

A. Meta Reviewer

Thank you very much for your comments. We have addressed the reviewers' concerns and revision items in the revised manuscript. Please find our response to your comments as follows.

**Required Changes.**

R1: O1-O4. (It is noted that the  $O_i$  are in some cases formulated as statements about results. The authors may still elaborate on these aspects in the revision.)

R2: O1, O2, O4

R3: O1-O4

Reply:

R1:

For O1, we explain the reason of proposing the baseline algorithm BOS-V, and emphasize the technical depth and the novelty of the advanced algorithm BOS-B.

For O2, we further optimize the compression time of the BOS-B and BOS-M algorithms, and report the evaluation in Figure 16, Section ??, Page 13.

For O3, we prove the theoretical guarantee on the quality of  $x_l$  and  $x_u$  in Proposition 4 in Section VI-C, Page 8.

For O4, we improve decompression time by arranging the storage layout of compressed data, in Section VII-A, Page 9.

R2:

For O1, we clarify the generality of the outlier issue in Section VIII-A2, Page 10.

For O2, we demonstrate the advantage of employing the operator on the storage and query processing costs, in Figure 12 in Section VIII-C1, Page 12.

For O4, we shorten the equations in the proof of Proposition 2, and omit some derivations in Proposition 3 that are similar.

R3:

For O1, we discuss compression techniques in signal processing/speech processing/data compression fields in Section II-B, Page 2, together with more references in recent years.

For O2, we add an experiment in Figure 14, Section VIII-D1, Page 12, to compare with the compression techniques in signal processing/speech processing/data compression fields.

For O3, we address the issues of referring to later sections, and simplify the long equations to enhance readability.

For O4, we explain necessity to divide the space in Figure 1 into three parts, and evaluate how the number of parts affects compression in Figure 15 in Section VIII-D2, Page 13.

**Optional Changes. R2: O3, O5**

Reply: R2:

For O3, we show how lower outliers affect compression if not separated, in Figure 13 in Section VIII-C2, Page 12.

For O5, we discuss the issue of BOS-M for other distributions such as skew, in Section VIII-B1, Page 12.

B. Reviewer 1

Thank you very much for reading our paper carefully and the helpful suggestions. Below is our response to your comments (in magenta).

**O1.** The first optimum algorithm for finding the boundaries  $(x_l, x_u)$  is straightforward and runs in quadratic time so it is not very practical. The improvement  $O(n \log n)$  follows almost straightforward. The main idea is that the value  $x_u$  is not necessary to be one from the input time series  $X$ . Given  $x_l$ , we can always compute the optimum value of  $x_u$ . It is a cute result, but I am not completely sure about the technical depth or the novelty of the submission.

Reply: It is true that the first optimum algorithm BOS-V running in quadratic time is not very practical. The reason of introducing this baseline is as follows. (1) It illustrates the rationale of traversing the values in  $X$  for the optimal solution in Proposition 1, which motivates the following algorithms. (2) It introduces some notations such as cumulative count in Definition 6, which are used in the following algorithms as well. (3) It is used to verify the correctness of the following advanced algorithm BOS-B, showing exactly the same compression results in Figure 11a. Referring to the comment, we add the aforesaid explanation at the beginning of Section IV, Page 4, before introducing the baseline.

While the improved  $O(n \log n)$  algorithm BOS-B is still concise, the foundation behind however is not-trivial. To find the optimal  $x_u$ , we need to prove that it is not necessary to traverse all the values in  $X$  in  $O(n)$  time for each  $x_l$ . The novelty of the proposal is to give the solution determined by the bit-width  $\beta$ , which takes only  $O(\log n)$  time. The technical depth roots in the existence of another better solution based on bit-width, for each solution  $(x_l, x_u)$  formed by values of  $X$ . The conclusion needs to be proved for two different cases as presented in Propositions 2 and 3, respectively. The complicated cost functions in Formulas 5 and 7, for center values, lower outliers and upper outliers, respectively, make the derivation difficult. Again, to better clarify the technical depth and the novelty of the submission, we add the aforesaid discussion at the beginning of Section V, Page 5, before introducing the advanced algorithm.

**O2.** Even an  $O(n \log n)$  or  $O(n)$  overhead might be too much in some cases. As shown in the experiments, the compression time is increasing but not significantly in almost all cases. In fact the  $O(n \log n)$  algorithm increase the compression time a lot. On the other hand, the BOS-M algorithm with  $O(n)$  execution time, even though it does not find the optimal  $x_l, x_u$ , only slightly increases the compression time.

Reply: We agree that even  $O(n \log n)$  or  $O(n)$  overhead of BOS-B and BOS-M algorithms can be significant in some cases, and thus propose to further optimize the performance.

(1) We improve BOS-B by traversing the bit-width  $\beta$  first, so that the cumulative counts for  $x_l$  and  $x_u$  can be more efficiently fetched. In the previous version, in Line 8 of Algorithm 2, for each  $x_l$ , BOS-B enumerates each  $x_u$  to find the responding cumulative count of  $x_l + 2^\beta$ , with varying  $\beta$ . In the revised BOS-B, we traverse the bit-width  $\beta$  first. That is, for each  $\beta$ , the cumulative counts for  $x_l$  and  $x_u = x_l + 2^\beta$  can be more efficiently fetched, given the fixed difference  $2^\beta$ . Although the time complexity is still  $O(n \log n)$ , Figure

double  
check  
the  
revision

revise  
all the  
replies  
below  
follow-  
ing this

A below shows that the average compression time over all datasets in Table III is significantly reduced by the improved BOS-B, compared to the original version, without losing compression ratio. We revise the aforesaid algorithm changes in Section V-B, Page 7. The corresponding compression time in Figure 11c is also updated.

(2) To compress all values, it requires to scan them at least once, and thus the  $O(n)$  time complexity of BOS-M is inevitable. Nevertheless, we replace QuickSelect [15] by a faster approximate median implementation [16] to improve the time cost of BOS-M. Again, Figure A illustrates the reduction of compression time by the improved BOS-M compared to the original version, without losing compression ratio. We update the algorithm details in Section VI-B, Page 8, and the corresponding compression time in Figure 11c as well.

Finally, we conduct an experiment on the average compression time over all datasets of improved BOS-V, BOS-B and BOS-M, by varying block size  $n$ , in Figure 16, Section ??, Page 13. All the methods increase almost linearly owing to the existence of duplicate values in the datasets. The advanced BOS-B increases much slower than BOS-V, while the approximate BOS-M is the most efficient. (The corresponding decompression time is presented for O4 below.)

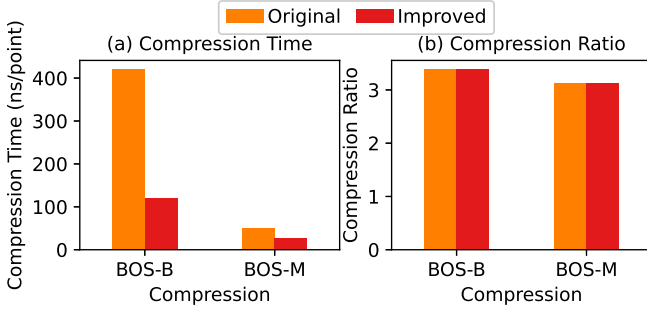


Fig.A: Improved compression time without losing compression ratio.

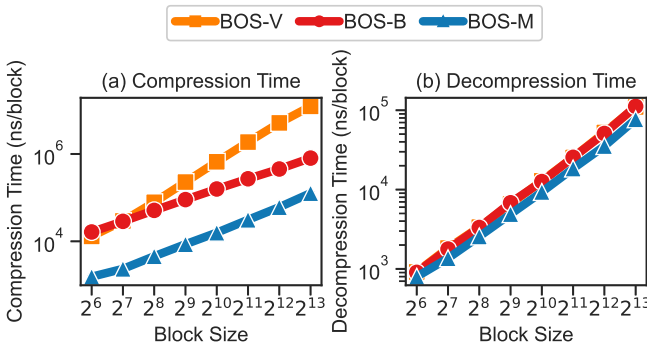


Fig.16: Compression and decompression time by varying block size  $n$ .

**O3.** Generally, it is good to mention that BOS-V and BOS-B are optimal algorithms, while BOS-M is a heuristic. While BOS-M is quite fast, it does not have any theoretical guarantee on the quality of  $x_l$  and  $x_u$ .

Reply: Thanks for your suggestions about the theoretical guarantee on the quality of  $x_l$  and  $x_u$  returned by the heuristic BOS-M. While it is difficult to bound the approximation ratio in general, given the various data distributions, we do obtain some theoretical guarantee for the special case of normal distribution, in Proposition 4 in Section VI-C, Page 8. The full proof procedure can be found in the appendix [1].

Let  $C_{\text{opt}}$  be the storage cost of the optimal solution for outlier separation problem, and  $C_{\text{approx}}$  be the storage cost of the solution  $x_l$  and  $x_u$  returned by the heuristic BOS-M. Since many real-world datasets follow the normal distribution as illustrated in Figure 9, we study the theoretical bound of approximation ratio  $\rho = \frac{C_{\text{approx}}}{C_{\text{opt}}}$  under the normal distribution.

**Proposition 4.** For normal distribution  $X \sim N(\mu, \sigma^2)$ , the approximation ratio  $\rho$  of BOS-M satisfies

$$\rho \leq \begin{cases} 2 & \text{if } \sigma \leq \frac{5}{3}, \\ \lceil \log(3\sigma - 1) \rceil & \text{otherwise.} \end{cases}$$

*Proof.* (1) Firstly, we prove the upper bound of the storage cost  $C_{\text{approx}}$  for BOS-M,

$$C_{\text{approx}} \leq \begin{cases} \lceil \log(6\sigma + 1) \rceil n & \text{if } \sigma < \frac{1}{2}, \\ 2n & \text{if } \frac{1}{2} \leq \sigma \leq \frac{5}{3}, \\ \lceil \log(3\sigma - 1) \rceil n & \text{otherwise.} \end{cases}$$

For normal distribution  $X \sim N(\mu, \sigma^2)$ , the median is  $\mu$ , the maximum value  $x_{\max}$  is approximately  $\mu + 3\sigma$ , and the minimum value  $x_{\min}$  is approximately  $\mu - 3\sigma$ , which correspond to the 99.7% confidence interval under the empirical rule (within three standard deviations from the mean).

The storage cost  $C_\beta$  of BOS-M with bit-width  $\beta$  is

$$\begin{aligned} C_\beta &= C(\mu - 2^\beta, \mu + 2^\beta) \\ &= n_l \lceil \log(\mu - 2^\beta - x_{\min} + 1) \rceil \\ &\quad + n_u \lceil \log(x_{\max} - (\mu + 2^\beta) + 1) \rceil \\ &\quad + (n - n_l - n_u) \lceil \log((\mu + 2^\beta) - (\mu - 2^\beta) + 1) \rceil \\ &= n_l \lceil \log(3\sigma - 2^\beta + 1) \rceil + n_u \lceil \log(3\sigma - 2^\beta + 1) \rceil \\ &\quad + (n - n_l - n_u)(\beta + 1) \\ &= (n_l + n_u) \lceil \log(3\sigma - 2^\beta + 1) \rceil \\ &\quad + (n - n_l - n_u)(\beta + 1). \end{aligned}$$

According to  $\beta$  from  $\lceil \log(6\sigma + 1) \rceil$  to 1,  $C_\beta$  first decreases and then increases. The upper bound of  $C_{\text{approx}}$  is thus

$$C_{\text{approx}} \leq \min\{C_{\lceil \log(6\sigma + 1) \rceil}, C_1\},$$

where

$$C_{\lceil \log(6\sigma + 1) \rceil} = \lceil \log(6\sigma + 1) \rceil n,$$

and

$$C_1 = (n_l + n_u) \lceil \log(3\sigma - 1) \rceil + 2(n - n_l - n_u).$$

Considering 4 different cases below, we rewrite  $C_1$  as

$$C_1 = 2n + (\lceil \log(3\sigma - 1) \rceil - 2)(n_l + n_u).$$

approxim where is this used below?

By increasing  $\beta$ ?

why

a) When  $\sigma \leq \frac{1}{2}$  ( $\mu - 2 < \mu - 3\sigma$ ), i.e., there are no upper and lower outliers with  $\beta = 1$ , we have  $n_l + n_u = 0$  and  $C_1 = 2n \geq \lceil \log(6\sigma + 1) \rceil n = C_{\lceil \log(6\sigma + 1) \rceil}$ .

b) When  $\frac{1}{2} \leq \sigma \leq \frac{2}{3}$  ( $\mu - 2 < \mu - 3\sigma$ ), i.e., there are no upper and lower outliers with  $\beta = 1$ , we have  $n_l + n_u = 0$  and  $C_1 = 2n < C_{\lceil \log(6\sigma + 1) \rceil}$ .

c) When  $\frac{2}{3} \leq \sigma \leq \frac{5}{3}$ , we have  $0 \leq \lceil \log(3\sigma - 1) \rceil \leq 2$  and  $2(n - n_l - n_u) \leq C_1 \leq 2n < C_{\lceil \log(6\sigma + 1) \rceil}$ .

d) When  $\sigma > \frac{5}{3}$ , we have  $\lceil \log(3\sigma - 1) \rceil \geq 2$  and  $C_1$  increases with  $\sigma$  growing. Then, when  $\sigma$  tends to positive infinity, we have there are no center values and  $C_1 = \lceil \log(3\sigma - 1) \rceil n < C_{\lceil \log(6\sigma + 1) \rceil}$ .

Therefore, we conclude that

$$C_{\text{approx}} \leq \begin{cases} \lceil \log(6\sigma + 1) \rceil n & \text{if } \sigma < \frac{1}{2}, \\ 2n & \text{if } \frac{1}{2} \leq \sigma \leq \frac{5}{3}, \\ \lceil \log(3\sigma - 1) \rceil n & \text{otherwise.} \end{cases}$$

(2) Moreover, we can prove that  $C_{\text{opt}} \geq n$ .

The storage cost is larger than the sum of bit-width for each value, thus the optimal cost  $C_{\text{opt}}$  has

$$C_{\text{opt}} = \sum_{i=1}^n b_i,$$

where

$$b_i = \begin{cases} 1 & \text{if } x_i = x_{\min}, \\ \lceil \log(x_i - x_{\min} + 1) \rceil & \text{otherwise.} \end{cases}$$

Thus, we have  $C_{\text{opt}} \geq n$ , even when all values are the same.

(3) Finally, we derive that

$$\rho = \frac{C_{\text{approx}}}{C_{\text{opt}}} \leq \begin{cases} \lceil \log(6\sigma + 1) \rceil & \text{if } \sigma < \frac{1}{2}, \\ 2 & \text{if } \frac{1}{2} \leq \sigma \leq \frac{5}{3}, \\ \lceil \log(3\sigma - 1) \rceil & \text{otherwise.} \end{cases}$$

a) When  $\sigma < \frac{1}{2}$ , we have  $\rho \leq \lceil \log(6\sigma + 1) \rceil \leq 2$ .

b) When  $\frac{1}{2} \leq \sigma \leq \frac{5}{3}$ , we have  $\rho \leq 2$ .

c) When  $\sigma > \frac{5}{3}$ , we have  $\rho \leq \lceil \log(3\sigma - 1) \rceil$ .

To sum up, we conclude that

$$\rho = \frac{C_{\text{approx}}}{C_{\text{opt}}} \leq \begin{cases} 2 & \text{if } \sigma \leq \frac{5}{3}, \\ \lceil \log(3\sigma - 1) \rceil & \text{otherwise.} \end{cases}$$

□

**O4.** The decompression time of BOS is generally larger than the decompression time without using BOS.

Reply: We acknowledge that the decompression time of BOS is generally larger than the decompression time without using BOS, since the outliers need to be processed separately. Nevertheless, we can improve the decompression time by arranging the storage layout in Figure 7, so that the data only needs to be scanned once. In the previous submission, the center values, lower outliers and upper outliers are stored separately. They need to be merged in decompression, by scanning all the values twice. In the revised layout, the lower outliers, center values and upper outliers are stored together in

use fractional data order. Their corresponding bit-widths,  $\alpha, \beta, \gamma$ , to avoid by a bitmap. Consequently, the decompression confuses only needs to scan the data once. We revise the storage layout in Section VII-A, Page 9.

same shown in Figure B below, the average decompression below over all datasets for BOS-V, BOS-B and BOS-M is improved by the new layout, compared to the original version, without losing compression ratio. We also update the decompression time of BOS in Figure 11c. It is not surprising that the decompression time increases linearly with the block size  $n$ , as illustrated in Figure 16b. BOS-M has less decompression time, since it separates fewer outliers. We explain the results in Section VIII-B2, Page 13.

double check

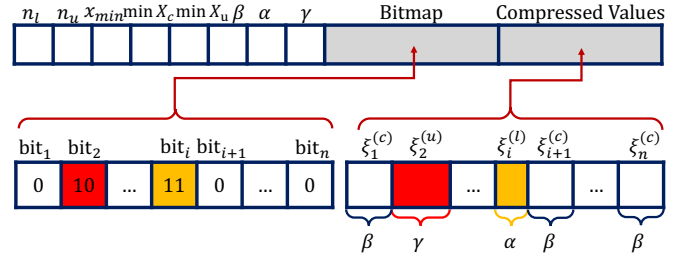


Fig. 7: (Improved) storage layout of bit-packing with outlier separation (BOS).

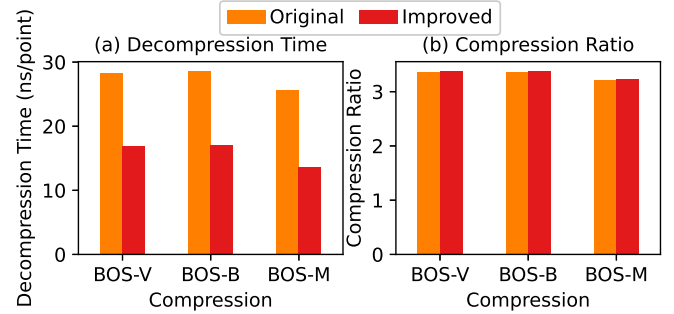


Fig.B: Improved decompression time without losing compression ratio.

### C. Reviewer 2

Thank you very much for reading our paper carefully and the helpful suggestions. Below is our response to your comments (in blue).

**O1.** The issue is a little bit special or the generality needs to be clearly pointed out.

Reply: Following the suggestion, we clarify the generality of the outlier issue in Section VIII-A2, Page 10. As the data distribution illustrated in Figure 9, outliers commonly exist in real datasets. We count the corresponding number of lower and upper outliers separated by BOS-V in each dataset in Figure 10. As discussed in O3 below, even for those datasets with a relatively small proportion of outliers, by separating them, the compression ratio could still be significantly improved. Therefore, the outlier issue is general and worthwhile to address in compression.

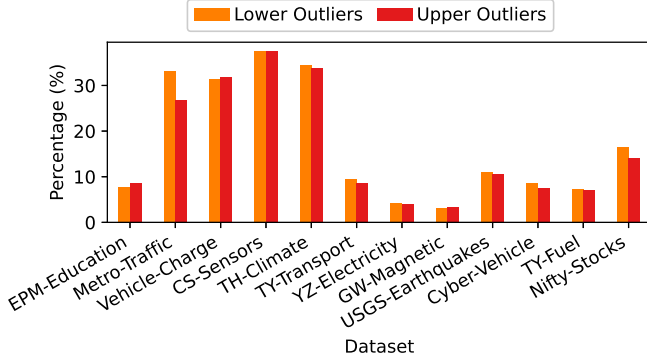


Fig.10: Percentage of lower and upper outliers separated by BOS-V

**O2.** To better motivate the issue, the advantage of employing the operator should be explicitly pointed out such as the storage and query processing costs. Is the scenario (upper and lower outliers) general? One should explicitly indicate that the issue is a common case.

Reply: Thanks again for the valuable suggestions to better motivate the outlier issue. In addition to illustrating that the scenario (upper and lower outliers) is general in the above O1, we further demonstrate the advantage of employing the operator, in Section VIII-C1, Page 12. Figure 12 reports the average storage and query processing cost over all datasets. As shown, with a better compression ratio in Figure 11a, our BOS operator yields lower storage costs. It leads to lower IO costs and thus query processing time comparable to the simple bit-packing operator (BP).

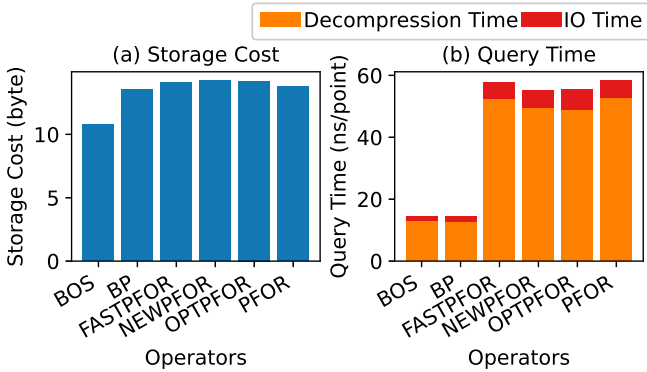


Fig. 12: Storage and query cost by various bit-packing operators in TS2DIFF

**O3.** Regarding the idea of considering the cumulative count, we may particularly focus on the range of a large number of values. This is because the number of lower outliers could be small. Then, the overall storage cost for them may not be significant. Can we terminate the loop early instead of enumerating all possible values.

Reply: As discussed in O1, it is true that the number of lower outliers could be small in some datasets, such as GW-Magnetic and YZ-Electricity, illustrated in Figure 10. While the overall storage cost for them may not be significant,

they could affect the storage of other center values if not separated. The reason is that as illustrated in Figure 1 and presented in Formula 5, the storage cost is determined by the minimum value of a set, i.e., lower outliers if not separated. Figure 13 reports the results of BOS by terminating the loop early without enumerating possible values for separating lower outliers, i.e., considering upper outliers only. As shown, even for those datasets with a relatively small proportion of lower outliers, such as the aforementioned GW-Magnetic and YZ-Electricity, considering both upper and lower outliers could have better compression ratio than separating upper outliers only (without considering lower outliers). We add the result discussion in Section VIII-C2, Page 12.

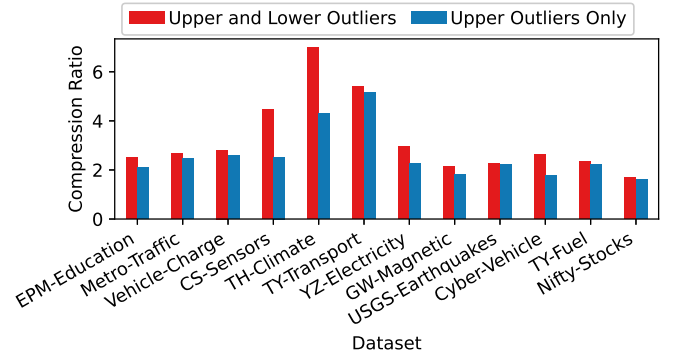


Fig. 13: Evaluating BOS terminating early without enumerating lower outliers.

**O4.** The proof procedure is solid but could be shortened or put into appendix.

Reply: We shorten the equations in the proof of Proposition 2, and omit some derivations in Proposition 3 that are similar to Proposition 2. For instance, we simplify Formula 13 in Proposition 2, Page 6. Moreover, in Proposition 3, we remove the proof of two conditions and Formula 12, in Page 6. The full proof procedure can be found in the appendix [1].

**O5.** The approximate median separation is based on the assumption that the data sets follow a normal distribution. One should also discuss the issue for other distributions such as skew.

Reply: As suggested, we discuss the issue of BOS-M for other distributions such as skew, in addition to the normal distribution, in Section VIII-B1, Page 12. First, for a normal distribution (after TS2DIFF), e.g., Figure 9(c) Vehicle-Charge, the approximate median separation works well, i.e., the compression ratio of (TS2DIFF+)+BOS-M is similar to that of BOS-V/B in Figure 11a (datasets VC). However, for other distributions such as skew, e.g., Figure 9(e) TH-Climate, there are a large number of low outliers in a very small range. It is difficult for BOS-M to find the proper separation of lower outliers by only enumerating bit-width  $\beta$ . Consequently, (TS2DIFF+)+BOS-M is much worse than BOS-V/B in Figure 11a (datasets TC).



Thank you very much for reading our paper carefully and the helpful suggestions. Below is our response to your comments (in red).

**O1.** Time series compression (or data compression) is a long-standing direction, which have been studied for few decades. Many research studies in signal processing/speech processing/data compression fields must propose many effective and efficient algorithms/software (e.g., 7-zip) for handling this compression task. However, the authors only provide those related studies that are in the database field. Therefore, I wonder whether those techniques in the aforementioned fields can be applied for this task. The authors should provide comprehensive review. It is quite shocking for me to see that there are only 27 references for this long-standing topic. Moreover, some of them are quite old.

Reply: Thanks for suggesting compression techniques in the signal processing/speech processing/data compression fields. We discuss them in Section II-B, Page 2, together with more other references in recent years.

Many research studies in signal processing/speech processing/data compression fields can be applied for the time series compression task. For example, 7-Zip [25] is a highly effective and efficient method for handling data compression. It is based on the LZMA (Lempel-Ziv-Markov chain algorithm) [24] compression algorithm, using dictionary compression and range encoding. LZ4 [7], derived from the LZ77 algorithm [37], searches for the longest matching string using a sliding window on the input stream. These data compression techniques for byte stream can be directly applied over the data encoded by bit-packing, i.e., complementary to our proposal, known as BOS+7-Zip or BOS+LZ4. For signal and speech processing, frequency-based methods are often employed [30], e.g., DCT [3] to compress speech data and FFT [14] to compress signal data. Since time-frequency transform could be lossy, to enable lossless compression, the corresponding residuals need to be stored. Again, our proposal BOS can be applied to improve the storage of the residuals often with outliers, known as BOS+DCT or BOS+FFT, i.e., again complementary to the existing methods. (The experiments below show that with our BOS, the compression ratio of these existing methods can be further improved.)

**O2.** As mentioned in O1, I also would like to see the comparison between those compression techniques in signal processing/speech processing/data compression fields and the proposed solution by the authors.

Reply: Following the suggestion, we add an experiment in Section VIII-D1, Page 12, to compare with the compression techniques in signal processing/speech processing/data compression fields. As discussed in O1, BOS as a fundamental bit-packing operator is complementary to these existing compression methods. Therefore, we compare compression ratio and time of 7-Zip [25], LZ4 [7], DCT [3], FFT [14] with and without our BOS in Figure 14. As shown, by combining these four compression algorithms with our BOS, the compression

ratios are all improved, of course with some extra overhead.

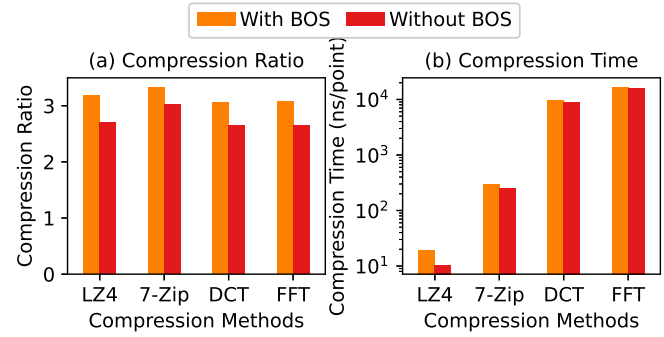


Fig. 14: Combining BOS with general data compression methods.

**O3.** The presentation of this paper should be further improved.

(1) The authors sometimes refer to later sections. As an example, Section I-B also refers to the later sections.(even for some tiny details). As another example, Section VII-B also refers to the Related Work section in Section VIII. Therefore, the flow is quite hard for readers to follow.

(2) The authors tend to write long equations (which spans for 6 lines in page 6). Therefore, I suggest the authors to simplify these equations.

Reply: Thanks for the detailed suggestions on improving the presentation of this paper.

(1) We address the issues of referring to later sections, in order to make the flow easy to follow. For example, we remove the references to later sections (for tiny details) in Section I-B in Page 2. Moreover, we move the Related Work section to Section II in Page 2 to avoid the problem of referring to Section VIII in the previous version.

(2) We simply the proof procedure and long equations to enhance readability. For instance, we shorten Formula 13 in Proposition 2, Page 6. Moreover, in Proposition 3, we remove the proof of two conditions and Formula 12, in Page 6, which is similar to Proposition 2. The full proof procedure can be found in the appendix [1].

**O4.** I do not know why it is necessary to divide the space in Figure 1 into three parts. Is it possible to use more parts (e.g., 6 parts)?

Reply: We divide the space in Figure 1 into three parts, owing to the existence of lower outliers, center values and upper outliers. It is possible to use more parts (e.g., 6 parts), which however may not further improve the compression ratio given the very close center values. Indeed, we can view the original bit-packing (BP) as a special case where the number of parts is 1. Likewise, the method considering only upper outliers is a special case with 2 parts. The BOS algorithms, in contrast, divide all values into 3 parts: center values, lower outliers and upper outliers. By further dividing the center values into 2 sets, it leads to a 4-part approach, and into 3 sets for a 5-part approach. The 6 parts are divided similarly.

As shown in Figure 15, when the number of parts increases from 1 to 3, the compression ratio improves significantly.

more  
new ref-  
erences,  
vldbj,  
icde24,  
...

many  
are url?  
check  
ency-  
clopae-  
dia in  
database  
or other  
fields

It verifies the intuition of our proposal in dividing the data into 3 parts, lower outliers, center values and upper outliers. However, the improvement is marginal by further dividing from 3 to 6 parts, as analyzed above. Unfortunately, the corresponding compression time increases considerably. Hence, we recommend to divide the space into 3 parts as in Figure 1. We add the above discussion in Section VIII-D2, Page 13.

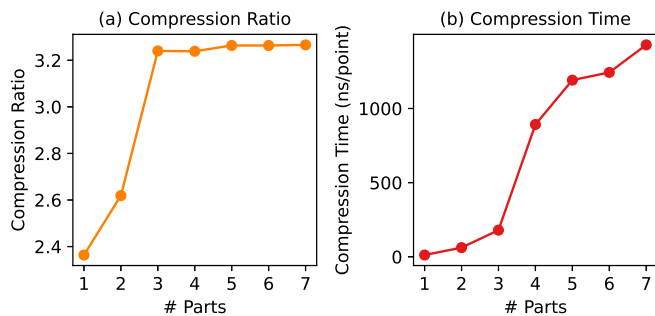


Fig. 15: Varying the number of divided value parts.

# BOS: Bit-packing with Outlier Separation

Jinzhao Xiao  
Tsinghua University  
xjc22@mails.tsinghua.edu.cn

Zihan Guo  
Tsinghua University  
gzh23@mails.tsinghua.edu.cn

Shaoxu Song\*  
Tsinghua University  
sxsong@tsinghua.edu.cn

**Abstract**—Bit-packing serves as the fundamental operator in various data encoding and compression methods. The idea is to use a fixed bit-width to represent all the (processed) values in a sequence. Some extremely large values, known as outliers, obviously amplify the bit-width, and thus lead to wasted bits for most other small values. We notice that not only the large values (upper outliers) but also the small ones (lower outliers) could incur wasted bit-width. In this paper, we propose to store both the upper and lower outliers separately, namely Bit-packing with Outlier Separation (BOS). While the remaining center values have a narrow spread, i.e., condensed bit-width, the separated outliers need some extra cost to denote their positions. The problem is thus how to determine better thresholds for separating the upper and lower outliers, yielding smaller storage cost. Rather than enumerating all the possible values as upper and lower outlier separators, in  $O(n^2)$  time, we consider bit-width as the separators, with  $O(n \log n)$  search time. Theoretical analysis illustrates all the possible cases such that the bit-width separation still returns the optimal solution as the value separation, and further leads to an approximate separation strategy with both median and bit-width, in  $O(n)$  time. Remarkably, our BOS is compatible to any existing compression methods using Bit-packing, and has replaced Bit-packing in Apache IoTDB and Apache TsFile. The extensive experiments on many real-world datasets demonstrate that by replacing Bit-packing with the proposed BOS in various compression methods, the compression ratio is significantly improved from about 2.75 to 3.25.

**Index Terms**—series, outlier, compression

## I. INTRODUCTION

There are many algorithms proposed to compress series data [38], [13], [26], [29], [19], [2]. Among them, many algorithms [38], [13], [2] employ Bit-packing [20] to improve storage by using the same bit-width for storing values in a block and removing leading zeros. Take a series of values  $X = (3, 2, 4, 5, 3, 2, 0, 8)$  as an example. Its maximum value is 8. The bit-width of 8 is 4 after removing leading zero. Thus, these 8 values can be stored with 4 bits respectively in bit-packing.

### A. Motivation

Note that in the above example series, only the large value 8, an outlier, needs 4 bits to store, while a bit-width 3 is sufficient for all the remaining values. That is, the outlier 8 incurs all the other values wasting 1 bit in Bit-packing.

1) *Outlier Separation Strategy*: A natural idea is thus to store the outliers separately, so that the remaining values could use a smaller bit-width. PFOR [38] and its variations, NewPFOR [34], OptPFOR [34], FastPFOR [18] and SimplePFOR

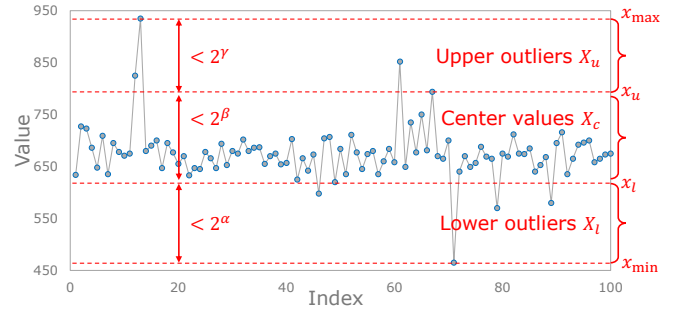


Fig. 1: A series  $X$  could be divided into 3 parts to compress separately the lower outliers  $X_l$ , center values  $X_c$  and upper outliers  $X_u$ , where  $\alpha, \beta, \gamma$  are the bit-widths for storing the corresponding values.

[18], propose to use  $b$  bits to store a part of the values, while separately storing others exceeding  $2^b - 1$ . For example, the aforesaid series can use 3 bits to store all the values except 8, which is processed separately.

We further notice that not only the large values, such as the above 8 known as upper outliers, but also the extremely small ones could amplify the bit width, e.g., value 0 in the series. By further separating the lower outlier 0, the remaining values (3, 2, 4, 5, 3, 2) need only a bit-width 2 to store, by subtracting the minimum value 2 from each value yielding (1, 0, 2, 3, 1, 0) in storage.

Figure 1 illustrates the outlier separation strategies over a series of values. The original Bit-packing needs  $\lceil \log(x_{\max} + 1) \rceil$  bits to store all the values, or  $\lceil \log(x_{\max} - x_{\min} + 1) \rceil$  bits by subtracting the minimum value  $x_{\min}$  from each value. PFOR separates the upper outliers greater than  $x_u$ , and uses  $\lceil \log(x_u + 1) \rceil$  bits (or  $\lceil \log(x_u - x_{\min} + 1) \rceil$  bits with  $x_{\min}$  subtraction) to store the remaining values. We propose to further separate the lower outliers smaller than  $x_l$ . The remaining values thus occupies only  $\lceil \log(x_u - x_l + 1) \rceil = \beta$  bits.

2) *Outlier Separation Determination*: Note that the separated outliers need some extra costs to indicate their indexes in the series, e.g., index 6 for value 0 and index 7 for value 8 in the example series. In other words, separating outliers may decrease the space of center values but introduce extra costs of indicating outlier positions. It thus needs a proper separation of outliers that would lead to lower total storage cost.

Unfortunately, existing methods use simple heuristics to determine outliers without considering the actual compression ratio performance. For example, NewPFOR [34] simply considers top 10% of values as outliers, and thus the storage

\*Shaoxu Song (<https://sxsong.github.io/>) is the corresponding author.

cost of values is not necessarily small. Other algorithms, such as OptPFOR [34], try to find the outliers with bit-width distribution of values. Again, the estimation by distribution does not lead to lower storage cost either.

Even worse, we need to determine both thresholds  $x_u$  and  $x_l$  for separating the upper and lower outliers, respectively, denoted by two red dashed lines in Figure 1, which are not considered in the existing studies. Note that a bit-width  $\beta$  can represent  $2^\beta$  distinct values. Intuitively, rather than enumerating all the values as possible  $x_l$  and  $x_u$ , we may search the separation in a small space of possible bit-widths  $\alpha, \beta, \gamma$  for lower, center and upper values, respectively.

## B. Contribution

In this paper, we propose *Bit-packing with Outlier Separation* (BOS) to improve Bit-packing, by separating both upper outliers and lower outliers. It is worth noting that BOS is complementary to the existing compression methods, such as RLE [13], SPRINTZ [2] and TS2DIFF [33]. By replacing their used Bit-packing operator with BOS, we have RLE+BOS, SPRINTZ+BOS, TS2DIFF+BOS, etc. Our major contributions in this paper are as follows.

(1) We formalize the optimization problem of outlier separation for bit-packing (BOS). Rather than simple heuristics, lower and upper outliers are formally defined with the corresponding data and position storage costs.

(2) We introduce the optimal separation based on values (BOS-V). It considers all the values from  $X$  as possible  $x_l$  and  $x_u$  for separating lower and upper outliers, with  $O(n^2)$  search time. Proposition 1 ensures the correctness of the separation, i.e., with the minimum storage cost. The solutions based on values from  $X$  lead to a more efficient method on bit-width below.

(3) We devise a bit-width separation method (BOS-B). Note that we only need to consider the bit-widths  $\beta$  and  $\gamma$  for center values and upper outliers, in  $O(n \log n)$  search time. Propositions 2 and 3 illustrate all the possible cases such that the bit-width separation still returns the optimal solution as the value separation. The correctness is ensured by transforming the solutions based on values from  $X$ .

(4) We propose an approximate median separation (BOS-M). Observing the normal distribution of values, especially after delta processing, we consider median in separation. Together with the aforesaid bit-width, the approximate separation can be determined in  $O(n)$  time.

(5) We conduct extensive experiments on many real world datasets. As summarized in Figure 11b, the compression ratio is significantly improved from about 2.75 by existing methods to 3.25 by our BOS-B on average. It is thus highly suggested to replace bit-packing by the proposed BOS in practice.

Our method BOS has been adopted in Apache IoTDB [29] and Apache TsFile [36], to replace Bit-packing. The code of the compression algorithm is included in the official GitHub repository of Apache IoTDB [5] and Apache TsFile [6] by system developers. The experiment related code and data are available in [12] for reproducibility.

## II. RELATED WORK

There are many algorithms to compress series data, including some lossy compression algorithms [10], [9], [17], and more importantly lossless algorithms [31], [32], [33].

### A. Compression of Integer and Float

*Integer Compression:* Lossless compression algorithms of integer include run-length-based [13] and differential-based algorithms. The storage cost of run-length-based algorithm, such as RLE [13], is better than that of other algorithms when values have high repeatability. However, these algorithms perform worse on values with small consecutive repeat. Differential-based algorithms, including TS2DIFF [33] and SPRINTZ [2], perform better on the series of values with small delta. These algorithms subtract the previous data from the current data and remove redundant leading zeros with bit-packing to reduce the storage cost of values. However, when there are several outliers leading to larger bit-width of values, the storage cost of these algorithms is very high.

*Float Compression:* GORILLA [26], CHIMP [21], Elf [19], and BUFF [22] are compression algorithms designed for floating-point numbers. GORILLA [26] computes a XOR of the current and previous float values, and then compresses these XOR values. CHIMP [21] improves GORILLA with distribution of leading and trailing zeros, and Elf [19] eases trailing zeros with precision of floating-point before computing XOR. However, if there are several larger outliers in float datasets, these algorithms have to store larger XOR values. BUFF [22] uses sparse encoding to handle outliers of floats. Nevertheless, BUFF [22] only splits values into two parts, outliers and normal values according to frequency, and does not optimize the outlier separation.

### B. Compression in Various Fields

Many research studies in signal processing/speech processing/data compression fields can be applied for the time series compression task. For example, 7-Zip [25] is a highly effective and efficient method for handling data compression. It is based on the LZMA (Lempel-Ziv-Markov chain algorithm) [24] compression algorithm, using dictionary compression and range encoding. LZ4 [7], derived from the LZ77 algorithm [37], searches for the longest matching string using a sliding window on the input stream. These data compression techniques for byte stream can be directly applied over the data encoded by bit-packing, i.e., complementary to our proposal, known as BOS+7-Zip or BOS+LZ4. For signal and speech processing, frequency-based methods are often employed [30], e.g., DCT [3] to compress speech data and FFT [14] to compress signal data. Since time-frequency transform could be lossy, to enable lossless compression, the corresponding residuals need to be stored. Again, our proposal BOS can be applied to improve the storage of the residuals often with outliers, known as BOS+DCT or BOS+FFT, i.e., again complementary to the existing methods.



### C. Compression with Outliers

Several previous compression schemes attempt to optimize the bit-packing algorithm by additionally handling outliers.

*Patched Frame-of-Reference, PFOR*: Zukowski et al. [38] propose the compression method to use a small bit-width  $b$  to bit-pack the center value and store outlier separately, but it does not compress additional outliers. PFOR stores the positions of outliers by organizing their indexes into lists. The offset to the next outlier index is stored in  $b$  bits of the center value position. This solution may introduce a large number of compulsory outliers. Therefore, PFOR still has a lot of room for improvement in general. Zhang et al. [35] propose to store the outlier values using either 8, 16, or 32 bits, which slightly improve the performance of PFOR.

*NewPFOR and OptPFOR*: Two other algorithms are proposed by Yang et al. [34] to obtain better storage. They use a bit-width for 128 integers, and store low  $b$  bits of the outlier value, so that the compulsory outlier can be avoided. Then they compress the 32 -  $b$  high bits and the index of outliers by Simple16. The difference between these two compression schemes lies in the strategy to determine  $b$ . NewPFOR chooses to search for the smallest  $b$  such that the number of outliers does not exceed 10% of the total data, while OptPFOR looks for the  $b$  that can compress the data best.

*FastPFOR and SimplePFOR*: To improve the compression effect of NewPFOR and accelerate it at the same time, Lemire and Boytsov [18] propose two algorithms called FastPFOR and SimplePFOR. They also use a bit-width per block of 128 integers, but store the outlier value and outlier index on each block. SimplePFOR compresses them together using Simple-8b, and FastPFOR classifies outliers according to the length of their high bits. Then they are packed according to the corresponding length bits one by one. Since FastPFOR also separates outliers rather than SimplePFOR, we compare FastPFOR in Section VIII.

However, this family of PFOR algorithms still had many problems. The first is that all of these algorithms only consider upper outliers are shown in Figure 1. In this case, the  $b$  used to pack most of the center values will be greatly affected. Secondly, bitmap is not considered to store index of outliers. In some cases, bitmap could save the index storage. Finally, the storage of outlier values is still not optimal. The value of each outlier point requires at least  $b$  bits to store the low bits. In fact, in our solution, it is very likely that less than  $b$  bits are needed to store the outlier value.

### III. PROBLEM STATEMENT

In this section, we give some basic definitions about storage cost of bit-packing and outlier separation. The optimization problem of outlier separation is then formalized. Table I lists the frequently used notations.

#### A. Bit-packing Encoding

Bit-packing [20] specifies a fixed bit-width for all the values in a series. The corresponding storage cost is given as follows.

TABLE I: Notations

Notation	Description
$X$	a series
$n$	the number of values in a block of series
$x_l, x_u$	floor value and ceiling value in center values
$X_c, X_l, X_u$	center values, lower outliers and upper outliers
$n_l, n_u$	the number of lower outliers and upper outliers
$\alpha, \beta, \gamma$	the bit-widths of lower, center and upper values
$C(x_l, x_u)$	storage cost with outlier separation
$c_i, c'_i$	the cumulative count

$x$	634	727	723	686	648	...	465	640	770	...	675
bitmap	0	10	10	0	0	...	11	0	10	...	0

Fig. 2: Example of using bitmap to indicate the positions of outliers.

**Definition 1** (Storage Cost). For a series  $X = (x_1, \dots, x_n)$ , its storage cost by Bit-packing is

$$C(X) = n \lceil \log(x_{\max} - x_{\min} + 1) \rceil \quad (1)$$

where  $x_{\max} = \max X$  and  $x_{\min} = \min X$  are the maximum and minimum values in the series  $X$ .

#### B. Outlier Separation

As shown in Figure 1, some large or small values increase storage cost in Definition 1. We propose to separate the outliers of both large and small values to store them separately, and thus reduce bit-widths of the remaining center values.

Specifically, we define lower bound of center values as  $x_l$ , and upper bound of center values as  $x_u$ . Based on  $x_l$  and  $x_u$ , all the values are split into 3 parts, including lower outliers, center values and upper outliers.

**Definition 2** (Center Values). Center values  $X_c$  are a set of values which are in the range of spread  $(x_l, x_u)$ ,

$$X_c = \{x_i \in X \mid x_l < x_i < x_u\}. \quad (2)$$

Center values are neither too larger nor too smaller with reduced bit-width  $\lceil \log(\max X_c - \min X_c + 1) \rceil$ .

**Definition 3** (Lower Outliers). Lower outliers  $X_l$  are a set of values which are less than center values,

$$X_l = \{x_i \in X \mid x_i \leq x_l\}. \quad (3)$$

The bit-width of lower outliers is reduced from  $\lceil \log(x_{\max} - x_{\min} + 1) \rceil$  to  $\lceil \log(\max X_l - x_{\min} + 1) \rceil$ . Thus, the storage cost of lower outliers is improved.

**Definition 4** (Upper Outliers). Upper outliers  $X_u$  are a set of values which are larger than center values,

$$X_u = \{x_i \in X \mid x_i \geq x_u\}. \quad (4)$$

The bit-width of upper outliers is decreased from  $\lceil \log(x_{\max} - x_{\min} + 1) \rceil$  to  $\lceil \log(x_{\max} - \min X_u + 1) \rceil$ . Again, the storage cost of upper outliers is improved.

Let  $n_l$  and  $n_u$  be the number of the lower outliers and upper outliers in the series  $X$ , i.e.,  $n_l = |X_l|$  and  $n_u = |X_u|$ . To

may  
simplify

store lower outliers and upper outliers individually, we need to record the positions of outliers in the original series. Figure 2 gives an example of storing outlier index with bitmap. We write '0' for the index of center values, '10' for lower outliers, and '11' for upper outliers. In this case, the storage cost of index is  $n + n_l + n_u$  bits.

### C. Separation Problem

In the following, we formulate the outlier separation problem. Let us first introduce the storage cost with outlier separation. The storage of index for outliers incurs extra storage cost. The total cost of values contains index cost and value cost of lower outliers, upper outliers and center values.

**Definition 5** (Storage Cost with Outlier Separation). *The cost  $C(x_l, x_u)$  of storing series  $X$  based on outlier separation by  $(x_l, x_u)$  is*

$$C(x_l, x_u) = n_l(\lceil \log(\max X_l - x_{\min} + 1) \rceil + 1) + n_u(\lceil \log(x_{\max} - \min X_u + 1) \rceil + 1) + (n - n_l - n_u)\lceil \log(\max X_c - \min X_c + 1) \rceil + n, \quad (5)$$

where  $x_{\min}$  and  $x_{\max}$  are the minimum and maximum values in the series  $X$ , having  $x_{\min} < \max X_l < \min X_c < \max X_c < \min X_u < x_{\max}$ .

If  $\max X_l = x_{\min}$ , the first term of  $C(x_l, x_u)$  is  $2n_l$ . If  $\min X_u = x_{\max}$ , the second term of  $C(x_l, x_u)$  is  $2n_u$ . If  $\max X_c = \min X_c$ , the third term of  $C(x_l, x_u)$  is  $(n - n_l - n_u)$ . When  $x_l < x_{\min}$  or  $x_u > x_{\max}$ , the number and bit-width of lower outliers or upper outliers are zero.

The outlier separation problem is to find the optimal range of center values with the minimum storage cost.

**Problem 1** (Outlier Separation Problem). *For a given series  $X$ , the outlier separation problem is to find the best  $(x_l, x_u)$  that minimizes the cost  $C(x_l, x_u)$ ,*

$$\arg \min_{x_l, x_u} C(x_l, x_u). \quad (6)$$

**Example 1.** Take the series in Figure 1 as an example. In the series, we set  $x_l$  as 620 and  $x_u$  as 794. Then,  $n_l$  and  $n_u$  are 5 and 4. Hence, the value cost is 698 and the cost of bitmap is 109. As a result, the storage cost is 807.

## IV. EXACT VALUE SEPARATION

Since the storage cost with outlier separation only depends on  $x_l$  and  $x_u$ , we could obtain the optimal solution by considering all the possible  $x_l$  and  $x_u$ . However, it takes too much time to consider each value from  $x_{\min}$  to  $x_{\max}$  for  $x_l$  and  $x_u$ . Thereby, we propose a separation algorithm by investigating only a set of values (BOS-V), still finding the optimal solution of outlier separation problem. First, we prove that an optimal solution  $(x_u, x_l)$  can always be found by considering the values  $x_l \in X$  and  $x_u \in X$  in Section IV-A. Moreover, the storage cost of different solutions can be efficiently calculated by maintaining the corresponding cumulative count, in Section IV-B. Furthermore, the pseudo

code of BOS-V algorithm is then presented in Section IV-C, together with the time complexity analysis in Section IV-D.

The reason of introducing this baseline is as follows. (1) It illustrates the rationale of traversing the values in  $X$  for the optimal solution, which motivates the following algorithms. (2) It introduces some notations such as cumulative count, which are used in the following algorithms as well. (3) It is used to verify the correctness of the following advanced algorithm BOS-B, showing exactly the same compression results.

### A. Optimal Separation with Values

According to Definition 5, since the cost of storing series  $X$  only depends on the values in the series, an optimal solution  $(x_l, x_u)$  must exist such that  $x_l$  and  $x_u$  are in the series  $X$ .

**Proposition 1.** *There must exist an optimal solution of outlier separation problem  $(x_u, x_l)$ , where  $x_l \in X$  and  $x_u \in X$ .*

*Proof.* For any optimal solution  $(x_l, x_u)$ , we can always construct another solution  $(\max X_l, \min X_u)$ , which has the same cost as  $(x_l, x_u)$ .

$$\begin{aligned} C(x_l, x_u) &= n_l(\lceil \log(\max X_l - x_{\min} + 1) \rceil + 1) \\ &\quad + n_u(\lceil \log(x_{\max} - \min X_u + 1) \rceil + 1) \\ &\quad + (n - n_l - n_u)\lceil \log(\max X_c - \min X_c + 1) \rceil + n \\ &= C(\max X_l, \min X_u). \end{aligned}$$

Note that the solution  $(\max X_l, \min X_u)$  has  $\max X_l \in X$  and  $\min X_u \in X$ . The conclusion is proved.  $\square$

### B. Cumulative Count

To calculate the storage cost  $C(x_l, x_u)$  for each solution, traversing all the values of the series  $X$  to obtain  $n_l$  and  $n_u$  in Definition 5 is very costly. Hence, we maintain a cumulative count to reduce times of traversing. The definition of cumulative count of values is as follows.

**Definition 6** (Cumulative Count). *The cumulative count  $c_i$  or  $c'_i$  of a value is the number of values less than and equal to it*

$$\begin{aligned} c_i &= |\{x_j \mid x_j \leq x_i, 1 \leq j \leq n\}|, \\ c'_i &= |\{x_j \mid x_j < x_i, 1 \leq j \leq n\}|. \end{aligned}$$

We present an example of cumulative count of values in Figure 3 for the series  $X$  from Figure 1. It is easy to see that lower outliers are in the left of the red line  $x_l$ , upper outliers are in the right of the red line  $x_u$ . Thus, we could get  $n_l$  and  $n_u$  with cumulative count efficiently.

Then, according to Definition 5, the value cost could be derived by cumulative count,

$$\begin{aligned} C(x_l, x_u) &= c_l(\lceil \log(\max X_l - x_{\min} + 1) \rceil + 1) \\ &\quad + (n - c'_u)(\lceil \log(x_{\max} - \min X_u + 1) \rceil + 1) \\ &\quad + (c'_u - c_l)\lceil \log(\max X_c - \min X_c + 1) \rceil + n. \end{aligned} \quad (7)$$

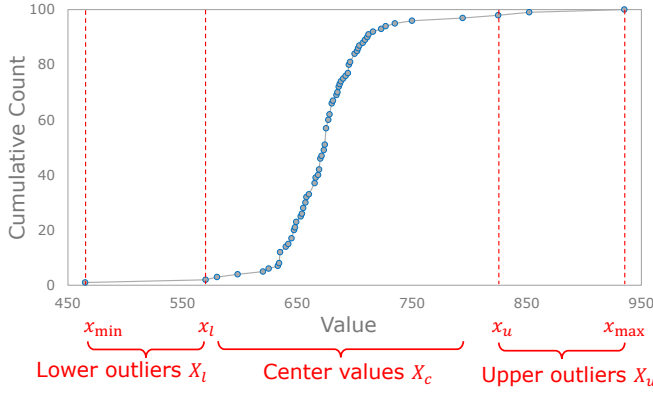


Fig. 3: Cumulative distribution function of values for outlier separation, where  $X_l, X_c, X_u$  denote the value separation for the series  $X$  in Figure 1.

### C. Value Separation Algorithm

Algorithm 1 presents the pseudo code of finding the optimal separation  $(x_l, x_u)$  with the minimum storage cost  $C_{\min}$ . First, we sort values in the series  $X$  in Line 1. Then, the cumulative count of each value in the series is calculated in Lines 2 and 3. Lines 4-10 get storage cost of each solution  $x_l, x_u$  with Formula 7, and find the one with the minimum storage cost.

---

#### Algorithm 1: Value Separation (BOS-V)

---

**Input:** Series  $X = (x_1, x_2, \dots, x_n)$   
**Output:** Optimal Solution  $(x'_l, x'_u)$

```

1  $X = \text{Sort}(X)$  ;
2 for  $x_i \leftarrow x_{\min}$  to  $x_{\max}$  do
3   | Get  $c_i$  with Definition 6 ;
4  $C_{\min} = n * \lceil \log(x_{\max} - x_{\min} + 1) \rceil$  ;
5 for  $x_i \leftarrow x_{\min}$  to  $x_{\max}$  do
6   | for  $x_j \leftarrow x_{\max}$  to  $x_i$  do
7     |  $C_i = C(x_i, x_j)$  with Formula 7 ;
8     | if  $C_i < C_{\min}$  then
9       |  $x'_l = x_i$  ;
10      |  $x'_u = x_j$  ;
11      |  $C_{\min} = C_i$  ;
12 return  $(x'_l, x'_u)$  ;
```

---

**Example 2.** Consider the series in Figure 1. First, we sort the series in ascending order and get cumulative count as shown in Figure 3. In the series, Algorithm 1 enumerates  $x_l$  from the minimum value 465 to the maximum 935, and  $x_u$  from the next value of  $x_l$  to the maximum value 935. Lastly, the algorithm finds the optimal solution (632, 696) with the minimum cost.

### D. Complexity Analysis

Algorithm 1 takes a time cost of  $O(n \log n)$  to sort values, where  $n$  represents the number of values in  $X$ . After sorting values, the time cost of getting cumulative count is  $O(n)$ . The search of solution  $(x_l, x_u)$  with the minimum cost enumerates pairs of values in the series  $X$  in  $O(n^2)$  time. In summary, the time complexity of Algorithm 1 is  $O(n^2)$ .

## V. EXACT BIT-WIDTH SEPARATION

The quadratic time complexity of Algorithm 1 is still costly. Rather than values from  $X$ , we propose to use bit-width as the separation (BOS-B), reducing the time complexity to  $O(n \log n)$ . In Section V-A, we prove that for each solution  $(x_l, x_u)$  with  $x_l \in X$  and  $x_u \in X$ , it can always find another no worse solution  $(x_l, x'_u)$ , where  $x'_u$  is determined by bit-width. Furthermore, we present the pseudo code and complexity analysis of BOS-B and in Sections V-B and V-C.

While the improved  $O(n \log n)$  algorithm BOS-B is still concise, the foundation behind however is not-trivial. To find the optimal  $x_u$ , we need to prove that it is not necessary to traverse all the values in  $X$  in  $O(n)$  time for each  $x_l$ . The novelty of the proposal is to give the solution determined by the bit-width  $\beta$ , which takes only  $O(\log n)$  time. The technical depth roots in the existence of another better solution based on bit-width, for each solution  $(x_l, x_u)$  formed by values of  $X$ . The conclusion needs to be proved for two different cases. The complicated cost functions in Formulas 5 and 7, for center values, lower outliers and upper outliers, respectively, make the derivation difficult.

### A. Optimal Separation with Bit-width

For any solution  $(x_l, x_u)$  with  $x_l \in X$  and  $x_u \in X$ , let

$$\beta = \lceil \log(\max X_c - \min X_c + 1) \rceil, \quad (8)$$

$$\gamma = \lceil \log(x_{\max} - \min X_u + 1) \rceil, \quad (9)$$

denote the bit-widths of center values and upper outliers.

**Proposition 2.** For any solution  $(x_l, x_u)$  with  $\beta \leq \gamma$ ,  $x_l \in X$  and  $x_u \in X$ , there always exists another solution  $(x_l, x'_u)$  having  $C(x_l, x'_u) \leq C(x_l, x_u)$ , where  $x'_u = \min X_c + 2^\beta$ .

*Proof.* According to  $\beta = \lceil \log(\max X_c - \min X_c + 1) \rceil$  in Formula 8, we have

$$\begin{aligned} \log(\max X_c - \min X_c + 1) &\leq \beta \\ \max X_c - \min X_c + 1 &\leq 2^\beta \\ \max X_c &\leq \min X_c + 2^\beta - 1 \\ \max X_c &< x'_u. \end{aligned}$$

(1) For  $x_u > x'_u$ , it follows  $\max X_c < x'_u < x_u = \min X_u$ . Since there is no value between  $\max X_c$  and  $\min X_u$  in  $X$ , according to Definitions 2 and 4, we have  $\min X'_u = \min X_u$ .

(2) For  $x_u \leq x'_u$ , referring to Definition 4, we have  $\max X'_u \geq \max X_u$ .

Combining the above two cases, we can conclude that

$$\min X'_u \geq \min X_u.$$

For  $n_u = |X_u|$  and  $n'_u = |X'_u|$  introduced after Definition 4, it follows  $n_u \geq n'_u$ . Let  $n_\Delta = |X_u \setminus X'_u|$  be the size of the increment, having  $n_\Delta = n_u - n'_u \geq 0$ .

Given the same  $x_l$  and the corresponding identical  $X_l, n_l$ , we could get the difference  $C_\Delta$  between  $C(x_l, x'_u)$  and  $C(x_l, x_u)$  defined in Formula 5,

$$\begin{aligned} C_\Delta &= C(x_l, x'_u) - C(x_l, x_u) \\ &= n_l(\lceil \log(\max X_l - x_{\min} + 1) \rceil + 1) \\ &\quad + n'_u(\lceil \log(x_{\max} - \min X'_u + 1) \rceil + 1) \\ &\quad + (n - n_l - n'_u)\lceil \log(\max X'_c - \min X'_c + 1) \rceil \\ &\quad - n_l(\lceil \log(\max X_l - x_{\min} + 1) \rceil + 1) \\ &\quad - n_u(\lceil \log(x_{\max} - \min X_u + 1) \rceil + 1) \\ &\quad - (n - n_l - n_u)\lceil \log(\max X_c - \min X_c + 1) \rceil. \end{aligned} \quad (10)$$

The same  $x_l$  also infers  $\min X'_c = \min X_c$ . Together with  $n_u = n'_u + n_\Delta$ , we have

$$C_\Delta = C_1 - C_2, \quad (11)$$

where

$$\begin{aligned} C_1 &= (n - n_l - n_u)\lceil \log(\max X'_c - \min X_c + 1) \rceil \\ &\quad + n_\Delta\lceil \log(\max X'_c - \min X_c + 1) \rceil \\ &\quad + n'_u(\lceil \log(x_{\max} - \min X'_u + 1) \rceil + 1) \end{aligned}$$

and

$$\begin{aligned} C_2 &= (n - n_l - n_u)\lceil \log(\max X_c - \min X_c + 1) \rceil \\ &\quad - n_\Delta(\lceil \log(x_{\max} - \min X_u + 1) \rceil + 1) \\ &\quad - n'_u(\lceil \log(x_{\max} - \min X_u + 1) \rceil + 1) \\ &= (n - n_l - n_u)\beta - n_\Delta(\gamma + 1) - n'_u(\gamma + 1). \end{aligned}$$

(i) Referring to Definition 2, we have  $\max X'_c < x'_u = \min X_c + 2^\beta$ . It follows

$$\begin{aligned} \log(\max X'_c - \min X_c + 1) &\leq \log(2^\beta) \\ \lceil \log(\max X'_c - \min X_c + 1) \rceil &\leq \beta. \end{aligned}$$

(ii) With the aforesaid proved  $\min X'_u \geq \min X_u$ , we infer  $\lceil \log(x_{\max} - \min X'_u + 1) \rceil \leq \lceil \log(x_{\max} - \min X_u + 1) \rceil = \gamma$ .

Applying the above two conditions, we further derive

$$\begin{aligned} C_\Delta &\leq (n - n_l - n_u)\beta + n_\Delta\beta + n'_u(\gamma + 1) \\ &\quad - (n - n_l - n_u)\beta - n_\Delta(\gamma + 1) - n'_u(\gamma + 1) \\ &= n_\Delta(\beta - \gamma - 1) \leq 0. \end{aligned}$$

Given  $\beta \leq \gamma$  and  $n_\Delta \geq 0$ , the conclusion is proved.  $\square$

Intuitively, as illustrated in Figure 4, all the points could be divided into 4 parts, lower outliers  $X_l$ , center values  $X_c$ , upper outliers moved from  $X_u$  to center values  $X'_c$ , and remaining upper outliers  $X'_u$ . During moving points, the cost of lower outliers  $X_l$  does not change, the bit-width of center values  $X_c$  is still  $\beta$ , and the bit-width of remaining upper outliers  $X'_u$  does not get larger. Moreover, the bit-width of upper outliers  $X_u$  moved to center values  $X'_c$  changes from  $\gamma$  to  $\beta$ , i.e., getting no larger given  $\beta \leq \gamma$ . In summary, the cost of all the points becomes no greater, having  $C(x_l, x'_u) \leq C(x_l, x_u)$ .

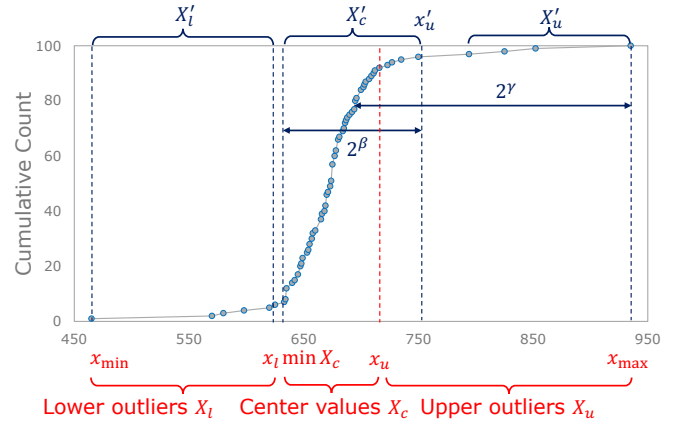


Fig. 4: Improving value separation  $(x_l, x_u)$  by bit-width separation  $(x_l, x'_u)$  using Proposition 2 with  $\beta \leq \gamma$ .

**Proposition 3.** For any solution  $(x_l, x_u)$  with  $\beta > \gamma$ ,  $x_l \in X$  and  $x_u \in X$ , there always exists another solution  $(x_l, x'_u)$  having  $C(x_l, x'_u) \leq C(x_l, x_u)$ , where  $x'_u = x_{\max} - 2^\gamma + 1$ .

*Proof.* According to  $\gamma = \lceil \log(x_{\max} - \min X_u + 1) \rceil$  in Formula 9, we have

$$\begin{aligned} \log(x_{\max} - \min X_u + 1) &\leq \gamma \\ x_{\max} - 2^\gamma + 1 &\leq \min X_u \\ x'_u &\leq \min X_u = x_u. \end{aligned}$$

(1) For  $x'_u = x_u = \min X_u$ , it is exactly the  $(x_l, x_u)$  solution, having  $\min X'_u = x'_u = \min X_u$ ,  $\max X'_c = \max X_c$ .

(2) For  $\max X_c < x'_u < x_u = \min X_u$ , since there is no value between  $\max X_c$  and  $\min X_u$  in  $X$ , according to Definitions 2 and 4, we have  $\min X'_u = \min X_u$ ,  $\max X'_c = \max X_c$  as well.

(3) For  $x'_u \leq \max X_c < x_u$ , referring to Definitions 2 and 4, it follows  $\max X'_c < \min X'_u \leq \max X_c < \min X_u$ .

Combining the above three cases, we can infer that

$$\begin{aligned} \min X'_u &\leq \min X_u \\ \max X'_c &\leq \max X_c. \end{aligned}$$

For  $n_u = |X_u|$  and  $n_u = |X'_u|$  introduced after Definition 4, it follows  $n'_u \geq n_u$ . Let  $n_\Delta = |X'_u \setminus X_u|$  be the size of the increment, having  $n_\Delta = n'_u - n_u \geq 0$ .

Given the same  $x_l$  and the corresponding identical  $X_l, n_l$ , we could get the difference  $C_\Delta$  according to Formula (12). The same  $x_l$  also infers  $\min X'_c = \min X_c$ . Together with  $n'_u = n_u + n_\Delta$ , we have

$$C_\Delta = C_1 - C_2,$$

where

$$\begin{aligned} C_1 &= (n - n_l - n'_u)\lceil \log(\max X'_c - \min X_c + 1) \rceil \\ &\quad + n_\Delta\lceil \log(x_{\max} - \min X'_u + 1) \rceil + 1 \\ &\quad + n_u(\lceil \log(x_{\max} - \min X'_u + 1) \rceil + 1) \end{aligned}$$



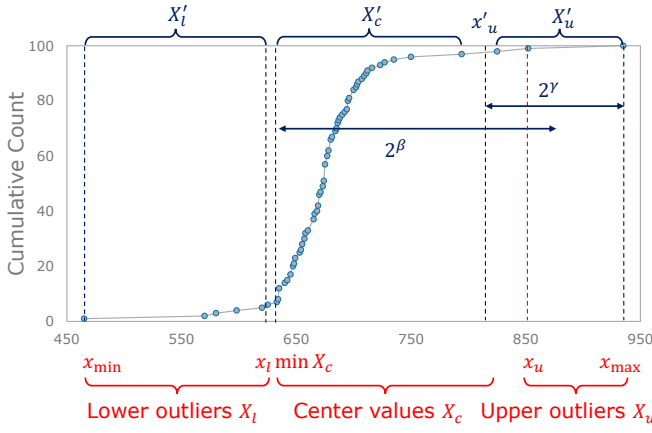


Fig. 5: Improving value separation  $(x_l, x_u)$  by bit-width separation  $(x_l, x'_u)$  using Proposition 3 with  $\beta > \gamma$ .

TABLE II: All possible cases of pruning by separation with bit-width

Proposition	Condition	Solution
Proposition 2	$\beta \leq \gamma$	$(x_l, \min X_c + 2^\beta)$
Proposition 3	$\beta > \gamma$	$(x_l, x_{\max} - 2^\gamma + 1)$

and

$$\begin{aligned}
 C_2 &= (n - n_l - n'_u) \lceil \log(\max X_c - \min X_c + 1) \rceil \\
 &\quad - n_\Delta \lceil \log(\max X_c - \min X_c + 1) \rceil \\
 &\quad - n_u (\lceil \log(x_{\max} - \min X_u + 1) \rceil + 1) \\
 &= (n - n_l - n'_u) \beta - n_\Delta \beta - n_u (\gamma + 1).
 \end{aligned}$$

Applying two conditions similar to condition (i) and (ii) in Proposition 2, we further derive

$$\begin{aligned}
 C_\Delta &\leq (n - n_l - n'_u) \beta + n_\Delta (\gamma + 1) + n_u (\gamma + 1) \\
 &\quad - (n - n_l - n'_u) \beta - n_\Delta \beta - n_u (\gamma + 1) \\
 &= n_\Delta (\gamma + 1 - \beta) \leq 0.
 \end{aligned}$$

Given  $\beta > \gamma$  and  $n_\Delta \geq 0$ , the conclusion is proved.  $\square$

Intuitively, as illustrated in Figure 5, all the points could be divided into 4 parts, lower outliers  $X_l$ , upper outliers  $X_u$ , center values moved from  $X_c$  to upper outliers  $X'_u$ , and the remaining center values  $X'_c$ . During moving points, the set of lower outliers  $X_l$  does not change, the bit-width of upper outliers  $X_u$  is still  $\gamma$ , and the bit-width of remaining center values  $X'_c$  does not get larger. Moreover, the bit-width of center values  $X_c$  moved to upper outliers  $X'_u$  changes from  $\beta$  to  $\gamma$ , i.e., become smaller given  $\gamma < \beta$ . To sum up, the cost of all the points becomes no greater, having  $C(x_l, x'_u) \leq C(x_l, x_u)$ .

### B. Bit-width Separation Algorithm

According to Propositions 2 and 3, we could get a solution no worse than value separation, including the optimal solution, by traversing each value as  $x_l$  and the corresponding bit-width  $\beta$  or  $\gamma$  for  $x_u$ . Table II summarizes all the possible cases of  $\beta \leq \gamma$  and  $\beta > \gamma$ , as well as their solutions to consider.

Algorithm 2 presents the pseudo code of bit-width separation (BOS-B). Firstly, same as Algorithm 1, we calculate

cumulative count of values in Lines 1 - 3. Then, the algorithm enumerates the cost of each  $x_l$  and each corresponding  $\beta$  with  $\beta \leq \gamma$  in Lines 5 - 12. The solution to consider is  $(x_l, \min X_c + 2^\beta)$ , according to the first case in Table II. **Note that we traverse the bit-width  $\beta$  first. That is, for each  $\beta$ , the cumulative counts for  $x_l$  and  $x_u = x_l + 2^\beta$  can be more efficiently fetched, given the fixed difference  $2^\beta$ .** For the second case in Table II, the algorithm enumerates the cost of each  $x_l$  and each  $\gamma$ , under the solution  $(x_l, x_{\max} - 2^\gamma + 1)$  in Lines 15 - 21.

### Algorithm 2: Bit-width Separation (BOS-B)

**Input:** Series  $X = (x_1, x_2, \dots, x_n)$

**Output:** Optimal Solution  $(x'_l, x'_u)$

```

1  $X = \text{Sort}(X)$  ;
2 for  $x_i \leftarrow x_{\min}$  to  $x_{\max}$  do
3   | Get  $c_i$  with Definition 6 ;
4  $C_{\min} = n * \lceil \log(x_{\max} - x_{\min} + 1) \rceil$  ;
5 for  $\beta \leftarrow 1$  to  $\lceil \log(x_{\max} - x_{i+1} + 1) \rceil - 1$  do
6   | for  $x_i \leftarrow x_{\min}$  to  $x_{\max}$  do
7     |  $x_u = x_{i+1} + 2^\beta$  ;
8     |  $C_i = C(x_l, x_u)$  with Formula 7 ;
9     | if  $C_i < C_{\min}$  then
10      |  $x'_l = x_l$  ;
11      |  $x'_u = x_u$  ;
12      |  $C_{\min} = C_i$  ;
13 for  $x_i \leftarrow x_{\min}$  to  $x_{\max}$  do
14   |  $x_l = x_i$  ;
15   | for  $\gamma \leftarrow 1$  to  $\lceil \log(x_{\max} - x_{i+1} + 1) \rceil - 1$  do
16     |  $x_u = x_{\max} - 2^\gamma + 1$  ;
17     |  $C_i = C(x_l, x_u)$  with Formula 7 ;
18     | if  $C_i < C_{\min}$  then
19       |  $x'_l = x_l$  ;
20       |  $x'_u = x_u$  ;
21       |  $C_{\min} = C_i$  ;
22 return  $(x'_l, x'_u)$  ;

```

**Example 3.** Consider the series in Figure 1. First, we sort the series in ascending order and get cumulative count, similar to Algorithm 1. In the series, Algorithm 2 enumerates  $x_l$  from the minimum value 465 to the maximum 935. For each  $x_l = x_i$  and  $\beta \leq \gamma$ , i.e.,  $\beta \leq \lceil \log(x_{\max} - x_{i+1} + 1) / 2 \rceil = \lceil \log(935 - x_{i+1} + 1) \rceil - 1$ , we consider the cost of  $x_u = x_{i+1} + 2^\beta$  as Figure 4. For each  $x_l = x_i$  and  $\beta > \gamma$ , i.e.,  $\gamma \leq \lceil \log(x_{\max} - x_{i+1} + 1) / 2 \rceil = \lceil \log(935 - x_{i+1} + 1) \rceil - 1$ , we consider the cost of  $x_u = x_{\max} - 2^\gamma + 1$  as Figure 5. Finally, the algorithm finds the optimal solution with  $x_l = 632$  and  $\beta = 6$ .

### C. Complexity Analysis

In Algorithm 2, the time cost of sorting values and getting cumulative count is  $O(n \log n)$ , similar to Algorithm 1. Then, it takes  $O(n \log n)$  time to calculate cumulative count of  $x_{i+1} + 2^\beta$  with cumulative count  $c_i$ , for each  $x_i$  and  $\beta$ . With the fixed  $x_{\max}$ , the calculation for cumulative count of  $x_{\max} - 2^\gamma + 1$  takes  $O(n)$  time. Finally, it takes  $O(n \log n)$  time

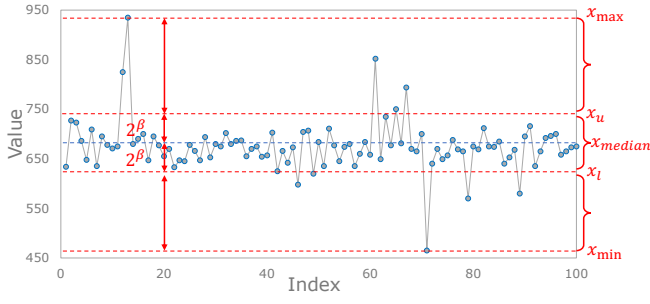


Fig. 6: Separation by  $x_{\text{median}}$  with bit-width  $\beta$  in both sides.

to find the minimum cost by enumerating  $x_l$  and  $\beta$  as well as  $\gamma$ . The overall complexity of Algorithm 2 is  $O(n \log n)$ .

## VI. APPROXIMATE MEDIAN SEPARATION

Algorithm 2 of bid-width separation still needs to traverse all possible values as  $x_l$ . In this section, we further narrow down the search space of  $x_l$  to the candidates determined by the median of  $X$  and bit-width  $\beta$ . This is motivated by the observation that many datasets (after pre-processing) follow a normal distribution, as illustrated in Figure 9 below.

### A. Approximate Separation with Median

Let  $x_{\text{median}}$  be the median of  $X$ . As illustrated in Figure 6, the center values are heuristically determined by  $x_l = x_{\text{median}} - 2^\beta$  and  $x_u = x_{\text{median}} + 2^\beta$ , for possible bit-width  $\beta$ . To efficiently calculate the storage cost in Formula 5, instead of the cumulative count, we define the count of buckets divided by median  $x_{\text{median}}$  and bit-width  $\beta$  as follows.

**Definition 7** (Bucket Count). *The bucket count  $h(\beta)$  is the number of values exceeding  $x_{\text{median}}$  in bit-width  $\beta$ ,*

$$h(\beta) = |\{x_i \in X \mid x_{\text{median}} + 2^{\beta-1} \leq x_i < x_{\text{median}} + 2^\beta\}|.$$

*The bucket count  $h(-\beta)$  is the number of values less than  $x_{\text{median}}$  with bit-width  $\beta$ ,*

$$h(-\beta) = |\{x_i \in X \mid x_{\text{median}} - 2^\beta < x_i \leq x_{\text{median}} - 2^{\beta-1}\}|.$$

*The special bucket count is  $h(0) = |\{x_i \in X \mid x_i = x_{\text{median}}\}|$ .*

### B. Median Separation Algorithm

We present the approximate median separation algorithm (BOS-M) in Algorithm 3. First, we use a faster approximate median implementation [16] of QuickSelect algorithm [15] to find median in Line 1. Then, we divide all the values for bucket count  $h(\beta)$  in Lines 2-10. Finally, it computes the storage cost of solution  $(x_{\text{median}} - 2^\beta, x_{\text{median}} + 2^\beta)$  for various bid-width  $\beta$  and finds the minimum.

**Example 4.** Consider the series in Figure 1. First, Algorithm 3 finds the median  $x_{\text{median}} = 674$ . Given the maximum  $\beta = 9$ , it divides values into 19 buckets. Then, the algorithm searches the bit-width  $\beta$  with the minimum cost by enumerating  $\beta$  from 9 to 1. It returns the solution (610, 738) with  $\beta = 6$ .

### Algorithm 3: Median Separation (BOS-M)

---

**Input:** Series  $X = (x_1, x_2, \dots, x_n)$   
**Output:** Approximate Solution  $(x'_l, x'_u)$

---

```

1  $x_{\text{median}} = \text{FindMedian}(X)$  ;
2 for  $i \leftarrow 1$  to  $n$  do
3   if  $x_i < x_{\text{median}}$  then
4      $\beta = \lceil \log(x_{\text{median}} - x_i + 1) \rceil$  ;
5      $h(-\beta) = h(-\beta) + 1$  ;
6   else if  $x_i > x_{\text{median}}$  then
7      $\beta = \lceil \log(x_i - x_{\text{median}} + 1) \rceil$  ;
8      $h(\beta) = h(\beta) + 1$  ;
9   else
10     $h(0) = h(0) + 1$  ;
11  $C_{\min} = n * \lceil \log(x_{\max} - x_{\min} + 1) \rceil$  ;
12 for  $\beta \leftarrow \lceil \log(x_{\max} - x_{\min} + 1) \rceil$  to 1 do
13    $n_l = n_l + h(-\beta)$  ;
14    $n_u = n_u + h(\beta)$  ;
15    $x_i = x_{\text{median}} - 2^\beta$  ;
16    $x_j = x_{\text{median}} + 2^\beta$  ;
17    $C_\beta = C(x_i, x_j)$  with  $n_l, n_u$  and Formula 5 ;
18   if  $C_\beta < C_{\min}$  then
19      $x'_l = x_i$  ;
20      $x'_u = x_j$  ;
21      $C_{\min} = C_\beta$  ;
22 return  $(x'_l, x'_u)$  ;
```

---

### C. Theoretical Guarantee

While it is difficult to bound the approximation ratio in general, given the various data distributions, we do obtain some theoretical guarantee for the special case of normal distribution.

Let  $C_{\text{opt}}$  be the storage cost of the optimal solution for outlier separation problem, and  $C_{\text{approx}}$  be the storage cost of the solution  $x_l$  and  $x_u$  returned by the heuristic BOS-M. Since many real-world datasets follow the normal distribution as illustrated in Figure 9, we study the theoretical bound of approximation ratio  $\rho = \frac{C_{\text{approx}}}{C_{\text{opt}}}$  under the normal distribution.

**Proposition 4.** For normal distribution  $X \sim N(\mu, \sigma^2)$ , the approximation ratio  $\rho$  of BOS-M satisfies

$$\rho \leq \begin{cases} 2 & \text{if } \sigma \leq \frac{5}{3}, \\ \lceil \log(3\sigma - 1) \rceil & \text{otherwise.} \end{cases}$$

### D. Complexity Analysis

In Algorithm 3, it takes  $O(n)$  amortized time complexity for the faster approximate median implementation [16] of QuickSelect algorithm [15] to find median [15]. Then, the algorithm takes  $O(n)$  time to divide all values into buckets, and  $O(\log n)$  time to calculate the storage cost of each solution with bit-width  $\beta$ . In summary, the time complexity of approximate median separation is  $O(n)$ .

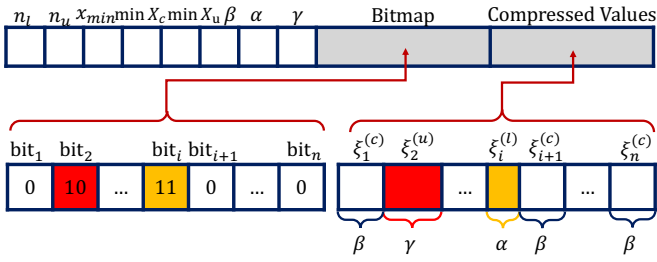


Fig. 7: Storage layout of bit-packing with outlier separation (BOS).

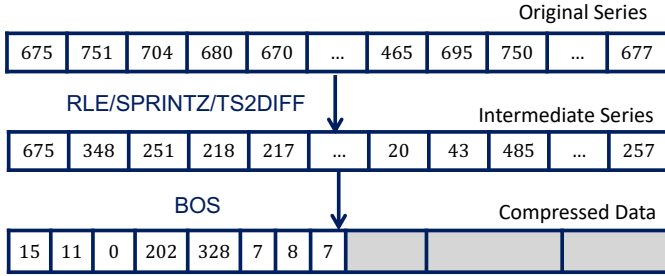


Fig. 8: Compression pipeline using bit-packing with outlier separation (BOS).

## VII. SYSTEM DEPLOYMENT

We implement BOS in Apache IoTDB [29] and Apache TsFile [36], and the code is available in the GitHub repository of the systems [5] and [6]. In the section, we introduce the storage layout of data compressed by BOS in the file format, and the compression pipeline with BOS in the database.

### A. Storage Layout

Figure 7 presents the storage structure of BOS in the file format. First, a block of values starts with some meta data of the series, including the number of outliers  $n_l$  and  $n_u$ , the minimum value  $x_{\min}$ , the minimum center value  $\min X_c$  and the minimum upper outlier  $\min X_u$ . It follows the bit-width of center values  $\beta = \lceil \log(\max X_c - \min X_c + 1) \rceil$ , bit-width of lower outliers  $\alpha = \lceil \log(\max X_l - x_{\min} + 1) \rceil$ , and bit-width of upper outliers  $\gamma = \lceil \log(x_{\max} - \min X_u + 1) \rceil$ . Then, we store the index of outliers with bitmap as shown in Figure 2, where  $\text{bit}_i$  is the indicator of  $i$ -th value. Next, the lower outliers, center values and upper outliers are stored together in the original data order. Their corresponding bit-widths,  $\alpha, \beta, \gamma$ , are marked by a bitmap. Consequently, the decompression process only needs to scan the data once. In Figure 7, center values  $\xi_i^{(c)}$  stores  $x_i - \min X_c$  in the blue boxes, lower outliers and upper outliers are stored as  $\xi_i^{(l)} = x_i - x_{\min}$  and  $\xi_i^{(u)} = x_i - \min X_u$  in the red and yellow boxes, respectively.

### B. Compression Pipeline

In addition to the original series data, the BOS algorithm could further compress the series output by other compression methods, such as RLE [13], SPRINTZ [2], TS2DIFF [29], etc. It is known as compression pipeline in databases. For example, for the series in Figure 1, the output of TS2DIFF is (675, 348, 251, 218, 217, ..., 20, 43, 485, ..., 257). BOS can

TABLE III: Real-world datasets

Dataset	Abbr.	Public	# Values	Data Type
EPM-Education	EE	[11]	900,000	Integer
GW-Magnetic	GM	[23]	933,984	Float
Metro-Traffic	MT	[28]	48,204	Integer
Nifty-Stocks	NS	[27]	295,193,088	Float
USGS-Earthquakes	UE	[8]	683,290	Float
Vehicle-Charge	VC	[4]	3,396	Integer
CS-Sensors	CS		100,000	Integer
Cyber-Vehicle	CV		35,676,900	Float, Integer
TH-Climate	TC		131,747	Integer
TY-Fuel	TF		183,556,352	Float, Integer
TY-Transport	TT		16,596,252	Integer
YZ-Electricity	YE		10,108	Float

be applied to further compress this intermediate series by TS2DIFF, as shown in Figure 8.

## VIII. EXPERIMENT

In this section, we experimentally compare the compression ratio and time of our BOS with other algorithms. In Section VIII-A, we introduce real-world datasets and metrics of the following experiments. In Section VIII-B, we compare the performance of our algorithms with baselines. Section ?? shows how the compression ratio of BOS changes with block size, and Section ?? evaluates BOS under various data distribution.

### A. Experimental Setting

The experiments were conducted on an Apple M1 Pro chip, featuring 8 CPU cores and 14 GPU cores, complemented by 16GB of unified memory.

1) *Baselines*: According to Section II of related work, we select several state-of-the-art algorithms in comparison, including floating-point compression algorithms (Gorilla [26], Chimp [21], Elf [19] and BUFF [22]) and integer encoding algorithms (RLE [13], SPRINTZ [2] and TS2DIFF [29]). Note that RLE, SPRINTZ and TS2DIFF use bit-packing, and thus are also denoted as RLE+BP, SPRINTZ+BP and TS2DIFF+BP.

We compare our algorithms, BOS with value separation (BOS-V), bit-width separation (BOS-B) and approximate median separation (BOS-M) with PFOR [38], NEWPFOR [34], OPTPFOR [34] and FASTPFOR [18], which also handle outliers in bit-packing. They can cooperate with other compression methods as well, by replacing BOS in the compression pipeline in Figure 8, e.g., RLE+BOS-V vs RLE+PFOR.

2) *Datasets*: The real world datasets utilized in our experimental evaluation encompass both publicly available data and the data acquired by our partners in various industries. The full dataset names in Table III indicate the corresponding diverse domains.

EPM-Education [11] is derived from the recorded activities of subjects engaged in educational simulations, facilitated by a logging application. GW-Magnetic [23] represents a multi-source and multivariate dataset tailored for indoor localization methodologies, drawing from WLAN and Geo-Magnetic fields. Metro-Traffic [28] comprises hourly traffic volume data for a specific thoroughfare. Nifty-Stocks [27] encompasses

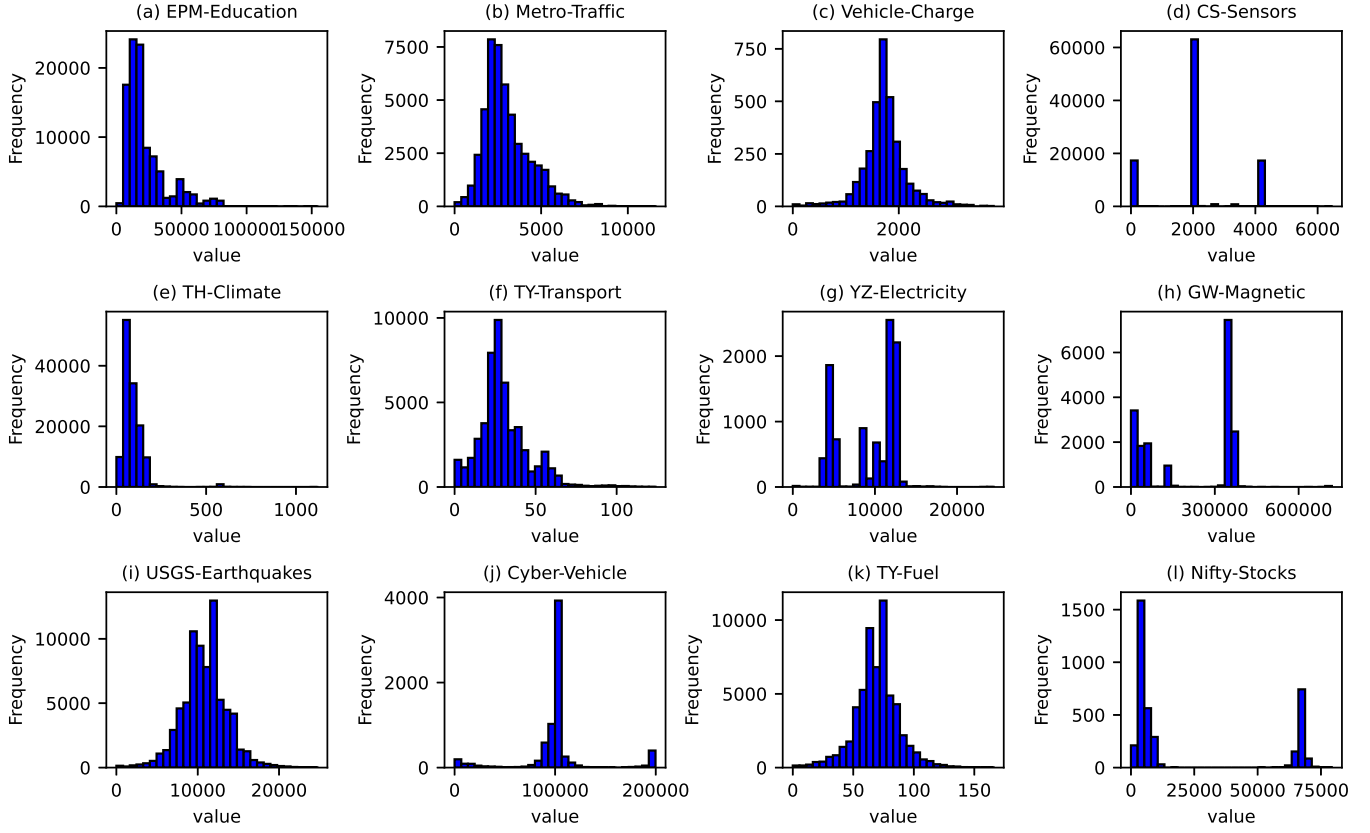


Fig. 9: Value distribution of all datasets after TS2DIFF.

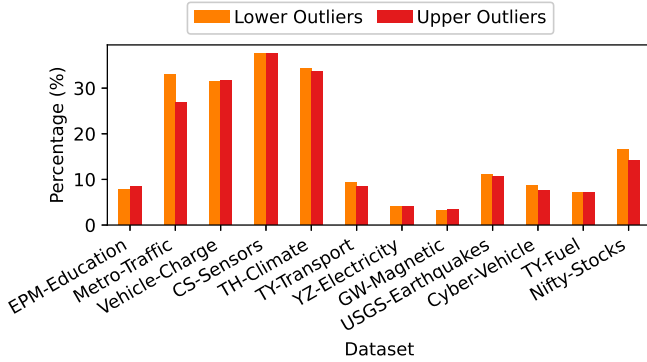


Fig. 10: Percentage of lower and upper outliers separated by BOS-V

comprehensive information including opening prices, high and low prices, closing prices, and volume data for Nifty 50 Stocks. USGS-Earthquakes [8] is an extensive dataset sourced from the United States Geological Survey (USGS), capturing detailed information on the magnitude and locations of earthquakes across the United States and its adjoining regions. Vehicle-Charges [4] is a specialized dataset tailored for electric vehicle drivers, collating usage data from multiple charging stations.

Additionally, our study integrates real-world datasets obtained through collaborations with industrial partners. CS-

Sensors is an assemblage of data obtained through the monitoring systems of maritime vessels. Cyber-Vehicle tracks series data pertaining to the operational status of concrete mixer trucks. TH-Climate represents a compilation of weather station data collected through wind sensors. TY-Fuel comprises data concerning the fuel consumption of various vehicle engines. TY-Transport encompasses data associated with the transportation duration of trucks. YZ-Electricity includes data derived from electrical power equipment systems.

Data types and the number of values in these datasets are shown in Table III. Some of the datasets contain only integers, where all the compression algorithms can be applied directly. There are also some datasets that contain floating-point numbers. Algorithms designed for integers, such as RLE, SPRINTZ and TS2DIFF, first convert float into integer by scaling  $10^p$ , where  $p$  is the precision of the original floating-point data [22].

To explain the results with compression pipeline in Figure 8, we draw the value distribution of all datasets after TS2DIFF in Figure 9. Since TS2DIFF removes trend by differencing values, most datasets after TS2DIFF follow normal distribution including datasets EE, MT, VC, TC, TT, UE, CV and TF. Moreover, owing to the existence of outliers, there are still some extreme delta values in the intermediate series by TS2DIFF, e.g., in TH.

As the data distribution illustrated in Figure 9, outliers

may  
simplify  
above

R2O1



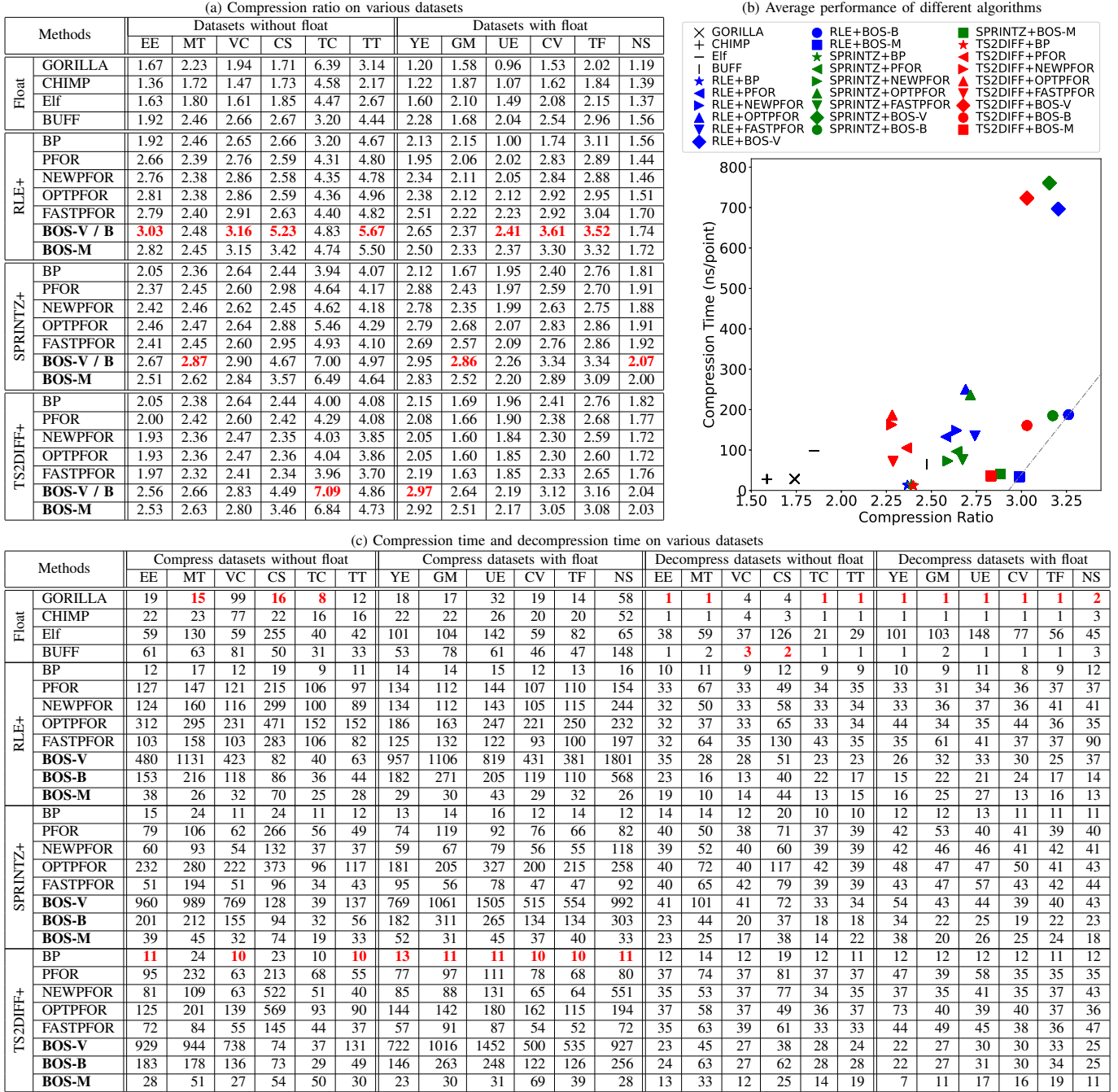


Fig. 11: Compression ratio and time of applying bit-packing with outlier separation (BOS) in different compression methods.

commonly exist in real datasets. We count the corresponding number of lower and upper outliers separated by BOS-V in each dataset in Figure 10. Even for those datasets with a relatively small proportion of outliers, by separating them, the compression ratio could still be significantly improved. Therefore, the outlier issue is general and worthwhile to address in compression.

3) *Metric*: We compare the compression ratio with other methods, which measures the ratio of uncompressed data size to compressed data size,  $compressionRatio =$

$$\frac{uncompressedSize}{compressedSize}.$$

We also evaluate the compression and decompression time per value (ns/points) by different algorithms. Each experiment is conducted 500 times and report the average.

### B. Comparison with Existing Methods

In the section, we compare performance of our proposals combined and compared with others. The compression ratio of algorithms is shown in Figure 11a. The corresponding compression time and decompression time are presented in Figure

11c. Figure 11b presents a summary of average compression ratio and time of each algorithm on all the datasets.

1) *Compression Ratio*: In Figure 11a, the red compression ratio is the best for the dataset in each column. As shown, the compression ratio of algorithms combined with BOS-V or BOS-B is always the best on all the datasets. In Figure 11b, BOS-B shows exactly the same compression ratio as BOS-V, verifying its correctness of returning the optimal solution. When combined with RLE or SPRINTZ, BOS-M has an overall performance better than the PFOR baseline and its variations.

Although TS2DIFF+BOS-M might not outperform some others, its compression ratio is still better than the PFOR baselines combined with TS2DIFF. The reason is that the output of TS2DIFF follows normal distribution as described in Section VIII-A2, where our median separation performs.

For a normal distribution (after TS2DIFF), e.g., Figure 9(c) Vehicle-Charge, the approximate median separation works well, i.e., the compression ratio of (TS2DIFF+)BOS-M is similar to that of BOS-V/B in Figure 11a (datasets VC). However, for other distributions such as skew, e.g., Figure 9(e) TH-Climate, there are a large number of low outliers in a very small range. It is difficult for BOS-M to find the proper separation of lower outliers by only enumerating bit-width  $\beta$ . Consequently, (TS2DIFF+)BOS-M is much worse than BOS-V/B in Figure 11a (datasets TC).

2) *Compression Time and Decompression Time*: As shown in Figure 11c, compression with value separation is very slow, since the time cost is high to sort all the values and enumerate value pairs as possible solutions. BOS-B with bit-width separation has lower time cost than BOS-V. The result is not surprising, given the time complexity reduced from  $O(n^2)$  to  $O(n \log n)$ . Finally, BOS-M with approximate median separation in  $O(n)$  time has comparable compression time cost as other baselines, while its compression ratio is better, as illustrated in Figure 11b.

As for decompression time, there is no clear difference observed between our BOS and the PFOR baselines with outlier separation. It is due to the same  $O(n)$  time cost in decompression.

3) *Trade-off between Compression Ratio and Time*: As illustrated in Figure 11b, the optimal solution BOS-B such as RLE+BOS-B has much better compression ratio than other algorithms, but is a bit slower in compression time. The linear time approximation BOS-M, e.g., RLE+BOS-M, achieves significantly lower compression time, and slightly weaker compression ratio (still outperforming baselines), i.e., a practical trade-off.

### C. Motivation Validation

1) *Storage and Query Cost*: To demonstrate the advantage of employing the operator, we perform an experiment to report the average storage and query processing cost over all datasets. With a better compression ratio in Figure 11a, our BOS operator yields lower storage costs as shown in Figure 12.

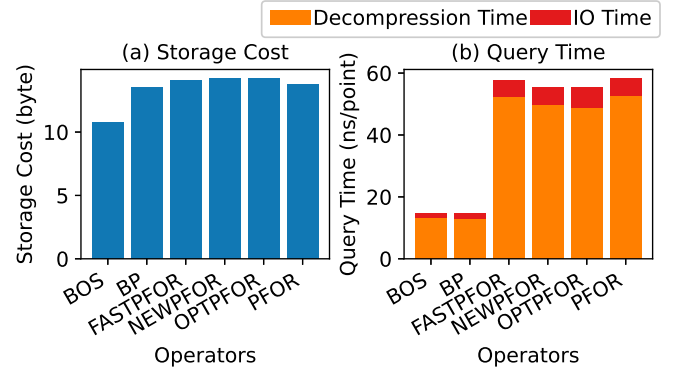


Fig. 12: Comparing query cost of BP and BOS-M

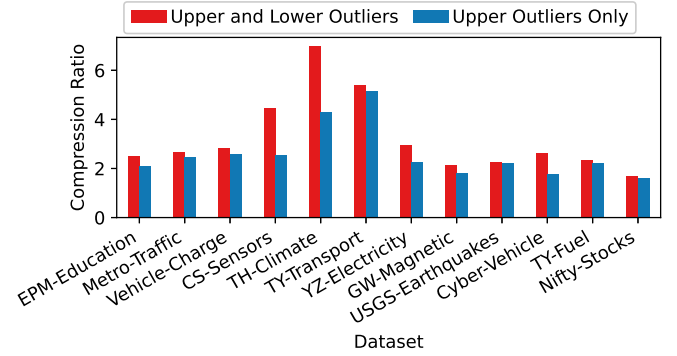


Fig. 13: Evaluating BOS terminating early without enumerating lower outliers.

It leads to lower IO costs and thus query processing time comparable to the simple bit-packing operator (BP).

2) *Comparison BOS with and without Lower Outliers*: It is true that the number of lower outliers could be small in some datasets, such as GW-Magnetic and YZ-Electricity, illustrated in Figure 10. While the overall storage cost for them may not be significant, they could affect the storage of other center values if not separated. The reason is that as illustrated in Figure 1 and presented in Formula 5, the storage cost is determined by the minimum value of a set, i.e., lower outliers if not separated. Figure 13 reports the results of BOS by terminating the loop early without enumerating possible values for separating lower outliers, i.e., considering upper outliers only. As shown, even for those datasets with a relatively small proportion of lower outliers, such as the aforementioned GW-Magnetic and YZ-Electricity, considering both upper and lower outliers could have better compression ratio than separating upper outliers only (without considering lower outliers).

### D. Variation Evaluation

1) *Comparison with Compressions in other Fields*: We perform an experiment to compare with the compression techniques in signal processing/speech processing/data compression fields. BOS as a fundamental bit-packing operator is complementary to these existing compression methods.

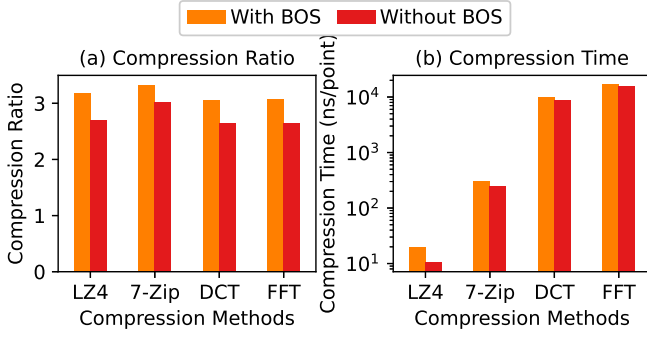


Fig. 14: Comparison with compression ratio of BOS with and without compression

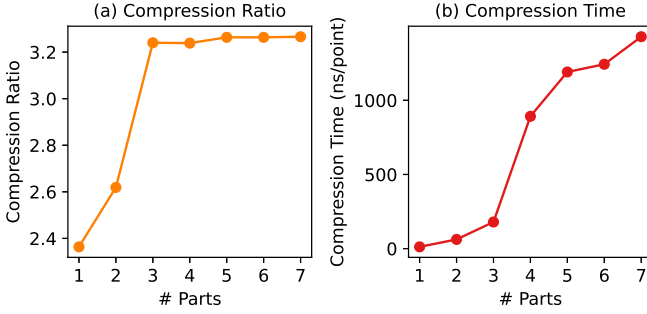


Fig. 15: Compression ratio and time varying part

Therefore, we compare compression ratio and time of 7-Zip [25], LZ4 [7], DCT [3], FFT [14] with and without our BOS in Figure 14. As shown, by combining these four compression algorithms with our BOS, the compression ratios are all improved, of course with some extra overhead.

2) *Varying the Part*: We perform an experiment about compression ratio and compression time varying the number of parts divided in Figure 15. when the number of parts increases from 1 to 3, the compression ratio improves significantly. It verifies the intuition of our proposal in dividing the data into 3 parts, lower outliers, center values and upper outliers. However, the improvement is marginal by further dividing from 3 to 6 parts, as analyzed above. Unfortunately, the corresponding compression time increases considerably. Therefore, we recommend to divide the space into 3 parts as shown in Figure 1.

#### E. Scalability

We conduct an experiment on the average compression time and decompression time over all datasets of BOS-V, BOS-B and BOS-M, by varying block size  $n$ , in Figure 16. All the methods increase almost linearly owing to the existence of duplicate values in the datasets. The advanced BOS-B increases much slower than BOS-V, while the approximate BOS-M is the most efficient. It is not surprising that the decompression time increases linearly with the block size  $n$ , as illustrated in Figure 16b. BOS-M has less decompression time, since it separates fewer outliers.

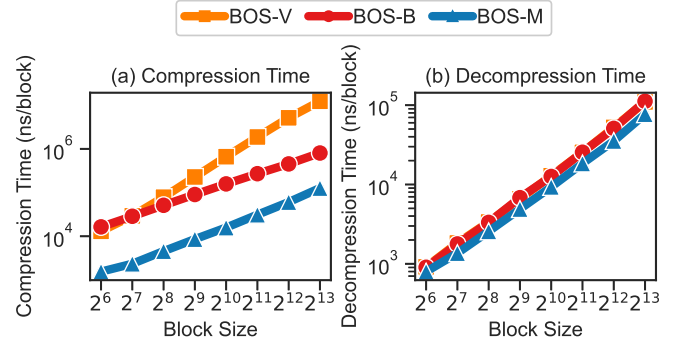


Fig. 16: Compression time varying block size

## IX. CONCLUSION

In this paper, we propose Bit-packing with Outlier Separation (BOS), which improves compression ratio of algorithms using bit-packing, by storing the outliers separately. It separates not only the upper outliers, occupying a large bit-width, but also the lower outliers, which waste the bit-width of center values as well. In order to determine a proper separation of outliers for better compression ratio, we devise an optimal separation strategy by enumerating the values in  $O(n^2)$  time, known as the value separator (BOS-V). With Propositions 2 and 3, the efficiency is improved by considering bit-width as the separator (BOS-B), still returning the optimal solution but taking only  $O(n \log n)$  search time. To further reduce the time cost, we propose an approximate median separation (BOS-M) in  $O(n)$  time. Experiments on real world datasets demonstrate that BOS-B with bit-width separation shows significantly higher compression ratio than existing methods, and lower compression time than the value separation BOS-V. As summarized in Figure 11b, together with RLE, BOS-M with approximate median separation achieves relatively high compression ratio and low compression time. In short, our proposal BOS is highly suggested to replace bit-packing in compression algorithms, which indeed has been adopted in Apache IoTDB and Apache TsFile.

## ACKNOWLEDGMENT

This work is supported in part by the National Natural Science Foundation of China (92267203, 62021002, 62072265, 62232005), the National Key Research and Development Plan (2021YFB3300500), and Beijing Key Laboratory of Industrial Big Data System and Application. Shaoxu Song (<https://xsong.github.io/>) is the corresponding author.

## REFERENCES

- [1] Appendix. <https://github.com/thssdb/encoding-outlier/blob/main/appendix.pdf>, 2024.
- [2] Davis W. Blalock, Samuel Madden, and John V. Guttag. Sprintz: Time series compression for the internet of things. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.*, 2(3):93:1–93:23, 2018.
- [3] Din-Yuen Chan, Jar-Ferr Yang, and Chun-Chin Fang. Fast implementation of MPEG audio coder using recursive formula with fast discrete cosine transforms. *IEEE Trans. Speech Audio Process.*, 4(2):144–148, 1996.

- [4] Charge. <https://www.kaggle.com/datasets/michaelbryantds/electric-vehicle-charging-dataset>, 2024.
- [5] Code. <https://github.com/apache/iotdb/tree/research/encoding-outlier>, 2024.
- [6] Code1. <https://github.com/apache/tsfile/tree/research/encoding-outlier>, 2024.
- [7] Yann Collet. Lz4: Extremely fast compression algorithm. <https://lz4.github.io/lz4/>, 2013. Available online.
- [8] Earthquakes. <https://www.kaggle.com/datasets/thedevastator/uncovering-geophysical-insights-analyzing-usgs-e>, 2024.
- [9] Frank Eichinger, Pavel Efros, Stamatis Karnouskos, and Klemens Böhm. A time-series compression technique and its application to the smart grid. *VLDB J.*, 24(2):193–218, 2015.
- [10] Hazem Elmeleegy, Ahmed K. Elmagarmid, Emmanuel Cecchet, Walid G. Aref, and Willy Zwaenepoel. Online piece-wise linear approximation of numerical streams with precision guarantees. *Proc. VLDB Endow.*, 2(1):145–156, 2009.
- [11] EPM. <https://doi.org/10.24432/C5NP5K>, 2024.
- [12] Experiment. <https://github.com/thssdb/encoding-outlier>, 2024.
- [13] Solomon W. Golomb. Run-length encodings (corresp.). *IEEE Trans. Inf. Theory*, 12(3):399–401, 1966.
- [14] Jinmoo Heo, Yongchul Jung, Seongjoo Lee, and Yunho Jung. FPGA implementation of an efficient FFT processor for FMCW radar signal processing. *Sensors*, 21(19):6443, 2021.
- [15] C. A. R. Hoare. Algorithm 65: find. *Commun. ACM*, 4(7):321–322, 1961.
- [16] C. A. R. Hoare. Quicksort. *Comput. J.*, 5(1):10–15, 1962.
- [17] Iosif Lazaridis and Sharad Mehrotra. Capturing sensor-generated time series with quality guarantees. In Umeshwar Dayal, Krithi Ramamritham, and T. M. Vijayaraman, editors, *Proceedings of the 19th International Conference on Data Engineering, March 5-8, 2003, Bangalore, India*, pages 429–440. IEEE Computer Society, 2003.
- [18] Daniel Lemire and Leonid Boytsov. Decoding billions of integers per second through vectorization. *Softw. Pract. Exp.*, 45(1):1–29, 2015.
- [19] Ruiyuan Li, Zheng Li, Yi Wu, Chao Chen, and Yu Zheng. Elf: Erasing-based lossless floating-point compression. *Proc. VLDB Endow.*, 16(7):1763–1776, 2023.
- [20] Yanan Li and Jignesh M. Patel. Bitweaving: fast scans for main memory data processing. In Kenneth A. Ross, Divesh Srivastava, and Dimitris Papadias, editors, *Proceedings of the ACM SIGMOD International Conference on Management of Data, SIGMOD 2013, New York, NY, USA, June 22-27, 2013*, pages 289–300. ACM, 2013.
- [21] Panagiotis Liakos, Katia Papakonstantinou, and Yannis Kotidis. Chimp: Efficient lossless floating point compression for time series databases. *Proc. VLDB Endow.*, 15(11):3058–3070, 2022.
- [22] Chunwei Liu, Hao Jiang, John Paparrizos, and Aaron J. Elmore. Decomposed bounded floats for fast compression and queries. *Proc. VLDB Endow.*, 14(11):2586–2598, 2021.
- [23] Magnetic. <https://doi.org/10.24432/C5DW43>, 2024.
- [24] Igor Pavlov. Lzma sdk (software development kit). <https://www.7-zip.org/sdk.html>, 2008. Available online.
- [25] Igor Pavlov. <https://www.7-zip.org/>, 2024.
- [26] Tuomas Pelkonen, Scott Franklin, Paul Cavallaro, Qi Huang, Justin Meza, Justin Teller, and Kaushik Veeraraghavan. Gorilla: A fast, scalable, in-memory time series database. *Proc. VLDB Endow.*, 8(12):1816–1827, 2015.
- [27] Stocks. <https://www.kaggle.com/datasets/tadakasuryateja/nifty-50-stocks>, 2024.
- [28] UCI. <https://archive.ics.uci.edu>, 2024.
- [29] Chen Wang, Jialin Qiao, Xiangdong Huang, Shaoxu Song, Haonan Hou, Tian Jiang, Lei Rui, Jianmin Wang, and Jiaguang Sun. Apache iotdb: A time series database for iot applications. *Proc. ACM Manag. Data*, 1(2):195:1–195:27, 2023.
- [30] Haoyu Wang and Shaoxu Song. Frequency domain data encoding in apache iotdb. *Proc. VLDB Endow.*, 16(2):282–290, 2022.
- [31] Tianrui Xia, Jinzhao Xiao, Yuxiang Huang, Changyu Hu, Shaoxu Song, Xiangdong Huang, and Jian-min Wang. Time series data encoding in apache iotdb: comparative analysis and recommendation. *VLDB J.*, 33(3):727–752, 2024.
- [32] Jinzhao Xiao, Wendi He, Shaoxu Song, Xiangdong Huang, Chen Wang, and Jianmin Wang. REGER: reordering time series data for regression encoding. In *40th IEEE International Conference on Data Engineering, ICDE 2024, Utrecht, The Netherlands, May 13-16, 2024*, pages 1242–1254. IEEE, 2024.
- [33] Jinzhao Xiao, Yuxiang Huang, Changyu Hu, Shaoxu Song, Xiangdong Huang, and Jianmin Wang. Time series data encoding for efficient storage: A comparative analysis in apache iotdb. *Proc. VLDB Endow.*, 15(10):2148–2160, 2022.
- [34] Hao Yan, Shuai Ding, and Torsten Suel. Inverted index compression and query processing with optimized document ordering. In Juan Quemada, Gonzalo León, Yoëlle S. Maarek, and Wolfgang Nejdl, editors, *Proceedings of the 18th International Conference on World Wide Web, WWW 2009, Madrid, Spain, April 20-24, 2009*, pages 401–410. ACM, 2009.
- [35] Jiangong Zhang, Xiaohui Long, and Torsten Suel. Performance of compressed inverted list caching in search engines. In Jinpeng Huai, Robin Chen, Hsiao-Wuen Hon, Yunhao Liu, Wei-Ying Ma, Andrew Tomkins, and Xiaodong Zhang, editors, *Proceedings of the 17th International Conference on World Wide Web, WWW 2008, Beijing, China, April 21-25, 2008*, pages 387–396. ACM, 2008.
- [36] Xin Zhao, Jialin Qiao, Xiangdong Huang, Chen Wang, Shaoxu Song, and Jianmin Wang. Apache tsfile: An iot-native time series file format. *Proc. VLDB Endow.*, 17(12):4064–4076, 2024.
- [37] Jacob Ziv and Abraham Lempel. A universal algorithm for sequential data compression. *IEEE Trans. Inf. Theory*, 23(3):337–343, 1977.
- [38] Marcin Zukowski, Sándor Héman, Niels Nes, and Peter A. Boncz. Super-scalar RAM-CPU cache compression. In Ling Liu, Andreas Reuter, Kyu-Young Whang, and Jianjun Zhang, editors, *Proceedings of the 22nd International Conference on Data Engineering, ICDE 2006, 3-8 April 2006, Atlanta, GA, USA*, page 59. IEEE Computer Society, 2006.

[4] Code.  
<https://github.com/apache/iotdb/tree/research/encoding-outlier>, 2024. [5] Code1.  
<https://github.com/apache/tsfile/tree/research/encoding-outlier>, 2024.



## APPENDIX

**Proposition 2.** For any solution  $(x_l, x_u)$  with  $\beta \leq \gamma$ ,  $x_l \in X$  and  $x_u \in X$ , there always exists another solution  $(x_l, x'_u)$  having  $C(x_l, x'_u) \leq C(x_l, x_u)$ , where  $x'_u = \min X_c + 2^\beta$ .

*Proof.* According to  $\beta = \lceil \log(\max X_c - \min X_c + 1) \rceil$  in Formula 8, we have

$$\begin{aligned} \log(\max X_c - \min X_c + 1) &\leq \beta \\ \max X_c - \min X_c + 1 &\leq 2^\beta \\ \max X_c &\leq \min X_c + 2^\beta - 1 \\ \max X_c &< x'_u. \end{aligned}$$

(1) For  $x_u > x'_u$ , it follows  $\max X_c < x'_u < x_u = \min X_u$ . Since there is no value between  $\max X_c$  and  $\min X_u$  in  $X$ , according to Definitions 2 and 4, we have  $\min X'_u = \min X_u$ .

(2) For  $x_u \leq x'_u$ , referring to Definition 4, we have  $\max X'_u \geq \max X_u$ .

Combining the above two cases, we can conclude that

$$\min X'_u \geq \min X_u.$$

For  $n_u = |X_u|$  and  $n'_u = |X'_u|$  introduced after Definition 4, it follows  $n_u \geq n'_u$ .

Let  $n_\Delta = |X_u \setminus X'_u|$  be the size of the increment, having  $n_\Delta = n_u - n'_u \geq 0$ .

Given the same  $x_l$  and the corresponding identical  $X_l, n_l$ , we could get the difference  $C_\Delta$  between  $C(x_l, x'_u)$  and  $C(x_l, x_u)$  defined in Formula 5,

$$\begin{aligned} C_\Delta &= C(x_l, x'_u) - C(x_l, x_u) \\ &= n_l(\lceil \log(\max X_l - x_{\min} + 1) \rceil + 1) \\ &\quad + n'_u(\lceil \log(x_{\max} - \min X'_u + 1) \rceil + 1) \\ &\quad + (n - n_l - n'_u)\lceil \log(\max X'_c - \min X'_c + 1) \rceil \\ &\quad - n_l(\lceil \log(\max X_l - x_{\min} + 1) \rceil + 1) \\ &\quad - n_u(\lceil \log(x_{\max} - \min X_u + 1) \rceil + 1) \\ &\quad - (n - n_l - n_u)\lceil \log(\max X_c - \min X_c + 1) \rceil. \end{aligned} \quad (12)$$

The same  $x_l$  also infers  $\min X'_c = \min X_c$ . Together with  $n_u = n'_u + n_\Delta$ , we have

$$C_\Delta = C_1 - C_2, \quad (13)$$

where

$$\begin{aligned} C_1 &= (n - n_l - n_u)\lceil \log(\max X'_c - \min X_c + 1) \rceil \\ &\quad + n_\Delta \lceil \log(\max X'_c - \min X_c + 1) \rceil \\ &\quad + n'_u(\lceil \log(x_{\max} - \min X'_u + 1) \rceil + 1) \end{aligned}$$

and

$$\begin{aligned} C_2 &= (n - n_l - n_u)\lceil \log(\max X_c - \min X_c + 1) \rceil \\ &\quad - n_\Delta(\lceil \log(x_{\max} - \min X_u + 1) \rceil + 1) \\ &\quad - n'_u(\lceil \log(x_{\max} - \min X_u + 1) \rceil + 1) \\ &= (n - n_l - n_u)\beta - n_\Delta(\gamma + 1) - n'_u(\gamma + 1). \end{aligned}$$

(i) Referring to Definition 2, we have  $\max X'_c < x'_u = \min X_c + 2^\beta$ . It follows

$$\begin{aligned} \log(\max X'_c - \min X_c + 1) &\leq \log(2^\beta) \\ \lceil \log(\max X'_c - \min X_c + 1) \rceil &\leq \beta. \end{aligned}$$

(ii) With the aforesaid proved  $\min X'_u \geq \min X_u$ , we infer

$$\lceil \log(x_{\max} - \min X'_u + 1) \rceil \leq \lceil \log(x_{\max} - \min X_u + 1) \rceil = \gamma.$$

Applying the above two conditions, we further derive

$$\begin{aligned} C_\Delta &\leq (n - n_l - n_u)\beta + n_\Delta\beta + n'_u(\gamma + 1) \\ &\quad - (n - n_l - n_u)\beta - n_\Delta(\gamma + 1) - n'_u(\gamma + 1) \\ &= n_\Delta(\beta - \gamma - 1) \leq 0. \end{aligned}$$

Given  $\beta \leq \gamma$  and  $n_\Delta \geq 0$ , the conclusion is proved.  $\square$

**Proposition 3.** For any solution  $(x_l, x_u)$  with  $\beta > \gamma$ ,  $x_l \in X$  and  $x_u \in X$ , there always exists another solution  $(x_l, x'_u)$  having  $C(x_l, x'_u) \leq C(x_l, x_u)$ , where  $x'_u = x_{\max} - 2^\gamma + 1$ .

*Proof.* According to  $\gamma = \lceil \log(x_{\max} - \min X_u + 1) \rceil$  in Formula 9, we have

$$\begin{aligned} \log(x_{\max} - \min X_u + 1) &\leq \gamma \\ x_{\max} - 2^\gamma + 1 &\leq \min X_u \\ x'_u &\leq \min X_u = x_u. \end{aligned}$$

(1) For  $x'_u = x_u = \min X_u$ , it is exactly the  $(x_l, x_u)$  solution, having  $\min X'_u = x'_u = \min X_u$ ,  $\max X'_c = \max X_c$ .

(2) For  $\max X_c < x'_u < x_u = \min X_u$ , since there is no value between  $\max X_c$  and  $\min X_u$  in  $X$ , according to Definitions 2 and 4, we have  $\min X'_u = \min X_u$ ,  $\max X'_c = \max X_c$  as well.

(3) For  $x'_u \leq \max X_c < x_u$ , referring to Definitions 2 and 4, it follows  $\max X'_c < \min X'_u \leq \max X_c < \min X_u$ .

Combining the above three cases, we can infer that

$$\begin{aligned} \min X'_u &\leq \min X_u \\ \max X'_c &\leq \max X_c. \end{aligned}$$

For  $n_u = |X_u|$  and  $n'_u = |X'_u|$  introduced after Definition 4, it follows  $n'_u \geq n_u$ . Let  $n_\Delta = |X'_u \setminus X_u|$  be the size of the increment, having  $n_\Delta = n'_u - n_u \geq 0$ .

Given the same  $x_l$  and the corresponding identical  $X_l, n_l$ , we could get the difference  $C_\Delta$  between  $C(x_l, x'_u)$  and  $C(x_l, x_u)$  defined in Formula 5,

$$\begin{aligned} C_\Delta &= C(x_l, x'_u) - C(x_l, x_u) \\ &= n_l(\lceil \log(\max X_l - x_{\min} + 1) \rceil + 1) \\ &\quad + n'_u(\lceil \log(x_{\max} - \min X'_u + 1) \rceil + 1) \\ &\quad + (n - n_l - n'_u)\lceil \log(\max X'_c - \min X'_c + 1) \rceil \\ &\quad - n_l(\lceil \log(\max X_l - x_{\min} + 1) \rceil + 1) \\ &\quad - n_u(\lceil \log(x_{\max} - \min X_u + 1) \rceil + 1) \\ &\quad - (n - n_l - n_u)\lceil \log(\max X_c - \min X_c + 1) \rceil. \end{aligned}$$

The same  $x_l$  also infers  $\min X'_c = \min X_c$ . Together with  $n'_u = n_u + n_\Delta$ , we have

$$C_\Delta = C_1 - C_2,$$

where

$$\begin{aligned} C_1 &= (n - n_l - n'_u) \lceil \log(\max X'_c - \min X_c + 1) \rceil \\ &\quad + n_\Delta (\lceil \log(x_{\max} - \min X'_u + 1) \rceil + 1) \\ &\quad + n_u (\lceil \log(x_{\max} - \min X'_u + 1) \rceil + 1) \end{aligned}$$

and

$$\begin{aligned} C_2 &= (n - n_l - n'_u) \lceil \log(\max X_c - \min X_c + 1) \rceil \\ &\quad - n_\Delta \lceil \log(\max X_c - \min X_c + 1) \rceil \\ &\quad - n_u (\lceil \log(x_{\max} - \min X_u + 1) \rceil + 1) \\ &= (n - n_l - n'_u) \beta - n_\Delta \beta - n_u (\gamma + 1). \end{aligned}$$

(i) Referring to Definition 2, we have  $x'_u = x_{\max} - 2^\gamma + 1 \leq \min X'_u$ . It follows

$$\begin{aligned} \log(x_{\max} - \min X'_u + 1) &\leq \log(2^\gamma) \\ \lceil \log(x_{\max} - \min X'_u + 1) \rceil &\leq \gamma. \end{aligned}$$

(ii) With the aforesaid proved  $\max X'_c \leq \max X_c$ , we infer  $\lceil \log(\max X'_c - \min X_c + 1) \rceil \leq \lceil \log(\max X_c - \min X_c + 1) \rceil = \beta$ .

Applying the above two conditions, we further derive

$$\begin{aligned} C_\Delta &\leq (n - n_l - n'_u) \beta + n_\Delta (\gamma + 1) + n_u (\gamma + 1) \\ &\quad - (n - n_l - n'_u) \beta - n_\Delta \beta - n_u (\gamma + 1) \\ &= n_\Delta (\gamma + 1 - \beta) \leq 0. \end{aligned}$$

Given  $\beta > \gamma$  and  $n_\Delta \geq 0$ , the conclusion is proved.  $\square$

**Proposition 4.** For normal distribution  $X \sim N(\mu, \sigma^2)$ , the approximation ratio  $\rho$  of BOS-M satisfies

$$\rho \leq \begin{cases} 2 & \text{if } \sigma \leq \frac{5}{3}, \\ \lceil \log(3\sigma - 1) \rceil & \text{otherwise.} \end{cases}$$

*Proof.* (1) Firstly, we prove the upper bound of the storage cost  $C_{\text{approx}}$  for BOS-M,

$$C_{\text{approx}} \leq \begin{cases} \lceil \log(6\sigma + 1) \rceil n & \text{if } \sigma < \frac{1}{2}, \\ 2n & \text{if } \frac{1}{2} \leq \sigma \leq \frac{5}{3}, \\ \lceil \log(3\sigma - 1) \rceil n & \text{otherwise.} \end{cases}$$

For normal distribution  $X \sim N(\mu, \sigma^2)$ , the median is  $\mu$ , the maximum value  $x_{\max}$  is approximately  $\mu + 3\sigma$ , and the minimum value  $x_{\min}$  is approximately  $\mu - 3\sigma$ , which correspond to the 99.7% confidence interval under the empirical rule (within three standard deviations from the mean).

The storage cost  $C_\beta$  of BOS-M with bit-width  $\beta$  is

$$\begin{aligned} C_\beta &= C(\mu - 2^\beta, \mu + 2^\beta) \\ &= n_l \lceil \log(\mu - 2^\beta - x_{\min} + 1) \rceil \\ &\quad + n_u \lceil \log(x_{\max} - (\mu + 2^\beta) + 1) \rceil \\ &\quad + (n - n_l - n_u) \lceil \log((\mu + 2^\beta) - (\mu - 2^\beta) + 1) \rceil \\ &= n_l \lceil \log(3\sigma - 2^\beta + 1) \rceil + n_u \lceil \log(3\sigma - 2^\beta + 1) \rceil \\ &\quad + (n - n_l - n_u) (\beta + 1) \\ &= (n_l + n_u) \lceil \log(3\sigma - 2^\beta + 1) \rceil \\ &\quad + (n - n_l - n_u) (\beta + 1). \end{aligned}$$

With  $\beta$  increasing from 1 to  $\lceil \log(6\sigma + 1) \rceil$ , bit-widths of 3 parts of values decrease firstly, and then center values become smaller. Thus,  $C_\beta$  firstly decreases and then increases. The upper bound of  $C_{\text{approx}}$  is thus

$$C_{\text{approx}} \leq \min\{C_{\lceil \log(6\sigma + 1) \rceil}, C_1\},$$

where

$$C_{\lceil \log(6\sigma + 1) \rceil} = \lceil \log(6\sigma + 1) \rceil n,$$

and

$$C_1 = (n_l + n_u) \lceil \log(3\sigma - 1) \rceil + 2(n - n_l - n_u).$$

Considering 4 different cases below, we rewrite  $C_1$  as

$$C_1 = 2n + (\lceil \log(3\sigma - 1) \rceil - 2)(n_l + n_u).$$

a) When  $\sigma \leq \frac{1}{2}$  ( $\mu - 2 < \mu - 3\sigma$ ), i.e., there are no upper and lower outliers with  $\beta = 1$ , we have  $n_l + n_u = 0$  and  $C_1 = 2n \geq \lceil \log(6\sigma + 1) \rceil n = C_{\lceil \log(6\sigma + 1) \rceil}$ .

b) When  $\frac{1}{2} \leq \sigma < \frac{5}{3}$  ( $\mu - 2 < \mu - 3\sigma$ ), i.e., there are no upper and lower outliers with  $\beta = 1$ , we have  $n_l + n_u = 0$  and  $C_1 = 2n < C_{\lceil \log(6\sigma + 1) \rceil}$ .

c) When  $\frac{2}{3} \leq \sigma \leq \frac{5}{3}$ , we have  $0 \leq \lceil \log(3\sigma - 1) \rceil \leq 2$  and  $2(n - n_l - n_u) \leq C_1 \leq 2n < C_{\lceil \log(6\sigma + 1) \rceil}$ .

d) When  $\sigma > \frac{5}{3}$ , we have  $\lceil \log(3\sigma - 1) \rceil \geq 2$  and  $C_1$  increases with  $\sigma$  growing. Then, when  $\sigma$  tends to positive infinity, we have there are no center values and  $C_1 = \lceil \log(3\sigma - 1) \rceil n < C_{\lceil \log(6\sigma + 1) \rceil}$ .

Therefore, we conclude that

$$C_{\text{approx}} \leq \begin{cases} \lceil \log(6\sigma + 1) \rceil n & \text{if } \sigma < \frac{1}{2}, \\ 2n & \text{if } \frac{1}{2} \leq \sigma \leq \frac{5}{3}, \\ \lceil \log(3\sigma - 1) \rceil n & \text{otherwise.} \end{cases}$$

(2) Moreover, we can prove that  $C_{\text{opt}} \geq n$ .

The storage cost is larger than the sum of bit-width for each value, thus the optimal cost  $C_{\text{opt}}$  has

$$C_{\text{opt}} = \sum_{i=1}^n b_i,$$

where

$$b_i = \begin{cases} 1 & \text{if } x_i = x_{\min}, \\ \lceil \log(x_i - x_{\min} + 1) \rceil & \text{otherwise.} \end{cases}$$

Thus, we have  $C_{\text{opt}} \geq n$ , even when all values are the same.

(3) Finally, we derive that

$$\rho = \frac{C_{\text{approx}}}{C_{\text{opt}}} \leq \begin{cases} \lceil \log(6\sigma + 1) \rceil & \text{if } \sigma < \frac{1}{2}, \\ 2 & \text{if } \frac{1}{2} \leq \sigma \leq \frac{5}{3}, \\ \lceil \log(3\sigma - 1) \rceil & \text{otherwise.} \end{cases}$$

a) When  $\sigma < \frac{1}{2}$ , we have  $\rho \leq \lceil \log(6\sigma + 1) \rceil \leq 2$ .

b) When  $\frac{1}{2} \leq \sigma \leq \frac{5}{3}$ , we have  $\rho \leq 2$ .

c) When  $\sigma > \frac{5}{3}$ , we have  $\rho \leq \lceil \log(3\sigma - 1) \rceil$ .

To sum up, we conclude that

$$\rho = \frac{C_{\text{approx}}}{C_{\text{opt}}} \leq \begin{cases} 2 & \text{if } \sigma \leq \frac{5}{3}, \\ \lceil \log(3\sigma - 1) \rceil & \text{otherwise.} \end{cases}$$

□