# GreedyGD: Enhanced Generalized Deduplication for Direct Analytics in IoT

**Aaron Hurst, Daniel E. Lucani and Qi Zhang**
DIGIT Centre and the Department of Electrical and Computer Engineering
Aarhus University, Aarhus, Denmark
E-mail: {ah,daniel.lucani,qz}@ece.au.dk

April 13, 2023

## Abstract

Exponential growth in the amount of data generated by the Internet of Things currently pose significant challenges for data communication, storage and analytics and leads to high costs for organisations hoping to leverage their data. Novel techniques are therefore needed to holistically improve the efficiency of data storage and analytics in IoT systems. The emerging compression technique Generalized Deduplication (GD) has been shown to deliver high compression and enable direct compressed data analytics with low storage and memory requirements. In this paper, we propose a new GD-based data compression algorithm called GreedyGD that is designed for analytics. Compared to existing versions of GD, GreedyGD enables more reliable analytics with less data, while running $11.2\times$ faster and delivering even better compression.

*Keywords* Data Compression, Data Analytics, Generalized Deduplication, Clustering, IoT

## 1 Introduction

Sensor data collected via the Internet of Things (IoT) are increasingly relied upon to make critical business decisions, improve efficiency and keep people safe [1]. However, managing these data in a traditional, centralised manner becomes costly not only in terms of storage and compute, but also latency, as data generation outpaces improvements in hardware. Likewise, Edge servers, while able to reduce communication costs and latency [2], are also limited by their lower computing power and higher storage unit-cost. Overall, the cost of transmitting, storing and analysing IoT data may limit the utility and sustainability of IoT systems by forcing organisations to discard data or retain it only for a short time [3, 4]. Thus, there is a need for consolidated frameworks that reduce the cost of IoT data storage, while improving analytics latency.

The emerging compression technique Generalized Deduplication (GD) [5–7] provides an effective solution to this need. Namely, it offers high, lossless end-to-end compression in IoT systems, as well as efficient approximate analytics without the need for decompression [7–9]. Moreover, the structure of GD-compressed data allows data to be split between Edge and Cloud servers to optimise for storage cost and analytics latency [9]. GD performs compression by splitting data chunks into *bases*, which are deduplicated, and *deviations*, which are stored verbatim (see Fig. 1). Deviations are stored together with an ID that links them to their associated base, allowing for lossless decompression. Compressing a dataset with GD involves two stages: 1) configuration, i.e., deciding which bits to allocate to the base, and 2) compression, i.e., splitting data chunks into bases and deviations.

In this paper, we propose GreedyGD, which is a GD-based data compression algorithm designed for running analytics directly on compressed data. Compared to previous versions of GD, it addresses a number of critical pitfalls for compression speed and analytics quality, including configuration time and base-sample order preservation. GreedyGD even delivers better compression than existing versions of GD. The main contributions of this paper include:

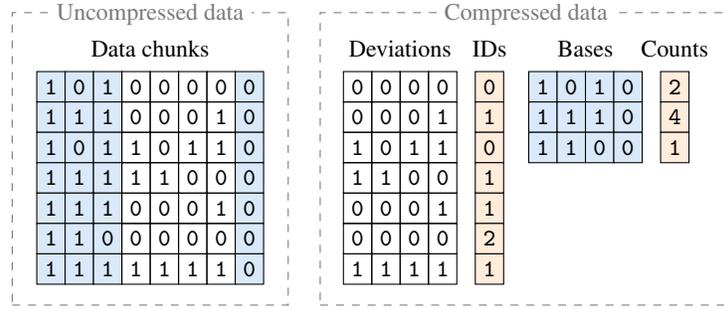1. The GreedyGD compression algorithm, which is designed for direct compressed data analytics,

Figure 1: GD applied to 8-bit data chunks. The first three bits and last bit (highlighted in blue) are allocated to the base and deduplicated, while the remaining bits are the deviations and are stored verbatim alongside an ID linking them to the relevant bases.

2. Data preprocessing for GD that improves compression of floating point data,

3. Comprehensive evaluation GreedyGD's performance in comparison to existing GD versions and universal lossless compressors in terms of compression ratio, runtime and analytics, and

4. Evaluation of configuring GD using a data subset.

The paper is structured as follows. Section 2 discusses related work. Section 3 details our nomenclature. Section 4 describes GreedyGD. Section 5 evaluates the performance of GreedyGD. Finally, Section 6 concludes the paper.

## 2   Related Work

Data compression algorithms can be either lossy or lossless. Lossy algorithms typically offer more compression, but at the cost of degraded data fidelity, which may be unsuitable for critical IoT applications. We therefore focus on lossless compression. A wide variety of universal lossless compression algorithms are available [10–14], as well as powerful IoT-specific lossless time series compression algorithms, such as Sprintz [15]. However, existing algorithms are limited in terms of random access (important for analytics), small-chunk performance (relevant for IoT sensors) and their ability to run on resource-limited IoT devices [16].

GD-based data compression, on the other hand, performs well in these areas and also provides high compression. To achieve this, the choice of base bits, i.e., the set of bit indices mapped to the bases, is critical and should depend on both data and application. The version of GD described in [7], which we refer to as GD-INFO, selects base bits according to inter-bit correlation. It first maps all bits to the base and computes the compressed data size. Next, the bit with lowest inter-bit correlation is moved from the base to the deviation and compressed data size is re-computed. This continues until the first local minimum for compressed data size. The base bits at this point are then used for compression. To limit runtime, configuration is run on at most the first $10^6$ samples.

While GD-INFO delivers good performance, it has some significant limitations, including: 1) no efficient method for computing the compressed data size, 2) inter-bit correlation is not always an accurate indicator of deduplication potential, 3) terminating at the first local minimum is often sub-optimal, 4) no standard approach for multidimensional data, and 5) poor compressed data analytics performance, as shown in [9]. The effect of using a data subset (i.e., the first $10^6$ samples) for configuration has also not been investigated.

In [9], an analytics-tailored version of GD was proposed, which we refer to as GD-GLEAN. Different to GD-INFO, this compressor concedes some compression for improved analytics performance, but otherwise suffers from the same limitations as GD-INFO.

Two related algorithms that bear notable similarities to GD are ALACRITY [17] and DigitHist [18]. While both algorithms perform compression by extracting base-like "bins" from data, they are limited with respect to GD in that they place significant restrictions on which bits can be assigned to the these bins. Moreover, in contrast to GD, DigitHist uses lossy compression.

In this paper, we address the aforementioned limitations by proposing GreedyGD, which delivers significantly more accurate analytics on compressed data, while running faster and providing even better compression.

Table 1: Nomenclature

| Notation | Meaning | Notation | Meaning |
|---|---|---|---|
| $n$ | No. of data samples | $S$ | Compressed data size |
| $d$ | No. of dimensions | $\Delta$ | Maximum deviation |
| $n_b$ | No. of bases | $\lambda$ | Balancing factor |
| $B$ | Set of base bits | $\alpha$ | Exploration factor |
| $l_c$ | Bits per chunk | $\mathcal{C}$ | `GreedySelect` cost |
| $l_b$ | Bits per base | CR | Compression ratio |
| $l_d$ | Bits per deviation | ADR | Analytics data ratio |
| $l_{id}$ | Bits per base ID | AR | Approximation ratio |
| $l_{bc}$ | Bits per base count | AMI | Adjusted mutual info. |

## 3   Nomenclature

See Table 1 for a summary of our nomenclature. We consider a dataset to contain $n$ samples and $d$ dimensions. The set of bit indices that are allocated to bases by GD is referred to as the *base bits* and denoted $B$. The total number of (deduplicated) bases is denoted $n_b$. A binary data *chunk* is formed by concatenating a data sample across its dimensions and has length $l_c$ bits. The number of bits per base and deviation are denoted $l_b$ and $l_d$, respectively. Hence, $l_c = l_b + l_d$. GD stores base IDs and base counts using $l_{id} = \lceil \log_2(n_b) \rceil$ and $l_{bc} = \lceil \log_2(n) \rceil$ bits each, respectively.

The *maximum deviation*, $\boldsymbol{\Delta}$, is the maximum difference between samples mapped to the same base. For example, the minimum value mapped to first base (all zero deviation bits) in Fig. 1 is 10100000, or 160, while the maximum value (all one deviation bits) is 10111110, or 190. This gives a maximum deviation of 00011110, or 30, which is simply all deviation bits set to 1. Note, $\boldsymbol{\Delta}$ is a vector for multidimensional data and varies between bases for floating point data since the values of mantissa bits depend on the exponent.

The total compressed data size for GD, $S$, is equal to:

$$S = n_b \left( l_b + l_{bc} \right) + n \left( l_{id} + l_d \right) + S_{params}, \tag{1}$$

where $S_{params}$ is the size of GD's parameters (typically negligible). Compression performance is measured using the *compression ratio* (CR), defined as follows:

$$\text{CR} = \frac{\text{size of compressed data}}{\text{size of uncompressed data}}. \tag{2}$$

As such, lower values (below 1) are better. Additionally, the amount of (compressed) data that must be accessed in order to run analytics directly on GD-compressed data is measured using the *analytics data ratio* (ADR), which is equal to:

$$\text{ADR} = \frac{\text{size of data used for analytics}}{\text{size of uncompressed data}} \approx \frac{n_b(l_b + l_{bc})}{n l_c}. \tag{3}$$

We use $k$-means clustering [19] as a test case for the analytics performance of GreedyGD. Three performance metrics are used. The first, approximation ratio (AR), is defined in terms of the sum of squared errors (SSE) as follows:

$$\text{AR} = \frac{\text{clustering SSE using compressed data}}{\text{clustering SSE using uncompressed data}}, \tag{4}$$

As such, $\text{AR} \geq 1$ and lower values indicate better performance. A value of 1 indicates that no loss in clustering quality due to running clustering on compressed data instead of uncompressed data.

The second metric we use is adjusted mutual information (AMI), which is an information theoretic measure of the similarity between two clusterings. In this case, AMI is always between 0 and 1 and higher values are better, i.e., indicate more similar clusterings.

The third metric is the Silhouette coefficient, which measures how well data points fit with those in the same cluster compared to those in the next most appropriate cluster. More formally, for a given data point $\boldsymbol{x}$, if $a(\boldsymbol{x})$ is the mean distance between $\boldsymbol{x}$ and other points in the same cluster and $b(\boldsymbol{x})$ is the mean distance between $\boldsymbol{x}$ and points in the next best cluster, then the silhouette of $\boldsymbol{x}$ is:

$$s(\boldsymbol{x}) = \frac{b(\boldsymbol{x}) - a(\boldsymbol{x})}{\max\left(a(\boldsymbol{x}), b(\boldsymbol{x})\right)}. \tag{5}$$

---

**Algorithm 1** GreedyGD

---

**Inputs:** dataset $\mathcal{D}$, exploration factor $\alpha$, balancing factor $\lambda$
**Outputs:** compressed dataset $\mathcal{D}_{GD}$

  1: $\mathcal{D}' \leftarrow$ preprocess $\mathcal{D}$                                                     *// Subsection 4.3*
  2: $D \leftarrow$ extract subset from $\mathcal{D}'$                                     *// Subsection 4.4 (optional)*
  3: $B \leftarrow \texttt{GreedySelect}(D', \alpha, \lambda)$                               *// Subsection 4.2*
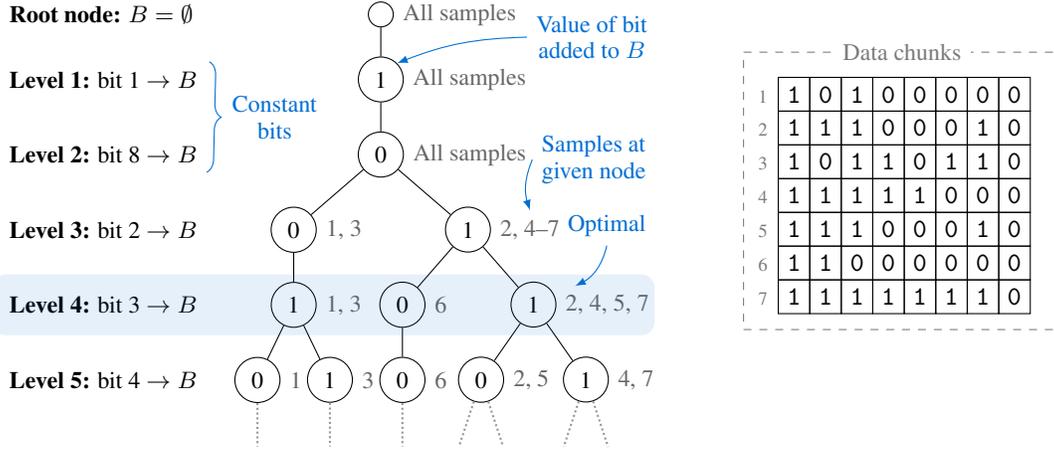  4: Apply GD to $\mathcal{D}'$ according to $B$

---



Figure 2: `BaseTree` partially applied to the data chunks from Fig. 1, which are repeated at the bottom of the figure for convenience. The highlighted optimal configuration, level 4, corresponds to the allocate of base and deviation bits in Fig. 1.

The Silhouette coefficient for an entire clustering is the mean of the silhouettes of all data points. The silhouette coefficient ranges between $\pm 1$ such that higher values are better and values close to 1 indicate dense, well-separated clusters.

## 4 GreedyGD Algorithm

The key to GD-based compression is selecting an effective set of base bits. This is challenging due to the number of possible combinations. For example, d-dimensional 32-bit data has $2^{32d}$ possible configurations, i.e., over 4 billion *per dimension*. Enumerating all of these is infeasible since one must count the number of unique bases for each possible combination, which is highly time consuming.

GD-INFO solves this problem by ordering bits according to inter-bit correlation and terminating at the first local minima for compressed data size. This limits the number of configurations enumerated to at most $l_c$. However, while this makes GD-INFO computationally feasible, the aforementioned limitations make it suboptimal.

The proposed GreedyGD compression algorithm, outlined at a high level in Algorithm 1, aims to find better base bit sets and evaluate them more quickly. In a nutshell, GreedyGD works by 1) applying novel data preprocessing, 2) optionally selecting a subset of the data to accelerate configuration, 3) selecting the base bits using a greedy algorithm and 4) compressing the data using GD.

The remainder of this section is structured as follows. Subsection 4.1 describes the `BaseTree` algorithm, which is used to quickly count bases. Subsection 4.2 describes the `GreedySelect`, which builds on `BaseTree` to select high quality base bits. Subsection 4.3 details how preprocessing can improve compression, while Subsection 4.4 details how using data subsets can improve runtime. Finally, Subsection 4.5 analyses the time complexity of GreedyGD.

### 4.1 Counting bases

To accelerate counting bases for GD configuration, we propose the `BaseTree` algorithm, which is illustrated in Fig. 2. The idea is to store the distribution of bases for a given base bits set in a tree structure. As shown in Fig. 2, the root

---

**Algorithm 2** `GreedySelect`

---

**Inputs:** dataset $\mathcal{D}$, exploration factor $\alpha$, balancing factor $\lambda$
**Outputs:** base bits $B$

 1: *// Initialisation*
 2: $B \leftarrow$ constant bits in $\mathcal{D}$
 3: Expand `BaseTree` with $B$
 4: $B_{\text{best}} \leftarrow B$
 5: $\mathcal{C}_{\text{best}} \leftarrow \infty$
 6: *// Iteratively add base bits until termination*
 7: **while** not all bits in $B$ **do**
 8:     $b_{\text{loc}} \leftarrow -1$                                                    *// reset lowest cost bit for local optimisation*
 9:     $\mathcal{C}_{\text{loc}} \leftarrow \infty$                                       *// reset lowest cost for local optimisation*
10:     **for** $i = 1$ to number of columns **do**
11:         $b_i \leftarrow$ most significant bit from column $i$ not in $B$
12:         $n_{b,i} \leftarrow n_b$ after adding $b_i$ to $B$                             *// use `BaseTree`*
13:         $S_i \leftarrow$ compressed data size after adding $b_i$ to $B$                *// Eq. 1*
14:         $\Delta_i' \leftarrow \Delta_i$ after adding $b_i$ to $B$                      *// Eq. 6*
15:         $\mathcal{C}_i \leftarrow$ cost after adding $b_i$ to $B$                      *// Eq. 7*
16:         **if** $\mathcal{C}_i < \mathcal{C}_{\text{loc}}$ **then**                     *// lowest cost bit in current iteration*
17:             $b_{\text{loc}} \leftarrow b_i$
18:             $\mathcal{C}_{\text{loc}} \leftarrow \mathcal{C}_i$
19:         **end if**
20:     **end for**
21:     **if** $\mathcal{C}_{\text{loc}} > (1 + \alpha) \cdot \mathcal{C}_{\text{best}}$ **then**    *// early termination*
22:         **return** $B_{\text{best}}$
23:     **end if**
24:     Append $b_{\text{loc}}$ to $B$
25:     Expand `BaseTree` with $b_{\text{loc}}$
26:     **if** $\mathcal{C}_{\text{loc}} < \mathcal{C}_{\text{best}}$ **then**             *// update global best*
27:         $B_{\text{best}} \leftarrow B$
28:         $\mathcal{C}_{\text{best}} \leftarrow \mathcal{C}_i$
29:     **end if**
30: **end while**
31: **return** $B_{\text{best}}$

---

node corresponds to $B = \emptyset$ and contains all data samples. Each subsequent layer corresponds to an additional bit added to $B$ and potentially additional nodes, depending on the values of the selected bit. For example, in level 2, only one node is created because bit 8 is always 0 (i.e. constant). In contrast, in level 3, the data are split between two child nodes because bit 3 contains both 0s and 1s. In level 4, only one child is spawned from the left parent node because samples 1 and 3 both have a 1 in bit 3, while the right parent node spawns two child nodes because bit 3 contains both 0s (left child node) and 1s (right child node) among samples 2 and 4–7.

This structure can be used to quickly count the number of bases for a given set of base bits $B$ since tree width corresponds to $n_b$. Counting bases, therefore, simply requires 1) adding the desired base bits to the tree and 2) counting the number of leaf nodes. Note that tree height corresponds to the number of base bits $l_b$ and root-to-leaf paths correspond to individual bases.

## 4.2   Selecting base bits

As noted previously, selecting appropriate base bits is critical to the performance of GD. For this task, we propose the `GreedySelect` algorithm, which is designed to optimise compressed data analytics performance while maintaining good compression. As outlined in Algorithm 2, `GreedySelect` begins by adding all constant bits (i.e. bits that have the same value for all samples) to $B$ and expanding `BaseTree` accordingly (lines 2–3). It then iteratively adds bits to $B$ while minimising a cost function until termination. Each iteration includes a local optimisation process followed by a global optimisation step.

In the local optimisation, `GreedySelect` considers each column of the data and identifies the most significant bit, $b_i$, that is currently mapped to the deviation (i.e. $b_i \notin B$, line 11). The number of bases, $n_{b,i}$, resulting from adding $b_i$ to $B$ is then computed using `BaseTree` (line 12) as well as the resulting compressed data size, $S_i$ (line 13). The

maximum deviation in column $i$ if $b_i$ is moved to the base, $\Delta'_i$, is then computed (line 14) as:

$$\Delta'_i = \Delta_i \oplus 2^{b_i}, \tag{6}$$

where $\Delta_i$ is the current maximum deviation in column $i$ and $\oplus$ is the binary XOR operator. Finally, the cost, $\mathcal{C}_i$, after adding $b_i$ to $B$ is computed (line 15) using the following cost function:

$$\mathcal{C}_i = \left(1 - \lambda \left(\Delta'_i / \Delta_i^{(0)}\right)^2\right) \cdot S_i, \tag{7}$$

where $\Delta_i^{(0)}$ is the maximum deviation in dimension $i$ after adding only the constant bits to $B$, and $\lambda$ is known as the *balancing factor* and defined such that $0 \leq \lambda < 1$.

The cost function in Eq. 7 is primarily dependent on $S$, which favours high compression. However, it is also scaled based on the maximum deviation such that it balances the number of bits added to $B$ from each dimension to improve analytics performance. In the first iteration, scaling has no effect since $\Delta'_i / \Delta_i^{(0)}$ is $1/2$ for all dimensions. In subsequent iterations, the ratio $\Delta'_i / \Delta_i^{(0)}$ will be lower for dimensions which have had bits added to $B$. Thus, the term $(1 - \lambda(\Delta'_i / \Delta_i^{(0)})^2)$ will be larger (resulting in higher cost) for dimensions that have had more bits added to $B$ and smaller (lower cost) for dimensions that have had fewer bits added to $B$. The magnitude of this effect is controlled by the balancing factor, $\lambda$. We found setting $\lambda = 0.02$ provided good results for most datasets. At this setting, once 3–4 bits are added to $B$ from each dimension, the balancing effect becomes negligible and $\mathcal{C}_i \approx S_i$. The local optimisation process is completed by keeping track of the lowest cost, $\mathcal{C}_{\text{loc}}$, and the corresponding bit, $b_{\text{loc}}$ (lines 16–19).

The first step in the global optimisation stage of `GreedySelect` is to check the termination condition (line 21). Termination is controlled by the *exploration factor*, $\alpha > 0$, which limits the amount of exploration beyond local minima. The algorithm terminates if $\mathcal{C}_{\text{loc}}$ exceeds $1 + \alpha$ times the current global best cost, $\mathcal{C}_{\text{best}}$. We found setting $\alpha = 0.1$ performed well for most datasets. If not terminated, `GreedySelect` adds $b_{\text{loc}}$ to $B$, updates `BaseTree` and, if required, updates $\mathcal{C}_{\text{best}}$ (lines 24–29). Overall, up to $l_c d$ base bits combinations are examined ($d$ per iteration and up to $l_c$ iterations).

By selecting base bits from most to least significant, `GreedySelect` guarantees *order preservation* [8]. This means that for any two data samples, $x_1$ and $x_2$, and their corresponding bases, $\text{base}(x_1)$ and $\text{base}(x_2)$,

$$x_1 < x_2 \Rightarrow \text{base}(x_1) \leq \text{base}(x_2). \tag{8}$$

Order preservation is critical for analytics performance. Without it, results can be orders of magnitude worse. Overall, by providing order preservation, balancing the number of base bits from each dimension and basing the cost function on compressed data size, `GreedySelect` effectively balances the objectives of compressed data analytics and compression.

### 4.3 Data preprocessing

Transforming data into a more amenable form is a typical step in compression algorithms. In the case of GD, floating point data can be compressed significantly better by applying scaling and then converting to integers. For example, consider the data illustrated in Fig. 3, which is from the *Aarhus CityLab* dataset [20]. The original data, Fig. 3(a), has just two decimal places but only one constant bit. After scaling to remove the decimal places, Fig. 3(b), many more bits are constant, meaning that they can be assigned to the base by GD, improving compression performance. Finally, converting the data to integer results in even more constant bits.

### 4.4 Configuration on a data subset

For large and complex datasets, configuring GD can still take substantial time. In cases where this is prohibitive, we follow the suggestion of [7] to use a subset of the data. However, we note that this introduces two important risks: 1) the number of decimal places in the dataset could by greater than in the data subset, resulting in inaccurate preprocessing, and 2) constant bits in the data subset may vary elsewhere in the dataset, which could result in violating order preservation. In light of these risks, we propose using the full dataset to select preprocessing and determine constant bits and then run the remainder of `GreedySelect` on a subset of the data.

### 4.5 Complexity

The time complexity for configuring GreedyGD (i.e. running `GreedySelect`) is as follows. First, the constant bits must be identified, which requires examining all $n l_c$ bits. Next, the base tree must be extended with these bits, but

0.39     00111110110001111010111000010100

37.83    01000010000101110101000111101100

98.92    01000010110001011101011100001010

(a) Original data (1 constant bit)

39.      01000010000111000000000000000000

3,783.   01000101011011000111000000000000

9,892.   01000110000110101001000000000000

(b) After scaling (14 constant bits)

39       00000000000000000000000000100101

3,783    00000000000000000000111011000111

9,892    00000000000000000010011010100100

(c) Integer (18 constant bits)

Figure 3:   Selected data points from the first dimension of the *Aarhus Citylab* dataset [20] in decimal and IEEE-754 [21] floating point representation (a) before and (b) after scaling and (c) as integers. Bits that are constant across the entire dataset are highlighted in blue.

since they are constant, there is no need to split the tree and only nodes for the constant bits must be added. In the worst case this is $l_c$ (all bits constant). The main while loop in GreedySelect runs up to $l_c$ times before all bits are in $B$. For each iteration, the for loop runs $d$ times. Within the for loop, all steps are $O(1)$ except for finding $n_{b,i}$, which requires checking whether to split each node in BaseTree. In the worse case, this is $O(n)$ since one bit needs to be checked across potentially every sample. The global update step is constant time, except for updating BaseTree, which is $O(n)$ in the worse case. Overall then, the while loop has a worst case time complexity of $O(ndl_c)$. Therefore, the total worst case time complexity for GreedyGD configuration is $O(nl_c + l_c + ndl_c)$ or $O(ndl_c)$. Noting that $l_c$ is a multiple of $d$ ($32d$ for single precision data and $64d$ for double precision data), we can simplify this to $O(nd^2)$. In practice, the while loop runs substantially fewer than $l_c$ times, with the number of iterations reduced by both the number of constant bits and early termination. The local optimisation loop may also have fewer iterations if all bits from one or more dimensions are added to the base.

Configuring GreedyGD on a data subset can reduce its time complexity. Assuming the full dataset is used for preprocessing and identification of constant bits and a data subset is used thereafter, then the time complexity of GreedySelect is $O(nl_c + l_c + ndrl_c)$ or $O(nl_c(1 + dr))$, where $0 < r < 1$ is the fraction of the data contained in the subset. If $r < 1/d$, then the worst case time complexity reduces to $O(nd)$.

Compression is comparatively straightforward as it only involves splitting the base bits from the deviation bits and adding base IDs and counts. Nonetheless, it requires accessing the full binary data and thus has time complexity $O(nd)$.

## 5   Performance Evaluation

We evaluated the performance of GreedyGD in terms of compression (Subsection 5.1), compressed data analytics (Subsection 5.2), configuration runtime (Subsection 5.3) and configuration using data subsets (Subsection 5.4). In total, 18 datasets were used for evaluation (Table 2), which include a range of sizes, dimensionalities, data types and precisions. Note that non-numerical attributes and rows with missing values were omitted. The *Chicago beach weather* dataset was also split between float and integer columns and some derived fields were omitted from the *Chicago taxi trips* dataset. All evaluations were performed on a laptop with an Intel Core i7-10510U 1.80 GHz CPU and 16 GB of RAM using Python 3.8.10.

First and foremost, GreedyGD was compared to previous versions of GD, namely GD-INFO [7] and GD-GLEAN [9]. Note that, to apply GD-INFO to multidimensional data, it was necessary to modify its termination criteria such that it explores slightly beyond local minima (similar to GreedyGD). We also compared to enhanced version of GD-INFO and GD-GLEAN that used preprocessing and counted bases using BaseTree. We refer to these compressors as GD-INFO+ and GD-GLEAN+, respectively. To use BaseTree, it was necessary to reverse the order of iteration so that GD-INFO+ starts with $B = \emptyset$ and iteratively adds bits, rather than starting with all bits in $B$ and iteratively removing bits. Throughout our evaluation, we set $\alpha = 0.1$ and $\lambda = 0.02$ for GreedyGD.

Table 2: Datasets used for evaluating GreedyGD.

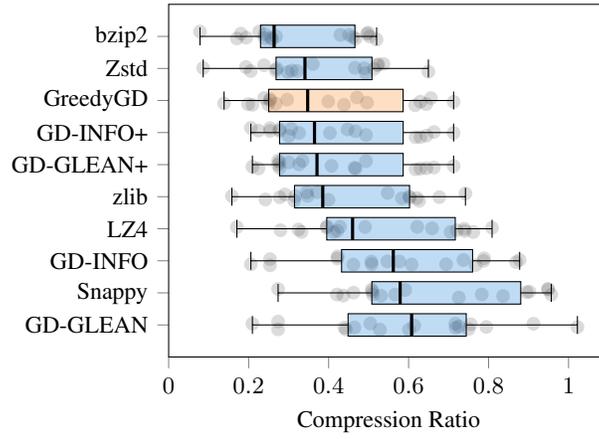| Dataset | Type | Precision | $n$ | $d$ | Size (kB) |
|---|---|---|---|---|---|
| Aarhus Citylab [20] | float | single | 26,387 | 4 | 422 |
| Aarhus pollution 172156 [22] | int | single | 17,568 | 5 | 351 |
| Aarhus pollution 204273 [22] | int | single | 17,568 | 5 | 351 |
| Chicago beach water I [23] | float | single | 39,829 | 5 | 797 |
| Chicago beach water II [23] | float | single | 10,034 | 6 | 241 |
| Chicago beach weather [24] | float | single | 86,694 | 9 | 3,121 |
| Chicago beach weather [24] | int | single | 86,763 | 5 | 1,735 |
| Chicago taxi trips [25] | float | double | 3,466,498 | 10 | 277,320 |
| CMU IMU acceleration [26] | float | single | 134,435 | 3 | 1,613 |
| CMU IMU velocity [26] | float | single | 134,435 | 3 | 1,613 |
| CMU IMU magnetic [26] | float | single | 134,435 | 3 | 1,613 |
| CMU IMU position [26] | float | single | 134,435 | 4 | 2,151 |
| CMU IMU all [26] | float | single | 134,435 | 13 | 6,991 |
| COMBED mains power [27] | float | double | 82,888 | 3 | 995 |
| COMBED UPS power [27] | float | double | 86,199 | 3 | 1,035 |
| Melbourne city climate [28] | float | single | 56,570 | 3 | 679 |
| Gas turbine emissions [29] | float | single | 36,733 | 11 | 1,616 |
| Household power usage [30] | float | single | 2,049,280 | 7 | 57,380 |



Figure 4: Box plots of CR for each compressor ordered from best (top) to worst (bottom) based on median CR.
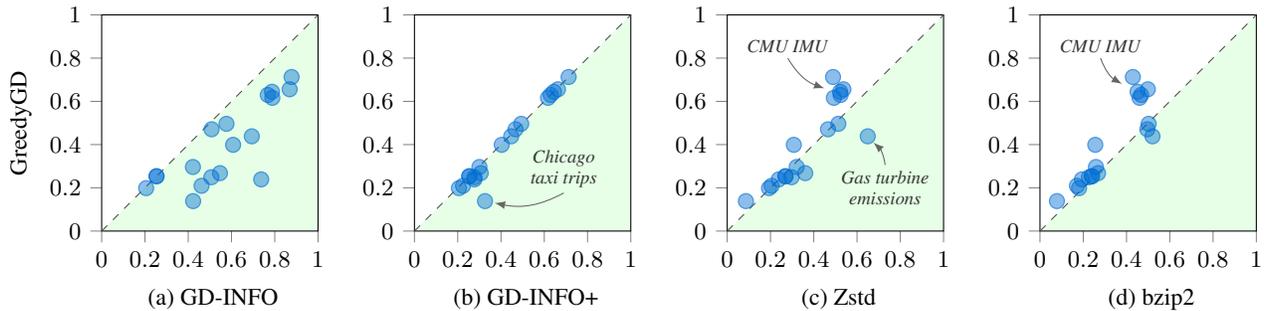


Figure 5: Pairwise comparison of compression ratios for GreedyGD and selected compressors across all 18 datasets with noteworthy datasets annotated. Points below (above) the diagonal indicate better (worse) compression with GreedyGD.

## 5.1 Compression Ratio

The compression performance of GreedyGD was compared to the other GD versions mentioned above, as well as five widely-used universal compressors, namely zlib [10], Zstd [11], Snappy [12], bzip2 [13] and LZ4 [14]. All universal
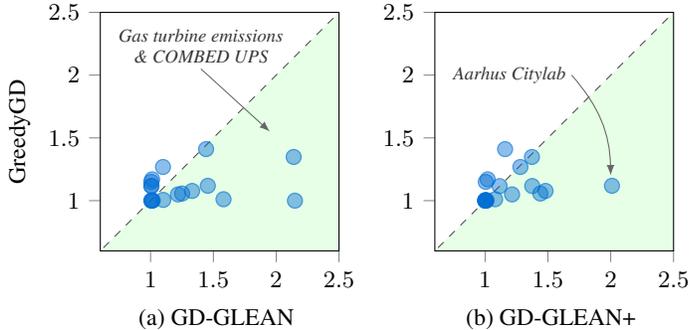
Figure 6: Pairwise comparison of AR for $k$-means clustering on compressed data using GreedyGD and (a) GD-GLEAN and (b) GD-GLEAN+ across all 18 datasets. Points below (above) the diagonal indicate better (worse) analytics performance with GreedyGD.
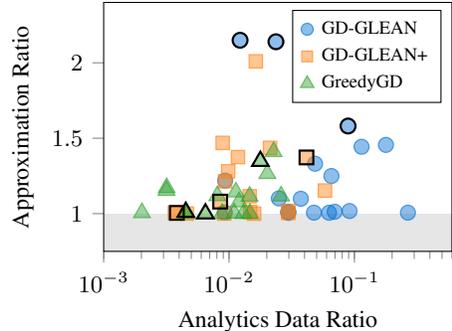
Figure 7: Joint distribution of analytics data ratio (ADR) and clustering approximation ratio (AR) for GreedyGD and both GD-GLEAN versions across all datasets. The three datasets with the worst performance under GD-GLEAN are highlighted with a black border for all three compressors.

Table 3: Summary of compression & analytics results. Best (second best) performance shown in bold (italics).

| Compressor | CR | ADR | AR | AMI | Silhouette |
|---|---|---|---|---|---|
| GreedyGD | **0.348** | *0.011* | 1.109 | **0.758** | *0.309* |
| GD-INFO+ | *0.365* | **0.009** | 1.246 | 0.630 | 0.294 |
| GD-GLEAN+ | 0.371 | 0.013 | **1.097** | *0.753* | 0.295 |
| GD-INFO | 0.562 | 0.014 | 1.806 | 0.368 | 0.162 |
| GD-GLEAN | 0.608 | 0.048 | *1.099* | 0.639 | **0.313** |

compressors were run in one-shot mode at maximum compression. The distributions of compression ratios (CRs) resulting from applying these compressors to all 18 datasets are illustrated in Fig. 4, where grey dots indicate CR values for individual datasets. As can be seen, GreedyGD is competitive with the best of the universal compressors and outperforms all other GD versions.

Fig. 5 provides more granulated comparisons. Each point in these figures represents the CR for a single dataset with GreedyGD (y-axis) and another compressor (x-axis). Points below the diagonal therefore represent better performance with GreedyGD. As can be seen in Fig. 5(a), GreedyGD substantially outperforms GD-INFO, which is the existing state-of-the-art version of GD in terms of compression. Indeed, with GreedyGD having a median CR of 0.348 compared to GD-INFO's 0.562, GreedyGD provides $1.6\times$ more compression on average. Fig. 5(b), which compares GreedyGD to GD-INFO+, shows that, even with data preprocessing, GD-INFO never outperforms GreedyGD. Fig. 5(c) and (d) compare GreedyGD to Zstd and bzip2, respectively. As can be seen, GreedyGD has very similar CR for most datasets compared to these compressors, with the exception of the IMU data. In particular, GreedyGD has very similar performance to Zstd, which has a median CR of 0.341. Note that GD can achieve even more compression if applying entropy coding to the base IDs, however, this compromises random access and thereby analytics efficiency [7].

## 5.2 Analytics on compressed data

Approximations to various analytics tasks can be obtained directly from GD-compressed data without decompressing it. This is done by performing analytics on the bases weighted by their counts [8]. A significant benefit of this approach is that only a fraction as many points need to be analysed, since the number of bases is typically far less than the number of samples, i.e. $n_b << n$. This translates into savings in runtime and memory [9]. Thus, to evaluate the effectiveness of compressed data analytics, we consider both the *quality* of approximate analytics results and the *amount* of data that needs to be accessed to perform analytics.
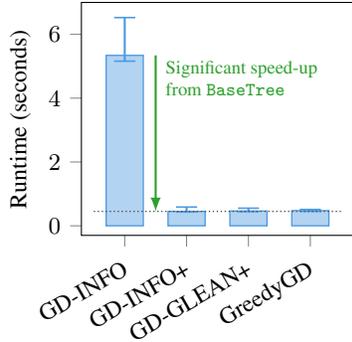
Figure 8: Configuration runtime for GD-based compressors for the *COMBED mains power* dataset.
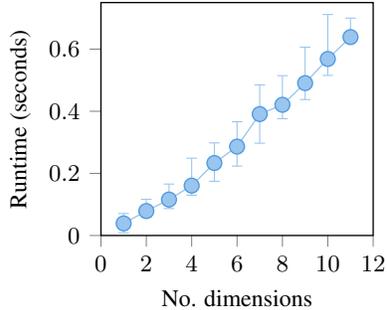


Figure 9: GreedyGD configuration runtime versus dimensionality averaged over subsamples of the dimensions of the *Gas turbine emissions* dataset.
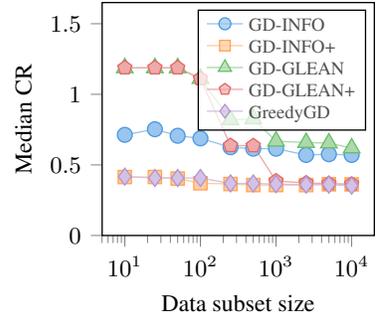


Figure 10: Median full dataset CR versus size of the training set used to configure each compressor.

The median amount of data needed for analytics for each GD-based compressor, expressed via the analytics data ratio (ADR), is shown in Table 3. As can be seen, GreedyGD requires a mere 1.1% of the uncompressed data to perform analytics, which is less than a quarter of that required by GD-GLEAN and 15% less than GD-GLEAN+.

Existing GD compressors have been shown to support high clustering performance [8, 9]. We have thus thus used $k$-means clustering [19] to evaluate GreedyGD's analytics performance. Our methodology mimicked that of [9], such that we computed cluster centres directly from the bases and used these to cluster the original data points. Clusterings of the original data were then used to compute Silhouette coefficients and then compared to clustering performed on uncompressed data to compute AR and AMI. Note that each clustering was repeated 10 times with 100 initialisations each to minimise variance and sampling with $n = 10,000$ was used to compute the Silhouette coefficient to avoid excessive runtime. The median results across all 18 datasets are presented in Table 3. As can be seen, GreedyGD provides best or near-best performance in each clustering performance metric.

The performance of GreedyGD in comparison to GD-GLEAN and GD-GLEAN+ was also examined more closely. Fig. 6(a) and (b) show the AR for each dataset using GreedyGD (y-axis) compared to GD-GLEAN and GD-GLEAN+, respectively (x-axis). As can be seen, GreedyGD generally delivers similar or better performance and has a more reliable distribution of AR values with no outliers.

Overall, GreedyGD provides solid performance across all analytics-related metrics we tested. To illustrate this more holistically, we also present Fig. 7, which plots AR against ADR for GreedyGD, GD-GLEAN and GD-GLEAN+ across all datasets. In this graph, lower and further to the left corresponds to better performance. As can be seen, the distribution of AR/ADR values for GreedyGD is the most optimal. The three datasets with the worts performance under GD-GLEAN are also highlighted with black borders for each compressor. Through this, it can also be seen that both AR and ADR improve from GD-GLEAN to GD-GLEAN+ and from GD-GLEAN+ to GreedyGD.

## 5.3 Configuration runtime

The proposed `BaseTree` algorithm significantly improves configuration runtime for GD. Fig. 8 shows the median configuration runtime for the *COMBED mains power* dataset over 50 trials for selected compressors. Error bars indicate the minimum and maximum runtimes across the 50 trials. As can be seen, applying `BaseTree` results in a significant speed-up. Indeed, comparing GD-INFO (5.341 s) with GD-INFO+ (0.452 s) indicates that `BaseTree` provides a speed-up of 11.8×. GD-GLEAN+ (0.463 s) and GreedyGD (0.475 s) are slightly slower than GD-INFO+ due to their optimisations for analytics, but still deliver speed-ups of 11.5× and 11.2×, respectively. Given the dataset size of 995 kB, GreedyGD's runtime corresponds to a throughput of 2.1 MB/s. Note that our implementation is not optimised for speed and is written in Python, so the absolute runtime values should not be taken as indicative of the algorithms' capabilities.

As stated in Subsection 4.5, GreedyGD's worst case configuration time complexity is $O(nd^2)$. However, as was noted, typical runtimes should be substantially faster. To examine the relationship between configuration runtime and dimensionality, we measured the time to configure GreedyGD on random subsets of the columns of the *Gas turbine emissions* dataset. For each dimensionality (1 to 11), we selected 50 random combinations of the dataset's columns (apart from $d = 1$, 10 and 11 where there are fewer than 50 combinations) and ran 10 trials on each. The median runtime across all combinations and trials for each dimensionality is shown in Fig. 9. Error bars indicate the 5$^{th}$ and

95[th] percentiles. As can be seen, GreedyGD scales much better than quadratic in practice. That is, the runtime for $d = 11$ (0.638 s) is only 16.4 times that for $d = 1$ (0.039 s), which is close to linear instead of quadratic.

### 5.4  Configuration on data subsets

The effect of configuring GD on a data subset was evaluated in terms of CR for a range of subset sizes from 10 to 10,000 samples. Subsets were chosen randomly and, as proposed in Subsection 4.4, data preprocessing and constant bits were determined from the entire dataset. The median CR across all datasets for each subset size is shown in Fig. 10. As can be seen, GreedyGD and GD-INFO+ exhibit very stable compression ratios across all subset sizes. Indeed, even with a mere 250 samples, GreedyGD delivers a median CR of 0.368, which is only 5.7% worse than if the entire dataset is used for configuration. If 10,000 samples are used, the median CR for GreedyGD is only 1.4% worse compared to configuration on the full dataset at 0.353.

## 6  Conclusion

In this paper, we have proposed a new GD-based compression algorithm, GreedyGD, that is designed for high accuracy direct compressed data analytics, while maintaining good compression. The proposed scheme shows improvement in analytics performance compared to existing GD methods, while also running considerably faster and delivering even better compression. We have also demonstrated that GreedyGD can be reliably configured on very small data subsets without compromising significantly on compression. In future work, we intend to investigate using GreedyGD for online compression in an IoT network and additional analytics tasks.

## 7  Acknowledgements

## References

[1] S. Kumar, P. Tiwari, and M. Zymbler, "Internet of Things is a revolutionary approach for future technology enhancement: a review," *Journal of Big Data*, vol. 6, no. 1, 2019.

[2] W. Yu, F. Liang, X. He, W. G. Hatcher, C. Lu, J. Lin, and X. Yang, "A survey on the edge computing for the internet of things," *IEEE Access*, vol. 6, pp. 6900–6919, 2018.

[3] E. Sisinni, A. Saifullah, S. Han, U. Jennehag, and M. Gidlund, "Industrial internet of things: Challenges, opportunities, and directions," in *IEEE Transactions on Industrial Informatics*, vol. 14, no. 11, 2018, pp. 4724–4734.

[4] Y. Lee, E. Hwang, and J. Choi, "A unified approach for compression and authentication of smart meter reading in AMI," in *IEEE Access*, vol. 7, 2019, pp. 34 383–34 394.

[5] R. Vestergaard, Q. Zhang, and D. E. Lucani, "Generalized deduplication: Lossless compression for large amounts of small IoT data," in *European Wireless*, 2019.

[6] ——, "Lossless compression of time series data with generalized deduplication," in *IEEE Global Communications Conference*, 2019.

[7] R. Vestergaard, D. E. Lucani, and Q. Zhang, "A randomly accessible lossless compression scheme for time-series data," in *IEEE INFOCOM*, 2020.

[8] A. Hurst, Q. Zhang, D. E. Lucani, and I. Assent, "Direct analytics of generalized deduplication compressed IoT data," in *IEEE Global Communications Conference*, 2021.

[9] A. Hurst, D. E. Lucani, I. Assent, and Q. Zhang, "GLEAN: Generalized deduplication enabled approximate edge analytics," *IEEE Internet of Things Journal*, 2022.

[10] P. Deutsch and J.-L. Gailly, "ZLIB compressed data format specification version 3.3," 1996, RFC 1950. [Online]. Available: https://datatracker.ietf.org/doc/html/rfc1950

[11] Facebook, "Zstandard – real-time data compression algorithm," accessed: Apr. 7, 2022. [Online]. Available: https://facebook.github.io/zstd/

[12] Google, "Snappy: A fast compressor/decompressor," 2021. [Online]. Available: https://github.com/google/snappy

[13] J. Seward, "bzip2," 2019. [Online]. Available: https://sourceware.org/bzip2/

[14] Y. Collet, F. Handte, I. Rosen, and R. Odaira, "LZ4—extremely fast compression," accessed: May. 4, 2022. [Online]. Available: https://lz4.github.io/lz4/

[15] D. Blalock, S. Madden, and J. Guttag, "Sprintz: Time series compression for the internet of things," *ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, vol. 2, no. 3, pp. 1–23, 2018.

[16] R. Vestergaard, Q. Zhang, M. Sipos, and D. E. Lucani, "Titchy: Online time-series compression with random access for the internet of things," *IEEE Internet of Things Journal*, vol. 8, no. 24, pp. 17 568–17 583, 2021.

[17] J. Jenkins, I. Arkatkar, S. Lakshminarasimhan, D. A. Boyuka, E. R. Schendel, N. Shah, S. Ethier, C.-S. Chang, J. Chen, H. Kolla, S. Klasky, R. Ross, and N. F. Samatova, "ALACRITY: Analytics-driven lossless data compression for rapid in-situ indexing, storing, and querying," in *Trans. on Large-Scale Data- and Knowledge-Centered Systems X*, 2013, pp. 95–114.

[18] M. Shekelyan, A. Dignös, and J. Gamper, "DigitHist: A histogram-based data summary with tight error bounds," in *Proceedings of the VLDB Endowment*, vol. 10, no. 11, 2017, pp. 1514–1525.

[19] J. A. Hartigan and M. A. Wong, "Algorithm AS 136: A k-means clustering algorithm," *Applied Statistics*, vol. 28, no. 1, p. 100, 1979.

[20] Aarhus Kommune, "Sensordata," 2017. [Online]. Available: https://www.opendata.dk/city-of-aarhus/sensordata

[21] "IEEE standard for floating-point arithmetic," *IEEE Standard 754-2019 (Revision of IEEE 754-2008)*, 2019.

[22] M. I. Ali, F. Gao, and A. Mileo, "CityBench: A configurable benchmark to evaluate RSP engines using smart city datasets," in *The Semantic Web - ISWC 2015*, 2015, pp. 374–389.

[23] City of Chicago, "Beach water quality - automated sensors," 2022. [Online]. Available: https://data.cityofchicago.org/Parks-Recreation/Beach-Water-Quality-Automated-Sensors/qmqz-2xku

[24] ——, "Beach weather stations - automated sensors," 2022. [Online]. Available: https://data.cityofchicago.org/Parks-Recreation/Beach-Weather-Stations-Automated-Sensors/k7hf-8y75

[25] ——, "Taxi trips - 2020," 2022. [Online]. Available: https://data.cityofchicago.org/Transportation/Taxi-Trips-2020/r2u4-wwk3

[26] F. D. la Torre, J. Hodgins, A. Bargteil, X. Martin, J. Macey, A. Collado, and P. Beltran, "Guide to the carnegie mellon university multimodal activity (cmu-mmac) database," Carnegie Mellon University, Tech. Rep., 2009, tech. report CMU-RI-TR-08-22.

[27] N. Batra, O. Parson, M. Berges, A. Singh, and A. Rogers, "A comparison of non-intrusive load monitoring methods for commercial and residential buildings," 2014, arXiv:1408.6595.

[28] City of Melbourne, "Sensor readings, with temperature, light, humidity every 5 minutes at 8 locations (trial, 2014 to 2015)," 2020. [Online]. Available: https://data.melbourne.vic.gov.au/Environment/Sensor-readings-with-temperature-light-humidity-ev/ez6b-syvw

[29] H. Kaya and P. Tüfekci, "Gas turbine CO and NOx emission data set data set," 2019. [Online]. Available: https://archive.ics.uci.edu/ml/datasets/GasTurbineCOandNOxEmissionDataSet

[30] G. Hebrail and A. Berard, "Individual household electric power consumption data set," 2012. [Online]. Available: https://archive.ics.uci.edu/ml/datasets/Individual+household+electric+power+consumption