

Evaluating Explanation Methods for Multivariate Time Series Classification

Davide Italo Serramazza, Thu Trang Nguyen, Thach Le Nguyen, and
Georgiana Ifrim

School of Computer Science, University College Dublin, Ireland
{davide.serramazza,thu.nguyen}@ucdconnect.ie
{thach.lenguyen,georgiana.ifrim}@ucd.ie

Abstract. Multivariate time series classification is an important computational task arising in applications where data is recorded over time and over multiple channels. For example, a smartwatch can record the acceleration and orientation of a person’s motion, and these signals are recorded as multivariate time series. We can classify this data to understand and predict human movement and various properties such as fitness levels. In many applications classification alone is not enough, we often need to classify but also understand what the model learns (e.g., why was a prediction given, based on what information in the data). The main focus of this paper is on analysing and evaluating explanation methods tailored to Multivariate Time Series Classification (MTSC). We focus on saliency-based explanation methods that can point out the most relevant channels and time series points for the classification decision. We analyse two popular and accurate multivariate time series classifiers, ROCKET and dResNet, as well as two popular explanation methods, SHAP and dCAM. We study these methods on 3 synthetic datasets and 2 real-world datasets and provide a quantitative and qualitative analysis of the explanations provided. We find that flattening the multivariate datasets by concatenating the channels works as well as using multivariate classifiers directly and adaptations of SHAP for MTSC work quite well. Additionally, we also find that the popular synthetic datasets we used are not suitable for time series analysis.

Keywords: Time Series Classification · Explanation · Evaluation

1 Introduction

Real-world time series data are often multivariate, i.e., data collected over a period of time on different channels. An example is human motion data collected from participants wearing a tri-axial accelerometer on their dominant wrist. The tri-variate data can be examined to identify epilepsy convulsions in everyday life [25]. Another example is traffic data where multiple sensors are set up at different locations to measure the traffic occupancy in a city¹.

¹ <https://pems.dot.ca.gov/>

While univariate time series have been the main research focus, there is a steadily growing interest in multivariate time series (MTS), in particular for the classification task (MTSC). The release of the MTSC benchmark [2], a collaborative effort by researchers from multiple institutions, is an important milestone that has accelerated studies of MTSC methods.

Explainable AI is another important topic due to the explosion of interest in complex machine learning models and deep learning methods. Pioneers in this field have been working mostly on text and image data and, as a result, a number of explanation frameworks including LIME [20], DeepLift [14], Shapley [15] have been introduced. The similarity between image and time series data allows such techniques to be adapted to time series models [26]. Nevertheless, there are some notable differences between images and time series. Firstly, images are usually represented using RGB encoding and all the 3 channels contain necessary information, while for time series it is common to have channels that do not contribute to, or even hinder, the classification decision. Secondly, in images there is a lot of homogeneity in the pixel values while moving between pixels belonging to the same objects and a sharp difference when moving between pixels belonging to different objects. In time series, it is less common to find such a strong locality, especially across all the channels. Furthermore, the data magnitude and pre-processing, such as normalisation, are important factors for time series, but less so for images.

In this work, we focus on methods for explaining MTSC as this is an important open problem that is often as important as the classification itself. In a scenario in which people wear accelerometers on their body while executing a physical exercise, other than classifying the exercise as correctly executed or not, it is also important to provide feedback to users, e.g., an explanation of why the exercise was incorrectly executed by pointing out the relevant data.

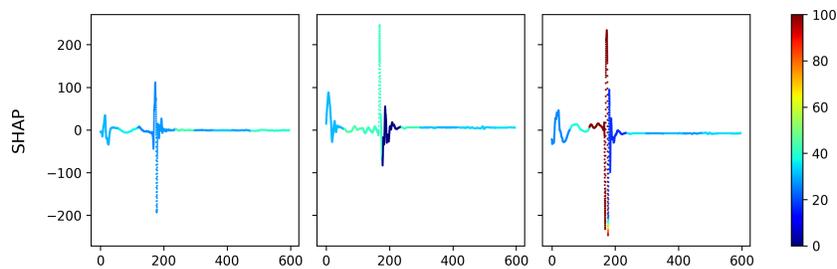


Fig. 1: Sample multivariate time series and explanation heat map. The 3 plots show the x, y, z channels for a jump sample.

In this paper, a *multivariate time series explanation* is a 2D saliency map [3] highlighting the importance of each time series channel and each time point for the classification decision, as illustrated in Figure 1. A proper MTSC explanation should be able to point out for each channel the relevant time points that may

be located at different parts of the time series. For example, CAM [27] was designed for explaining univariate time series thus it can not identify important time points which vary across channels.

In this work we aim to analyse and evaluate a few MTSC explanation methods we found in the literature. Throughout our literature research, the only bespoke MTS explanation methods found are all tailored for deep learning methods (especially CNN), while few others are able to provide a 2D heat map by adapting *univariate time series explanation* to work in a multivariate scenario (most of the time by flattening the dataset and reshaping the 1D heat map into a matrix).

The lack of bespoke multivariate time series explanations, combined with the lack of explanation evaluation methods, is an important gap in the scientific literature. The main aim of this work is to study and evaluate existing MTSC explanation methods in order to start addressing this gap.

Our main contributions in this paper are:

- We analyse the literature on saliency-based explanation methods for MTSC and find very few bespoke methods, all of which are designed for deep learning models. Among these, we select dCAM [3] which extends CAM, a very popular method for time-series and image explanations.
- We conduct experiments using state-of-the-art multivariate time series classifiers ROCKET [6] and dResNet [3] and explanation methods SHAP [16] and dCAM [3]. We study ways to adapt SHAP to work with multivariate time series and compare it to the bespoke MTSC explanation method dCAM.
- We use 3 synthetic datasets and 2 real-world datasets to compare the classifiers and the explanations. We evaluate the explanations both quantitatively, using the methodology proposed in [18], as well as qualitatively. We find that for truly multivariate datasets (i.e., where multiple channels are needed for the correct classification), ROCKET-SHAP works better than dCAM, but is also more computationally expensive. We also find that flattening the datasets by concatenating the channels and using univariate classifiers works as well as using multivariate classifiers directly.

In the rest of the paper, in Section 2 we discuss prior work addressing the MTSC explanation task. In Section 3 we formally define the problem addressed, the classifiers and the explanation methods used in the experiments. In Section 4 we describe the datasets used in our study, in Section 5 we describe our experiments and in Section 6 we summarise our main findings.

2 Related Work

Explanation Methods adapted from Univariate to Multivariate TSC.

Some multivariate time series explanation methods are simple adaptations of methods developed for univariate data. In [1], the authors explain the adapted classifiers by applying the timeXplain [17] framework on each channel *independently*. The result is a multivariate explanation that highlights the important

segments in each channel of the multivariate sample. Nonetheless, it is arguable whether this approach is appropriate since the explained model(s) (univariate) and the model that needs to be explained (multivariate) are not the same. Additionally, it is not clear if the accuracy of the channel-wise univariate model is similar or worse than that of the multivariate model, and this is not discussed in the paper.

Bespoke Explanation Methods for MTSC. Most of the previous explanation methods designed for MTSC are tailored to deep learning methods, which are not state-of-the-art with regard to classification accuracy. In [3], the authors discussed the drawbacks of the CAM explanation method for MTS data. CAM can only produce a univariate saliency map, thus it is unable to identify the important channels. Features that depend on more than one channel are also not detectable. dCAM, proposed in the same paper, addressed these limitations by rearranging the input time series with all the permutations of the channels. The paper shows that this technique can be applied to any architecture with a Global Average Pooling layer (GAP) such as ResNet or InceptionTime. dCAM limitations are discussed by comparing this method with other deep learning explanation methods, as for instance it was shown that dCAM is not the best option when dealing with multivariate datasets that can be classified focusing on just one channel, but there is no comparison against model agnostic methods such as SHAP [15] or LIME [20].

Evaluation of Explanation Methods for MTSC. While explanation methods for MTSC are few, works on evaluating such methods are even fewer. For univariate time series, several approaches have been proposed to compare explanation methods from different angles. The work in [5,11] benchmarks the methods with controllable synthetic datasets. The work of [8] attempted to extract "ground-truth" explanations with a white-box classifier. The "ground-truth" explanation is then used to evaluate post-hoc explanations. AMEE [18] is a recent framework to quantitatively compare explanation methods on a dataset by perturbing input time series and measuring the impact on the classification accuracy of several classifiers. For multivariate time series, recently [24] designed an evaluation framework that is also based on the idea of perturbation, but the work is only limited to evaluating deep learning classifiers and associated explanations. The paper also proposed a synthetic multivariate time series dataset to benchmark explanation methods.

3 Background

A multivariate time series X can be represented as a $d \times L$ matrix, where the d rows are also called channels and the L columns store the values associated with each channel at every time point. Hence X_i^j is the value of the time series at time point i and channel j , with $0 \leq i < L$ and $0 \leq j < d$. We also refer to X^j as the univariate time series at channel j , therefore X can be written as $X = [X^0, X^1, \dots, X^{d-1}]$.

An explanation of a time series X is a saliency map W that provides an importance weight for each data point (at every time point i and every channel

j) in the time series. Hence the saliency map can also be represented by a $d \times L$ matrix. A common visualisation method for the saliency map is a heat map where more important data points are highlighted with warmer colours.

An explanation method for MTSC is a method that, given the input MTS, can produce a saliency map highlighting the relevance of each time point to the classifier decision. Intrinsically explainable models such as Ridge Classifier can also be an explanation method while black-box models such as ResNet (dResNet) and ROCKET need a post-hoc explanation method.

In our experiments we compare three different classifiers and explanation methods: ROCKET [6] coupled with SHAP [16], dResNet coupled with dCAM [3] and the Ridge Classifier [10] which is an intrinsically explainable model. We also use a random explanation (a matrix of random weights) as a sanity check.

3.1 Classification Methods

The first classifier we used is **ROCKET** [6] which was originally designed for UTS, but also has an adaptation for MTS: it applies a large set of random convolution kernels to the time series in order to transform it into tabular data with 20,000 features. It introduced some key concepts such as dilation, proportion of positive values (PPV), etc., starting an algorithm family in which recent members such as Minirocket [7], MultiRocket [23] are improvements of the original idea. All the hyper-parameters for ROCKET were learned from the UCR archive. The authors selected 40 random datasets from the archive and used them as the development set to find the best values for the hyper-parameters. Finally, all the kernel weights are sampled from a distribution $\mathcal{N}(0, 1)$. After the transformation step, the authors use classic linear classifiers Ridge or Logistic Regression.

The second classifier is **dResNet** [3] which is a variation of ResNet [9]. This last one, originally designed for image classification, was used for the first time in TSC in [26]. It introduced the key concept of *shortcut connections* to mitigate the gradient vanishing problem. The main architecture of the network is composed of three consecutive blocks which in turn contain three different convolutional layers. These three blocks are followed by a GAP layer and a softmax layer for classification.

The dResNet version uses the same architecture with two differences specifically designed to work alongside dCAM. Firstly, for a multivariate time series X with d channels, i.e., a matrix $X = [X^0, X^1, \dots, X^{d-1}]$, the input $C(X)$ of the network will be a 3D tensor:

$$C(X) = \begin{bmatrix} X^{d-1} & X^0 & \dots & X^{d-3} & X^{d-2} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ X^1 & X^2 & \dots & X^{d-1} & X^0 \\ X^0 & X^1 & \dots & X^{d-2} & X^{d-1} \end{bmatrix}$$

In other words, the input is turned from a 2D matrix into a 3D one in which each row contains the d channels in different positions. The second change was

to turn the convolution shapes from $1D$ to $2D$ to have the same output shape as ResNet. These changes were made so that the network is able to capture patterns depending on multiple channels while still learning on individual channels.

The third model we used is the well-known **Ridge Classifier** [10], meant to be a baseline in the experiments: we used the scikit-learn [19] package RidgeCV using Cross Validation, leaving the other solver parameters as default. This classifier disregards the time ordering in each time series as it treats each time series as a tabular vector of features.

3.2 Explanation Methods

The first explanation method considered in this paper is **SHAP** [15] which measures feature importance using Shapley values borrowed from game theory. SHAP quantifies the contribution of each feature by examining the differences in the model output when a specific feature is masked, i.e., it is replaced with a specific value and when it is not. SHAP considers every possible masking configuration, thus is computationally expensive. The timeXplain library [17] applies SHAP on the UTSC task by dividing the time series into segments, each is treated as a feature. The segmentation exploits locality in time series and significantly reduces the number of features before applying SHAP. As SHAP is a model-agnostic method, it works with any TSC model. We couple it with ROCKET due to its efficiency and accuracy.

The second explanation method (used along dResNet), is **dCAM** [3]. It computes *CAM* [27] for each row of the input (described in Section 3.1), resulting in a $2D$ matrix \mathcal{M} where all channels are brought back to their original positions to evaluate their contribution. Since the network is trained to compute meaningful predictions regardless of the order in which the channels are provided, dCAM computes k different matrices \mathcal{M} each of them obtained by a different random permutation of the channel order: all these k matrices are then averaged into $\hat{\mathcal{M}}$. The final step to retrieve the explanation W consists in filtering out uninformative time points and uninformative channels using respectively the average value of $\hat{\mathcal{M}}$ in each channel and the variance of all positions for a single channel. dCAM can tell how important a time point was for the classification by taking the differences in $\hat{\mathcal{M}}$ when the time point is present in different positions.

The third explanation method is **Ridge**. As mentioned before, this method is intrinsically explainable because the explanation weights are the weights learned by the classifier. The model is basically a vector of coefficients for each feature, i.e., data point in the time series.

The final explanation method **Random** is a baseline that generates the saliency map W by sampling values randomly from a continuous uniform distribution. The idea is that any good explanation method should provide a better explanation than the random one.

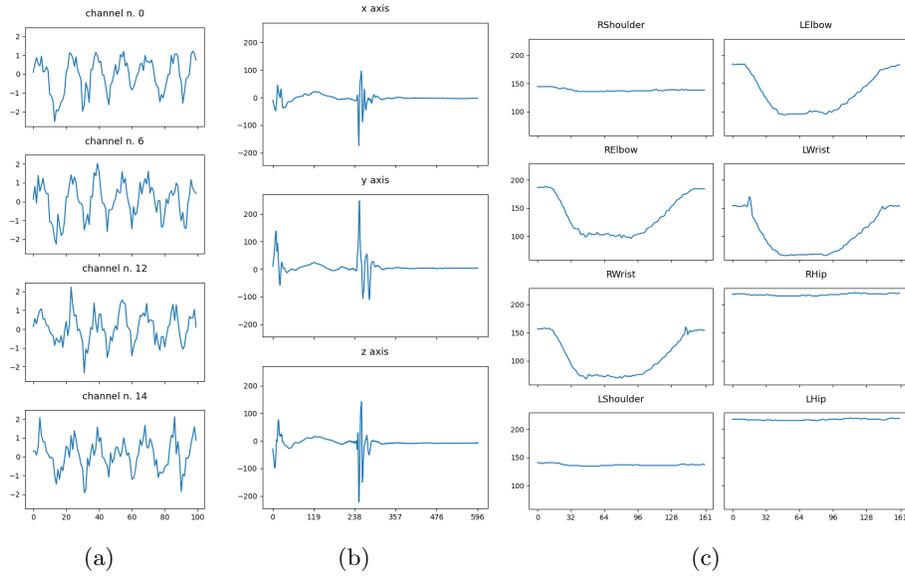


Fig. 2: Sample time series: Fig 2a PseudoPeriodic negative sample. Fig 2b one instance from CMJ Bend. Fig 2c one instance from MP Normal.

4 Datasets

We work with 3 synthetic multivariate time series classification datasets and 2 real-world ones. In Figure 2 we present one sample from one synthetic dataset and one sample each for the real-world datasets.

4.1 Synthetic Datasets

For the synthetic datasets, we use the multivariate time series classification benchmark by Ismail et al. [12]. We generated three different datasets, using the *Pseudo Periodic*, *Gaussian* and *Auto Regressive* distributions. Each has 100 samples in both train and test sets, with $L = 100$ and $d = 20$. The two classes for classification are *positive* and *negative*. The discriminative data points are stationary and within a small box, i.e., X_i^j is discriminative if and only if $10 \leq i < 20$ and $0 \leq j < 10$. In other words, 50% of the channels and 10% of the time steps are relevant. Overall, only 5% of the time series matter for predicting the class.

4.2 Real-World Datasets

The first real-world dataset is **Counter Movement Jump (CMJ)** [13]. The data were collected using accelerometer sensors attached on the participants while performing the counter-movement jump exercise. The three classes are: jumps with acceptable form (class 1), with their legs bending during the flight

(class 2), and with a stumble upon landing (class 3). The training set has 419 samples while the test set has 179 samples. Each time series has 3 channels ($d = 3$) that record the acceleration in x , y , and z axis. The original data is variable-length thus we resampled every time series to the same length ($L = 596$). From the domain experts, we know that the distinctions between classes are more observable on channel y , thus it makes this channel the most important one.

The second real-world dataset is **Military Press** (MP) [22]. To collect the data, 56 participants were asked to perform the Military Press strength-and-conditioning exercise. Each of them completed 10 repetitions in the normal form and another 30 in induced forms, with 10 repetitions each (simulating 3 types of errors). The time series were extracted from video using the OpenPose library [4]. The dataset has 1452 samples in the training set and 601 in the test set, each time series has 161 time points and 50 channels corresponding to the x, y coordinates of 25 body parts. From the original dataset we have selected 8 channels representing the y coordinates of both left and right Shoulder, Elbow, Wrist and Hip. This dataset has 4 different classes representing the kind of exercise done, namely Normal (N), Asymmetrical (A), Reduced Range (R) and Arch (Arch). We know from domain experts that the importance of channels for this dataset is in decending order: Elbows, Wrists, Shoulders, Hips. High accuracy can be obtained only by using the Elbows and Wrists while it is not possible to achieve a high accuracy by only using one channel. We later show experiments both in Section 5 and in the Appendix to document this behaviour.

5 Experiments

In our experiments we aim to understand the strengths and weaknesses of existing methods for explaining multivariate time series classification. As summarised in Table 1, we compared one of the bespoke multivariate method found (dResNet), the popular SHAP, which has the downside of being adapted to provide a 2D heatmap, and Ridge as a sanity check baseline. Some different coupling such as ROCKET paired with dCAM or dResNet paired with SHAP are not possible respectively because dCAM can only explain models having a GAP layer and the timeXplain library (used for ROCKET-SHAP concatenated) is implemented only for 1D-vector instances (univariate time series).

Classifier	Explanation Method	MTS Approach
dResNet	dCAM	Bespoke MTSC
ROCKET	SHAP	Concatenated
ROCKET	SHAP	Channel by Channel
Ridge Classifier	Ridge Classifier	Concatenated
n/a	Random	n/a

Table 1: Summary of the explanation methods tested in this paper.

To make the timeXplain library work with MTS, we apply the following two strategies (Figure 3): (1) **Concatenated**: Concatenating all the channels

to a single univariate time series. As a result, the output saliency map is also univariate and thus needs to be reshaped. (2) **Channel by Channel**: Train and explain one model for each channel independently. The MTSC model in this case is an ensemble of per-channel UTSC models.

For SHAP-channel-by-channel, we assign the number of segments to 10 while, for SHAP-concatenated, the number of segments is set to $d \times 10$. Since Ridge can only work using univariate datasets, we only used the dataset concatenation strategy for this classifier. The output of all explanation methods is a saliency map in the form of either $d \times L$ or $d \times 10$ matrix (reshaped if necessary).

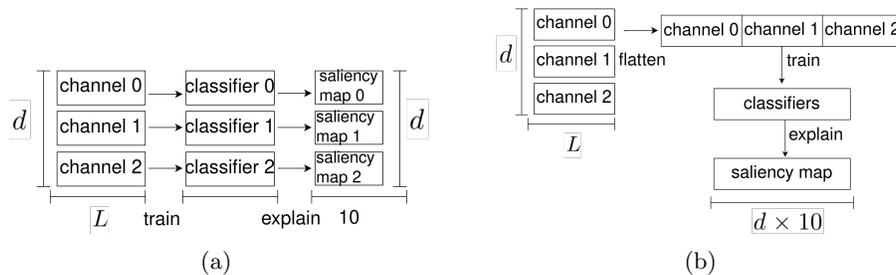


Fig. 3: Strategies to use the timeXplain library in a multivariate scenario, for $d = 3$. In Fig 3a, a classifier is trained for each channel: for explaining each classifier, d heat maps of length 10 are produced: stacking these vectors together results in a matrix of dimension $d \times 10$. In Fig 3b the time series are concatenated and one single classifier is trained. We explain the classifier using a number of segments $d \times 10$ and reshape the resulting vector into a 2D matrix having the same shape as in the previous case.

It is important to note that we have only one bespoke method for multivariate time series, dCAM, that computes a saliency map of the same shape as the original time series instance.

All the experiments were done using a machine with 32GB RAM, Intel i7-12700H CPU and an Nvidia GeForce RTX 350 Ti GPU (the GPU was used only for dResNet/dCAM). All the code used to perform the experiments is available on a Github repository².

5.1 Classification Accuracy Analysis

Before diving into the explanations, we first take a look at the accuracy of the classifiers used for producing the explanations. All the classifiers listed in Table 2 were trained 5 different times (for ROCKET we also tried to either normalize the data or not). In this Table are reported the most accurate ones i.e., the models

² <https://github.com/mlgig/Evaluating-Explanation-Methods-for-MTSC>

used in the experiments as well as the accuracy for the univariate concatenated datasets.

Having a look at the table we can notice that all the times both ROCKET and dResNet have high accuracy (with some exceptions for the synthetic datasets): this is an important pre-requisite when comparing explanations methods applied to different classifiers as we did.

We note that RidgeCV does particularly well on the synthetic datasets. On Military Press, the multivariate models are more accurate than the univariate ones (on concatenated data). This is expected since it is difficult to achieve a high accuracy with a single channel for this dataset, so this is a trully multivariate dataset. Concatenating all the channels for Military Press hurts more dResNet which loses 9 percentage points accuracy, while ROCKET loses only 4. For CMJ, the behaviour is reversed, with univariate models being more accurate than the multivariate ones. dResNet has a noticeable 9 percentage points improvement on the concatenated dataset, while ROCKET gains 1 percentage point.

Classifier/Dataset	PseudoPeriodic	Gaussian	AutoRegressive	CMJ	MilitaryPress
dResNet multivariate	1.0	0.83	0.82	0.82	0.79
dResNet concatenated	1.0	0.89	0.81	0.91	0.68
ROCKET multivariate	1.0	0.93	0.87	0.87	0.87
ROCKET concatenated	1.0	0.72	0.73	0.88	0.83
ROCKET ch-by-ch	0.99	0.72	0.95	0.85	0.65
RidgeCV	1.0	1.0	1.0	0.44	0.61

Table 2: Accuracy for the models listed in Table 1 plus dResNet concatenated and ROCKET multivariate: using this table it is possible to appreciate the differences when using multivariate vs univariate datasets.

5.2 Synthetic Data

For the synthetic data, we performed 5-fold cross-validation to train a Logistic Regression classifier for ROCKET, allowing up to 1000 iterations. For dResNet we used 64 filters, and we trained using the Cross-Entropy Loss and Adam optimizer with a learning rate set to 0.0001. Finally for RidgeCV we used the standard scikit-learn parameters for cross-validation using 5 folds.

Regarding the explanation methods we used 10 segments for ROCKET concatenated in the channel-by-channel scenario and 200 segments in the concatenated one; for dCAM the number of permutations to evaluate k was set to 200 (this is the maximum recommended in [3]).

The steps done for syntethic data evaluation are illustrated in Figure 4. The first step is to reshape all the explanations so that they all have the same dimension. Specifically, the saliency maps we obtained from dCAM and Ridge have a shape of $d \times L = 20 \times 100$ while the ones from SHAP concatenated and channel-by-channel have a shape of $d \times \text{n. segments} = 20 \times 10$. We chose

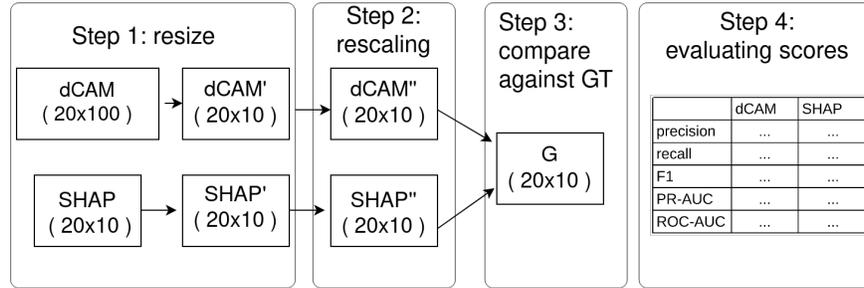


Fig. 4: Steps performed in the synthetic data evaluation when comparing dCAM and SHAP. In Step 1, dCAM is reshaped into $(20, 10)$ averaging 10 consecutive elements in each channel, while SHAP is untouched. In Step 2, the reshaped matrices are rescaled in the range $[0, 100]$. In Step 3, both the explanations achieved so far are compared against the ground truth matrix G and finally in Step 4 the scores computed in the previous step are evaluated.

to average 10 consecutive elements for dCAM and Ridge explanation as we empirically verified that all the metrics had slight improvements. The alternative was to repeat 10 times the same item in SHAP explanations.

After this stage, all explanations have the shape of 20×10 .

The second step rescales the explanation weights as they can have different magnitudes among different instances and different methods. First of all, we take the absolute value of each explanation (to also take into consideration variables that have a negative contribution for the classification) and then we rescale by min-max normalization in the range $[0, 100]$.

The third step is to instantiate a ground truth matrix G and compare each explanation against it. For the settings described before, this is a binary matrix having shape 20×10 (same dimension of explanations after Step 2), all the elements are set to 0 except for the ones in G_i^j with $i = 1, 0 \leq j < 10$ that are set to 1. In other words, this is a binary matrix describing whether or not a segment is important for the classification.

To be noted that synthetic dataset parameters such as the number and range of informative time points and channels, and explanation method parameters such as the number of segments were chosen so that the resulting segments are made up either by only informative time points or by only uninformative time points.

The last step is simply to compute the metrics used for the evaluation i.e., Precision, Recall, F1-score, PR-AUC and ROC-AUC [3]. For Precision and Recall we had to fix a threshold dividing the values considered uninformative from the ones considered informative: we have chosen 50 as the medium value between 0 and 100. On the other hand PR-AUC and ROC-AUC do not fix any threshold as they average multiple scores using different thresholds into one single value. All these metrics computed for the 3 synthetic datasets are reported in Table 3.

Looking at the table it is possible to note that all the time Ridge has perfect metrics but for Recall (and consequently F1 score) with the Gaussian dataset.

Dataset	XAI Method	Precision	Recall	F1	PR-AUC	ROC-AUC	Time
Pseudo-Periodic	SHAP ch-by-ch	0.73	0.94	0.82	0.99	0.99	6.2 h
Pseudo-Periodic	SHAP concatenated	0.92	0.66	0.77	0.99	0.99	3.5 h
Pseudo-Periodic	dCAM	0.50	0.50	0.50	0.63	0.98	50 m
Pseudo-Periodic	Ridge	1.0	1.0	1.0	1.0	1.0	0 s
Gaussian	SHAP ch-by-ch	0.88	0.63	0.73	0.91	0.99	6.2 h
Gaussian	SHAP concatenated	0.34	0.18	0.24	0.16	0.71	3.5 hr
Gaussian	dCAM	0.36	0.15	0.21	0.35	0.94	50 m
Gaussian	Ridge	0.83	1.0	0.9	1.0	1.0	0 s
Auto-Regressive	SHAP ch-by-ch	0.85	0.60	0.71	0.49	0.77	6.2 h
Auto-Regressive	SHAP concatenated	0.27	0.13	0.18	0.29	0.57	3.5 h
Auto-Regressive	dCAM	0.34	0.15	0.21	0.06	0.57	50 m
Auto-Regressive	Ridge	1.0	1.0	1.0	1.0	1.0	0 s
All	Random	0.05	0.15	0.08	0.05	0.5	0 s

Table 3: Scores and runtime of each XAI method for synthetic datasets: h stands for hours, m stands for minutes and s stand for seconds; ch-by-ch stands for channel by channel.

These results along with the one provided in Table 2 (perfect accuracies of Ridge for the 3 synthetic datasets), are very strong evidence that these commonly used benchmarks are not ideal for time series analyses at least using the parameters described before. We think this is the case due to the way the benchmarks are created, by adding or subtracting a single value to consecutive time points. This means that a simple tabular classifier such as Ridge is enough to perfectly classify these datasets. In conclusion, we recommend against the use of these synthetic benchmarks for analysing time series classification or explanation methods.

Comparing the other method, most of the time SHAP channel by channel is the second best model, while comparing dCAM with SHAP concatenated there is no clear winner as in some metrics the first one has better results while in some others is the opposite.

The two last points to be noted are that some methods have metrics close to random, especially for Recall, and the time required for computing the explanations is high, taking into account that these are small datasets: 50 minutes for dCAM, 3.5 hours for SHAP concatenated, and more than 6 hours for SHAP channel by channel.

5.3 Real-world Data

In this section we used some different hyper-parameters: for dResNet the number of filters is now set to 128 as we found better classification results, the number of dCAM permutations k was set to 6 for dCAM (this dataset has 3 channels so the number of possible channel permutations is just 6) while it is still 200, i.e. the maximum recommended, for MP which has 8 channels. We set the number of timeXplain segments using the concatenated dataset to 30 for CMJ and 80 for MP so that they are still equal to $d \times 10$.

Looking at the classifier accuracy in Table 2 we notice how for the two real-world datasets, the accuracies achieved by dResNet and ROCKET are comparable or even better when using the concatenated dataset versions. This means that analysing the explanation methods for MTSC by turning the multivariate problems into univariate ones could be useful.

The close accuracy between original multivariate and concatenated univariate datasets can arise some questions whether these datasets are truly multivariate (i.e., the necessary information for correct classification is spread among different channels). This seems to be the case for Military Press, but less so for CMJ. We plan to investigate further this point in future work.

In this work we decided to use the concatenated datasets and the methodology developed by [18] to evaluate the explanation methods. For the case of dCAM which produces a matrix as an explanation, we flatten the matrix to a vector by concatenating the rows and using it as any other univariate explanation. So dCAM is obtained in a truly multivariate setting (dResNet is a multivariate classifier and dCAM a multivariate explanation), but reshaped to look like a univariate explanation. The explanations obtained from SHAP and Ridge, on the other hand, are univariate explanations obtained by first concatenating the channels and then running univariate classifiers.

For the real-world datasets we do not have precise explanation ground truth as for the synthetic datasets, but we do have domain knowledge about which channels and parts of the time series are important for the classification. Finally in this section we didn't include SHAP channel by channel in the MP dataset experiment as the accuracy is low (Table 2) therefore it does not make sense to derive an explanation.

Evaluation of Explanation Methods. We apply AMEE [18], an explanation evaluation method for the univariate time series classification task, on the CMJ and MP univariate datasets obtained through concatenating all channels. This method aims to measure the faithfulness of an explanation by estimating its impact on a set of independent classifiers (the *referee classifiers*). If an explanation correctly points out the important areas of a univariate time series, perturbation of such areas will lead to a drop in accuracies of the referee classifiers. The faithfulness of the explanation is then calculated using the Area Under the Curve (AUC) of the accuracy drop curves of each of the referee classifiers. AMEE is designed to be robust by employing various perturbation strategies (i.e. how an important area is perturbed and replaced with a new value) and a diverse set of high-performing referee classifiers. The main idea is that masking away important parts of the data as pointed out by the explanation, should affect a set of high performing classifiers leading to a drop in accuracy across the board.

For our task, we use the default perturbation strategies with three classifiers included in the standard referees set: MrSEQL [13] WEASEL 2.0 [21] and ROCKET (for more information and results about this methodology we invite the readers to have a look to the original publication [18]). Table 4 shows the accuracy of these referee classifiers on the evaluated datasets.

Dataset	MrSEQL	ROCKET	WEASEL 2.0
CMJ-concat	0.76	0.88	0.92
MP-concat	0.82	0.84	0.80

Table 4: Accuracy of referee classifiers for the AMEE evaluation of explanation methods for univariate time series classification.

The result of the explanation evaluation is presented in Table 5 as well as the methodology and the evaluation running time. The methodology running time is dependent on the number of both perturbation strategies and employed referees. It is specific to our choice of the three mentioned referees and four perturbation types using Mean and Gaussian sample from both time-point dependent (local) and time-point independent (global) statistics of the test samples. Looking at the second one (time for running the explanation methods) we notice the high SHAP computational complexity: this was the main reason why we used only 2 real-world datasets for the experiments. We focused on human motion data because in this case we can rely on domain expertise.

From the quantitative evaluation with AMEE, we note that for the CMJ dataset, SHAP concat is the best method, although it is close to a random explanation. dCAM ranks third for this dataset. We note that this dataset is quite noisy due to quiet parts after the jump, and this could explain why SHAP and Random are so close in ranking.

For the MP dataset, SHAP concatenated is by far the best method, significantly better than dCAM, as well as Random and Ridge. This is an interesting finding considering that dCAM was proposed to deal with datasets where there are clear dependencies between channels, but for MP this method does not seem to perform so well.

We supplement the quantitative ranking with a more detailed qualitative analysis in the Appendix. In short we find that for CMJ, the importance rankings of channels given by SHAP concat and dCAM are the same, while for MP, SHAP provides a ranking more in line with domain knowledge, while dCAM places the least informative channels at the top of the ranking.

6 Conclusion

In this paper we have investigated explanation methods for MTSC. We studied two very popular explanation methods, dCAM and SHAP, and have provided a quantitative and qualitative analysis of their behavior on synthetic as well as real-world datasets. We found that adaptations of SHAP for MTSC work quite well, and they outperform the recent bespoke MTSC explanation method dCAM. We have also pointed out that a very popular synthetic MTSC benchmark does not seem suitable for MTSC evaluation, since a simple Ridge classifier outperforms all other methods both in classification accuracy and in explanation quality. Finally, while SHAP seems to work effectively to point out important time series channels and time points, we highlighted the time required to run SHAP and

Dataset	XAI method	Explanation Power	Rank	Evaluation Time	Explanation Time
CMJ-concat	SHAP concat	1.00	1	2h	7.15h
	Random	0.99	2	2h	0s
	dCAM	0.39	3	2h	30s
	SHAP ch-by-ch	0.05	4	2h	7.5h
	Ridge	0.0	5	2h	0s
MP-concat	SHAP concat	1.00	1	4.8h	24h
	dCAM	0.33	2	4.8h	15m
	Random	0.07	3	4.8h	0s
	Ridge	0.0	4	4.8h	0s

Table 5: Results of AMEE to rank XAI methods on CMJ and MP datasets concatenated.

pointed out the open problem of excessive time requirements for this method. In future work we plan to investigate the computation time for SHAP, as well as other frameworks for evaluating bespoke explanation methods for MTSC.

Acknowledgments

This publication has emanated from research supported in part by a grant from Science Foundation Ireland under Grant number 18/CRT/6183. For the purpose of Open Access, the author has applied a CC BY public copyright licence to any Author Accepted Manuscript version arising from this submission.

References

1. Babayev, R., Wiese, L.: Interpreting decision-making process for multivariate time series classification. In: *New Trends in Database and Information Systems: AD-BIS 2021 Short Papers, Doctoral Consortium and Workshops: DOING, SIMPDA, MADEISD, MegaData, CAoNS, Tartu, Estonia, August 24-26, 2021, Proceedings*. pp. 146–152. Springer (2021)
2. Bagnall, A., Dau, H.A., Lines, J., Flynn, M., Large, J., Bostrom, A., Southam, P., Keogh, E.: The uea multivariate time series classification archive, 2018. arXiv preprint arXiv:1811.00075 (2018)
3. Boniol, P., Meftah, M., Remy, E., Palpanas, T.: dcam: Dimension-wise class activation map for explaining multivariate data series classification. In: *Proceedings of the 2022 International Conference on Management of Data*. pp. 1175–1189 (2022)
4. Cao, Z., Hidalgo Martinez, G., Simon, T., Wei, S., Sheikh, Y.A.: Openpose: Real-time multi-person 2d pose estimation using part affinity fields. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2019)
5. Crabbé, J., Van Der Schaar, M.: Explaining Time Series Predictions with Dynamic Masks. In: Meila, M., Zhang, T. (eds.) *Proceedings of the 38th International Conference on Machine Learning*. *Proceedings of Machine Learning Research*, vol. 139, pp. 2166–2177. PMLR (1 2021), <https://proceedings.mlr.press/v139/crabbe21a.html>

6. Dempster, A., Petitjean, F., Webb, G.I.: Rocket: exceptionally fast and accurate time series classification using random convolutional kernels. *Data Mining and Knowledge Discovery* **34**(5), 1454–1495 (2020)
7. Dempster, A., Schmidt, D.F., Webb, G.I.: Minirocket: A very fast (almost) deterministic transform for time series classification. In: *Proceedings of the 27th ACM SIGKDD conference on knowledge discovery & data mining*. pp. 248–257 (2021)
8. Guidotti, R.: Evaluating local explanation methods on ground truth. *Artificial Intelligence* **291**, 103428 (2021)
9. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. arXiv 2015. arXiv preprint arXiv:1512.03385 **14** (2015)
10. Hoerl, A.E., Kennard, R.W.: Ridge regression: Biased estimation for nonorthogonal problems. *Technometrics* **12**(1), 55–67 (1970)
11. Ismail, A.A., Gunady, M., Corrada Bravo, H., Feizi, S.: Benchmarking Deep Learning Interpretability in Time Series Predictions. In: Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M.F., Lin, H. (eds.) *Advances in Neural Information Processing Systems*. vol. 33, pp. 6441–6452. Curran Associates, Inc. (2020), <https://proceedings.neurips.cc/paper/2020/file/47a3893cc405396a5c30d91320572d6d-Paper.pdf>
12. Ismail, A.A., Gunady, M., Corrada Bravo, H., Feizi, S.: Benchmarking deep learning interpretability in time series predictions. *Advances in neural information processing systems* **33**, 6441–6452 (2020)
13. Le Nguyen, T., Gsponer, S., Ilie, I., O’reilly, M., Ifrim, G.: Interpretable time series classification using linear models and multi-resolution multi-domain symbolic representations. *Data mining and knowledge discovery* **33**, 1183–1222 (2019)
14. Li, J., Zhang, C., Zhou, J.T., Fu, H., Xia, S., Hu, Q.: Deep-lift: deep label-specific feature learning for image annotation. *IEEE Transactions on Cybernetics* **52**(8), 7732–7741 (2021)
15. Lundberg, S.M., Lee, S.I.: A Unified Approach to Interpreting Model Predictions. In: Guyon, I., Luxburg, U.V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., Garnett, R. (eds.) *Advances in Neural Information Processing Systems* 30, pp. 4765–4774. Curran Associates, Inc. (2017), <http://papers.nips.cc/paper/7062-a-unified-approach-to-interpreting-model-predictions.pdf>
16. Lundberg, S.M., Lee, S.I.: A unified approach to interpreting model predictions. In: Guyon, I., Luxburg, U.V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., Garnett, R. (eds.) *Advances in Neural Information Processing Systems* 30, pp. 4765–4774. Curran Associates, Inc. (2017), <http://papers.nips.cc/paper/7062-a-unified-approach-to-interpreting-model-predictions.pdf>
17. Mujkanovic, F., Doskoč, V., Schirneck, M., Schäfer, P., Friedrich, T.: timexplain—a framework for explaining the predictions of time series classifiers. arXiv preprint arXiv:2007.07606 (2020)
18. Nguyen, T.T., Nguyen, T.L., Ifrim, G.: Amee: A robust framework for explanation evaluation in time series classification. arXiv preprint arXiv:2306.05501 (2023)
19. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., Duchesnay, E.: Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research* **12**, 2825–2830 (2011)
20. Ribeiro, M.T., Singh, S., Guestrin, C.: ” why should i trust you?” explaining the predictions of any classifier. In: *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*. pp. 1135–1144 (2016)

21. Schäfer, P., Leser, U.: WEASEL 2.0—A Random Dilated Dictionary Transform for Fast, Accurate and Memory Constrained Time Series Classification. arXiv preprint arXiv:2301.10194 (2023)
22. Singh, A., Bevilacqua, A., Nguyen, T.L., Hu, F., McGuinness, K., O’Reilly, M., Whelan, D., Caulfield, B., Ifrim, G.: Fast and robust video-based exercise classification via body pose tracking and scalable multivariate time series classifiers. *CoRR* **abs/2210.00507** (2022). <https://doi.org/10.48550/arXiv.2210.00507>, <https://doi.org/10.48550/arXiv.2210.00507>
23. Tan, C.W., Dempster, A., Bergmeir, C., Webb, G.I.: Multirocket: multiple pooling operators and transformations for fast and effective time series classification. *Data Mining and Knowledge Discovery* **36**(5), 1623–1646 (2022)
24. Turbé, H., Bjelogrić, M., Lovis, C., Mengaldo, G.: Evaluation of post-hoc interpretability methods in time-series classification. *Nature Machine Intelligence* **5**(3), 250–260 (Mar 2023). <https://doi.org/10.1038/s42256-023-00620-w>, <https://doi.org/10.1038/s42256-023-00620-w>
25. Villar, J., VERGARA, P., Menéndez González, M., Marín, E., González, V., Sedano, J.: Generalized models for the classification of abnormal movements in daily life and its applicability to epilepsy convulsion recognition. *International Journal of Neural Systems* **26** (04 2016). <https://doi.org/10.1142/S0129065716500374>
26. Wang, Z., Yan, W., Oates, T.: Time series classification from scratch with deep neural networks: A strong baseline. In: 2017 International joint conference on neural networks (IJCNN). pp. 1578–1585. IEEE (2017)
27. Zhou, B., Khosla, A., Lapedriza, A., Oliva, A., Torralba, A.: Learning deep features for discriminative localization. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2921–2929 (2016)

A Supplementary Material

channel	importance
LElbow	0.60
RElbow	0.59
RWrist	0.58
LWrist	0.57
LShoulder	0.52
RShoulder	0.49
LHip	0.39
RHip	0.36

(a) Rocket ranking

channel	importance
LHip	1.0
RHip	0.99
RWrist	0.74
LWrist	0.73
RElbow	0.54
LElbow	0.53
RShoulder	0.50
LShoulder	0.50

(b) dCAM ranking

channel	importance
RWrist	1.0
LWrist	0.98
LElbow	0.93
RElbow	0.86
LShoulder	0.82
RShoulder	0.80
RHip	0.77
LHip	0.76

(c) SHAP ranking

Table 6: Ranking of MP columns using different methods. In Table 6a ordered accuracy of single-channel classifiers in ROCKET channel-by-channel scenario: we take this as a channel importance baseline. In Table 6b and 6c respectively dCAM and SHAP ranking achieved by averaging all time points in the single channels among all the time series (for readability the values were rescaled such that the most important channel has a value 1.0).

Table 6c is closer to 6a than Table 6b. The major difference is in the RHip and LHip ranking: while SHAP places them as the least important ones, agreeing with Rocket, dCAM ranks them as the most important channels.

method	LElbow	RElbow	RWrist	LWrist	LShoulder	RShoulder	RHip	LHip
ROCKET	1	2	3	4	5	6	7	8
dCAM	6	5	3	4	8	7	2	1
SHAP	3	4	1	2	5	6	7	8

Table 7: MP channels importance: same analysis as in Table 6 but showing the ranking rather than the raw values.

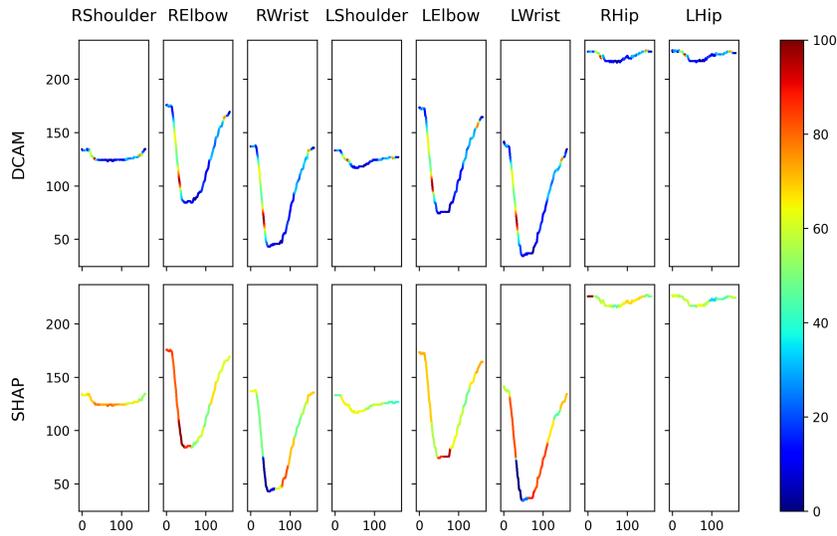


Fig. 5: Different saliency maps for the same MP instance correctly classified as Normal: SHAP focuses on LWrist, RElbow, RWrist and LElbow. dCAM has a similar behavior in terms of channel importance but it focuses on smaller sections: this is just partially explainable due to division in segments by SHAP since sometimes there is no overlapping between highlighted regions in the same channel.

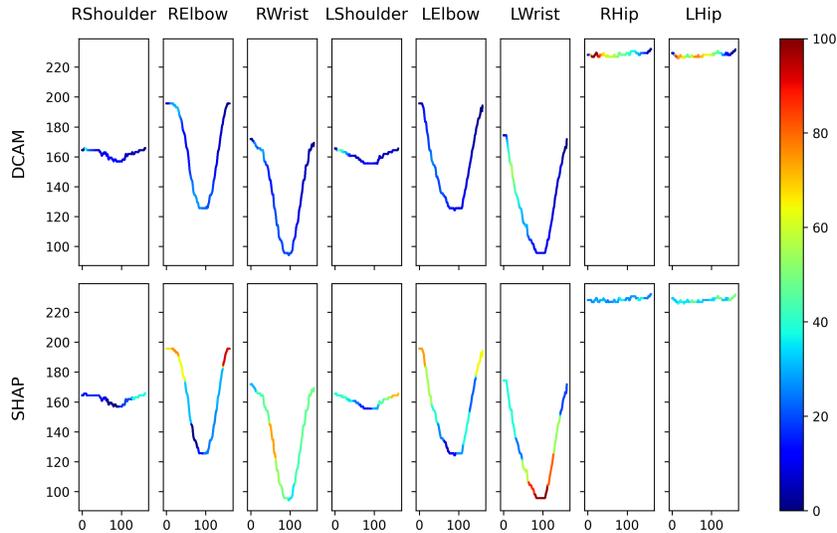


Fig. 6: Different saliency maps for the same MP instance correctly classified as Asymmetric. dCAM focuses on Hips channels while SHAP focuses on Wrists and Elbows in accordance to Table 7.

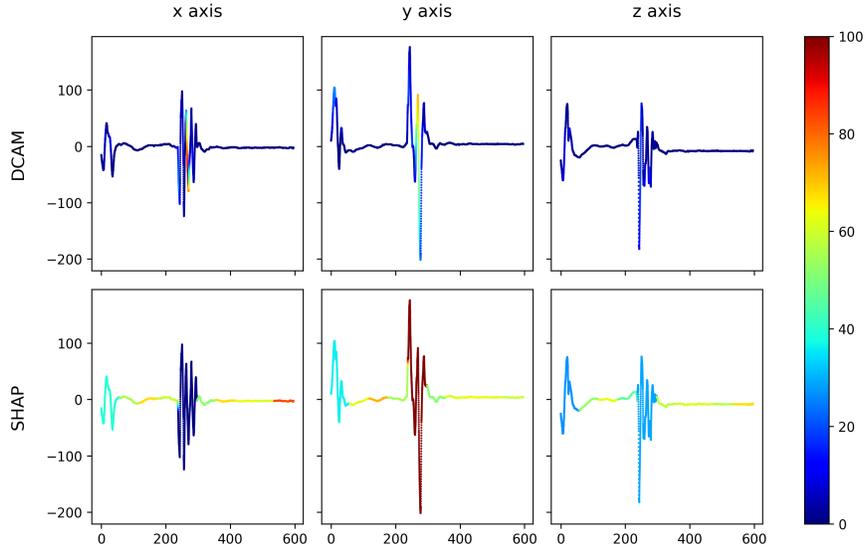


Fig. 7: Different saliency maps for the same CMJ instance correctly classified as Acceptable form: SHAP focuses more on the y axis according to what is shown in Table 8 while dCAM highlights more the x axis.

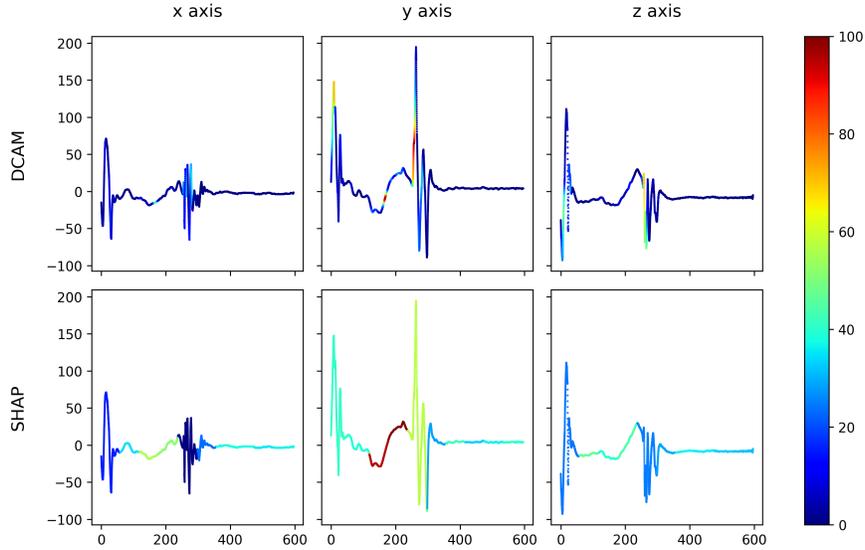


Fig. 8: Different saliency maps for the same CMJ instance correctly classified as Legs bending: both explanation methods focus on y axis, but different parts. SHAP focuses on the small peak, around time step 200 while dCAM focuses more on the beginning of the following higher peak.

Method	y	z	x
ROCKET	0.85	0.81	0.79
dCAM	1.0	0.92	0.89
SHAP	1.0	0.85	0.83

Table 8: Ranking of CMJ columns using different methods. In this case all the methods agree.