposed bounded floats for fast compression and queries," *PVLDB*, vol. 14, no. 11, pp. 2586–2598, jul 2021.
- [75] D. W. Blalock, S. Madden, and J. V. Guttag, "Sprintz: Time series compression for the internet of things," *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.*, vol. 2, no. 3, pp. 93:1–93:23, 2018.
- [76] P. Lindstrom and M. Isenburg, "Fast and efficient compression of floating-point data," *IEEE Transactions on Visualization and Computer Graphics*, vol. 12, no. 5, pp. 1245–1250, 2006.
- [77] F. Knorr, P. Thoman, and T. Fahringer, "ndzip: a high-throughput parallel lossless compressor for scientific data," in *Proc. 31st Data Compression Conference (DCC)*, 2021, pp. 103–112.
- [78] N. J. Larsson and A. Moffat, "Off-line dictionary-based compression," *Proceedings of the IEEE*, vol. 88, no. 11, pp. 1722–1732, 2000.
- [79] Y. Liu, J. Li, H. Gao, and X. Fang, "Enabling  $\epsilon$ -approximate querying in sensor networks," *Proc. VLDB Endow.*, vol. 2, no. 1, pp. 169–180, 2009.
- [80] S. K. Jensen, T. B. Pedersen, and C. Thomsen, "ModelarDB: modular model-based time series management with spark and cassandra," *PVLDB*, vol. 11, no. 11, pp. 1688–1701, jul 2018.
- [81] P. Ferragina and G. Vinciguerra, "The PGM-index: a fully-dynamic compressed learned index with provable worst-case bounds," *PVLDB*, vol. 13, no. 8, pp. 1162–1175, 2020.
- [82] X. Kitsios, P. Liakos, K. Papakonstantinou, and Y. Kotidis, "Sim-Piece: Highly accurate piecewise linear approximation through similar segment merging," *PVLDB*, vol. 16, no. 8, pp. 1910–1922, 2023.
- [83] B. Barbarioli, G. Mersy, S. Sintos, and S. Krishnan, "Hierarchical residual encoding for multiresolution time series compression," *Proc. ACM Manag. Data*, vol. 1, no. 1, may 2023.
- [84] X. Kitsios, P. Liakos, K. Papakonstantinou, and Y. Kotidis, "Sim-piece: Highly accurate piecewise linear approximation through similar segment merging," *PVLDB*, vol. 16, no. 8, pp. 1910–1922, jun 2023.
- [85] T. Kraska, A. Beutel, E. H. Chi, J. Dean, and N. Polyzotis, "The case for learned index structures," in *Proc. 44th International Conference on Management of Data (SIGMOD)*, 2018, pp. 489–504.
- [86] R. Marcus, E. Zhang, and T. Kraska, "CDFShop: Exploring and optimizing learned index structures," in *Proc. International Conference on Management of Data (SIGMOD)*, 2020, pp. 2789–2792.
- [87] A. Kipf, R. Marcus, A. van Renen, M. Stoian, A. Kemper, T. Kraska, and T. Neumann, "RadixSpline: a single-pass learned index," in *Proc. 3rd International Workshop on Exploiting Artificial Intelligence Techniques for Data Management (aiDM)*, 2020.
- [88] P. Ferragina and G. Vinciguerra, "Learned data structures," in *Recent Trends in Learning From Data*, L. Oneto, N. Navarin, A. Sperduti, and D. Anguita, Eds. Springer International Publishing, 2020, pp. 5–41.
- [89] C. Wongkham, B. Lu, C. Liu, Z. Zhong, E. Lo, and T. Wang, "Are updatable learned indexes ready?" *PVLDB*, vol. 15, no. 11, pp. 3004–3017, 2022.
- [90] Z. Sun, X. Zhou, and G. Li, "Learned index: A comprehensive experimental evaluation," *PVLDB*, vol. 16, no. 8, pp. 1992–2004, 2023.
- [91] P. Ferragina, M. Frasca, G. C. Marinò, and G. Vinciguerra, "On nonlinear learned string indexing," *IEEE Access*, vol. 11, pp. 74021–74034, 2023.
- [92] D. Amato, G. Lo Bosco, and R. Giancarlo, "On the suitability of neural networks as building blocks for the design of efficient learned indexes," in *Proc. 23rd International Conference on Engineering Applications of Neural Networks (EANN)*, 2022, pp. 115–127.
- [93] D. Amato, G. L. Bosco, and R. Giancarlo, "Standard versus uniform binary search and their variants in learned static indexing: The case of the searching on sorted data benchmarking software platform," *Softw. Pract. Exp.*, vol. 53, no. 2, pp. 318–346, 2023.
- [94] D. Amato, R. Giancarlo, and G. Lo Bosco, "Learned sorted table search and static indexes in small-space data models," *Data*, vol. 8, no. 3, 2023.
- [95] M. H. Abrar and P. Medvedev, "PLA-index: A k-mer index exploiting rank curve linearity," in *Proc. 24th International Workshop on Algorithms in Bioinformatics (WABI)*, 2024, pp. 13:1–13:18.