### Granulation of Large Temporal Databases: An Allan Variance Approach

Lorina Sinanaj<sup>1</sup> a, Hossein Haeri<sup>2</sup> b, Satya Prasad Maddipatla<sup>3</sup> c, Liming Gao<sup>3</sup> d, Rinith Pakala<sup>1</sup> e, Niket Kathiriya<sup>1</sup> f, Craig Beal<sup>4</sup> e, Sean Brennan<sup>3</sup> h, Cindy Chen<sup>1</sup> i, and Kshitij Jerath\*<sup>2</sup> i <sup>1</sup> Computer Science Department, University of Massachusetts Lowell, 220 Pawtucket St, Lowell, USA

<sup>2</sup> Mechanical Engineering Department, University of Massachusetts Lowell, Lowell, USA

<sup>3</sup> Mechanical Engineering Department, The Pennsylvania State University, University Park, USA

<sup>4</sup> Mechanical Engineering Department, Bucknell University, Lewisburg, USA

{lorina\_sinanaj, hossein\_haeri, rinith\_pakala, niket\_kathiriya} @student.uml.edu, {cindy\_chen, kshitij\_jerath} @uml.edu, {lug358, szm888, snb10} @psu.edu, cbeal@bucknell.edu

Keywords: Big Data, Data Reduction, Temporal Granulation, Allan Variance.

Abstract:

As the use of Big Data begins to dominate various scientific and engineering applications, the ability to conduct complex data analyses with speed and efficiency has become increasingly important. The availability of large amounts of data results in ever-growing storage requirements and magnifies issues related to query response times. In this work, we propose a novel methodology for granulation and data reduction of large temporal databases that can address both issues simultaneously. While prior data reduction techniques rely on heuristics or may be computationally intensive, our work borrows the concept of Allan Variance (AVAR) from the fields of signal processing and sensor characterization to efficiently and systematically reduce the size of temporal databases. Specifically, we use Allan variance to systematically determine the temporal window length over which data remains relevant. Large temporal databases are then granulated using the AVAR-determined window length. Averaging over the resulting granules produces aggregate information for each granule, resulting in significant data reduction. The query performance and data quality are evaluated using existing standard datasets, as well as for two large datasets that include temporal information for vehicular and weather data. Our results demonstrate that the AVAR-based data reduction approach is efficient and maintains data quality, while leading to an order of magnitude improvement in query execution times compared to three existing clustering-based data reduction methods.

### 1 INTRODUCTION

Whether we are monitoring software systems, tracking applications, financial trading analytics, business intelligence tools, or other temporal systems, timeseries data flows through our data pipelines and applications at warp speed, enabling us to discover hidden and valuable information on how that data changes

<sup>a</sup> https://orcid.org/0000-0003-4687-5809

b https://orcid.org/0000-0002-6772-6266

over time. This has led to a growing interest in the development of data mining techniques capable in the automatic extraction of patterns, anomalies, trends and other useful knowledge from data (Johnston, 2001), (Liu and Motoda, 2002).

However, such hundreds of terabytes of data pose an I/O bottleneck—both while writing the data into the storage system and while reading the data back during analysis. Given this magnitude and faced with the curse of dimensionality which requires exponential running time to uncover significant knowledge patterns (Keogh and Mueen, 2017), much research has been devoted to the data reduction task (Januzaj et al., 2004). Dealing with a reduced set of representative and relevant data (instead of enormous amounts of raw and potentially redundant data), has the potential to significantly benefit data analysis.

Traditional data reduction approaches have proposed methods such as data compression (Rehman

<sup>&</sup>lt;sup>c</sup> https://orcid.org/0000-0002-5785-3579

d https://orcid.org/0000-0002-0159-4010

e https://orcid.org/0000-0001-7442-4170

f https://orcid.org/0000-0002-4146-3402

g https://orcid.org/0000-0001-7193-9347

h https://orcid.org/0000-0001-9844-6948 https://orcid.org/0000-0002-8712-8108

j https://orcid.org/0000-0001-6356-9438

et al., 2016), data cube aggregation (Gray et al., 1997), sampling (Madigan and Nason, 2002), clustering (Kile and Uhlen, 2012), among others. One of the most notable techniques to reduce data while limiting loss of important information is to use cluster representations of data instead of the original instances. Unfortunately, many widely-used clustering methods such as K-means (MacQueen et al., 1967), K-medoids (Kaufmann, 1987), Fuzzy C-means (Bezdek et al., 1984), etc., can be computationally expensive and often rely on heuristics for choosing the appropriate number of clusters to use.

To overcome these drawbacks, we propose a granulation and data reduction method based on Allan Variance (AVAR) (Allan, 1966), (Haeri et al., 2021) for large temporal databases. In time-series data, order and time are fundamental elements that are crucial to the meaning of the data. However, there is a trade-off between how much past information should be used to make predictions about the future. Specifically, this trade-off directly relates to the window size of past time-series data. For example, using the most recent information (i.e. the smallest time window) for prediction may be useful, but it may also contain noisy data which does not fully capture the true trends of the time-series data. On the other hand, using a large amount of historical data (i.e. large time windows) may lead to the incorporation of old data that captures older time-series trends that have become irrelevant, and hence not useful for making predictions about the near future (Lu et al., 2018). Allan variance provides a systematic way of resolving this trade-off and identifying the appropriate granule size, as discussed in Section 3. The time complexity of the proposed method is O(n) where n is the number of input data points. After segmenting the time-series dataset into granules according to the characteristic timescale given by AVAR, we use the average value of the granules as partition representatives instead of the original data points. As a result, a reduced representation of the data is produced without significantly losing important information, while still being computationally efficient. While time-series data is a primary focus of this work, we also demonstrate that the AVARbased granulation approach can be applied to other non-temporal attributes as well. Moreover, the data sets may have both time-series and non-time-series data. A comparative evaluation of standard data sets that have non-time-series data is included in Section

The rest of the paper is organized as follows. In Section 2, we discuss related work. Section 3 describes the AVAR approach for finding the characteristic timescale over which measurements are relevant.

The proposed algorithms and theoretical analysis of the algorithms are presented in Section 4. Section 5 discusses the various experiments on a set of standard data sets, as well as on large data sets pertaining to vehicular and weather data. Section 5 also discusses the associated data reduction and query performance results. The conclusions and future work follow in Section 6.

#### 2 RELATED WORK

Granular computing is an information processing paradigm to represent and process data into chunks or clusters of information called information granules (Pedrycz, 2001). Information granules are a collection of entities grouped together by similarity, proximity, indistinguishability and functionality (Zadeh, 1997). The process of forming information granules is called granulation. In this section, we briefly review the related work which investigate data reduction based on granules or clusters.

Lumini and Nanni (Lumini and Nanni, 2006) present a data reduction method based on clustering (CLUT). The CLUT approach adopts the fuzzy C-means clustering algorithm to divide the original dataset into granules, then use the centroid of each granule as the representative instance to achieve the reduced dataset. The authors use the Hartigan's greedy heuristic (Hartigan, 1975) to select the optimal number of clusters. The time complexity for the fuzzy C-means method (Bezdek et al., 1984) is  $O(n^2)$  with the increase in the size of the original data, as a number of successive iterations need to be completed with the intention to converge on an optimal set of partitions. In addition, requiring a priori specification of the number of clusters by the Hartigan's method, adds more computational complexity to the overall cost.

Olvera et al. (Olvera-López et al., 2010) achieve data reduction by using a Prototype Selection based Clustering method (PSC). PSC first divides the original dataset into clusters using the C-means algorithm, then checks each cluster if it is homogeneous, such that all instances belong to the same class, or not. For the final reduced prototype set, PSC selects the set of the mean prototypes from each homogeneous cluster and the border prototypes from each non-homogeneous cluster.

Sun et al. (Sun et al., 2019) propose to achieve fast data reduction using granulation-based instances importance labeling (FDR-GIIL). The approach uses K-means to generate the granules and then labels the importance of each instance in each granule using the Hausdorff distance (Henrikson, 1999). Data reduc-

tion is achieved by eliminating those instances which have the lowest importance labels, until a user-defined reduction ratio is reached. However, K-means algorithm (MacQueen et al., 1967) is computationally expensive with high volume datasets. Similar to the fuzzy C-means algorithm, when the number of input data points n increases, it is observed that the time complexity becomes  $O(n^2)$ .

For Big Data, the previous related works are time consuming because the used clustering algorithms often rely on heuristics such as Hartigan's statistics, rule of thumb, elbow method, cross-validation, etc. (Kodinariya and Makwana, 2013) to choose the appropriate number of clusters or granules. Additionally, they need to process many iterations in order to converge to optimal cluster centers. This leads to a quadratic time complexity which is very prohibitive for large datasets. To the best of our knowledge, this paper is the first to present a systematic granular data reduction method for temporal databases using Allan Variance, which does not rely on heuristics and successive iterations to process data into information granules.

#### 3 ALLAN VARIANCE

Allan Variance (AVAR) was first proposed to characterize the time drift or instability in atomic clocks (Allan, 1966), but it later became a practical method for sensor noise characterization (Jerath and Brennan, 2011). In its original context, AVAR was developed to quantify the drift and its correlation with time. In the context of signal processing, Allan variance characterizes the noise in a sensor by quantifying the variance observed in measurements across various timescales (Jerath et al., 2018). If the temporal data stream from a sensor is correlated across time, then AVAR can help identify the timescale over which such correlations are most stable. By identifying this stable correlation time, AVAR can systematically quantify the trade-off between aggregating enough data to remove noise, but not aggregating so much data such that signals whose characteristics have drifted are incorporated. Drawing inspiration from these classical applications of AVAR in signal processing and sensor characterization, the authors recently proposed a novel method which utilizes AVAR to identify the characteristic timescale of any given temporal data set with numerical entries (Haeri et al., 2021)(Sinanaj. et al., 2021). Given noisy temporal data that follows a certain unknown pattern, the characteristic timescale determines the time horizon over which measurements yield a near-optimal moving average estimate. Specifically, our method proposes that estimation (or averaging) tasks should not use any measurements that are older than the characteristic timescale (Haeri et al., 2021). This forms the basis of our granulation strategy.

## 3.1 AVAR-informed Characteristic Timescale

Allan variance (or two-sample variance) of a given temporal data set  $\mathbf{y} = \{y_1, y_2, ..., y_n\}$  is mathematically defined as the expected variance of two successive averaged groups of measurements, i.e. data blocks, at a given timescale or window size m (Allan, 1966), (Jerath et al., 2018). Specifically, the Allan variance is given by:

$$\sigma_A^2(m) = \frac{1}{2} \mathbb{E}\left[ (\bar{y}_k - \bar{y}_{k-m})^2 \right] \tag{1}$$

where  $\bar{y}_k$  is simple moving average of the measurements at time k, averaged over the window length m, and is given by:

$$\bar{y}_k = \frac{1}{m} \sum_{i=k-m}^k y_i \tag{2}$$

To numerically estimate AVAR, we can compute average of the term  $(\bar{y}_k - \bar{y}_{k-m})^2$  across all possible time steps k as follows (Sesia and Tavella, 2008):

$$\hat{\sigma}_A^2(m) = \frac{1}{2(n-2m)} \sum_{k=2m+1}^n (\bar{y}_k - \bar{y}_{k-m})^2$$
 (3)

Figure 1(a) shows uncorrelated temporal data (more specifically Gaussian white noise), and Figure 1(d) shows the corresponding Allan variance curve evaluated at several different window lengths or timescales. It is observed that for the uncorrelated Gaussian white noise, Allan variance decreases as the timescale or averaging window length m increases. This is intuitive because as we average across time blocks of larger size m, these averages  $\bar{y}_k$  include more time-series data and also tend closer to the mean value of the Gaussian white noise.

On the other hand, if data points are correlated in time, there is a possibility of the Allan variance increasing with increasing window lengths. For example, Figure 1(b) shows a random walk process with its prototypical drift over time. It is intuitively understood that averages  $\bar{y}_k$  taken across increasingly larger window lengths m will also be larger due to the underlying drift behavior of the random walk. Thus, as seen in Figure 1(e), the AVAR increases as block size increases since the data is correlated in time and the most relevant measurements to any point at k is its immediate neighboring measurement at k-1 and k+1.

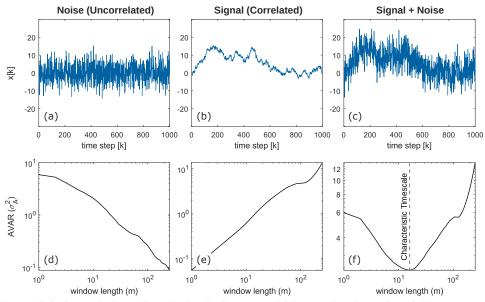


Figure 1: Characteristic timescale determines the time horizon over which averaging data points yields a near optimal results. (a) Gaussian white noise (b) Random walk signal. (c) Random walk corrupted with Gaussian white noise. Figure (d), (e), and (f) show AVAR of the signals on top calculated at various window lengths  $(1 \le m \le n/4)$ .

A lower value of the AVAR at a given timescale indicates lower bias and reduced variance between data aggregated at that timescale.

In this work, however, we are more interested in cases like the one in Figure 1(c) where a correlated signal is corrupted with an uncorrelated noise. In this case, we face a trade-off between averaging over more data points to help eliminate noise, and averaging over fewer data points to help avoid drift as seen in the random walk behavior. Allan variance can help systematically resolve this trade-off by explicitly identifying the characteristic timescale at which data remains relevant, as shown in Figure 1(f). This AVAR-determined window length (or characteristic time scale) will be used to separate the original dataset into granules to achieve data reduction.

#### 3.2 Fast AVAR Calculation

Although AVAR yields valuable information regarding the characteristic timescale of the temporal data, it still needs to be efficiently computed for large data sets. To accomplish this, we use the algorithm suggested in (Maddipatla et al., 2021) and constrain the potential window lengths to powers of 2, i.e  $m=2^p$  where  $p \in \mathbb{Z}^+$ . Then we can quickly estimate the expression in (1) using a simple dynamic programming method explained in Algorithm 1, which has O(n) running time, where n is the number of input data points. The algorithm first constructs a list of exponentially growing window length candidates

 $T = \{2^0, 2^1, 2^2, ...\}$ . Then, for each window length, starting from  $\tau = 1$ , the algorithm calculates the associated AVAR by averaging  $0.5(Y_k - Y_{k+1})^2$  across all valid k values. Meanwhile, it pre-computes the next temporal list Y' by averaging adjacent data values  $Y_k$  and  $Y_{k+1}$ . A detailed description of the algorithm can be found in our earlier work (Maddipatla et al., 2021).

```
Algorithm 1: Fast AVAR calculator
  Input: a set Y of raw temporal data points
  Output: a set \mathcal{A} of AVAR values and a set \mathcal{T}
              of associated timescales
  construct T = \{1, 2, 4, 8, ..., 2^p\} such that
   2^p < Y.length/2 < 2^{p+1}
  for \tau \in \mathcal{T} do
       Y' = \text{empty list}
      k = 1
       while k < Y.length do
           c = c + 0.5(Y_k - Y_{k+1})^2
           add 0.5(Y_k + Y_{k+1}) to Y'
       end
       add \frac{c}{Y'.length} to \mathcal{A}
       Y
  end
  return \mathcal A and \mathcal T
```

The theoretical analysis and the experiments in Figure 2 prove that the running time for the Algorithm

1 is O(n). In the first iteration where  $\tau=1$ , the algorithm will perform  $n\times 1$  computations, in the second iteration because  $\tau=2$ , the number of calculations will be  $n\times \frac{1}{2}$ , in the third iteration  $\tau=4$ , so the number of AVAR calculations is  $n\times \frac{1}{4}$  and so on until the m-th iteration where  $\tau=2^p$ . In the last iteration the total number of calculations is  $n\times (\frac{1}{2})^p$ . Adding all the calculations in each iteration yields a geometric series with the first term is n, the common ratio is  $\frac{1}{2}$ , and the number of terms is m. Using the sum formula for a finite geometric series we can calculate the total time complexity of Algorithm 1:

$$S = \sum_{p=0}^{\infty} n(\frac{1}{2})^p = \frac{n}{1 - \frac{1}{2}} \tag{4}$$

Since the common ratio in this geometric series is  $\frac{1}{2}$  (less than 1), the series converges with sum S = O(n) = 2n. The linear complexity time of the fast AVAR calculator is also shown in Figure 2, where the execution time is linear and proportional to the number of measurements in the data. Being able to efficiently compute the AVAR characteristic window length for large amounts of data is crucial, because the number of the final granules or clusters in the reduced dataset is determined by the temporal aggregation timescale given by the fast AVAR algorithm.

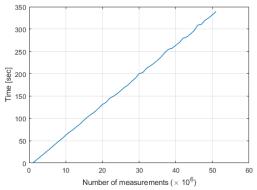


Figure 2: Fast AVAR calculator execution time.

#### 4 AVAR-BASED GRANULATION

This section presents the structure of the proposed AVAR-based granulation approach for temporal databases by describing the overall workflow of the technique and the theoretical analysis of the algorithm complexity, using the AVAR-informed characteristic timescale.

#### 4.1 System Architecture

The current state-of-the-art granulating temporal data algorithms require successive iterations and hence significant computational effort. In addition, they also often rely on heuristics for choosing how to partition the original data. These are computationally expensive operations, especially for temporal databases with a large number of rows. The main issue with using heuristics to specify the number of clusters in advance, is that the final result will be sensitive to the initialization of parameters. A practical approach is to compare the outcomes of multiple runs with different k-number of clusters and choose the best one based on a predefined criterion. However, this method is also very time consuming due to the large number of iterations that the algorithm needs to carry out (MacQueen et al., 1967).

In our study, we propose an AVAR-based granulation technique on large temporal databases, that in O(n) time complexity takes as input the n raw data points, systematically determines the time window over which data is relevant, and calculates the aggregated information of the relevant data for each time window. Compared to previous approaches where authors used heuristics (Kodinariya and Makwana, 2013) to determine the optimal number of partitions (granules), in our work the number of partitions is systematically determined by the AVAR characteristic timescale itself.

Figure 3 shows the overall workflow of the proposed method including four steps: (1) Preprocessing: A simple pre-processing step is applied on the raw dataset (if it is not originally sorted) to sort the data points in time prior to applying the AVAR algorithm, (2) Allan variance calculation: The AVAR algorithm takes as input the sorted temporal database, and outputs a characteristic timescale over which the measurements are relevant and should be averaged across, (3) Granulation: The granulation algorithm takes in input the AVAR timescale, partitions the data in a series of non-overlapping time interval segments and generates aggregated information for each interval, and (4) Data preparation: Aggregated information for each partition will then be used as representative instances for the new reduced dataset. As a result, at a later stage, data mining methods can analyze the data faster due to the decrease in volume without losing data quality.

Compared to the other methods in the related works section, the AVAR-based technique systematically identifies a characteristic timescale at which measurements are relevant, representative and stable without relying on iterations or heuristics. As

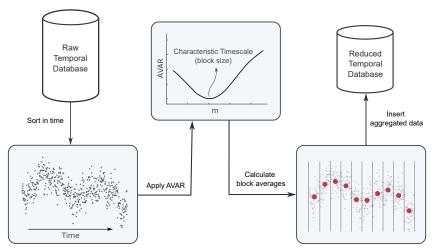


Figure 3: System Architecture—Schematic view of data granulation process based on Allan Variance.

a limitation to this method, input data should be sorted in time before getting processed. However, the worst-case time complexity when sorting the data is  $O(n\log(n))$  (Mishra and Garg, 2008). Even with sorting operations included in the assessment, the presented AVAR-based granulation methodology performs better compared to methods that compute a large number of iterations and use heuristics to achieve approximate cluster representatives in  $O(n^2)$ .

### 4.2 AVAR Granulation Approach

A challenge in large sized databases is to evaluate complex queries over a continuous stream of inputs. The key idea is to reduce the volume using moving window aggregations, i.e. the calculations of aggregates in a series of non-overlapping (tumbling) windows. Tumbling windows (Helsen et al., 2017) are a series of fixed-sized, non-overlapping, and contiguous time intervals where tuples are grouped in a single window based on time. A tuple in the database cannot belong to more than one tumbling window. In the proposed algorithm, AVAR method determines the size of the tumbling window over which the measurements are relevant.

The complete temporal granulation process defined according to the proposed methodology is described by Algorithm 2. The algorithm takes as an input the set of the original data points and the AVAR characteristic timescale  $\tau_{AVAR}$ , and returns the set of the granule representatives. Initially, it starts by sequentially scanning every data point in the original set associated with their respective timestamp values. It then partitions the data points into different time intervals, by checking which window interval the timestamps fall into (if statement). Before jumping to a

new time window interval (else statement), the algorithm calculates the average value for the current granule and inserts that value in the representatives set

In Algorithm 2, the data is scanned only once inside the while loop. Each step inside the "if" and "else" statements takes only O(1) time as it does not contain loops, recursion or call to any other nonconstant time function. As a result, time complexity of the AVAR-based granulation algorithm becomes O(n), where n is the number of input data points. This time complexity is linear and proportional to the size of the original data. If we combine the time complexity of Algorithm 1 with the time complexity of Algorithm 2, the total time complexity of the proposed approach is simply O(n). The more voluminous the original data is, the slower the reduction algorithm will be. However, our approach is a pre-processing step to prepare the data for the future data analysing techniques. Efficiency is achieved because data analytical methods can analyze the data faster due to the decrease in volume without losing data quality.

The space complexity of the proposed algorithm is the amount of memory space it needs to run according to its input size. If the number of Information Granules generated by the algorithm is q and the input size is n, the storage requirement for the algorithm to complete its task is O(nq).

#### 5 EXPERIMENTS

In the following we present experimental results that assess the performance of the AVAR method in the granulation algorithm. As the experimental basis, a simulation environment is used to produce a large

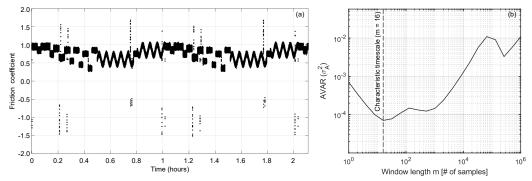


Figure 4: (a) Friction data. Negative values are outliers but we retained them to show that the technique is robust under different scenarios. (b) Allan variance for vehicular friction measurements containing 91 million data points, evaluated across various window lengths. The minimum point in the AVAR curve indicates the characteristic timescale of the data.

amount of temporal information for vehicular data. Each tuple contains information about road-tire friction measurements at a certain point in time.

On this dataset we perform three separate experiments: (1) we apply the AVAR method to determine the characteristic timescale over which measurements are relevant, (2) we apply the granulation algorithm based on the AVAR output, on data of different sizes to observe the algorithm execution time with respect to input data size, and (3) on the input dataset and on the reduced dataset, we run the same exact query to analyze the performance of our reduction method. In addition experimental analysis is conducted to compare the proposed method against the competitors in the related works.

The test results show that performing the same type of query in the reduced dataset compared to executing the query in the original data, not only drastically reduces the query execution time but it can also generate close results with a relatively low absolute error.

**Hardware:** The platform of our experimentation is a PC with a 2.60 GHz Single Core CPU, 64 GB RAM using PostgreSQL 10.12 on Linux kernel 3.10.0.

## 5.1 Characteristic Timescale of Vehicular Data

In this subsection, we evaluate the characteristic timescale of the vehicular friction measurements by calculating the AVAR across various window lengths using the Algorithm 1. Figure 4(a) shows the pattern of the friction values over time and Figure 4(b) shows the corresponding AVAR curve. In Figure 4(a) we can observe negative values in the form of outliers in the friction data. As a future work, outliers will be removed from the original dataset but for now we are retaining them to show that the proposed granulation

technique is robust under different scenarios. Figure 4(b) indicates that the size of the granule to be used is equal to the characteristic timescale of the data, which points to the minimum AVAR value in the graph. The reader may also note that the Allan variance curve indicates the presence of two local minima, which arise due to the presence of periodic sinusoid patterns in the temporal data. In Section 5.3, we also investigate the effectiveness of AVAR-based granulation of real-world weather datasets where the periodicity is inherent in the dataset, and where we also found multiple minima in the Allan variance curve.

**Algorithm 2:** AVAR-based temporal granulation algorithm

```
Input: a set Y of raw temporal data points
        and timescale \tau_{AVAR}
Output: a set Y' of granule representatives
initialize window<sub>start</sub> and window<sub>size</sub>=\tau_{AVAR}
q=1 // Information Granules counter
while y \in Y do
    if timestamp(y) \in window interval then
        insert y in Information Granule IG_q
    else
        set window<sub>start</sub> as
          (window_{start} + window_{size})
         shift window interval (windowstart,
          window_{start} + window_{size})
        calculate the average \bar{y}_q for all
          y \in IG_q
        add \bar{y}_q to Y'
        q=q+1
    end
end
return Y'
```

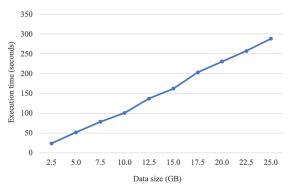


Figure 5: Execution time of the AVAR-based granulation algorithm.

# 5.2 Time Efficiency of AVAR-based Granulation Algorithm

To measure the time efficiency of the AVAR-based granulation algorithm, we will gradually increment the size of the data to be granulated. We initially start by applying the granulation algorithm on a temporal database of size 2.5 GB ( $\approx$  32 million data points), and we measure its execution time. We repeat these steps by applying the algorithm on data with different sizes up to 25 GB ( $\approx$  320 million data points) and present the results in Figure 5.

From these results, we observe that for larger amounts of input data, the algorithm takes more time to execute. However, the increase in execution time is linear to the size of the input data, once again proving our theoretical algorithm analysis that the time complexity for our proposed AVAR-based granulation technique is O(n), where n is the number of input data points.

# 5.3 Effectiveness of AVAR-based Granulation Technique

We evaluate the proposed approach on a vehicular information dataset of size 7216 MB, which consists of 91,475,050 tuples measured at a time range of around 25 hours. Each tuple holds road-tire friction values associated with their respective timestamps. The purpose of AVAR-based granulation is to reduce the query execution time without losing important information from data reduction. To show that this technique efficiently produces a reduced representation of the data without losing data quality, we run the same query on both datasets, the original and the reduced, and we analyze the query execution time as well as the query results. We use a query to output the average friction value for a given interval in time. The

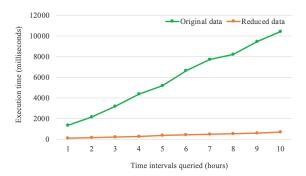


Figure 6: The runtime rising tendency of the query as more data is analyzed

query is run for different time intervals, to observe how the query execution time changes when an increasingly number of tuples have to be analyzed. We show that after using our approach, data can be analyzed faster due to the decrease in volume, without losing important information.

First, we show that the query execution time can be reduced drastically after applying the AVAR-based granulation algorithm. The size of the reduced dataset after the granulation technique with a characteristic timescale of 16 milliseconds is 426 MB, with the total number of tuples being 5,717,191. The storage requirement is efficiently reduced as shown by the calculated reduction rate of  $\approx 94$  %. Experimental results in Figure 6 show that for the same time interval, the query takes less time to execute in the reduced data compare to executing it in the original data. In addition, we observe that for the original dataset, the runtime growth of the query is higher with the growth of the amount of data queried. From these results we can observe that the benefits of data reduction processes are sometimes not evident when the data is small; they begin to become obvious when the datasets start growing in size and more instances have to be analyzed.

Second, we show that we do not lose data quality while reducing the temporal database in a representative subset. The advantage of using representatives is that besides improving query execution time, it also improves the model generalization for the use of data mining techniques in the future. We measured the average, minimum and maximum friction values for different time intervals on both datasets. Numerical results are shown in Table 1. In addition we calculated the percentage error rate for the average friction values in increasing time intervals of 1 hour and we show the graph results in Figure 7.

From Table 1 we can observe that for the average (AVG) query, the error is occurring at the  $5^{th}$  or  $6^{th}$  decimal place. This event mostly happens due

Table 1: Query performance on the friction data.

Interval of time	AVG(fr	ciction)	AVG que	ery error	MIN(fr	iction)	MIN que	ry error	MAX(fr	iction)	MAX que	ery error
interval of time	Original data	Reduced data	Percentage	Absolute	Original data	Reduced data	Percentage	Absolute	Original data	Reduced data	Percentage	Absolute
between '2020-06-15 05:00:00' and '2020-06-15 05:10:00'	0.754902	0.754902	1.14E-05	8.60E-08	-1.608	-0.929	42.214	0.679	1.559	1.326	14.968	0.233
between '2020-06-15 05:10:00' and '2020-06-15 05:20:00'	0.654777	0.654776	1.05E-04	6.87E-07	0.253	0.324	28.170	0.071	0.949	0.870	8.349	0.079
between '2020-06-15 05:20:00' and '2020-06-15 05:30:00'	0.600036	0.600037	1.97E-04	1.18E-06	0.349	0.431	23.460	0.082	0.839	0.770	8.246	0.069
between '2020-06-15 05:30:00' and '2020-06-15 05:40:00'	0.686079	0.686078	1.46E-04	1.00E-06	-0.723	-0.258	64.385	0.466	1.732	1.532	11.571	0.200
between '2020-06-15 05:40:00' and '2020-06-15 05:50:00'	0.849033	0.849026	8.27E-04	7.02E-06	-1.536	-1.011	34.193	0.525	1.413	1.084	23.317	0.329
between '2020-06-15 05:50:00' and '2020-06-15 06:00:00'	0.897936	0.897934	2.36E-04	2.12E-06	0.633	0.728	14.931	0.095	1.074	0.975	9.214	0.099

to the computer representation for binary floating-point numbers in the IEEE Standard for Floating-Point Arithmetic (IEEE 754)(IEEE, 2019). IEEE 754 standard, for floating point representation, allows 23 bits for the fraction. 23 bits is equivalent to  $\log_{10}(23) \approx 6$  decimal digits. Beyond those number of significant digits, accuracy is not preserved, hence round-off starts to occur as reported in Table 1. The absolute error values, mostly due to numerical accuracy round-off, are extremely small, proving that the AVAR-based granulation technique in temporal databases keeps the quality of the original data. Nevertheless, the vehicular friction values round to the  $4^{th}$  decimal place, already have enough accuracy for vehicular control application.

For the minimum (MIN) and maximum (MAX) queries, even though the absolute error is larger due to the presence of outliers (as observed in Figure 4(*Top*)), the results are still close in value to each other. A good approach in the future is to apply a filtering algorithm that removes the outliers prior to applying the AVAR granulation approach.

From Figure 7 we notice that the more data that is used, the smaller the percentage error is. This is often a characteristic of the round-off error (Goldberg, 1991). Using a finite number of bits to represent real numbers which can have infinitely many digits, requires an approximate representation. Rounding error occurs as the difference between the calculated approximation of a number and its exact mathematical value. The IEEE standard uses round-to-nearest method where results are rounded to the nearest representable value. As more data is calculated, rounding up and down many times causes the cumulative round-off error to decrease.

We further tested our proposed AVAR-based granulation technique on NASA's Solar and Meteorology dataset (NASA, ) for the Boston area with coordinates latitude 42.3601 and longitude -71.0589. These data were obtained from the NASA Langley Research

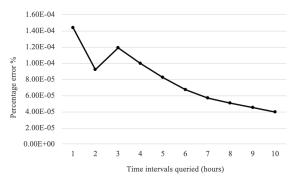


Figure 7: Percentage error rate for the average friction values

Center (LaRC) POWER Project funded through the NASA Earth Science/Applied Science Program. The datasets contain 177480 records for the time range from 01/01/1980 to 01/31/2021. The temperature dataset and humidity dataset provide hourly temperature and specific humidity at 2 meters high above sea level, respectively. The columns – year, month, day, and hour are converted into the sequential hour and the AVAR graph is plotted for hourly temperature and humidity data, shown in Figures 8 and 9, respectively.

The AVAR for temperature data interestingly shows that there are two local minima for the granulation window length: the first minimum occurs at 16 hours and the second occurs at 256 hours. We then generate two granulated temperature datasets with window length as 16 hours and one with window length as 256 hours. The query performance result for average (AVG), minimum (MIN) and maximum (MAX) for both granulated datasets are shown in Table 2 and Table 3. For the larger granulation window length, data are compacted at greater level, however, when answering min and max queries, the error rate is higher than using the smaller granulation window length.

Figure 9 shows that the AVAR for humidity data is minimized at a window length of 256 hours. The

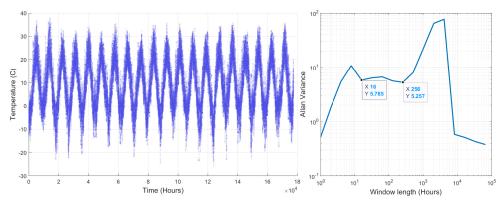


Figure 8: AVAR of temperature of the Solar and Meteorology dataset calculated at various window lengths. (a) Temperature data (b) AVAR curve showing the two characteristic window length = 16 hours and length = 256 hours.

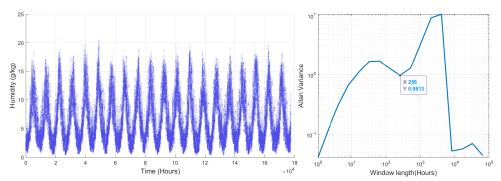


Figure 9: AVAR of humidity of the Solar and Meteorology dataset calculated at various window lengths. (a) Humidity data (b) AVAR curve showing the characteristic window length = 256 hours.

Table 2: Query performance on the Temperature data: Granulation using first minimum of Allan variance.

Interval of time	AVG(Tem	perature)	AVG que	ry error	MIN(Tem	perature)	MIN que	ry error	MAX(Tem	perature)	MAX que	ery error
interval of cline	Original data	Reduced data	Percentage	Absolute	Original data	Reduced data	Percentage	Absolute	Original data	Reduced data	Percentage	Absolute
between '2001-01-01 00:00:00' and '2005-01-19 00:00:00'	9.136	9.141	0.054	0.005	-23.02	-19.294	16.185	3.726	38.01	31.659	16.708	6.351
between '2005-01-19 00:00:00' and '2009-02-06 00:00:00'	9.141	9.135	0.065	0.006	-21.64	-18.025	16.705	3.615	35.39	30.203	14.656	5.187
between '2009-02-06 00:00:00' and '2013-02-24 00:00:00'	9.788	9.79	0.02	0.002	-24.31	-19.039	21.682	5.271	35.59	31.579	11.27	4.011
between '2013-02-24 00:00:00' and '2017-03-14 00:00:00'	9.473	9.471	0.021	0.002	-25.92	-21.368	17.561	4.552	35.74	30.58	14.437	5.16
between '2017-03-14 00:00:00' and '2021-03-31 23:00:00'	9.642	9.643	0.01	0.001	-22.74	-19.341	14.947	3.399	35.61	30.86	13.338	4.75

Table 3: Query performance on the Temperature data: Granulation using second minimum of Allan variance.

Interval of time	AVG(Tem	perature)	AVG que	ry error	MIN(Tem	perature)	MIN que	ry error	MAX(Ten	perature)	MAX que	ery error
interval of time	Original data	Reduced data	Percentage	Absolute	Original data	Reduced data	Percentage	Absolute	Original data	Reduced data	Percentage	Absolute
between '2001-01-01 00:00:00' and '2005-01-19 00:00:00'	9.136	9.082	0.591	0.054	-23.02	-12.426	46.02	10.594	38.01	26.821	29.436	11.189
between '2005-01-19 00:00:00' and '2009-02-06 00:00:00'	9.141	9.233	1.006	0.092	-21.64	-9.135	57.786	12.505	35.39	24.729	30.124	10.661
between '2009-02-06 00:00:00' and '2013-02-24 00:00:00'	9.788	9.747	0.418	0.041	-24.31	-7.143	70.617	17.167	35.59	24.427	31.365	11.163
between '2013-02-24 00:00:00' and '2017-03-14 00:00:00'	9.473	9.438	0.369	0.035	-25.92	-10.597	59.116	15.323	35.74	25.457	28.771	10.283
between '2017-03-14 00:00:00' and '2021-03-31 23:00:00'	9.642	9.666	0.248	0.024	-22.74	-11.002	51.618	11.738	35.61	26.142	26.588	9.468

Table 4: Query performance on the Humidity data.

Interval of time	AVG(Hu	ımidity)	AVG que	ry error	MIN(Hu	midity)	MIN que	ry error	MAX(H	umidity)	MAX que	ery error
interval of cline	Original data	Reduced data	Percentage	Absolute	Original data	Reduced data	Percentage	Absolute	Original data	Reduced data	Percentage	Absolute
between '2001-01-01 00:00:00' and '2005-01-19 00:00:00'	6.487	6.474	0.2	0.013	0.49	1.341	173.673	0.851	20.63	14.807	28.225	5.823
between '2005-01-19 00:00:00' and '2009-02-06 00:00:00'	6.786	6.81	0.353	0.024	0.61	1.699	178.524	1.089	20.51	14.678	28.434	5.832
between '2009-02-06 00:00:00' and '2013-02-24 00:00:00'	6.939	6.927	0.172	0.012	0.49	2.086	325.714	1.596	19.71	14.958	24.109	4.752
between '2013-02-24 00:00:00' and '2017-03-14 00:00:00'	6.673	6.662	0.164	0.011	0.43	1.532	256.279	1.102	19.84	15.916	19.778	3.924
between '2017-03-14 00:00:00' and '2021-03-31 23:00:00'	6.87	6.87	0	0	0.55	1.562	184	1.012	20.02	14.978	25.184	5.042

humidity is defined as the ratio of the mass of water vapor to the total mass of air at 2 meters, and has units of g/kg. The query performance results for average (AVG), minimum (MIN) and maximum (MAX) are shown in Table 4, and show comparable performance to the other data sets, with the query error for averaging being especially low.

#### **5.4** Comparative Evaluation

Comparison experiments are conducted by comparing the proposed AVAR-based granulation algorithm with the other clustering data reduction methods in the related works. Five datasets selected from the UCI Repository (Dua and Graff, 2019) with different sizes are reduced to demonstrate the effectiveness of the AVAR approach. The chosen datasets (Segmentation, Magic, Letter, Shuttle, Covertype) are considered as "large" datasets by the competitors FDR-GIIL (Sun et al., 2019), CLU (Lumini and Nanni, 2006) and PSC (Olvera-López et al., 2010). CLU and PSC are not applicable to run in the Covertype dataset (250,000 instances), because they are very expensive when large datasets are processed. A description of these datasets is given in Table 5.

Table 5: Description of datasets.

Dataset	Number of	Number of
Dataset	instances	attributes
Segmentation	2100	19
Magic	19,020	10
Letter	20,000	16
Shuttle	58,000	9
Covertype	250,000	54

There are two key experiments conducted in this section; (i) measurement of the query performance in each dataset after they have been reduced by the AVAR-based approach and (ii) evaluation of how fast the execution time of the proposed AVAR approach is

compared to the existing clustering based data reduction methods.

Figure 10, 11, 12, 13 and 14 show the AVAR calculated at various window lengths for each dataset. We can observe from the graphs in Figures 8-12(a), how each data has different characteristics and shapes. Some of them have a large number of outliers (Shuttle and Covertype), while some others have many repetitive data points (Letter). It is important to see how the AVAR-method performs against different kind of data, so that future work can be planned to make the method applicable in more general cases.

The first step in the proposed granulation process measures the characteristic window length at the minimum AVAR, shown in Figures 8-12(b). The next step is to use this window length to separate the data into granules and generate the aggregated information for each such granule. After the reduction step is performed, we observe the query performance for each reduced dataset.

For the Segmentation dataset, we measured the average, minimum and maximum rawblue values where rawred was between a given interval. For the Magic dataset, we measured the average, minimum and maximum fConc values where fSize was between a given interval. For the Letter dataset, we measured the average, minimum and maximum box\_width values where horizontal\_position was between a given interval. For the Shuttle dataset, we measured the average, minimum and maximum Column\_7 values where time was between a given interval. For the Covertype dataset, we measured the average, minimum and maximum hillshade\_index\_noon values where slope was between a given interval.

Numerical results for each data are shown in Tables 6, 7, 8, 9 and 10. For the AVERAGE query we can observe that the absolute error is almost always close to 0, indicating that there is little difference between the results in the original and the reduced data. There are special cases where the absolute error is large, as in certain intervals of Table 9 and 10 which

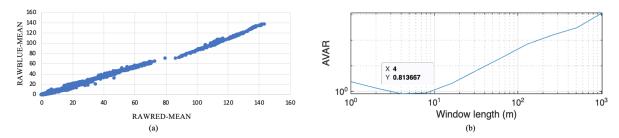


Figure 10: AVAR of the Segmentation dataset calculated at various window lengths. (a)Segmentation data (b)AVAR curve showing the characteristic window length=4 units.

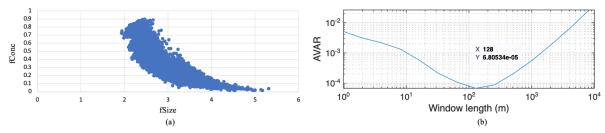


Figure 11: AVAR of the Magic dataset calculated at various window lengths. (a) Magic data (b) AVAR curve showing the characteristic window length = 128 units.

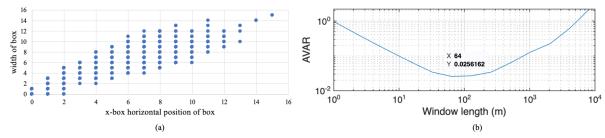


Figure 12: AVAR of the Letter dataset calculated at various window lengths. (a) Letter data (b) AVAR curve showing the characteristic window length = 64 units.

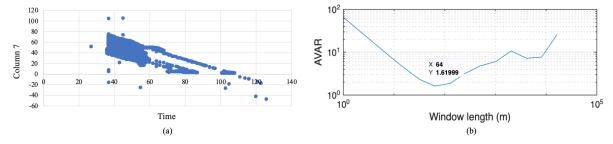


Figure 13: AVAR of the Shuttle dataset calculated at various window lengths. (a) Shuttle data (b) AVAR curve showing the characteristic window length = 64 units.

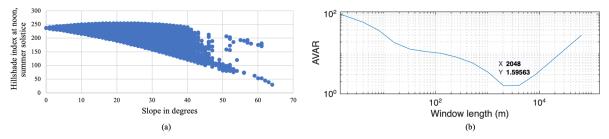


Figure 14: AVAR of the Covertype dataset calculated at various window lengths. (a) Covertype data (b) AVAR curve showing the characteristic window length = 2048 units.

Table 6: Query performance on the Segmentation data.

Interval of	AVG(RAWB	LUE-MEAN)	AVG que	ry error	MIN(RAWE	BLUE-MEAN)	MIN que	ry error	MAX(RAW	BLUE-MEAN)	MAX que	ery error
RAWRED-MEAN	Original data	Reduced data	Percentage	Absolute	Original data	Reduced data	Percentage	Absolute	Original data	Reduced data	Percentage	Absolute
[0, 29]	9.506	9.506	0.000	0.000	0.000	0.000	0.000	0.000	26.111	25.083	3.937	1.028
[29, 58]	39.698	39.778	0.202	0.080	18.333	24.417	33.186	6.084	54.111	51.639	4.568	2.472
[58, 87]	55.188	55.161	0.049	0.027	50.000	51.278	2.556	1.278	70.778	65.389	7.614	5.389
[87, 116]	92.963	92.877	0.093	0.086	72.000	72.250	0.347	0.250	107.444	103.639	3.541	3.805
[116, 145]	118.453	118.543	0.076	0.090	103.667	105.778	2.036	2.111	137.111	136.000	0.810	1.111

Table 7: Query performance on the Magic data.

Interval of fSize	AVG(	fConc)	AVG que	ry error	MIN(fConc)		MIN query error		MAX(fConc)		MAX query error	
micrvai of 1312e	Original data	Reduced data	Percentage	Absolute	Original data	Reduced data	Percentage	Absolute	Original data	Reduced data	Percentage	Absolute
[2.0, 2.5]	0.594	0.596	0.337	0.002	0.290	0.492	69.655	0.202	0.893	0.734	17.805	0.159
[2.5, 3.0]	0.377	0.378	0.265	0.001	0.116	0.265	128.448	0.149	0.885	0.49	44.633	0.395
[3.0, 3.5]	0.221	0.221	0.000	0.000	0.049	0.174	255.102	0.125	0.638	0.262	58.934	0.376
[3.5, 4.0]	0.153	0.154	0.654	0.001	0.025	0.120	380.000	0.095	0.292	0.177	39.384	0.115
[4.0, 4.5]	0.080	0.085	6.250	0.005	0.013	0.068	423.077	0.055	0.155	0.102	34.194	0.053

Table 8: Query performance on the Letter data.

Interval of	AVG(bo	AVG(box_width)		AVG query error		MIN(box_width)		MIN query error		MAX(box_width)		ery error
horizontal_position	Original data	Reduced data	Percentage	Absolute	Original data	Reduced data	Percentage	Absolute	Original data	Reduced data	Percentage	Absolute
[0, 3]	2.592	2.610	0.694	0.018	0.000	0.281	0.000	0.281	5.000	3.906	21.880	1.094
[3, 6]	5.235	5.243	0.153	0.008	2.000	4.156	107.800	2.156	9.000	6.516	27.600	2.484
[6, 9]	7.408	7.409	0.013	0.001	4.000	6.703	67.575	2.703	12.000	8.297	30.858	3.703
[9, 12]	8.567	8.567	0.000	0.000	6.000	7.797	29.950	1.797	14.000	9.703	30.693	4.297
[12, 15]	11.259	11.375	1.030	0.116	9.000	11.375	26.389	2.375	14.000	11.375	18.750	2.625

Table 9: Query performance on the Shuttle data.

Interval of time	AVG(cc	olumn_7)	AVG que	ry error	MIN(column_7)		MIN query error		MAX(column_7)		MAX query error	
interval of cline	Original data	Reduced data	Percentage	Absolute	Original data	Reduced data	Percentage	Absolute	Original data	Reduced data	Percentage	Absolute
[25, 45]	44.626	44.625	0.002	0.001	-16.000	37.531	334.569	53.531	104.000	50.750	51.202	53.250
[45, 65]	34.178	34.178	0.000	0.000	-26.000	28.000	207.692	54.000	105.000	48.016	54.270	56.984
[65, 85]	6.857	6.720	1.998	0.137	3.000	3.188	6.267	0.188	48.000	44.078	8.171	3.922
[85, 105]	2.156	2.039	5.427	0.117	-27.000	0.565	102.093	27.565	22.000	5.516	74.927	16.484
[105, 125]	0.321	-2.560	897.508	2.881	-43.000	-20.688	51.888	22.312	3.000	1.203	59.900	1.797

Table 10: Query performance on the Covertype data.

Interval o	of slope	AVG(hillshad	de_index_noon)	AVG que	ery error	MIN(hillsha	de_index_noon)	MIN que	ry error	MAX(hillsha	de_index_noon)	MAX qu	ery error
interval	or stope	Original data	Reduced data	Percentage	Absolute	Original data	Reduced data	Percentage	Absolute	Original data	Reduced data	Percentage	Absolute
[0, 1	10]	232.452	232.383	0.030	0.069	219.000	227.854	4.043	8.854	249.000	237.337	4.684	11.663
[10,	20]	222.861	222.766	0.043	0.095	194.000	213.032	9.810	19.032	254.000	229.754	9.546	24.246
[20,	30]	204.780	204.199	0.284	0.581	162.000	190.453	17.564	28.453	254.000	212.414	16.372	41.586
[30,	40]	184.955	182.371	1.397	2.584	126.000	172.000	36.508	46.000	252.000	190.098	24.564	61.902
[40,	50]	149.370	122.424	18.040	26.946	87.000	122.424	40.717	35.424	240.000	122.424	48.990	117.576

is explained by the presence of outliers in those intervals of the original data. The presence of outliers has a more obvious effect in the results of the MINI-MUM and MAXIMUM queries. In such queries there is a larger gap in the results between the original and the reduced data, hence the absolute error is worse. As a recommendation in the future, a filtering algorithm for the outliers removal will be used before the

AVAR-based granulation approach. As the number of outliers increases, so does the absolute error.

While the competitors pick K-means or fuzzy C-means to decide on the number of clusters, the AVAR approach sets the size of the cluster by using the characteristic window length on available data. The comparative experiments show that the AVAR-method can still give good performance results in preserving data

Table 11: Execution time of large datasets (seconds).

Dataset	Number of instances	Number of attributes	AVAR	Execution time (seconds)					
Dataset	Number of histances	Number of attributes	window size	AVAR	FDR-GIIL	CLU	PSC		
Segmentation	2100	19	4	0.286	3.9	6.0	7.0		
Magic	19,020	10	128	0.217	88.6	167.1	172.1		
Letter	20,000	16	64	0.210	14.5	217.2	226.2		
Shuttle	58,000	9	64	0.351	30.9	277.4	288.4		
Covertype	250,000	54	2048	0.321	649.7	-	-		

quality even in data streams that are not time-series. In the future we plan to extend this approach for both spatio-temporal data. Next, the execution time of the competitor algorithms FDR-GIIL, CLU and PSC are added for each dataset to show the improvement in the computational cost of the proposed AVAR-based granulation method. The corresponding results are recorded in Table 11. The '-' sign indicates that the execution time of the algorithm is more than 100 hours, as we can see for CLU and PSC which have an expensive computational cost when data increases in size (Covertype data). Nevertheless, we have shown that the AVAR approach can be executed for data up to  $\approx 90,000,000$  instances (friction data). For temporal data, while the other approaches perform two dimensional clustering, the AVAR approach generates clusters based on the characteristic timescale. The AVARbased method will offer fewer advantages for datasets that are more complex or for which averaging is less useful.

When the characteristic AVAR window size is small, the number of the representative prototypes is larger, hence there is a higher insertion cost of these instances in the new reduced table (Segmentation dataset). When the characteristic AVAR window size is large, the number of the representative prototypes is smaller, hence there is a lower insertion cost of these instances in the new reduced table (Magic dataset). When the characteristic AVAR window size is the same for two different datasets, the dataset with a larger number of instances will have a higher computational cost than the dataset with a lower number of instances (Letter and Shuttle datasets). From Table 11 it can be concluded that the execution time of our algorithm is much smaller than the compared algorithms, proving once again that the AVAR approach is fast when applied on Big data.

## 6 CONCLUSIONS AND FUTURE WORK

In this paper we have proposed a computationally efficient granulation algorithm for large temporal databases using Allan Variance. The proposed method systematically determines the characteristic timescale over which data is relevant and calculates the aggregated information for each time window. The total time complexity of this approach is O(n), which is an improvement from existing algorithms  $O(n^2)$ . Experimental results show that the proposed technique considerably reduces the query execution time by reducing the storage requirement, while preserving data integrity. Overall, the presented approach increases query efficiency due to the decrease in data volume while preserving the quality of the result of the queries.

In the future, we plan to incorporate an outlierremoval filtering algorithm to our method and further improve this approach for spatio-temporal databases with respect to both their time domain and spatial layouts. One interesting challenge is that in multidimensional data, different granularities may exist. Choosing the appropriate level of detail or granularity is crucial. To mitigate this challenge, we plan to extend the AVAR-based approach to represent different levels of resolution by creating a hierarchical structure of spatio-temporal data. In addition, we will focus whether and when to recalculate the characteristic timescale, as the AVAR estimator would be useful in an incremental scenario when data keeps coming in. Furthermore, we intend to construct dynamic AVAR estimators which can cope with local changes in the temporal and spatial characteristic sizes.

#### Acknowledgements

This material is based upon work supported by the National Science Foundation under Grant No. 1932138. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the National Science Foundation.

#### **Conflict of interest statement**

On behalf of all authors, the corresponding author states that there is no conflict of interest.

#### REFERENCES

- Allan, D. W. (1966). Statistics of atomic frequency standards. *Proceedings of the IEEE*, 54(2):221–230.
- Bezdek, J. C., Ehrlich, R., and Full, W. (1984). FCM: The fuzzy c-means clustering algorithm. *Computers & Geosciences*, 10(2-3):191–203.
- Dua, D. and Graff, C. (2019). UCI machine learning repository, 2017. URL http://archive.ics.uci.edu/ml.
- Goldberg, D. (1991). What every computer scientist should know about floating-point arithmetic. *ACM Computing Surveys (CSUR)*, 23(1):5–48.
- Gray, J., Chaudhuri, S., Bosworth, A., Layman, A., Reichart, D., Venkatrao, M., Pellow, F., and Pirahesh, H. (1997). Data cube: A relational aggregation operator generalizing group-by, cross-tab, and sub-totals. *Data Mining and Knowledge Discovery*, 1(1):29–53.
- Haeri, H., Beal, C. E., and Jerath, K. (2021). Nearoptimal moving average estimation at characteristic timescales: An Allan variance approach. *IEEE Con*trol Systems Letters, 5(5):1531–1536.
- Hartigan, J. A. (1975). Clustering algorithms. John Wiley & Sons, Inc.
- Helsen, J., Peeters, C., Doro, P., Ververs, E., and Jordaens, P. J. (2017). Wind farm operation and maintenance optimization using big data. In 2017 IEEE Third International Conference on Big Data Computing Service and Applications (BigDataService), pages 179–184.
- Henrikson, J. (1999). Completeness and total boundedness of the Hausdorff metric. *MIT Undergraduate Journal of Mathematics*, 1:69–80.
- IEEE (2019). IEEE standard for floating-point arithmetic. *IEEE Std 754-2019 (Revision of IEEE 754-2008)*, pages 1–84.
- Januzaj, E., Kriegel, H.-P., and Pfeifle, M. (2004). DBDC: Density-based distributed clustering. In *International Conference on Extending Database Technology*, pages 88–105. Springer.
- Jerath, K., Brennan, S., and Lagoa, C. (2018). Bridging the gap between sensor noise modeling and sensor characterization. *Measurement*, 116:350 366.
- Jerath, K. and Brennan, S. N. (2011). GPS-free terrainbased vehicle tracking performance as a function of inertial sensor characteristics. In *Dynamic Systems* and Control Conference, volume 54761, pages 367– 374.

- Johnston, W. (2001). Model Visualization, page 223–227. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA.
- Kaufmann, L. (1987). Clustering by means of medoids. In *Proc. Statistical Data Analysis Based on the L1 Norm Conference, Neuchatel, 1987*, pages 405–416.
- Keogh, E. and Mueen, A. (2017). Curse of dimensionality. In *Encyclopedia of Machine Learning and Data Mining*, pages 314–315.
- Kile, H. and Uhlen, K. (2012). Data reduction via clustering and averaging for contingency and reliability analysis. *International Journal of Electrical Power & Energy Systems*, 43(1):1435–1442.
- Kodinariya, T. M. and Makwana, P. R. (2013). Review on determining number of cluster in k-means clustering. *International Journal*, 1(6):90–95.
- Liu, H. and Motoda, H. (2002). On issues of instance selection. *Data Min. Knowl. Discov.*, 6:115–130.
- Lu, J., Liu, A., Dong, F., Gu, F., Gama, J., and Zhang, G. (2018). Learning under concept drift: A review. IEEE Transactions on Knowledge and Data Engineering, 31(12):2346–2363.
- Lumini, A. and Nanni, L. (2006). A clustering method for automatic biometric template selection. *Pattern Recognition*, 39(3):495–497.
- MacQueen, J. et al. (1967). Some methods for classification and analysis of multivariate observations. In *Proceedings of the Fifth Berkeley symposium on mathematical statistics and probability*, volume 1, pages 281–297. Oakland, CA, USA.
- Maddipatla, S. P., Haeri, H., Jerath, K., and Brennan, S. (2021). Fast Allan Variance (FAVAR) and Dynamic Fast Allan Variance (D-FAVAR) Algorithms for both Regularly and Irregularly Sampled Data. *IFAC-PapersOnLine*, 54(20):26–31.
- Madigan, D. and Nason, M. (2002). Data reduction: sampling. In *Handbook of data mining and knowledge discovery*, pages 205–208.
- Mishra, A. D. and Garg, D. (2008). Selection of best sorting algorithm. *International Journal of intelligent information Processing*, 2(2):363–368.
- NASA. Prediction of Worldwide Energy Resource (POWER) datasets. https://power.larc.nasa.gov/.
- Olvera-López, J. A., Carrasco-Ochoa, J. A., and Martínez-Trinidad, J. F. (2010). A new fast prototype selection method based on clustering. *Pattern Analysis and Applications*, 13(2):131–141.
- Pedrycz, W. (2001). Granular computing: An introduction. In *Proceedings joint 9th IFSA world congress and 20th NAFIPS international conference (Cat. No. 01TH8569)*, volume 3, pages 1349–1354. IEEE.
- Rehman, M. H., Liew, C. S., Abbas, A., Jayaraman, P. P., Wah, T. Y., and Khan, S. U. (2016). Big data reduction methods: a survey. *Data Science and Engineering*, 1(4):265–284.
- Sesia, I. and Tavella, P. (2008). Estimating the Allan variance in the presence of long periods of missing data and outliers. *Metrologia*, 45(6).

- Sinanaj., L., Haeri., H., Gao., L., Maddipatla., S., Chen., C., Jerath., K., Beal., C., and Brennan., S. (2021). Allan Variance-based Granulation Technique for Large Temporal Databases. In *Proceedings of the 13th International Joint Conference on Knowledge Discovery, Knowledge Engineering and Knowledge Management KMIS*,, pages 17–28. INSTICC, SciTePress.
- Sun, X., Liu, L., Geng, C., and Yang, S. (2019). Fast data reduction with granulation-based instances importance labeling. *IEEE Access*, 7:33587–33597.
- Zadeh, L. A. (1997). Toward a theory of fuzzy information granulation and its centrality in human reasoning and fuzzy logic. *Fuzzy sets and systems*, 90(2):111–127.