

Exploring Interpretable Features for Large Time Series with SE4TeC

Jingwei Zuo
DAVID Lab, University of Versailles
University of Paris-Saclay
Versailles, France
jingwei.zuo@uvsq.fr

Karine Zeitouni
DAVID Lab, University of Versailles
University of Paris-Saclay
Versailles, France
karine.zeitouni@uvsq.fr

Yehia Taher
DAVID Lab, University of Versailles
University of Paris-Saclay
Versailles, France
yehia.taher@uvsq.fr

ABSTRACT

Time Series (TS) data are ubiquitous in enormous application fields, such as medicine, multimedia, and finance. In this paper, we present the demonstration with **SE4TeC**: A Scalable Engine for efficient and expressive Time Series Classification, which is applicable to certain fields in Big Data context, where the TS features and their extraction process should be interpretable. **SE4TeC** improves the state of the art solutions by proposing a scalable and highly efficient method to classify TS based on characteristic subsequences (i.e., shapelets). We explain the techniques we adopt, and show how to use **SE4TeC** for exploring the real-life datasets in medical diagnosis and in industrial troubleshooting.

1 INTRODUCTION

Tony is informed unhealthy heart status on the basis of his electrocardiogram (ECG) during a medical diagnosis. An experienced doctor can easily correlate the abnormal ECG with the diseases, then explain to Tony the symbolic abnormality in the ECG and the relevant treatment. Nowadays, Machine Learning technique can partly replace the role of an experienced doctor and do the diagnosis very accurately. In Tony's case, the electrical activity of the heart from physiological sensor is collected as Time Series (TS). From the perspective of the machine, the diagnosis can be considered as a Time Series classification (TSC) problem.

The classical approaches [1] on TSC problems are usually based on the statistical features extracted from time series, such as mean, standard deviation of subsequences, which are assumed to represent the global characteristics of time series. Intuitively, they get very superficial information with low noise tolerance. On the basis of solving these limitations, the **Shapelet** approach [12] has attracted great interest over the past years, owing to its high discriminative feature and good interpretability. However, extracting shapelets from data series has a large computation cost. Even for small sized datasets, the algorithm can take days. This is mainly due to the repeated similarity search between a sub-sequence (i.e., a candidate shapelet) and TS instances in the database. Some typical speed-up techniques (i.e., indexing [11], lower-bounding [6] and early abandoning [12]), introduce always extra parameters, which is difficult to operate without prior knowledge. Some low-dimensional representation methods have also been proposed, such as Piece-wise Aggregate Approximation (PAA) [4], which computes the mean of each subsequence of time series in a given length, and transforms the raw data in coarse-grained sub-components. Symbolic Aggregate approXimation (SAX) [5], transforms subsequences of raw time series into value-characterized symbols, which is eligible for a hierarchic indexing

iSAX [2] to accelerate similarity search. Nevertheless, scalability remains the bottleneck.

Our work is biased towards raw time series processing which has a higher accuracy performance, but a relatively high time complexity [12]. Unlike some hardware-based implementations, such as using GPUs to accelerate the similarity calculation [3], we focus on the scalability of TSC based on shapelet extraction. Traditional TSC algorithms on raw TS data are not applicable for big data context, because of their low scalability. Here are our contributions in this paper:

- (1) We propose a novel method to assess the importance of shapelets in batches
- (2) We introduce a scalable engine to extract the shapelets
- (3) Based on the scalable engine, we propose an optimization strategy to speed-up the shapelets extraction.

The rest of this paper is organized as follows. In Section 2, we review the background and state the research problems. We present our scalable engine for Time Series Classification in Section 3. Section 4 shows an empirical evaluation of our method, as well as a guidance for the demonstration. Finally, we give our conclusions and perspectives for future work in Section 5.

2 BACKGROUND

2.1 Definitions and Notations

We start with defining the notions used in the paper:

Definition 1: A Time Series T is a sequence of real-valued numbers $T=(t_1, t_2, \dots, t_i, \dots, t_n)$, where n is the length of T .

Definition 2: A subsequence $T_{i,m}$ of Time Series T is a continuous subset of values from T of length m starting from position i . $T_{i,m} = (t_i, t_{i+1}, \dots, t_{i+m-1})$, where $i \in [0, n - m + 1]$.

Definition 3: Shapelet \hat{s} is a time series subsequence which is particularly representative of a class. As such, it shows a shape which can distinguish one class from the others.

Definition 4: A Dataset D is a collection of time series T_i , and its class label c_i . Formally, $D = \langle T_1, c_{j_1}, \dots, T_N, c_{j_N} \rangle$, where N is the number of instances in D . $C = c_1, c_2, \dots, c_{|C|}$ is a collection of class labels, where $|C|$ denotes the number of labels.

Definition 5: Z-Normalization Time Series is a formal representation of Time Series, which is defined as $ZNormal(T) = \frac{T - \mu}{\sigma}$, where μ is the sample mean, σ is the standard deviation:

$$\mu = \frac{1}{m} \sum_{i=1}^n t_i, \quad \sigma^2 = \frac{1}{m} \sum_{i=1}^n t_i^2 - \mu^2 \quad (1)$$

Z-Normalization allows us to focus on the structural feature of T , rather than its amplitude value. It addresses the problem of data stability. For instance, assume that the Euclidean Distance (ED) between two time series $T_{x,m}$, $T_{y,m}$ is expressed as follows:

$$ED_{x,y} = \sqrt{\sum_{i=1}^m (t_{x,i} - t_{y,i})^2} \quad (2)$$

Some little changes (e.g., the noise) will cause an evident bias for the result, Z-Normalization is a way of smoothing the bias value.

Definition 6: *Normalized Euclidean Distance (N-ED)*, is expressed by the formula $\sqrt{\frac{1}{m} \sum_{i=1}^m (t_{x,i} - t_{y,i})^2}$

Definition 7: *Distance Profile DP_i* is a vector which stores the Normalized Euclidean Distance between a given subsequence/query $T_{i,m}$ and every subsequences $T'_{j,m}$ of a target Time Series T' . Formally, $DP_{i,j}^m = \text{dist}(T_{i,m}, T'_{j,m}), \forall j \in [0, n' - m + 1]$

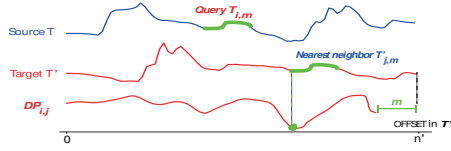


Figure 1: Distance Profile between Query $T_{i,m}$ and target time series T' , where n' is the length of T' . Obviously, $DP_{i,j}$ can be considered as a meta TS annotating target T'

Definition 8: *MASS*, namely Mueen’s ultra-fast Algorithm for Similarity Search, computes Distance Profile based on Fast Fourier Transform (FFT), which requires just $O(n \log n)$ time, other than $O(nm^2)$ time in classical *N-ED* similarity search.

Definition 9: *Matrix Profile MP* is a vector of distance between subsequence $T_{i,m}$ in source T and its nearest neighbor $T'_{j,m}$ in target T' . Formally, $MP_i^m = \min(DP_i^m)$, where $i \in [0, n - m + 1]$.

Unlike the distance profile, the matrix profile is a meta TS annotating the source time series. The highest point on *MP* corresponds to the TS discord, the lowest points correspond to the position of a query which has a similar matching in target TS.

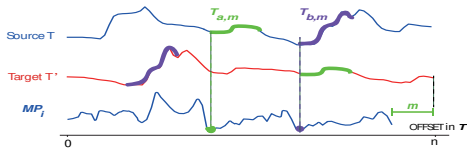


Figure 2: Matrix Profile between Source time series T and Target time series T' , where n is the length of T . Intuitively, MP_i shares the same offset as source T

2.2 Evaluation of Candidate Shapelets

The quality of a candidate shapelet \hat{s} , can be assessed by its ability to separate the instances of different class in the dataset D . A prerequisite of the quality measure for \hat{s} , is that a set of distance $D_{\hat{s}}$ must be calculated, where $D_{\hat{s}} = D_{\hat{s},1}, D_{\hat{s},2}, \dots, D_{\hat{s},n}$, n is the number of T in dataset D .

Information Gain, an evaluation method based on Decision Tree, is widely adopted in previous works [12]. An iteration test of $D_{\hat{s}}$ is conducted to extract the split distance which brings the highest Information Gain. The distance instance will be applied as a property of candidate Shapelet, to check the inclusion between the candidate and time series. Another simple approach, is to use the F-Statistic [7], based on the difference of means in an Analysis of Variance A(NOVA). The main idea of this statistic method, is to assess the difference in distributions of \hat{s} between the class distances.

2.3 Problem Statement

The high degree of coupling inside classical TSC algorithm [12] leads to the problem of not being able to parallelize. The speed-up method such as Early Abandoning [12] is based on classical Euclidean Distance measure, which has a time complexity of $O(N^2 n^4)$ with several orders of magnitude higher than

MASS [8]: $O(N^2 n^3 \log n)$. Another common trick played by previous work [12]: If we know \hat{s} is a low-quality candidate, then any similar subsequence \hat{s}' to \hat{s} must also result in a low quality and therefore, a costly computation of the distance set $D_{\hat{s}'}$ (evaluation of \hat{s}') can be skipped. However, a candidate shapelet is evaluated by its quality ranking among all candidates of the same length. Assume that the distributed nodes have generated from dataset a collection of candidates \hat{s}_j , an aggregation operation between nodes is required to extract the candidate with the best quality. Extra aggregations will be made along with the iteration of candidate length. Apparently, the acceleration from the classical pruning techniques can be easily offset by the communication cost caused by the aggregation.

3 SYSTEM OVERVIEW

The main idea of our system is that the calculation should be shared and executed independently, less communication between the nodes, more powerful the algorithm would be.

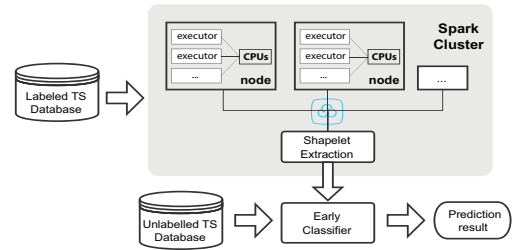


Figure 3: System overview

3.1 Main Structure

The conventional Time Series classification problems are tackled with nearest neighbor (kNN) algorithm [12] due to its easy-design feature. As shown in Figure 3, on the basis of kNN, an early classifier [10] adopted in the system allows to give the prediction result as earlier as possible without waiting for the entire sequence. The processing of labelled Time Series data requires to be flexibly arranged for nodes in the cluster, where the executors share the CPU/memory resource. To this end, a suitable algorithm is applied here allowing assignment of computing tasks which are relatively independent of each other.

3.2 SMAP: Shapelet Extraction on Matrix Profile

Matrix Profile provides a meta-data which facilitates the representation of a complex correlation between two time series. As shown in *Algorithm 1*, *SMAP* takes the time series as the smallest processing unit between nodes, and utilizes the normalized quality to extract the most important parts in each processing unit and then merges them by an aggregation process. For this reason, the number of candidate shapelet could be greatly reduced. Moreover, a single aggregation is required to get the global shapelet result of different class. In *line 5*, dataset is broadcast to distributed nodes in order to reduce the communication cost caused by accessing the common data. Then, each cluster partition shares the computing tasks for a set of time series. The function *computeDiscrimP* aims at computing in batches the quality of candidate shapelets. The visualized process is shown in **Figure 4**.

The batch quality of instances in a TS is defined by *Discrimination Profile*, which refers to the concept *Representative Profile*:

$$RP(T_i^C, D) = \text{avg}(MP_{T_i^C, T_j}) \quad (3)$$

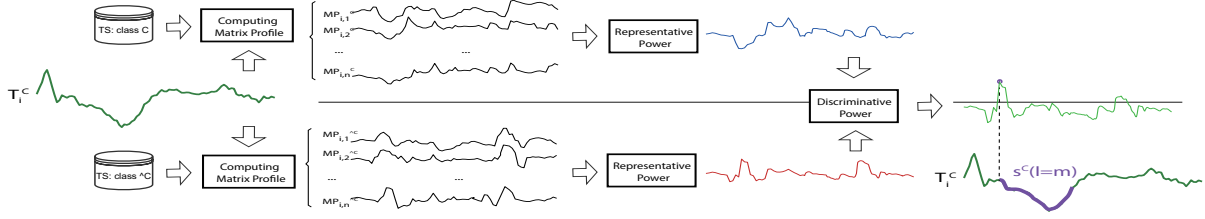


Figure 4: Discrimination Profile Extraction

Algorithm 1: SMAP(Shapelet on Matrix Profile)

```

Input: Dataset  $D$ , classSet  $\hat{C}$ ,  $k$ 
Output:  $\hat{S}$ 
1  $minLength \leftarrow 2, maxLength \leftarrow getMinLen(D)$ 
2  $double[] DiscmP \leftarrow [], double[] DistThresh \leftarrow [], \hat{S} \leftarrow \emptyset$ 
3  $D.cache();$  //cache all the dataset in the cluster, where each time series has a unique ID
4 MapPartition ( $Set\ of\ < ID, T >: T_{set}$ )
5   for  $< ID, T > \in T_{set}$  do
6     for  $m \leftarrow minLength\ to\ maxLength$  do
7        $DiscmP[m], DistThresh[m] \leftarrow computeDiscmP(T, D, m)$ 
8        $DiscmP[m] \leftarrow DiscmP[m] * \sqrt{1/l}$ 
9      $DiscmP \leftarrow pruning(DiscmP)$ 
10    emit( $DiscmP, DistThresh$ )
11 MapAggregation ( $class, (DiscmP, DistThresh)$ )
12   for  $c \in \hat{C}$  do
13      $\hat{S}' \leftarrow getTopK(DiscmP[c], DistThresh[c], k)$ 
14     for  $\hat{s} \in \hat{S}'$  do
15        $\hat{s}.matchingIndices \leftarrow getMatchingIndices(\hat{s}, D)$ 
16      $\hat{S} \leftarrow \hat{S} \cup \hat{S}'$ 
17 return  $\hat{S}$ 

```

where $T_j \in D^C$. Representative Profile targets thus on the minimal processing unit (i.e., time series) in SMAP, which shows a vector of Representative Power of each instance in the processing unit. To put it simply, the *Representative Power* of a subsequence in class C , is its normalized distance to the global instance cluster of class C . Intuitively, it represents the relevance between the subsequence (i.e., the candidate shapelet) and the class.

The fact that a subsequence is discriminative for its class towards others, can be expressed by the difference of Representative Power from class C to others (OVA, *one-vs-all*). *Discrimination Profile* is then defined as follows:

$$DiscmProfile(T_i^C, D) = -(RP(T_i^C, D^C) - RP(T_i^C, D^{1^C})) \quad (4)$$

A quality *Normalization* in line 8 is made which allows to assess the Discrimination Power for shapelet of different length in an uniform way. Similar as the concept *Information Gain*, but Discrimination Profile is a technique more interpretable serving to assess the candidate shapelets. Moreover, in this manner, a split distance can be given directly, other than iterating every possible distance and deciding the best one with the highest Information Gain, in time $O(N^2n^2)$. A strategy to check if T contains a shapelet can be defined as the following:

$$sInT(T, \hat{s}^C) = \begin{cases} true, & \text{if } dist(T, \hat{s}^C) \leq RP(\hat{s}^C, D^C) \\ false, & \text{otherwise} \end{cases} \quad (5)$$

3.3 Optimization Strategy

The pruning function in line 9 is capable of eliminating the number of candidate shapelet, and then reducing the communication cost during the aggregation process. We can simply take the "TopK" strategy, which extracts the biggest K values of the *Discrimination Profile*. However, such a technique is far from lightening the computation during *MapPartition* process.

Since each processing unit should be independent from each other, a tenable technique for updating the profile of a long range query could be adopted. The *Lower Bounding distance* [6] is defined to estimate a minimal possible Z-Normalized Euclidean Distance between two subsequences $T_{i,l+k}$ and $T_{j,l+k}$, based on the distance already computed between $T_{i,l}$ and $T_{j,l}$. Compared to a linear time complexity of computing the exact distance, LB Distance Profile can be calculated in a constant time, which can accelerate greatly the computation of Matrix Profile in Figure 4. For example, from shapelet length $l = m$ to $m + 1$, the time complexity of computing the distance $d_{i,j}^{l+1}$ is $O(n \frac{m(m-1)}{2})$, where $j \in [0, n \frac{m(m-1)}{2}]$ which represents the number of subsequences in D , n is the number of instance in D , m is the length of the longest instance in D . Accordingly, LB distance takes $O(n \frac{m(m-1)}{2})$ which shows an apparent advantage when the query length is relatively long. Lower Bounding distance [6] is defined as:

$$LB(d_{i,j}^{l+k}) = \begin{cases} \sqrt{l} \frac{\sigma_{j,l}}{\sigma_{j,l+k}}, & \text{if } q_{i,j} \leq 0 \\ \sqrt{l(1 - q_{i,j}^2)} \frac{\sigma_{j,l}}{\sigma_{j,l+k}}, & \text{otherwise} \end{cases} \quad (6)$$

where $q_{i,j} = \frac{\sum_{p=1}^l (t_{j+p-1} - t_{i+p-1})}{\sigma_{i,l} \sigma_{j,l}} - \mu_{i,l} \mu_{j,l}$

Empirically, the matching subsequence $T_{j,l}$ which is the nearest neighbor of $T_{i,l}$, can deduce a longer subsequence $T_{j,l+1}$, which is probably the nearest neighbor of $T_{i,l+1}$. Assume that the matching subsequence keeps in the same position in T_{target} when query $T_{i,l}$ length increases, then the time complexity for computing the minimal distance between $T_{i,l}$ and T_{target} is $O(l)$, other than $O(l(n-l+1))$. As mentioned in **Definition 9**, $MP_i^m = \min(DP_i^m)$, the main idea here is to utilize LB Distance to accelerate the computation of $\min(DP_i^m)$, other than computing the entire DP_i in a higher time complexity.

4 ABOUT THE DEMONSTRATION

This demonstration is intended to show a distributed approach to extract features from large-scale time series and to make the extraction process and extracted features very easy to understand and interpret, thanks to the visual power of shapelets. Through this demonstration, the attendees will have a general understanding of time series, as well as its application areas and the current challenges. With two real datasets, attendees will have the opportunity to experience and interact with SE4TeC from two aspects:

- (1) *Practical operation for distributing computation tasks*: The attendees are invited to connect to our elastic cluster, and will further explore the distribution mechanism for feature extraction. For instance, the relationship between performance and adjustable parallelism, the progress monitoring of parallel tasks on each distributed node, etc.
- (2) *Exploration of shapelet extraction process*: Based on a small-sized dataset, the attendees can interactively perform each intuitive step shown in Figure 4, and are invited to analyze the hidden meaning behind each intermediate features.

In the view of the attendees, the reliability of extracted shapelets can be proved by utilizing naked eyes or a smart classifier, on a test instance.

4.1 Demonstration Platform

SE4TeC is implemented by Python3.6, powered by Apache Spark. The program is executed on AWS EMR cluster. We provide also an 1-click cluster based on Docker, to facilitate the attendees to replay the distributed test offline. The baseline of the evaluation is USE in [9], which utilizes the traditional method for shapelet extraction based on Information Gain. 1NN classifier is applied for all accuracy tests, and 5 shapelets are extracted for each class.

4.2 Demonstration Scenario

Due to page limitations, here we show two examples, more details and videos can be found on the demonstration page¹.

4.2.1 ECG medical diagnosis. The dataset ECG200, from MIT-BIH Long-Term ECG Database (ltddb), is collected by two electrodes which record the brain activities in distinct body positions. Each heartbeat has an assigned label of normal or abnormal. All abnormal heartbeats are representative of a cardiac pathology known as supraventricular premature beat. ECG200 contains **100** labelled records with a fixed length of 96.

4.2.2 Wafer industrial troubleshooting. To put the evaluation into larger scale context, we choose the time series dataset **Wafer**, which contains **1000** training records with a fixed length of 152. The dataset, collected during the manufacture of semiconductor microelectronics, comprises a collection of sequences for the measurements recorded by one vacuum-chamber sensor during the etch process applied to one silicon wafer. The class of manufacture quality can be normal or abnormal.

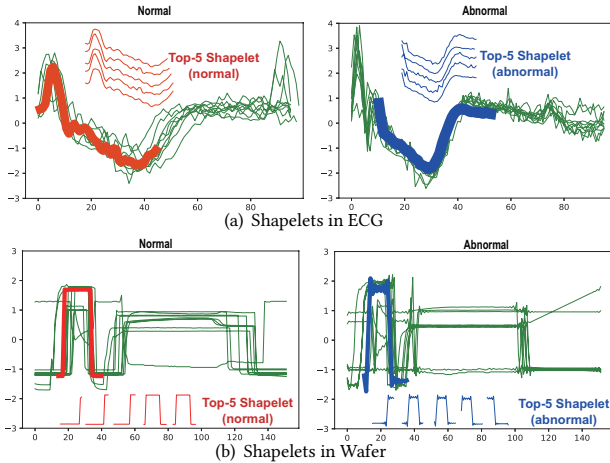


Figure 5: Interpretable Shapelet Feature Results

Reliability & Interpretability: Through the demonstration, the attendees are capable of extracting the shapelets in various manners, and comparing their difference. The shapelets extracted by SMAP are shown in Figure 5, which has a relatively higher prediction accuracy than USE: 84%/76% (ECG), 97%/92% (Wafer).

Scalability: The performance results are shown in Figure 6. SMAP_{LB} shows a gain in performance: **24.5X** (ECG), **587.3X** (Wafer) faster. As the size of ECG (i.e., 100) is lower than the parallelism power when we expand the cluster to 10 nodes, therefore, according to the time $O(N^2 n^3 \log n)$, the instance length n will be the decisive factor in execution time. We should know

¹<https://github.com/JingweiZuo/SE4TeC>

that the speed-up performance relies on the computing power of the cluster. We are more inclined to consider its parallel capability which can be assessed by the aggregation cost between distributed nodes. From 1 to 30 nodes on cluster mode, the aggregation cost for *Wafer* increases by **181%**, the total cost drops to **0.68%**. Obviously, considering the gain, the aggregation cost can be ignored when we expand the cluster to a larger scale.

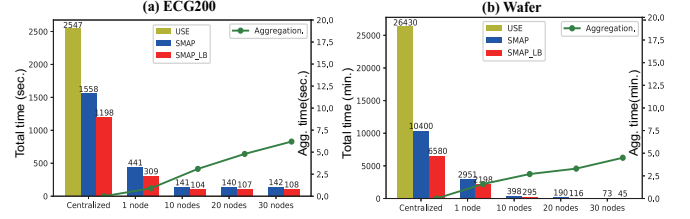


Figure 6: Scalability performance

5 CONCLUSION

In this paper, we propose a novel methodology, namely SMAP, for Time Series Classification. SMAP adopts the concept Matrix Profile, and extracts the shapelet features in a scalable and interpretable manner. Within SMAP, the Discrimination Profile is defined to assess in batches the quality of candidate shapelets. On the basis of Lower Bounding, we propose also an acceleration strategy, which works appropriately for distributed environment. The satisfactory results proved the efficiency of the scalable approach, and testified its competitiveness. Different optimizations are in our planning list. Specifically, to expand SMAP for longer time series by reducing the data dimension, but meanwhile to conserve its high accuracy and scalability.

ACKNOWLEDGEMENT

This research was supported by DATAIA convergence institute as part of the « Programme d'Investissement d'Avenir », (ANR-17-CONV-0003) operated by DAVID Lab, University of Versailles Saint-Quentin, and MASTER project that has received funding from the European Union's Horizon 2020 research and innovation programme under the Marie-Slodowska Curie grant agreement N. 777695.

REFERENCES

- [1] Ling Bao and Stephen S. Intille. 2004. Activity recognition from user-annotated acceleration data. Springer, 1–17.
- [2] Alessandro Camerra et al. 2010. iSAX 2.0: Indexing and Mining One Billion Time Series. In *Proc. ICDM 2010*. 58–67.
- [3] Kai Wei Chang et al. 2012. Efficient pattern-based time series classification on GPU. In *Proc. ICDM 2012*. 131–140.
- [4] Eamonn Keogh, Kaushik Chakrabarti, Michael Pazzani, and Sharad Mehrotra. 2001. Dimensionality Reduction for Fast Similarity Search in Large Time Series Databases. *Knowledge and Information Systems* (01 Aug 2001), 263–286.
- [5] Jessica Lin et al. 2003. A Symbolic Representation of Time Series, with Implications for Streaming Algorithms. In *Proc. 8th ACM SIGMOD (DMKD '03)*.
- [6] Michele Linardi and Yan et al. Zhu. 2018. VALMOD: A Suite for Easy and Exact Detection of Variable Length Motifs in Data Series. In *Proc. SIGMOD'18*.
- [7] Jason Lines, Luke M Davis, Jon Hills, and Anthony Bagnall. 2012. *A Shapelet Transform for Time Series Classification*. Technical Report.
- [8] Chin-Chia Michael Yeh et al. 2016. Matrix Profile I: All Pairs Similarity Joins for Time Series: A Unifying View That Includes Motifs, Discords and Shapelets. In *Proc. ICDM '16*. 1317–1322. <https://doi.org/10.1109/ICDM.2016.0179>
- [9] Raef Mousheimish, Yehia Taher, and Karine Zeitouni. 2017. Automatic Learning of Predictive CEP Rules: Bridging the Gap between Data Mining and Complex Event Processing. In *Proc. DEBS '17*. 158–169.
- [10] Zhengzheng Xing, Jian Pei, and Philip S. Yu. 2009. Early Prediction on Time Series: a Nearest Neighbor Approach. *21st IJCAI* (2009), 1297–1302.
- [11] Djamel-Edine Yagoubi et al. 2017. DPISAX: Massively Distributed Partitioned iSAX. In *Proc. ICDM 2017*. 1–6.
- [12] Lexiang Ye and Eamonn Keogh. 2009. Time series shapelets: A New Primitive for Data Mining. *Proc. 15th ACM SIGKDD* (2009).