

---

# Generalised Interpretable Shapelets for Irregular Time Series

---

Patrick Kidger\*

James Morrill\*

Terry Lyons

Mathematical Institute, University of Oxford  
 The Alan Turing Institute, British Library  
 {kidger, morrill, tlyons}@maths.ox.ac.uk

## Abstract

The shapelet transform is a form of feature extraction for time series, in which a time series is described by its similarity to each of a collection of ‘shapelets’. However it has previously suffered from a number of limitations, such as being limited to regularly-spaced fully-observed time series, and having to choose between efficient training and interpretability. Here, we extend the method to continuous time, and in doing so handle the general case of irregularly-sampled partially-observed multivariate time series. Furthermore, we show that a simple regularisation penalty may be used to train efficiently without sacrificing interpretability. The continuous-time formulation additionally allows for learning the length of each shapelet (previously a discrete object) in a differentiable manner. Finally, we demonstrate that the measure of similarity between time series may be generalised to a learnt pseudometric. We validate our method by demonstrating its performance and interpretability on several datasets; for example we discover (purely from data) that the digits 5 and 6 may be distinguished by the chirality of their bottom loop, and that a kind of spectral gap exists in spoken audio classification.

## 1 Introduction

Shapelets are a form of feature extraction for time series [1, 2, 3, 4]. Given some fixed hyperparameter  $K$ , describing how many shapelets we are willing to consider, then each time series is represented by a vector of length  $K$  describing how similar it is to each of the  $k$  selected shapelets.

We begin by recalling the classical definition of the shapelet transform [5].

### 1.1 Classical shapelet transform

Given  $N$  regularly sampled multivariate time series, with  $D$  observed channels, where the  $n$ -th time series is of length  $T_n$ , then the  $n$ -th time series is a matrix

$$f^n = (f_t^n)_{t \in \{0, \dots, T_n - 1\}} = (f_{t,d}^n)_{t \in \{0, \dots, T_n - 1\}, d \in \{1, \dots, D\}}, \quad (1)$$

with each  $f_{t,d}^n \in \mathbb{R}$  and  $n \in \{1, \dots, N\}$ .

Fix some hyperparameter  $K \in \mathbb{N}$ , which will describe the number of shapelets. Fix some  $S \in \{0, \dots, \min_{i \in \{1, \dots, N\}} T_i - 1\}$ , which will describe the length of each shapelet. Then the  $k$ -th shapelet is a matrix

$$w^k = (w_t^k)_{t \in \{0, \dots, S - 1\}} = (w_{t,d}^k)_{t \in \{0, \dots, S - 1\}, d \in \{1, \dots, D\}},$$

---

\*Equal contribution.

with each  $w_{t,d}^k \in \mathbb{R}$ .

Then the discrepancy between  $f^n$  and  $w^k$  is given by (sometimes without the square):

$$\sigma_S(f^n, w^k) = \min_{s \in \{0, \dots, T_n - S\}} \sum_{t=0}^{S-1} \|f_{s+t}^n - w_t^k\|_2^2, \quad (2)$$

where  $\|\cdot\|_2$  describes the  $L^2$  norm on  $\mathbb{R}^D$ . A small discrepancy implies that  $f^n$  and  $w^k$  are similar to one another. This corresponds to sweeping  $w^k$  over  $f^n$ , and finding the offset  $s$  at which  $w^k$  best matches  $f^n$ .

In this article we will refer to the map

$$f \mapsto (\sigma_S(f, w^1), \dots, \sigma_S(f, w^K)) \quad (3)$$

as the *classical shapelet transform*. The result is now a feature describing  $f$ , which may now be passed to some model to perform classification or regression.

## 1.2 Limitations

The classical shapelet method suffers from a number of limitations.

1. The technique only applies to regularly spaced time series.
2. The choice of shapelet length  $S$  is discrete and a hyperparameter. As such optimising it involves a relatively expensive hyperparameter search.
3. Learning the shapelets  $w^k$  by searching is expensive [1], whilst optimising differentially [2] typically sacrifices interpretability [6].

Besides this, the choice of  $L^2$  norm is ad-hoc and a general formulation should allow for other notions of similarity. It is these limitations that we seek to address here.

## 1.3 Contributions

We extend the method to continuous time rather than discrete time. This allows for the treatment of irregularly-sampled partially-observed multivariate time series on the same footing as regular time series. Additionally, this continuous-time formulation means that the length of each shapelet (previously a discrete value) takes its values in a continuous range, and may now be trained differentially.

Next, we demonstrate how simple regularisation is enough to achieve shapelets that resemble characteristic features of the data. This gives interpretability with respect to a classification result, and also offers pattern discovery for determining previously unknown information about the data. For example we discover – purely from data – that the digits 5 and 6 may be distinguished by the chirality of their bottom loop, and that a kind of spectral gap exists in spoken audio classification.

Finally, we generalise the discrepancy between a shapelet and a time series to be a learnt pseudometric. This is particularly useful for interpretability and pattern discovery, as doing so learns the importance of different channels.

Our code is available at [https://github.com/patrick-kidger/generalised\\_shapelets](https://github.com/patrick-kidger/generalised_shapelets).

## 2 Prior work

Shapelets may be selected as small intervals extracted from training samples [1]. However doing so is very expensive, requiring  $\mathcal{O}(N^2 \cdot \max_n T_n^4)$  work. Much work on shapelets has sought speedup techniques [7, 8, 9], for example via random algorithms [10, 11].

However [2] observe that the discrepancy  $\sigma_S$  of equation (2) is differentiable with respect to  $w^k$ , so that shapelets may be differentially optimised jointly with the subsequent model, as part of an end-to-end optimisation of the final loss function. (Although [2] include a ‘softmin’ procedure which we believe to be unnecessary, as the minimum function is already almost everywhere differentiable.) This costs only  $\mathcal{O}(N \cdot \max_n T_n^2)$  to train, and is the approach that we extend here.

This method is attractive for its speed and its ease of trainability via modern deep learning frameworks [12, 13, 14]. However, [6] observe that the predictive power of the distance between a shapelet and a time series need not correlate with a similarity between the two, so there is no pressure towards interpretability. [6] propose to solve this via adversarial regularisation; we will present a simpler alternative later. Without such procedures, then efficient training and interpretability become mutually exclusive.

The method may additionally be generalised by considering alternative notions of similarity between a shapelet and a time series; for example [15] replace the  $L^2$  norm with dynamic time warping.

The shapelet method is attractive for its normalisation of variable-length time series, and demonstration of typically good performance [4, 16]. Arguably its most important advantage is interpretability, as use of a particular feature corresponds to the importance of the similarity to the shapelet  $w^k$ . This may describe some shape that is characteristic of a particular class, and can discover previously unknown patterns in the data.

### 3 Method

#### 3.1 Continuous-time objects

We interpret a time series as a discretised sample from an underlying process, observed only through the time series. Similarly, a shapelet constructed as in Section 1.1 may be thought of as a discretisation of some underlying function. The first important step in our procedure is to construct continuous-time approximations to these underlying objects.

**Continuous-time path interpolants** Formally speaking, we assume that for  $n \in \{1, \dots, N\}$  indexing different time series, each of length  $T_n$ , we observe a collection of time series

$$f^n = (f_{t_\tau}^n)_{\tau \in \{1, \dots, T_n\}},$$

where  $t_\tau \in \mathbb{R}$  is the observation time of  $f_{t_\tau}^n \in (\mathbb{R} \cup \{*\})^D$ , where  $*$  denotes the possibility of a missing observation.

Next, interpolate to get a function  $\iota(f^n): [0, T_n - 1] \rightarrow \mathbb{R}^D$  such that  $\iota(f^n)(t_\tau) = f_{t_\tau}^n$  for all  $\tau \in \{0, \dots, T_n - 1\}$  such that  $f_{t_\tau}^n$  is observed. There are many possible choices for  $\iota$ , such as splines, kernel methods [17], or Gaussian processes [18, 19]. In our experiments, we use piecewise linear interpolation.

**Continuous-time shapelets** The shapelets themselves we are free to control, and so for  $k \in \{1, \dots, K\}$  indexing different shapelets, we take each  $w^{k,\rho}: [0, 1] \rightarrow \mathbb{R}^D$  to be some learnt function depending on learnt parameters  $\rho$ . For example, this could be an interpolated sequence of learnt points, an expansion in some basis functions, or a neural network. In our experiments we use linear interpolation of a sequence of learnt points.

Then for some learnt length  $S_k > 0$ , define  $w^{k,\rho,S_k}: [0, S_k] \rightarrow \mathbb{R}^D$  by

$$w^{k,\rho,S_k}(t) = w^{k,\rho} \left( \frac{t}{S_k} \right).$$

Taking the length  $S_k$  to be continuous is a necessary prerequisite to training it differentially. We will discuss the training procedure in a moment.

#### 3.2 Generalised discrepancy

The core of the shapelet method is that the similarity or discrepancy between  $f^n$  and  $w^{k,\rho,S_k}$  is important. In general, we approach this by defining a *discrepancy function* between the two, which will typically be learnt, and which we require only to be a pseudometric.

We denote this discrepancy function by  $\pi_S^A$ . It depends upon a length  $S$  and a learnt parameter  $A$ , consumes two paths  $[0, S] \rightarrow \mathbb{R}^D$ , and returns a real number describing some notion of closeness between them. We are being deliberately vague about the regularity of the domain of  $\pi_{S_k}^A$ , as it is a function space whose regularity will depend on  $\iota$ .

Given some  $\pi_S^A$ , then the discrepancy between  $f^n$  and  $w^{k,\rho,S_k}$  is defined as

$$\sigma_{S_k}^A(f^n, w^{k,\rho,S_k}) = \min_{s \in [0, T_n - S_k]} \pi_{S_k}^A(\iota(f^n)|_{[s, s+S_k]}(s + \cdot), w^{k,\rho,S_k}). \quad (4)$$

The collection of discrepancies  $(\sigma_{S_k}^A(f^n, w^{1,\rho,S_k}), \dots, \sigma_{S_k}^A(f^n, w^{K,\rho,S_k}))$  is now a feature describing  $f^n$ , and is invariant to the length  $T_n$ . Use of the particular feature  $\sigma_{S_k}^A(f^n, w^{k,\rho,S_k})$  corresponds to the importance of the similarity between  $f^n$  and  $w^{k,\rho,S_k}$ . In this way, the choice of  $\pi_{S_k}^A$  gives a great deal of flexibility.

**Existing shapelets fit into this framework** A simple example, in analogy to the classical shapelet method of equation (2), is to take

$$\pi_{S_k}^A(f, w) = \left( \int_0^{S_k} \|f(t) - w(t)\|_2^2 dt \right)^{\frac{1}{2}},$$

which in fact has no  $A$  dependence. If  $\iota$  is taken to be a piecewise constant ‘interpolation’ then this will exactly correspond to (the square root of) the classical shapelet approach.

**Learnt  $L^2$  discrepancies** The previous example may be generalised by taking our learnt parameter  $A \in \mathbb{R}^{D \times D}$ , and then letting

$$\pi_S^A(f, w) = \left( \int_0^S \|A(f(t) - w(t))\|_2^2 dt \right)^{\frac{1}{2}}. \quad (5)$$

That is, allowing some learnt linear transformation before measuring the discrepancy. In this way, particularly informative dimensions may be emphasised. In our experiments we take  $A$  to be diagonal. Allowing a general matrix was found during initial experiments to produce slightly worse performance.

**More complicated discrepancies** Moving on, we consider other more general choices of discrepancy, which may be motivated by the problem at hand. In particular we will discuss discrepancies based on the logsignature transform [20], and mel-frequency cepstrums (MFC) [21].

Our exposition on these two discrepancies will be deliberately brief, as the finer details on exactly when and how to use them is domain-specific. The point is that our framework has the flexibility to consider general discrepancies motivated by other disciplines, or which are known to extract information which is particular useful to the domain in question. An understanding of either logsignatures or mel-frequency cepstrums will not be necessary to follow the paper.

**Logsignature discrepancies** The logsignature transform is a transform on paths, known to characterise its input whilst extracting statistics which describe how the path controls differential equations [20, 22, 23, 24, 25]. Let  $\mu$  denote the Möbius function, and let

$$\beta_{D,R} = \sum_{r=1}^R \frac{1}{r} \sum_{\rho|r} \mu\left(\frac{r}{\rho}\right) D^\rho,$$

which is Witt’s formula [26]. Let

$$\text{LogSig}^R: \{f: [0, T] \rightarrow \mathbb{R}^D \mid T \in \mathbb{R}, f \text{ is of bounded variation}\} \rightarrow \mathbb{R}^{\beta_{D,R}}$$

be the depth- $R$  logsignature transform. Let  $A \in \mathbb{R}^{\beta_{D,R} \times \beta_{D,R}}$  be full or diagonal as before, and let  $\|\cdot\|_p$  be the  $L^p$  norm on  $\mathbb{R}^{\beta_{D,R}}$ . Then we define the  $p$ -logsignature discrepancy between two functions to be

$$\pi_S^A(f, w) = \|A(\text{LogSig}^R(f) - \text{LogSig}^R(w))\|_p. \quad (6)$$

**MFC discrepancies** The computation of an MFC is a function-to-function map derived from the short-time Fourier transform, with additional processing to focus on frequencies that are particularly relevant to human hearing [21]. Composing this with the  $L^2$  based discrepancy of equation (5) produces

$$\pi_S^A(f, w) = \left( \int_0^S \|A(\text{MFC}(f)(t) - \text{MFC}(w)(t))\|_2^2 dt \right)^{\frac{1}{2}}. \quad (7)$$

**The generalised shapelet transform** Whatever the choice of  $\pi_S^A$ , and in analogy to the classical shapelet transform [5], we call the map

$$f \mapsto (\sigma_{S_1}^A(f, w^{1,\rho,S_1}), \dots, \sigma_{S_K}^A(f, w^{K,\rho,S_K})) \quad (8)$$

the *generalised shapelet transform*.

### 3.3 Interpretable regularisation

As previously described, learning shapelets differentiably can sacrifice interpretability [6], as the learnt shapelets need not resemble the training data. We propose a novel regularisation penalty to solve this: simply add on

$$\sum_{k=1}^K \min_{n \in \{1, \dots, N\}} \sigma_S^A(f^n, w^{k, \rho, s}) \quad (9)$$

as a regularisation term, so that minimising the discrepancy between  $f^n$  and  $w^{k, \rho, S}$  is also important. Taking a minimum over  $n$  asks that every shapelet should be similar to a single training sample, as in the original approach of finding shapelets by searching through the dataset instead of training differentiably.

### 3.4 Minimisation objective and training procedure

Overall, suppose we have some differentiable model  $F^\theta$  parameterised by  $\theta$ , some loss function  $\mathcal{L}$ , and some observed time series  $f^1, \dots, f^N$  with targets  $y_1, \dots, y_N$ .

Then letting  $\gamma > 0$  control the amount of regularisation, we propose to seek a local minimum of

$$\frac{1}{N} \sum_{n=1}^N \mathcal{L}(y_n, F^\theta(\sigma_{S_1}^A(f^n, w^{1, \rho, S_1}), \dots, \sigma_{S_K}^A(f^n, w^{K, \rho, S_K}))) + \gamma \sum_{k=1}^K \min_{n \in \{1, \dots, N\}} \sigma_{S_k}^A(f^n, w^{k, \rho, S_k}) \quad (10)$$

over model parameters  $\theta$ , discrepancy parameters  $A$ , shapelet parameters  $\rho$ , and shapelet lengths  $S_k$ , via standard stochastic gradient descent based techniques.

**Differentiability** Some thought is necessary to verify that this construction is differentiable with respect to  $S_k$ . There are two operations that may seem to pose a problem, namely the minimum over a range  $\min_{s \in [0, T_n - S_k]}$ , and the restriction operator  $\iota(f^n) \mapsto \iota(f^n)|_{[s, s + S_k]}$ .

Practically speaking, however, it is straightforward to resolve both of these issues. For the minimum over a range, this may reasonably be approximated by a minimum over some collection of points  $s \in \{0, \varepsilon, 2\varepsilon, \dots, T_n - S_k - \varepsilon, T_n - S_k\}$ , for some  $\varepsilon > 0$  small and dividing  $T_n - S_k$ . This is now a standard piece of an autodifferentiation package. The error of this approximation may be controlled by the modulus of continuity of  $s \mapsto \pi_{S_k}^A(\iota(f^n)|_{[s, s + S_k]}(s + \cdot), w^{k, \rho, S_k})$ , but in practice we found this to be unnecessary, and simply took  $\varepsilon$  equal to the smallest gap between observations.

Next, the continuous-time paths  $\iota(f^n)$  and continuous-time shapelets  $w^{k, \rho, S_k}$  must both be represented by some parameterisation of function space, and it is thus sufficient to restrict to considering differentiability with respect to this parameterisation.

In our experiments we represent both  $\iota(f^n)$  and  $w^{k, \rho, S_k}$  as a continuous piecewise linear function stored as a collection of knots. In this context, the restriction operator is clearly differentiable, as a map from one collection of knots to a restricted collection of knots. Each knot is either kept (the identity function), thrown away (the zero function), or interpolated between to place a new knot at the boundary (a ratio of existing knots).

**Choice of  $F^\theta$**  Interpretability of the model will depend on an interpretable choice of  $F^\theta$ . In our experiments we thus used a linear model on the logarithm of every feature, so that a very negative coefficient corresponds to the importance of  $f^n$  and  $w^{k, \rho, S_k}$  being similar to each other.

## 4 Experiments

We compare the generalised shapelet transform to the classical shapelet transform, in terms of both performance and interpretability, on a large range of time series classification problems. The shapelets of the classical shapelet transform are learnt differentiably.

In every case the model is a linear map, for interpretability as previously described, on either the generalised (equation (8)) or classical (equation (3)) shapelet transforms. The learnt pseudometrics for the generalised shapelet transform scale each channel individually by taking  $A$  to be diagonal.

Precise experimental details (optimiser, training scheme, ...) may be found in Appendix A.

Table 1: Test accuracy (mean  $\pm$  std, computed over three runs) on UEA. A ‘win’ is the number of times each algorithm was within 1 standard deviation of the top performer for each dataset.

Dataset	Discrepancy		
	$L^2$	Logsignature	Classical
BasicMotions	90.8% $\pm$ 1.4%	80.8% $\pm$ 3.8%	<b>96.7% <math>\pm</math> 5.8%</b>
ERing	<b>82.6% <math>\pm</math> 6.3%</b>	43.3% $\pm$ 2.9%	67.2% $\pm$ 11.8%
Epilepsy	<b>88.4% <math>\pm</math> 3.0%</b>	<b>88.6% <math>\pm</math> 0.8%</b>	72.9% $\pm$ 5.4%
Handwriting	10.3% $\pm$ 2.6%	<b>11.8% <math>\pm</math> 1.2%</b>	6.5% $\pm$ 3.7%
JapaneseVowels	<b>97.2% <math>\pm</math> 1.1%</b>	53.9% $\pm$ 3.0%	91.5% $\pm$ 4.1%
Libras	<b>67.0% <math>\pm</math> 9.4%</b>	<b>67.8% <math>\pm</math> 5.5%</b>	62.2% $\pm$ 2.4%
LSST	<b>36.1% <math>\pm</math> 0.2%</b>	35.7% $\pm$ 0.4%	33.5% $\pm$ 0.5%
PenDigits	<b>97.3% <math>\pm</math> 0.1%</b>	96.7% $\pm$ 0.7%	<b>97.5% <math>\pm</math> 0.6%</b>
RacketSports	<b>79.6% <math>\pm</math> 0.7%</b>	61.2% $\pm$ 9.2%	<b>79.6% <math>\pm</math> 2.4%</b>
Wins	7	3	3

#### 4.1 The UEA Time Series Archive

This is a collection of 30 fully-observed regularly-sampled datasets with varying properties [27], see Appendix A. Evaluating on the full collection of datasets would take a prohibitively long time, and so we select 9 representing a range of difficulties.

We begin by performing hyperparameter optimisation for the classical shapelet transform, on each dataset. We then use the same hyperparameters for the generalised shapelet transform. For the generalised shapelet transform, the length hyperparameter is used to determine the initial length of the shapelet, but this may of course vary as it is learnt.

For the generalised shapelet transform, we consider two different discrepancy functions, specifically the learnt  $L^2$  and  $p$ -logsignature discrepancies of equations (5) and (6). For the latter, we take  $p = 2$  and the depth  $R = 3$ . We did not try to optimise  $p$  and  $R$ , as we use the logsignature discrepancy simply to highlight the possibility of using more unusual discrepancies if desired.

**Classification performance** The results are given in Table 1. We see that the generalised shapelet transform with  $L^2$  discrepancy function achieves within one standard deviation of the top performing algorithm on 7 of the 9 datasets, whilst the classical approach does so for only 3.

**Interpretability on PenDigits** We demonstrate interpretability by examining the PenDigits dataset. This is a dataset of handwritten digits 0–9, sampled at 8 points along their trajectory. We select the most informative shapelet for each of the ten classes (as in Section 3.4), for both the classical shapelet transform and the generalised shapelet transform, with  $L^2$  discrepancy. We then locate the training sample that it is most similar to, and plot an overlay of the two. See Figure 1.

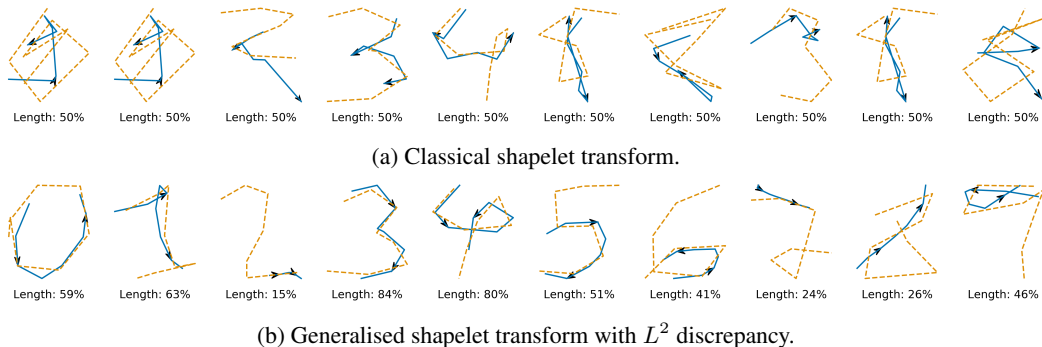


Figure 1: The most significant shapelet for each class (blue, solid), overlaid with the most similar training example (orange, dashed). Similarity is measured with respect to the (learnt) discrepancy function.

Table 2: Test accuracy (mean  $\pm$  std, computed over three runs) on three UEA datasets with missing data. A ‘win’ is defined as the number of times each algorithm was within 1 standard deviation of the top performer for each dataset.

Dataset	Dropped data	Lengths selected by	
		Differentiable optimisation	Hyperparameter searching
JapaneseVowels	10%	<b>93.2% <math>\pm</math> 2.1%</b>	<b>93.1% <math>\pm</math> 0.9%</b>
	30%	<b>91.2% <math>\pm</math> 4.1%</b>	<b>91.4% <math>\pm</math> 2.8%</b>
	50%	<b>93.5% <math>\pm</math> 1.1%</b>	92.0% $\pm$ 1.1%
Libras	10%	57.4% $\pm$ 4.2%	<b>59.3% <math>\pm</math> 1.6%</b>
	30%	<b>81.2% <math>\pm</math> 7.6%</b>	63.9% $\pm$ 8.2%
	50%	<b>62.5% <math>\pm</math> 14.8%</b>	<b>65.3% <math>\pm</math> 8.6%</b>
LSST	10%	40.2% $\pm$ 3.5%	<b>44.0% <math>\pm</math> 1.0%</b>
	30%	<b>38.1% <math>\pm</math> 0.3%</b>	<b>40.2% <math>\pm</math> 5.6%</b>
	50%	41.5% $\pm$ 2.7%	<b>44.4% <math>\pm</math> 0.5%</b>
Wins		6	7

We can clearly see multiple issues with the shapelets learnt with the classical approach. The most significant shapelet for the classes 0 and 1 is the same shapelet, and for classes 1, 5, 6, 7, 9, the most significant shapelet is not even closest to a member of that class. Visually, the shapelets for 3 and 4 seem to have identified distinguishing features of those classes, but the shapelets corresponding to the other classes appear to be little more than random noise.

In contrast, the results of the generalised shapelet approach are abundantly clear. Every class has a unique most significant shapelet, and every such shapelet is close to a member of the correct class. In the case of class 3, the shapelet has essentially reproduced the entire digit.

A point of interest is the difference between the shapelets for the digits 5 and 6, for the generalised shapelet transform. Whilst visually very similar, the difference between them is their direction. Whilst a 5 and a 6 may appear visually similar on the page (with a loop in the bottom half of the digit), they may clearly be distinguished by the direction in which they tend to be written. This is a nice example of discovering something about the data that was not necessarily already known!

Another such example is the shapelet corresponding to the class 7, for the generalised shapelet transform. This is perhaps surprising to see as a distinguishing feature of a 7. However it turns out that no other digit uses a stroke in that direction, in that place! (Figuring this out was a fun moment for the authors, sketching figures in the air.) A similar case can be made for the 2 shapelet.

For further details see Appendix A.2.

## 4.2 Learning lengths, with irregularly sampled partially observed time series

We now investigate the strategy of learning lengths differentially.

So as to keep things interesting, and to additionally provide benchmarks on irregularly sampled partially-observed datasets (to which the classical shapelet transform cannot be applied), for this test we drop either 10%, 30% or 50% of the data for each of the JapaneseVowels, Libras and LSST datasets, selected for representing a range of difficulties. The data dropped is independently selected for every channel of each time series, and is the same for every model and repeat.

We use the generalised shapelet transform with learnt  $L^2$  discrepancy, except we fix the lengths rather than learning them differentially. We then perform a hyperparameter search to determine the best and worst lengths for this model on each dataset. We then train a model with differentially learnt lengths initialised at the *worst* lengths, and compare it to the *best* performer from the hyperparameter search. See Table 2. We see that the performance is comparable! This demonstrates that lengths learnt differentially perform just as effectively as those selected by hyperparameter search, but without the relatively more expensive search.

For further details see Appendix A.3.

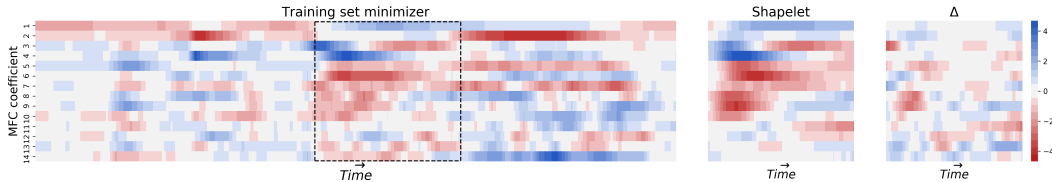


Figure 2: Generalised shapelet transform with learnt  $L^2$  discrepancy. First 15 MFC coefficients for the training set minimizer (left), shapelet (middle), and the difference between them (right). The dashed box in the minimizer plot indicates the position in the series that the shapelet corresponds to.

### 4.3 Speech Commands

Finally we consider the Speech Commands dataset [28]. This is comprised of one-second audio files, corresponding to words such as ‘yes’, ‘no’, ‘left’, ‘right’, and so on. We consider 10 classes so as to create a balanced classification problem.

For the generalised shapelet transform, we use the MFC discrepancy described in equation (7).

**Classification performance** For this more difficult dataset, the generalised shapelet transform substantially outperformed the classical shapelet transform. (To keep things fair, the classical shapelet transform is used in MFC-space; the performance gap is not due to this.) The classical shapelet transform produces a test accuracy of  $44.8\% \pm 8.6\%$ , whilst the generalised shapelet transform produces a test accuracy of  $91.9\% \pm 2.4\%$  (mean  $\pm$  std, averaged over three runs).

#### Interpretability

We examine interpretability in three different ways. First we consider MFC-space, see Figure 2. We see that the shapelets have learnt to resemble small segments of the training data, so that classification may be determined by the presence of different frequencies.

Furthermore, these shapelets may be listened to as audio! The audio files may be found at [https://github.com/patrick-kidger/generalised\\_shapelets/tree/master/audio](https://github.com/patrick-kidger/generalised_shapelets/tree/master/audio). The audio to MFC map is naturally lossy, so the shapelets are far from perfect, but the difference between them is nonetheless clear. The shapelet most strongly associated with ‘left’ captures the ‘eft’ sound, whilst the one for ‘stop’ actually sounds like the word itself. Much like the shapelet associated with class 7 in the PenDigits example, the sounds extracted need not resemble the word in isolation. Instead, they capture features that distinguishes that class from the others present.

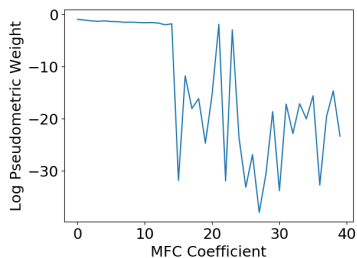


Figure 3: Pseudometric channel weighting identifies a spectral gap.

Finally, we examine the coefficients of the learnt  $L^2$  pseudometric, recalling that the matrix  $A$  of equation (7) is diagonal and thus weights the importance of each channel. See Figure 3. The coefficients of the pseudometric have learnt to be relatively large for the first 15 channels, and dramatically smaller for the later 25 channels. The pseudometric has learnt – purely from data – a quantitative description of the fact that lower frequencies are more important to distinguish words [29]. In short, it has discovered a kind of spectral gap!

See Appendix A.5 for further plots.

## 5 Conclusion

In this work we have generalised the classical shapelet method in several ways. We have generalised it from discrete time to continuous time, and in doing so extended the method to the general case of irregularly-sampled partially-observed multivariate time series. Furthermore this allows for the length of each shapelet to be treated as a parameter rather than a hyperparameter, and optimised differentially. We have introduced generalised discrepancies to allow for domain adaptation. Finally we have introduced a simple regularisation penalty that produces interpretable results capable of giving new insight into the data.



## Broader Impact

Interpretability is important in the application of many machine learning systems, often over and above raw performance, so that the reason for choices made on the basis of that system can be justified, seen to be made fairly, and without undue bias. Furthermore, methods which give new insight into the data are valuable for their ability to help the subsequent development of theory. The generalised shapelet transform, with interpretable regularisation, is capable of supporting both of these objectives, and so it is our hope that a substantial part of the broader impact of this work will be its contributions towards these strategic goals.

## Acknowledgments and Disclosure of Funding

PK was supported by the EPSRC grant EP/L015811/1. JM was supported by the EPSRC grant EP/L015803/1 in collaboration with Iterex Therapeutics. PK, JM, TL were supported by the Alan Turing Institute under the EPSRC grant EP/N510129/1.

## References

- [1] L. Ye and E. Keogh, “Time Series Shapelets: A New Primitive for Data Mining,” *KDD 2009*, 2009.
- [2] J. Grabocka, N. Schilling, M. Wistuba, and L. Schmidt-Thieme, “Learning Time-Series Shapelets,” *KDD 2014*, 2014.
- [3] L. Hou, J. Kwok, and J. Zurada, “Efficient Learning of Timeseries Shapelets,” *AAAI 2016*, 2016.
- [4] A. Bagnall, A. Bostrom, J. Large, and J. Lines, “The Great Time Series Classification Bake Off: An Experimental Evaluation of Recently Proposed Algorithms. Extended Version,” *arXiv:1602.01711*, 2016.
- [5] J. Hills, J. Lines, E. Baranauskas, J. Mapp, and A. Bagnall, “Classification of time series by shapelet transformation,” *Data Mining and Knowledge Discovery*, vol. 28, pp. 851–881, 2014.
- [6] Y. Wang, R. Emonet, E. Fromont, S. Malinowski, E. Menager, L. Mosser, and R. Tavenard, “Learning Interpretable Shapelets for Time Series Classification through Adversarial Regularization,” *CAP 2019*, 2019.
- [7] A. Mueen, E. Keogh, and N. Young, “Logical-Shapelets: An Expressive Primitive for Time Series Classification,” in *Proceedings of the 17th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 1154–1162, 2011.
- [8] J. Grabocka, M. Wistuba, and L. Schmidt-Thieme, “Scalable Discovery of Time-Series Shapelets,”
- [9] J. Grabocka, M. Wistuba, and L. Schmidt-Thieme, “Fast classification of univariate and multivariate time series through shapelet discovery,” *Knowl. Inf. Syst.*, vol. 49, pp. 429–454, 2016.
- [10] T. Rakthanmanon and E. Keogh, “Fast shapelets: A scalable algorithm for discovering time series shapelets,” in *Proceedings of the 13th SIAM International Conference on Data Mining*, pp. 668–676, 2013.
- [11] M. Wistuba, J. Grabocka, and L. Schmidt-Thieme, “Ultra-Fast Shapelets for Time Series Classification,” *arXiv:1503.05018*, 2015.
- [12] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser, M. Kudlur, J. Levenberg, D. Mané, R. Monga, S. Moore, D. Murray, C. Olah, M. Schuster, J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke, V. Vasudevan, F. Viégas, O. Vinyals, P. Warden, M. Wattenberg, M. Wicke, Y. Yu, and X. Zheng, “TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems,” 2015. Software available from tensorflow.org.
- [13] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimselshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, “PyTorch: An Imperative Style, High-Performance Deep Learning Library,” in *Advances in Neural Information Processing Systems 32*, pp. 8024–8035, Curran Associates, Inc., 2019.
- [14] J. Bradbury, R. Frostig, P. Hawkins, M. J. Johnson, C. Leary, D. Maclaurin, and S. Wanderman-Milne, “JAX: composable transformations of Python+NumPy programs,” 2018.
- [15] M. Shah, J. Grabocka, N. Schilling, M. Wistuba, and L. Schmidt-Thieme, “Learning DTW-Shapelets for Time-Series Classification,” *CODS 2016*, 2016.

- [16] A. Bostrom and A. Bagnall, “Binary shapelet transform for multiclass time series classification,” in *Big Data Analytics and Knowledge Discovery* (S. Madria and T. Hara, eds.), (Cham), pp. 257–269, Springer International Publishing, 2015.
- [17] S. N. Shukla and B. Marlin, “Interpolation-prediction networks for irregularly sampled time series,” in *International Conference on Learning Representations*, 2019.
- [18] S. C.-X. Li and B. M. Marlin, “A scalable end-to-end Gaussian process adapter for irregularly sampled time series classification,” in *Advances in Neural Information Processing Systems*, pp. 1804–1812, 2016.
- [19] J. Futoma, S. Hariharan, and K. Heller, “Learning to Detect Sepsis with a Multitask Gaussian Process RNN Classifier,” in *Proceedings of the 34th International Conference on Machine Learning*, pp. 1174–1182, 2017.
- [20] S. Liao, T. Lyons, W. Yang, and H. Ni, “Learning stochastic differential equations using RNN with log signature features,” *arXiv:1908.08286*, 2019.
- [21] M. Xu, L.-Y. Duan, J. Cai, L.-T. Chia, C. Xu, and Q. Tian, “Hmm-based audio keyword generation,” in *Advances in Multimedia Information Processing - PCM 2004* (K. Aizawa, Y. Nakamura, and S. Satoh, eds.), (Berlin, Heidelberg), pp. 566–574, Springer Berlin Heidelberg, 2005.
- [22] T. Lyons, M. Caruana, and T. Levy, *Differential equations driven by rough paths*. Springer, 2004. École d’Été de Probabilités de Saint-Flour XXXIV - 2004.
- [23] P. Bonnier, P. Kidger, I. Perez Arribas, C. Salvi, and T. Lyons, “Deep Signature Transforms,” in *Advances in Neural Information Processing Systems*, pp. 3099–3109, 2019.
- [24] P. Kidger and T. Lyons, “Signatory: differentiable computations of the signature and logsignature transforms, on both CPU and GPU,” *arXiv:2001.00706*, 2020. <https://github.com/patrick-kidger/signatory>.
- [25] S. Howison, A. Nevado-Holgado, S. Swaminathan, A. Kormilitzin, J. Morrill, and T. Lyons, “Utilisation of the signature method to identify the early onset of sepsis from multivariate physiological time series in critical care monitoring,” *Critical Care Medicine*.
- [26] M. Lothaire, “Combinatorics on words,” 1997.
- [27] A. Bagnall, H. A. Dau, J. Lines, M. Flynn, J. Large, A. Bostrom, P. Southam, and E. Keogh, “The uea multivariate time series classification archive, 2018,” *arXiv preprint arXiv:1811.00075*, 2018.
- [28] P. Warden, “Speech commands: A dataset for limited-vocabulary speech recognition,” *arXiv preprint arXiv:1804.03209*, 2018.
- [29] B. B. Monson and J. Caravello, “The maximum audible low-pass cutoff frequency for speech,” *The Journal of the Acoustical Society of America*, vol. 146, no. 6, pp. EL496–EL501, 2019.
- [30] D. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *ICLR*, 2015.

## A Experimental details

### A.1 General notes

Many details of the experiments are already specified in Section 4, and we do not repeat those details here.

**Code** Code to reproduce every experiment can be found at [https://github.com/patrick-kidger/generalised\\_shapelets](https://github.com/patrick-kidger/generalised_shapelets).

**Choice of  $\iota$**  The interpolation scheme  $\iota$  is taken to be piecewise linear interpolation. In particular efficient algorithms for computing the logsignature transform only exist for piecewise linear paths [24].

**Regularisation parameter** The parameter  $\gamma$  for the interpretable regularisation is taken to be  $10^{-4}$ . This was selected by starting at  $10^{-3}$  and reducing the value until test accuracy no longer improved, so as to ensure that it did not compromise performance.

**Optimisation** The loss was cross entropy, the optimiser was Adam [30] with learning rate 0.05 and batch size 1024. If validation loss stagnated for 20 epochs then the learning rate was reduced by a factor of 10 and training resumed, down to a minimum learning rate of 0.001. We note that these relatively large learning rates are (as is standard practice) proportional to the large batch size. If validation loss and accuracy failed to decrease over 60 epochs then training was halted. Once training was completed then the model parameters were rolled back to those which produced the highest validation accuracy.

**Computer infrastructure** Experiments were run on the CPU of a variety of different machines, all using Ubuntu 18.04 LTS, and running PyTorch 1.3.1.

### A.2 UEA

The datasets can be downloaded from <https://timeseriesclassification.com>.

The maximum number of epochs allowed for training was 250.

All UEA datasets were used unnormalised. Those samples which were shorter than the maximum length of the sequence were padded to the maximum length by repeating their final entry.

Hyperparameters were found by performing a grid search over 2, 3, 5 shapelets *per class*, with a maximum total number of shapelets of 30, and shapelets being set to (classical shapelet transform) / initialised at (generalised shapelet transform) 0.15, 0.3, 0.5, 1.0 times the maximum length of the time series.

The dataset comes with default train/test splits, which we respect here. The training data is split 80%/20% into train and validation sets, stratified by class label.

The details of each dataset are as below. We note that the train/test splits are sometimes of unusual proportion; we do not know the reason for this odd choice.

Table 3: UEA dataset details and hyperparameter choices

Dataset	Train size	Test size	Dimensions	Length	Classes	Shapelets	Shapelet length fraction
BasicMotions	40	40	6	100	4	12	0.5
ERing	30	30	4	65	6	12	0.5
Epilepsy	137	138	3	206	4	20	0.5
Handwriting	150	850	3	152	26	30	0.5
JapaneseVowels	270	370	12	29	9	18	0.5
Libras	180	180	2	45	15	30	1.0
LSST	2459	2466	6	36	14	28	1.0
PenDigits	7494	3498	2	8	10	30	0.5
RacketSports	151	152	6	30	4	12	0.5

### A.3 Learning lengths

Table 4 shows the hyperparameters used in Section 4.2. These were chosen to optimize the validation score for the generalised shapelet transform without learnt lengths, rather than the classical shapelet transform, and as such are different to those noted above. The hyperparameters were optimised for the 30% drop rate, and the same hyperparameters simply used for the 10% and 50% drop rate cases.

Table 4: Hyperparameter choices for study on learnt lengths

Dataset	Shapelets	Best length fraction	Worst length fraction
JapaneseVowels	27	0.15	0.5
Libras	30	1.0	0.15
LSST	28	0.3	1.0

### A.4 Full hyperparameter search results

For completeness we also give the full results from both hyperparameter searches in the preceding two sections, in Tables 5 and 6.

### A.5 Speech Commands

The dataset can be downloaded from

[http://download.tensorflow.org/data/speech\\_commands\\_v0.02.tar.gz](http://download.tensorflow.org/data/speech_commands_v0.02.tar.gz).

We began by selecting every sample from the ‘yes’, ‘no’, ‘up’, ‘down’, ‘left’, ‘right’, ‘on’, ‘off’, ‘stop’ and ‘go’ categories, and discarding the samples which were not of the maximum length (16000; nearly every sample is of this maximum length). This gives a total of 34975 samples.

The samples were preprocessed by computing the MFC with a Hann window of length 400, hop length 200, 400 frequency bins, 128 mels, and 40 MFC coefficients. Every sample is then of length 81 with 40 channels. Every channel was then normalised to have mean zero and variance one.

No hyperparameter searching was performed, due to the inordinately high cost of doing so - shapelets are an expensive algorithm that is primarily a ‘small data’ technique, and this represented the upper limit of problem size that we could consider! That said, this is in large part an implementation issue. An efficient GPU implementation should be possible. The problem is that current machine learning frameworks [12, 13, 14] typically parallelise by vectorising every operation, however this problem (which is embarrassingly parallel) is instead best handled via naïve parallelism at the top level. This is because of the need for different behaviour for different batch elements; they will in general have minimisers at different sections of the time series, meaning that a vectorised approach needs to keep track of the union of these points for every batch element. We highlight that not needing to perform hyperparameter optimisation on the length is an advantage of our generalised shapelet transform, thus reducing this kind of computational burden.

The maximum number of epochs allowed for training was 1000. The number of shapelets used per class was 4, for a total of 40 shapelets. The length of each shapelet (set to for the classical shapelet transform; initialised at for the generalised shapelet transform) was taken to be 0.3 of the full length of the dataset.

The data is combined into a single dataset and a 70%/15%/15% training/validation/test split taken, stratified by class label.

In Figure 4 we show all 40 MFC coefficients for the shapelet (blue) and the training set minimizer (orange, dashed) for the generalised shapelet transform with  $L^2$  discrepancy, for a particular (arbitrarily selected) run.

**Hyperparameter options.** Top row: shapelets per class. Bottom row: Shapelet length fraction.

Dataset	2		3			5			5				
	0.15	0.3	0.5	1.0	1.0	0.15	0.3	0.5	1.0	0.15	0.3	0.5	1.0
BasicMotions	100.0%	100.0%	100.0%	100.0%	100.0%	100.0%	<b>100.0%</b>	100.0%	75.0%	100.0%	100.0%	100.0%	100.0%
ERing	83.3%	100.0%	<b>100.0%</b>	100.0%	100.0%	100.0%	83.3%	100.0%	100.0%	100.0%	83.3%	83.3%	100.0%
Epilepsy	82.1%	75.0%	82.1%	82.1%	82.1%	82.1%	85.7%	82.1%	71.4%	82.1%	82.1%	<b>89.3%</b>	75.4%
Handwriting	13.3%	16.7%	23.3%	23.3%	23.3%	16.7%	23.3%	<b>26.7%</b>	23.3%	16.7%	20.0%	16.7%	23.3%
JapaneseVowels	90.7%	90.7%	<b>96.3%</b>	92.6%	92.6%	92.6%	92.6%	94.4%	92.6%	90.7%	92.6%	96.3%	90.7%
Libras	83.3%	80.6%	88.9%	83.3%	77.8%	77.8%	88.9%	86.1%	77.8%	77.8%	86.1%	80.6%	<b>91.7%</b>
LSST	34.8%	33.7%	33.3%	<b>35.2%</b>	33.7%	33.7%	33.9%	35.0%	34.1%	33.7%	33.7%	35.0%	34.6%
PenDigits	97.6%	98.0%	98.7%	96.7%	98.0%	98.0%	97.9%	98.4%	96.5%	97.6%	98.5%	<b>98.9%</b>	96.3%
RacketSports	61.3%	74.2%	83.9%	80.6%	58.1%	58.1%	67.7%	<b>87.1%</b>	77.4%	71.0%	77.4%	64.5%	80.6%

Table 5: Accuracy on the validation set for the hyperparameter runs performed to determine the hyperparameters used in Section 4.1. The upper row represents the number of shapelets per class, with the lower row being the shapelet length fraction. The best run is given in bold. When multiple options achieved the highest score, the hyperparameters were chosen randomly from that top performing set. Only one run was performed for each hyperparameter option.

**Hyperparameter options.** Top row: shapelets per class. Bottom row: Shapelet length fraction.

Dataset	2		3			5			5			
	0.15	0.3	0.5	1.0	0.15	0.3	0.5	1.0	0.15	0.3	0.5	1.0
JapaneseVowels	94.4%	94.4%	94.4%	92.6%	<b>96.3%</b> *	<b>94.4%</b> *	94.4%	94.4%	96.3%	94.4%	92.6%	92.6%
Libras	<b>63.9%</b> *	77.8%	77.8%	<b>88.9%</b> *	75.0%	80.6%	83.3%	83.3%	73.2%	69.4%	83.3%	86.1%
LSST	45.7%	<b>46.5%</b> *	42.5%	<b>37.0%</b> *	41.9%	43.3%	42.6%	39.0%	45.5%	43.5%	41.4%	39.4%

Table 6: Accuracy on the validation set for the hyperparameter runs performed to determine the hyperparameters used in Section 4.2. The upper row represents the number of shapelets per class, with the lower row being the shapelet length fraction. The best and worst hyperparameters – recall that worst is also used here – are denoted in bold. The best case additionally has a superscript \* and the worst case additionally has a superscript \*\*. When multiple options achieved the highest score, the hyperparameters were chosen randomly from that top performing set.

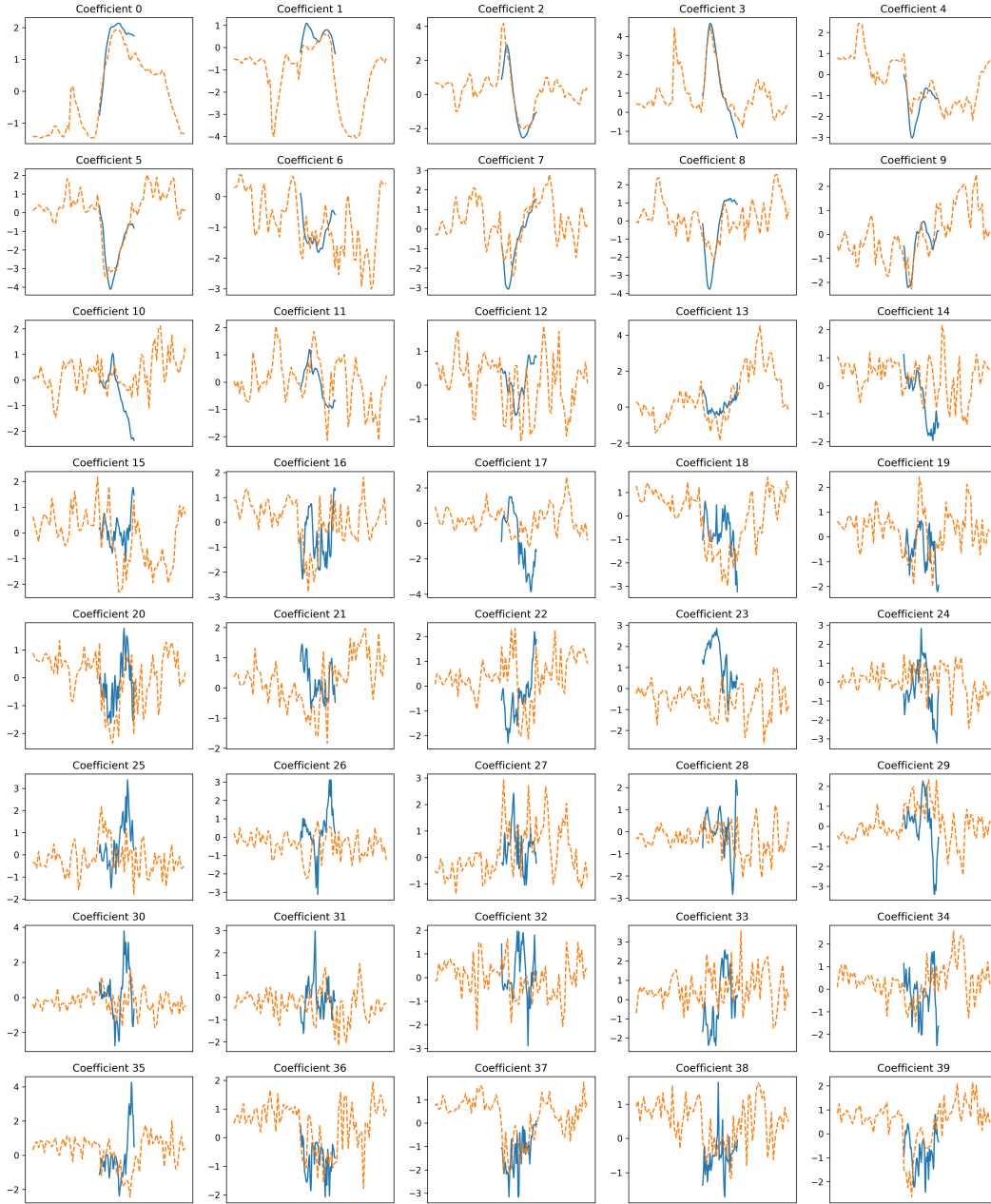


Figure 4: Generalised shapelet transform with  $L^2$  discrepancy. All MFC coefficients for the shapelet (blue) and the training set minimizer (orange, dashed).