# OTW: OPTIMAL TRANSPORT WARPING FOR TIME SERIES

*Fabian Latorre\*, Chenghao Liu†, Doyen Sahoo†, Steven C.H. Hoi†*

\*École polytechnique fédérale de Lausanne (EPFL)
†Salesforce Research Asia

## ABSTRACT

Dynamic Time Warping (DTW) has become the pragmatic choice for measuring distance between time series. However, it suffers from unavoidable quadratic time complexity when the optimal alignment matrix needs to be computed exactly. This hinders its use in deep learning architectures, where layers involving DTW computations cause severe bottlenecks. To alleviate these issues, we introduce a new metric for time series data based on the Optimal Transport (OT) framework, called Optimal Transport Warping (OTW). OTW enjoys linear time/space complexity, is differentiable and can be parallelized. OTW enjoys a moderate sensitivity to time and shape distortions, making it ideal for time series. We show the efficacy and efficiency of OTW on 1-Nearest Neighbor Classification and Hierarchical Clustering, as well as in the case of using OTW instead of DTW in Deep Learning architectures.

***Index Terms***— Time-series, Optimal Transport, Deep Learning, Optimization.

## 1. INTRODUCTION

Time-series data is ubiquitous in contemporary Machine Learning applications including Heart disease predictions based on ECG data [1], end-to-end text-to-speech synthesis [2] and sign-language recognition [3]. Common tasks like supervised classification and hierarchical clustering require a notion of *distance* or *similarity* between time series, and their performance strongly depends on such a choice. A key desired characteristic of a time series distance is the ability to identify commonly occurring patterns such as shape distortions and time delays. For example, two time series that differ by a slight distortion in shape or a slight delay in time should be considered *similar* or *close* [4].

The euclidean distance is unable to reason about shape and time distortions, and thus, is considered a poor choice for time series applications. Dynamic Time Warping (DTW) [5] can assign high similarity values to time series that have a closely matching shape, but which do not necessarily align perfectly in the time domain. For this reason, DTW has established itself as the most prominent similarity measure for time series [6, 7, 8, 9, 10].

However, DTW and all closely-related variants suffer from an unavoidable quadratic complexity with respect to the length of the time series [11]. This makes it highly expensive or outright unusable for time series of considerable length, especially for large datasets. This issue has been noted in literature, and researchers have developed faster variants of DTW, such as FastDTW[12]. Unfortunately, FastDTW has been dismissed as being slower than DTW in practice, despite its theoretical linear time complexity [13].

Recently, DTW has been incorporated inside Deep Learning pipelines for time-series data. It has been used either as a loss function [14, 15], or as a replacement for linear layers as in DTWNet [10]. In the latter case it replaces the inner product operation in the linear layers of the network. Typically, these layers have a linear time-complexity with regard to input size (e.g. convolutional layers, batch-norm, etc.), and replacing them with a quadratic complexity DTW layer introduces a significant computational bottleneck.

Moreover, DTW uses Dynamic Programming, an inherently sequential framework that does not lend itself to trivial paralellization in GPU. Even though there are available GPU implementations [16, 17, 18], they do not avoid the sequential nature. This limits the adoption of DTW in practice, as the cost of training and hyperparameter optimization increases considerably. Lastly, DTW is not a metric, as it violates the triangle inequality. Thus, it cannot exploit faster similarity search methods like the Approximating and Eliminating Search Algorithm (AESA) [19].

Thus, we need a distance notion that not only can be computed in linear time, but is also differentiable and can be parallelized on GPU, thus speeding up all training pipelines by a huge margin. Even though the euclidean distance enjoys such characteristics, it does not provide a good inductive bias for time-series data, as it ignores the time component. Precisely, the goal is to keep the theoretical properties and performance of DTW at a computational cost similar to that of the euclidean distance.

**Our contributions.** It is apparent that there is a need for new time series distance notions that overcomes some, if not all, the aforementioned drawbacks. We propose a new distance for time series data that we call **Optimal Transport Warping (OTW)**. Our distance is rooted in the theory of Optimal Transport [20] (OT), which is a well-known metric

---

used to compare the shape of two probability distributions.

We adapt Optimal Transport for the case of time-series data through an Unbalanced Optimal Transport formulation, given that two time series may not have the same mass (may not sum up to the same value). We also address the issue of negativity, i.e., while probability distributions are non-negative, time series may not be. Our final OTW formulation (1) can be computed in linear time and space, (2) is differentiable, (3) can be easily computed on massively parallel architectures (GPU/TPUs), (4) has moderate sensitivity to shape and time distortions, and (5) is a proper metric. A comparison of our method against several state-of-the-art time series distances is summarized in Table 1.

In experiments we observe that OTW runs up to 10-30x times faster than DTW, depending on the dataset. In the classification task it improves over DTW in 6 out of 7 types of datasets in the UCR time-series benchmark (table 2). In the clustering task, it improves over DTW in all 7 types of datasets considered (table 6). A total of 92 datasets were considered. Through synthethic and real experiments we show that by replacing DTW with OTW in DTW-Net we solve its computational bottleneck, and we can achieve a lower error in less than 50% of the time.

## 2. RELATED WORK

We summarize the most prominent time-series distances in table 1, and conclude that there are few low-complexity alternatives to DTW, which suffers from quadratic complexity. Only GDTW[21], FastDTW [12] and OTW (our proposed method) are fast alternatives that avoid computational bottlenecks. However, FastDTW has been dismissed by the community [13] while GDTW only provides an CPU implementation of a stochastic rather than a deterministic gradient. Because the activations appear at every layer, this prevents the computation of an unbiased estimator of the gradient of the loss function. Moreover, GDTW requires a Contiguous Sparse Matrix Data Structure that might not be efficient in GPU (the authors of the paper do not provide a GPU implementation). Hence, a priori we expect that training networks using GDTW activations on CPU is an unfeasible task given limited time.

In recent years, DTW has found popularity in several applications using Deep Learning architectures. DTW-like distances have been used either as a loss function for time series forecasting/regression, or as part of a feature-extracting module that can be used as a layer inside a neural network. When such distances are used to replace inner products, as in linear layers or transformers, the complexity is inevitably increased. Because layers are applied in a sequential fashion (as they cannot be parallelized), this forms a bottleneck that slows down the learning pipeline.

A few recent works developing Deep Learning pipelines that make use of a DTW-like module, potentially suffering from speed issues: DILATE [14] is a loss function module

| Method | Complexity | Gradient | GPU |
|---|---|---|---|
| TWED [22] | $O(n^2)$ | ✗ | ✗ |
| ERP [23] | $O(n^2)$ | ✗ | ✗ |
| EDR [24] | $O(n^2)$ | ✗ | ✗ |
| DTW [5] | $O(n^2)$ | ✗ | ✓ |
| SoftDTW [9] | $O(n^2)$ | ✓ | ✓ |
| DILATE [14] | $O(n^2)$ | ✓ | ✓ |
| SoftDTW div. [25] | $O(n^2)$ | ✓ | ✓ |
| FastDTW [12] | $O(n)$ | ✗ | ✗ |
| GDTW [21] | $O(n)$ | Stochastic | ✗ |
| OTW (This work) | $O(n)$ | ✓ | ✓ |

**Table 1**. Comparison of time series similarity measures. **Gradient:** differentiability with respect to its inputs. **GPU:** existence of a parallel GPU implementation.

computing Soft-DTW as a subroutine; DTWNet [10] proposes a feature-extraction module based on DTW; D3TW [15] is a discriminative Loss for Action Alignment and Segmentation based on Soft-DTW; SLR [26] is an Encoder-Decoder architecture with Soft-DTW alignment constraint for Continuous Sign Language Recognition; STRIPE [27] is similar to [14]; DTW-NN [28] is similar to [10]; SpringNet is a Transformer architecture with DTW replacing inner products; TTS [2] is an End-to-End text-to-speech architecture with a Soft-DTW-based prediction loss; [29] uses DTW to synchronize a resultant and a target sequence inside a text-to-video sign-language synthesizer network. Due to the soft-DTW computational demands, an alternative attention mechanism is studied.

**OT-based time series distances.** The potential of Optimal Transport for time-series has been explored in other works. For example, [30] proposes to compare positive time-series using sinkhorn divergences, a variant of OT. However, even in the one-dimensional case there exists no subquadratic algorithm for the problem. Time-adaptive Optimal Transport [31] is a similar approach. Nevertheless, their formulation is essentially different from ours: for a time series $a$, a uniform distribution is constructed, over pairs $(i, a_i)$, which are then compared using the traditional OT formulation. Because this is a two-dimensional distribution, the complexity of their proposed algorithm is at least quadratic and hence, as slow as DTW. In contrast our algorithms work in linear time. We recall that OT requires cubic time but it can be approximated in quadratic time using Sinkhorn iterations [32]. Only in the one-dimensional case it has a linear-time implementation [33]. Finally, [34] interprets the values of a time series as *set* of samples rather than a sequence. This has the effect of removing the time information of the sequence and hence OT, as used in this case, is oblivious to the temporal nature of the data. The resulting optimization problems are solved with the so-called Sinkhorn iterations which again lead to a quadratic complexity. In summary, ours is the first linear-time distance for time-series

based on OT.

## 3. OPTIMAL TRANSPORT WARPING

### 3.1. Problem Setting and Optimal Transport

Probability distributions are analogous to density functions in physics, which measure the concentration of mass along an infinitesimal unit of volume. We can think of certain time-series as measuring a concentration of mass along an infinitesimal length of *time*: if we let $a(t)$ be the amount of car traffic at time $t$ of the day, after normalization, we can think of the integral $a(t)$ over a time interval, as the probability that a car is in traffic during that interval.

Hence, a method that compares probability distributions can potentially compare time-series. Indeed, Optimal transport has been used to compare time-series or general sequences before [32, 34, 31]. Optimal transport provides a many-to-many alignment, instead of a one-to-many alignment in the case of DTW and its variants. A priori, it is reasonable to believe that some time-series data might be more amenable to alignment using one type of alignment over the other. Indeed, OT is also sometimes called the *Earth Mover's Distance*, because it captures the cost of *reshaping* one distribution into another. In this way it is sensitive to shape distortions, a fact that has been exploited and explained in depth in [30].

A discrete time series $a$ is a sequence of numbers $(a_1, \ldots, a_n)$, that can be understood as a function from the set $[n] = \{1, \ldots, n\}$ to the real numbers. Since Optimal Transport is a distance notion between probability measures, we first focus on time series that are positive and sum to one i.e., $\sum_{i=1}^{n} a_i = 1; a_i \geq 0$ and $\sum_{j=1}^{n} b_j = 1; b_i \geq 0$. Later, we will expand the notion of distance to arbitrary time series. To define optimal transport for probability distributions over the set $[n]$, we require a nonnegative function $d : [n] \times [n] \to \mathbb{R}_+$ defined over pairs of elements of $[n]$. We denote by $D$ the matrix with entries $D_{i,j} := d(i,j)$ and we refer to it as the *distance matrix*. Denoting the column vector of all-ones as $\mathbb{1} \in \mathbb{R}^n$, the Optimal Transport distance (with respect to the distance matrix $D$) between $a$ and $b$ is defined as:

$$W_D(a,b) := \min_T \left\{ \langle T, D \rangle : T \geq 0, T\mathbb{1} = a, \mathbb{1}^T T = b \right\} \quad (1)$$

For an introduction to Optimal Transport please refer to [20]. Note that we can extend definition (1) to sequences that sum up to the same value, not necessarily equal to 1: because $a$ is defined as the row-sums of $T$, and $b$ is defined as the column-sums of $T$, $a$ and $b$ need only sum up to the same value, equal to the sum of all the entries of $T$. However, we cannot remove the constraint that that $a, b \geq 0$, as the entries of $T$ are positive c.f. eq. (1).

In our case we will choose the matrix $D$ as the absolute value distance, i.e., $d(i,j) := |i - j|$. This choice means that transporting a unit of mass from the position $i$ to the position $j$ will be equal to the absolute difference between the two positions. In the following we will always assume that $D$ is chosen in this way and we let $W(a,b) := W_D(a,b)$ to simplify notation. This choice of $D$ ensures that there exists a closed form solution for the Optimal Transport problem eq. (1) that can be computed in linear time [33]:

$$W(a,b) = \sum_{i=1}^{n} |A(i) - B(i)|$$

$$A(i) := \sum_{j=1}^{i} a_j, \quad B(i) := \sum_{j=1}^{i} b_j \quad (2)$$

that is, $A$ and $B$ are the *cumulative distribution* functions of $a$ and $b$, respectively.

Because a time series is not necessarily a probability measure, we cannot directly use eq. (1). The two main issues are: (i) *unbalancedness*, as the time series may not sum up to the same value and (ii) *negativity*, as a time series may contain negative values. In order to extend the Optimal Transport framework to arbitrary time series we need to workaround the limitations of the definition in eq. (1), while trying to retain the linear space/time properties of the closed-form solution (2).

### 3.2. Unbalanced Optimal Transport for time-series

We will first assume that the sequences $a, b$ are nonnegative but that they do not necessarily sum up to the same value. In order to resolve the unbalancedness issue, we proceed by adding a *sink* as described in [35, Chapter 3]: we append one additional element to both sequences as follows:

$$\hat{a}_i := \begin{cases} a_i & \text{if } i \in [n] \\ \sum_{j=1}^{n} b_j & \text{if } i = n+1 \end{cases}$$

$$\hat{b}_i := \begin{cases} b_i & \text{if } i \in [n] \\ \sum_{j=1}^{n} a_i & \text{if } i = n+1 \end{cases} \quad (3)$$

In this way, we ensure that $\sum_i \hat{a}_i = \sum_i \hat{b}_i$. This step can be understood as balancing the total mass of the two sequences. Now, we also have to extend the distance matrix in some way, to account for the *sink* in the extended sequences $\hat{a}$ and $\hat{b}$. For some value $m \in \mathbb{R}_+$ called the *waste cost* we set:

$$D(m)_{i,j} = \begin{cases} |i - j| & \text{if } i, j \in [n] \\ m & \text{if } i = n+1, j \in [n] \\ m & \text{if } j = n+1, i \in [n] \\ 0 & \text{if } i = n+1, j = n+1 \end{cases} \quad (4)$$

the idea being that any excess can be transported to the $n+1$-th point (the sink) incurring a cost of $m$ per unit of mass. We define the *(nonnegative) unbalanced Optimal Transport Distance* problem as:

$$\widehat{W}_m(a,b) := W_{D(m)}(\hat{a}, \hat{b})$$

$$= \min_T \left\{ \langle T, D(m) \rangle : T \geq 0, T\mathbb{1} = \hat{a}, \mathbb{1}^T T = \hat{b} \right\} \quad (5)$$

note that if the sequences $a, b$ sum up to the same value, the problem reduces to the original *balanced* optimal transport problem eq. (1). With this modification, it is not clear if eq. (5) can also be solved in closed form in linear time and space as in eq. (2). However, we can compute an upper bound:

**Theorem 1.** *Let $D_{i,j} = |i - j|$ and let $D(m)$ be defined as in eq. (4). Define*

$$OTW_m(a, b) := m|A(n) - B(n)| + \sum_{i=1}^{n-1} |A(i) - B(i)| \quad (6)$$

*Then, $\widehat{W}_m(a, b) \leq OTW_m(a, b)$. Clearly, $OTW_m(a, b)$ can be computed in linear time/space.*

We defer the proof of theorem 1 to appendix C. Precisely, **we propose to use eq. (6) as the notion of distance** between time series of positive sign. In practice, the choice of $m < n$ works better. For large values of $n$, the choice $m = n$ would put too much weight on the first component $m|A(n) - B(n)|$, which strongly penalizes the total mass difference of the two sequences $a, b$.

As we show in lemma 1, the unbalanced OT distance (4) increases linearly when a time-shift is introduced. This makes it ideal for time-series applications like demand forecasting, where a shift in time can represent a change in the seasonality of a product, for example. In contrast, in speech recognition where time-shifts should be ignored, perhaps a distance like DTW might be more suitable.

**Lemma 1.** *Let $a \in \mathbb{R}_+^n$ and $b$ be two time series, and let $b'$ be a shifted version of $b$ by $t < n$ units, that is $b_i' = b_{i-t}$ for $i = t + 1, \ldots, n$. For simplicity assume that $b_{n-t+1} = \ldots = b_n = 0$ i.e., the sequence $b$ is zero-padded. In this way $\sum_i b_i = \sum_i b_i'$. It holds that $|\widehat{W}_m(a, b') - \widehat{W}_m(a, b)| \leq t \left( \sum_{i=1}^n a_i \right)$*

The proof of lemma 1 is deferred to appendix C. This property means that the distance increases linearly proportional to the magnitude of the time-shift. Note that Soft-DTW [9] has a quadratic lower bound on its sensitivity to time-shifts, and was presented as a main contribution in [30, Theorem 1]. In contrast we show an upper bound and our sensitivity is linear.

**Constraining the Transport map to be local.** For time series distances like DTW, it has been observed that in practice, considering all possible alignments might be detrimental to the performance in downstream tasks. Instead, the best choice is to constrain the alignment to perform only *local* matching i.e., map samples only inside a constrained window. This is precisely the so-called DTW distance with Sakoe-Chiba band constraint c.f., [36] for details.

Our proposed Optimal-Transport based distance is no different, as the transport map $T$ in eq. (5) does not have any constraint, and constitutes a many-to-many map between arbitrary positions in $[n] = \{1, \ldots, n\}$. To mitigate this issue we propose to use a windowed-cumulative-sum rather than the

full cumulative-sum $A(i) := \sum_{j=1}^i a_j$ as originally proposed in eq. (2). More precisely we let:

$$A_s(i) := \sum_{j=1}^i a_j - \sum_{j=1}^{i-s} a_i, \quad B_s(i) := \sum_{j=1}^i b_j - \sum_{j=1}^{i-s} b_j \quad (7)$$

and we define the local Optimal Transport Warping distance for $s \in \{1, \ldots, n\}$:

$$OTW_{m,s}(a, b) := m|A_s(n) - B_s(n)| \\ + \sum_{i=1}^{n-1} |A_s(i) - B_s(i)| \quad (8)$$

**Localness.** the parameter $s$ is akin to the window parameter of the Sakoe-Chiba constraint in DTW: it interpolates between the $\ell_1$-norm distance and the Unbalanced Optimal Transport. One the one hand, the $\ell_1$-norm compares the two sequences $a, b$ entry-wise: it does not allow for alignment between two different time positions. On the other hand, the Unbalanced Optimal Transport distance allows transport of mass between any two time positions in the time-series. The parameter $s$ interpolates between the two extremes. In practice not all datasets require the same level of localness, and it is important to choose $s$ by cross-validation, which usually results in better performance. We summarize this fact in the following lemma, with proof deferred to appendix C:

**Lemma 2.** *For simplicity assume $m = 1$. When $s = 1$ then $\overline{OTW}_{1,1}(a, b) = \|a - b\|_1$. When $s = n$ we recover the global OTW distance (6) i.e., $OTW_{1,n}(a, b) = OTW_m(a, b)$.*

Finally, note that $OTW_{m,s}$ is not differentiable when $A_s = B_s$, due to the presence of the absolute value function. In order to have a fully differentiable version, it suffices to use a smooth approximation of absolute value, like the well-known *smooth $\ell_1$-loss*:

$$L_\beta(x) = \begin{cases} x^2/(2\beta) & |x| < \beta \\ |x| - \beta/2 & |x| \geq \beta \end{cases} \quad (9)$$

We define:

$$OTW_{m,s}^\beta(a, b) = m L_\beta(A_s(b) - B_s(n)) \\ + \sum_{i=1}^{n-1} L_\beta(A_s(i) - B_s(i)) \quad (10)$$

and it holds that $OTW_{m,s}^\beta \to OTW_{m,s}$ as $\beta \to 0$.

### 3.3. Dealing with negative values

Up to this point, we have presented a way to compare two time series that have only positive entries, using Unbalanced Optimal Transport. However, time series can contain negative values. In order to deal with sequences of arbitrary sign we propose to choose one of the following strategies, using cross-validation:

1. Apply eq. (10) to arbitrary sequences. Note that this formula can be applied even in the case where some entries of the sequences $a$ or $b$ are negative.

2. Split arbitrary sequences $a, b$ into their positive and negative parts i.e., $a_+ = \max(a, 0)$ and $a_- = \max(-a, 0)$, and sum the unbalanced optimal transport distance between the parts of equal sign:

$$\overline{\mathrm{OTW}}^{\beta}_{m,s}(a,b) = \mathrm{OTW}^{\beta}_{m,s}(a_+, b_+) \\ + \mathrm{OTW}^{\beta}_{m,s}(a_-, b_-) \quad (11)$$

**Remark.** Because we only need to compute the windowed-cumulative-sums (7) and then apply the smooth $\ell_1$-loss to the differences $A_s(i) - B_s(i)$ for $i = 1, \ldots, n$, we only need a linear number of operations, as a function of $n$. Moreover, since we only use basic operations available in Deep Learning frameworks like PyTorch, the computation can be readily performed in GPU using optimized code. Finally, because the smooth $\ell_1$-loss is convex, this implies the triangle inequality and OTW becomes a true metric. Next, we perform experiments to validate the performance of OTW.

## 4. EXPERIMENTS

### 4.1. Comparing OT to DTW on 1-nearest-neighbors classification

In this experiment we train 1-nearest-neighbors classifiers on the UCR time series classification archive, which consists of a large number of univariate time series data. Please note that **the purpose is not to show state-of-the-art performance**. Rather, our goal is to perform an apples-to-apples comparison of the Dynamic Time Warping (DTW) distance and our proposed OTW distance.

To choose hyperparameters for the OTW distance, for each combination we follow a 80/20 random split of the training/validation sets, and choose the one that maximizes accuracy on the validation set. For the best hyperparameters found, we evaluate the method on the testing set. The accuracy for the 1-nearest-neighbors classification for the DTW distance with learned warping window is directly obtained from the UCR time series classification benchmark website[1].

For our method, we do 10 independent runs and obtain a 95% confidence interval for the test error. Because of the lower complexity of our method vs DTW (linear vs quadratic complexity), it runs considerably faster and hence we consider it a better method if it can attain the same performance as DTW. Because there is only a single testing observation for the DTW methods available in the reported benchmarks, we consider that our method is better if the confidence interval for

---

[1] https://www.cs.ucr.edu/~eamonn/time_series_data_2018/

**Table 2**. Average test error of 1-NN classification using Learned DTW or Learned OTW (OTW error) and number of datasets in the collection where OTW outperforms DTW

| Type | DTW error | OTW error | Improv. | Total | % |
|------|-----------|-----------|---------|-------|---|
| Image | 0.25 | **0.25 ± 0.02** | 19 | 31 | **61.29** |
| Spectro | 0.26 | **0.24 ± 0.03** | 6 | 8 | **75.00** |
| Sensor | 0.23 | 0.44 ± 0.03 | 9 | 28 | 32.14 |
| Device | 0.34 | 0.44 ± 0.03 | 5 | 9 | **55.56** |
| Medical | 0.38 | **0.36 ± 0.02** | 9 | 13 | **69.23** |
| Traffic | 0.12 | **0.06 ± 0.02** | 2 | 2 | **100.00** |
| Power | 0.08 | **0.04 ± 0.01** | 1 | 1 | **100.00** |

the test error of OTW contains or is below the reported test error for DTW.

In table 2 we summarize the results according to the type of dataset. We group the ECG, EOG, HRM and Hemodynamics types of datasets under the name *Medical*, as they contain a relatively small number of datasets. We observe a noticeable improvement in all the considered categories except the *sensor* category.

Note that we remove the datasets corresponding to the motion and trajectory type, as the OTW distance is not suited for such data. OT is suited to time series that can be interpreted as probability distributions or, more generally, measures. Time-series that track the position of an object through time, do not fall in this category. The results for each single dataset are collected in appendix A. Overall, we observe in table 2 that **OTW improves over DTW on 6 out of 7 types of datasets**. On some datasets like Medical and Traffic, the advantage is apparent. In total, we evaluated on 92 dataset from the UCR time series benchmark, only the ones with missing values or variable length of sequence were discarded.

### 4.2. Hierarchical Clustering

We compare the DTW and OTW distance for time series clustering. We use the Agglomerative Clustering algorithm [37], which only requires access to the distance matrix. We run such clustering algorithm on the UCR time series benchmark collection. We skip datasets having more than 500 samples, as the quadratic memory requirements of the algorithms imposes such constraint. We evaluate the quality of the clustering using the Rand Index (RI), given that the datasets in the UCR archive are labelled. We summarize the results in table 6. We observe that our proposed distance outperforms DTW on most datasets considered. This is striking, as the time required for our method to run is considerably less than the DTW-based clustering. Overall, we observe in table 6 that **OTW improves over DTW on 7 out of 7 types of datasets**. On some datasets like Image, Traffic and Sensor, the advantage is apparent.

**Table 3**. Average Rand Index (RI) and average cpu time (seconds) of OT-based and DTW-based Hierarchical Clustering methods.

| Type | DTW RI - (time) | OTW RI - (time) | Improv. | Total | % |
|------|-----------------|-----------------|---------|-------|---|
| Device | 0.40 - (2, 279) | **0.49** - (**81**) | 5 | 6 | **83** |
| Image | 0.62 - (277) | **0.67** - (**28**) | 25 | 29 | **86** |
| Medical | 0.79 - (884) | **0.82** - (**40**) | 7 | 9 | **78** |
| Power | 0.50 - (90) | **0.52** - (**20**) | 1 | 1 | **100** |
| Sensor | 0.57 - (384) | **0.67** - (**31**) | 15 | 17 | **88** |
| Spectro | 0.56 - (88) | **0.61** - (**11**) | 4 | 7 | **57** |
| Traffic | 0.62 - (61) | **0.72** - (**9**) | 2 | 2 | **100** |

### 4.3. Performance of DTW-Net vs OTW-Net in synthetic and real data

DTW distances have been employed to design neural network layers with inductive biases that are better suited for time series data. This is the case, for example, of DTW-Net [10]. In this network, the first hidden layer consists of DTW distances between the input and the rows of a matrix, which is the trainable parameter of the layer. If there are $k$ such rows, then the computational complexity of the layer is $O(kn^2)$, where $n$ is the length of the input. On top of such features an arbitrary network architecture is added, which outputs the class probabilities.

In contrast, in a vanilla multi-layer fully-connected neural network the first hidden-layer (a linear layer) consists of inner products between the input and the rows of a matrix. If there are $k$ such rows, the complexity of this linear layer is $O(kn)$. Hence, the complexity of DTW-Net [10] is higher than a regular fully-connected neural network: it suffers from a computational bottleneck. In this experiment, we replace the DTW distance in DTW-Net by our OTW distance. Because OTW can be computed in linear time, this restores the complexity to $O(kn)$, in line with the other layers of the network, getting rid of the computational bottleneck. We describe such layers in algorithms 1 and 2.

---

**Algorithm 1** OTW Feature Extraction Layer

---

**Input:** input $a \in \mathbb{R}^n$
**Parameters:** matrix $B \in \mathbb{R}^{k \times n}$
**for** $i = 1$ to $k$ **do**
$\quad b \leftarrow B_{[:,i]}$ $\qquad\qquad\qquad$ ▷ $i$-th row of $B$
$\quad z_i \leftarrow \text{OTW}_{m,s}^{\beta}(a, b)$
**return** $z \in \mathbb{R}^k$

---

**Algorithm 2** DTW Feature Extraction Layer

---

**Input:** $a \in \mathbb{R}^d$
**Parameters:** matrix $B \in \mathbb{R}^{k \times n}$
**for** $i = 1$ to $k$ **do**
$\quad b \leftarrow B_{[:,i]}$ $\qquad\qquad\qquad$ ▷ $i$-th row of $B$
$\quad z_i \leftarrow \text{DTW}(a, b)$
**return** $z \in \mathbb{R}^k$

---

**Synthetic data.** In order to illustrate the performance of both types of networks, we generate three synthetic datasets c.f. Figure 1. The synthetic datasets contain four classes which are determined by shape (triangle or square), location (left or right) and some added noise.

For the synthetic data experiment, the hidden layer sizes for both DTW-Net and OTW-Net are set as $[1, 128, 128]$. We train both networks for 500 epochs and we plot the test-error vs training time in Figure 2. To improve the readability, we plot the minimum achieved test error, up to time $t$. Due to the computational bottleneck in DTW-net, its training time is orders of magnitude larger than OTW-Net. However, DTW-Net converges in fewer epochs, which somewhat offsets its slower time-per-epoch. In any case, OTW-Net achieves zero-error in 50 to 60 percent of the time of DTW-Net. Note that even if the training time to convergence is only reduced by 50%-60%, **the inference is reduced by a larger margin** and is many times faster (see fig. 3 middle and right panes), which shows the true advantage of linear complexity + GPU usage.

**Real data.** We compare DTW-Net with OTW-Net on real datasets from the UCR time series archive. For the real data experiment, the hidden layer sizes for OTW-Net are set as $[500, 500, 500]$ and for DTW-Net as $[100, 500, 500]$. The smaller size of the first hidden layer of DTW-Net allows training in a reasonable amount of time. In fact, we estimate the time and memory required to train DTW-Net on the UCR time series archive, consisting of more than 100 datasets, and we conclude that DTW-Net can only be trained on a handful of them in less than 24 hours. We choose the MoteStrain dataset in this collection to demonstrate the speed and accuracy of both methods.

We train both networks for 5000 epochs and we plot the test-error vs training time in Figure 3. To improve the readability, we plot the minimum achieved test error, up to time $t$. For the MoteStrain dataset, we only show the time up to 400 seconds, when OTW-Net converges to around 14% test error. In contrast, DTW-Net only achieves 25% error and its training takes around 12 hours, compared to OTW-Net which only takes 30 minutes. We conclude that the quadratic complexity of the first layer in DTW-Net makes this approach unfeasible on realistic datasets, and that one way to solve this problem is to use our proposed architecture OTW-Net.

**Speed comparison on CPU/GPU.** We turn to understanding how the higher complexity of DTW affects the performance of Deep Neural Networks with DTW-based layers and our proposed OTW-based layer replacement. To this end, we compute the time that it takes to perform one forward/backward pass through the networks, for an input of increasing dimension.

In CPU, we compare our OTW-Net against DTW-Net. DTW does not have a GPU implementation available, so instead we use the state-of-the-art GPU implementation of Soft-DTW [17]. The results are plotted in Figure 3, in the middle and right panes. As expected we observe a stark difference

**Fig. 1**. The three synthetic labeled datasets considered, consisting of 4 different classes determined by a combination of shape (square/triangle) and time shift. Each color corresponds to a different class. 4 samples of each class are shown.



**Fig. 2**. Test error vs Wall clock time (in seconds) of DTW-Net and OTW-Net, trained on the synthetic datasets proposed.



**Fig. 3**. Left: Test error vs Wall clock time of DTW-Net and OTW-Net, trained on the MoteStrain dataset from the UCR time series archive. Center: Wall clock time of a forward/backward pass over the network (in CPU) as a function of the size of the input. Right: Wall clock time of a forward/backward pass over the network (in GPU) as a function of the size of the input.

in time, and our OTW-Net runs considerably faster than the DTW-based counterparts, both in CPU and GPU. This illustrates why the DTW based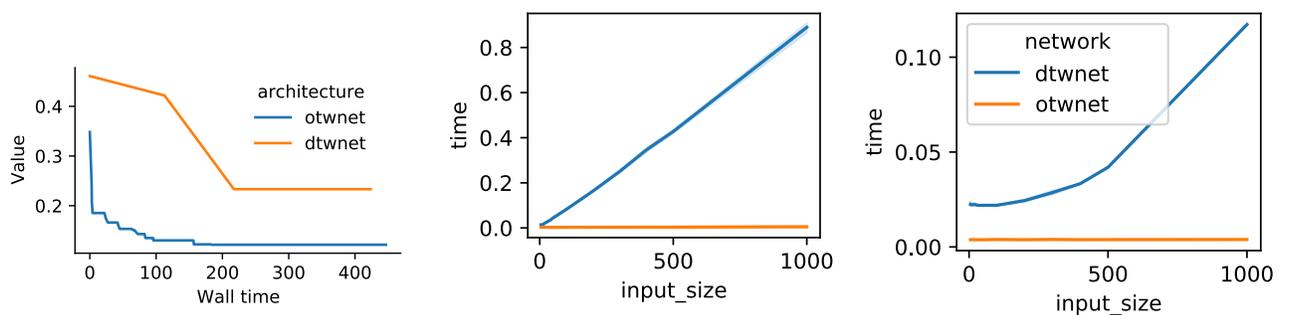 layers are not a feasible alternative already for moderate dimensions and faster alternatives, like our proposed method, have to be developed if practitioners are to adopt such kind of architectures.

## 5. CONCLUSION

We have introduced OTW, a distance for time-series that overcomes the main limitations of DTW while retaining or improving its performance in downstream tasks. In this way, it opens the path to the use of time-series distances inside Deep Learning pipelines, thanks to its easy GPU implementation and the absence of computational bottlenecks. One of the limitations of our work is the constraint that the time-series be one-dimensional. We leave possible extensions to the multivariate case as a promising direction of future research. Finally, we recall that despite the improvements, the purpose is not to completely replace DTW: for many datasets it is still the best performing distance. Rather, OTW is a complementary distance that should be preferred when there are resource constraints.

# Acknowledgements

## 6. REFERENCES

[1] Evanthia E. Tripoliti, Theofilos G. Papadopoulos, Georgia S. Karanasiou, Katerina K. Naka, and Dimitrios I. Fotiadis, "Heart failure: Diagnosis, severity estimation and prediction of adverse events through machine learning techniques," *Computational and Structural Biotechnology Journal*, vol. 15, pp. 26–47, 2017.

[2] Jeff Donahue, Sander Dieleman, Mikolaj Binkowski, Erich Elsen, and Karen Simonyan, "End-to-end adversarial text-to-speech," in *International Conference on Learning Representations*, 2021.

[3] Pat Jangyodsuk, Christopher Conly, and Vassilis Athitsos, "Sign language recognition using dynamic time warping and hand shape distance based on histogram of oriented gradient features," in *Proceedings of the 7th International Conference on PErvasive Technologies Related to Assistive Environments*, New York, NY, USA, 2014, PETRA '14, Association for Computing Machinery.

[4] Armir Bujari, Bogdan Licar, and Claudio E. Palazzi, "Movement pattern recognition through smartphone's accelerometer," in *2012 IEEE Consumer Communications and Networking Conference (CCNC)*, 2012, pp. 502–506.

[5] H. Sakoe and S. Chiba, "Dynamic programming algorithm optimization for spoken word recognition," *IEEE Transactions on Acoustics, Speech, and Signal Processing*, vol. 26, no. 1, pp. 43–49, 1978.

[6] Young-Seon Jeong, Myong K. Jeong, and Olufemi A. Omitaomu, "Weighted dynamic time warping for time series classification," *Pattern Recognition*, vol. 44, no. 9, pp. 2231–2240, 2011, Computer Analysis of Images and Patterns.

[7] Jiaping Zhao and Laurent Itti, "Shapedtw," *Pattern Recogn.*, vol. 74, no. C, pp. 171–184, feb 2018.

[8] Eamonn Keogh, "Chapter 36 - exact indexing of dynamic time warping," in *VLDB '02: Proceedings of the 28th International Conference on Very Large Databases*, Philip A. Bernstein, Yannis E. Ioannidis, Raghu Ramakrishnan, and Dimitris Papadias, Eds., pp. 406–417. Morgan Kaufmann, San Francisco, 2002.

[9] Marco Cuturi and Mathieu Blondel, "Soft-dtw: A differentiable loss function for time-series," in *Proceedings of the 34th International Conference on Machine Learning - Volume 70*. 2017, ICML'17, p. 894–903, JMLR.org.

[10] Xingyu Cai, Tingyang Xu, Jinfeng Yi, Junzhou Huang, and Sanguthevar Rajasekaran, "Dtwnet: a dynamic time warping network," in *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, Eds. 2019, vol. 32, Curran Associates, Inc.

[11] Amir Abboud, Arturs Backurs, and Virginia Vassilevska Williams, "Tight hardness results for lcs and other sequence similarity measures," in *2015 IEEE 56th Annual Symposium on Foundations of Computer Science*, 2015, pp. 59–78.

[12] Stan Salvador and Philip Ka-Fai Chan, "Fastdtw: Toward accurate dynamic time warping in linear time and space," 2004.

[13] Renjie Wu and Eamonn J. Keogh, "Fastdtw is approximate and generally slower than the algorithm it approximates," *ArXiv*, vol. abs/2003.11246, 2020.

---

[2] https://www.gofundme.com/f/ help-me-attend-icassp-2023

[14] Vincent LE GUEN and Nicolas THOME, "Shape and time distortion loss for training deep time series forecasting models," in *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, Eds. 2019, vol. 32, Curran Associates, Inc.

[15] C. Chang, De-An Huang, Yanan Sui, Li Fei-Fei, and Juan Carlos Niebles, "D3tw: Discriminative differentiable dynamic time warping for weakly supervised action alignment and segmentation," *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 3541–3550, 2019.

[16] Bertil Schmidt and Christian Hundt, "cudtw++: Ultrafast dynamic time warping on cuda-enabled gpus," in *Euro-Par 2020: Parallel Processing*, Maciej Malawski and Krzysztof Rzadca, Eds., Cham, 2020, pp. 597–612, Springer International Publishing.

[17] Mehran Maghoumi, Eugene M. Taranta II, and Joseph J. LaViola Jr, "DeepNAG: Deep Non-Adversarial Gesture Generation," 2020.

[18] Keon Lee, "Soft-dtw-loss," https://github.com/keonlee9420/Soft-DTW-Loss, 2021.

[19] Enrique Vidal Ruiz, "An algorithm for finding nearest neighbours in (approximately) constant average time," *Pattern Recognition Letters*, vol. 4, no. 3, pp. 145–157, 1986.

[20] Gabriel Peyré and Marco Cuturi, "Computational optimal transport: With applications to data science," *Foundations and Trends® in Machine Learning*, vol. 11, no. 5-6, pp. 355–607, 2019.

[21] Xiang Liu, Naiqi Li, and Shu-Tao Xia, "Gdtw: A novel differentiable dtw loss for time series tasks," in *ICASSP 2021 - 2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2021, pp. 2860–2864.

[22] Pierre-François Marteau, "Time warp edit distance with stiffness adjustment for time series matching," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 31, no. 2, pp. 306–318, 2009.

[23] Lei Chen and Raymond T. Ng, "On the marriage of lp-norms and edit distance," in *VLDB*, 2004.

[24] Lei Chen, M. Tamer Özsu, and Vincent Oria, "Robust and fast similarity search for moving object trajectories," in *Proceedings of the 2005 ACM SIGMOD International Conference on Management of Data*, New York, NY, USA, 2005, SIGMOD '05, p. 491–502, Association for Computing Machinery.

[25] Mathieu Blondel, Arthur Mensch, and Jean-Philippe Vert, "Differentiable divergences between time series," *CoRR*, vol. abs/2010.08354, 2020.

[26] Junfu Pu, Wengang Zhou, and Houqiang Li, "Iterative alignment network for continuous sign language recognition," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.

[27] Vincent LE GUEN and Nicolas THOME, "Probabilistic time series forecasting with shape and temporal diversity," in *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, Eds. 2020, vol. 33, pp. 4427–4440, Curran Associates, Inc.

[28] Brian Kenji Iwana, Volkmar Frinken, and Seiichi Uchida, "Dtw-nn: A novel neural network for time series recognition using dynamic alignment between inputs and weights," *Knowledge-Based Systems*, vol. 188, Jan. 2020.

[29] Jan Zelinka and Jakub Kanis, "Neural sign language synthesis: Words are our glosses," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, March 2020.

[30] Hicham Janati, Marco Cuturi, and Alexandre Gramfort, "Spatio-temporal alignments: Optimal transport through space and time," in *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, Silvia Chiappa and Roberto Calandra, Eds. 26–28 Aug 2020, vol. 108 of *Proceedings of Machine Learning Research*, pp. 1695–1704, PMLR.

[31] Zheng Zhang, Ping Tang, and Thomas Corpetti, "Time adaptive optimal transport: A framework of time series similarity measure," *IEEE Access*, vol. 8, pp. 149764–149774, 2020.

[32] Marco Cuturi, "Sinkhorn distances: Lightspeed computation of optimal transport," in *Advances in Neural Information Processing Systems*, C.J. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K.Q. Weinberger, Eds. 2013, vol. 26, Curran Associates, Inc.

[33] S. S. Vallender, "Calculation of the wasserstein distance between probability distributions on the line," *Theory of Probability and Its Applications*, vol. 18, pp. 435–435, 1974.

[34] Bing Su and Gang Hua, "Order-preserving optimal transport for distances between sequences," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 41, no. 12, pp. 2961–2974, 2019.

[35] Kevin Guittet, "Extended Kantorovich norms : a tool for optimization," Research Report RR-4402, INRIA, 2002.

[36] Chotirat Ratanamahatana and Eamonn J. Keogh, "Making time-series classification more accurate using learned constraints," in *SDM*, 2004.

[37] Marcel R. Ackermann, Johannes Blömer, Daniel Kuntze, and Christian Sohler, "Analysis of agglomerative clustering," *Algorithmica*, vol. 69, pp. 184–215, 2012.

**Table 4**. Test error of 1-NN classification using Learned DTW or Learned OT.

| Name | Type | DTW error | OT error |
|---|---|---|---|
| ACSF1 | Device | 0.38 | **0.34 ± 0.06** |
| Adiac | Image | 0.39 | **0.39 ± 0.01** |
| AllGestureWiimoteX | Sensor | 0.28 | 0.90 ± 0.00 |
| AllGestureWiimoteY | Sensor | 0.27 | 0.90 ± 0.00 |
| AllGestureWiimoteZ | Sensor | 0.35 | 0.90 ± 0.00 |
| ArrowHead | Image | 0.20 | **0.23 ± 0.03** |
| Beef | Spectro | 0.33 | 0.55 ± 0.04 |
| BeetleFly | Image | 0.30 | **0.26 ± 0.06** |
| BirdChicken | Image | 0.30 | **0.29 ± 0.08** |
| Car | Sensor | 0.23 | **0.27 ± 0.06** |
| Chinatown | Traffic | 0.05 | **0.07 ± 0.03** |
| ChlorineConcentration | Sensor | 0.35 | 0.39 ± 0.02 |
| Coffee | Spectro | 0.00 | 0.04 ± 0.03 |
| Computers | Device | 0.38 | **0.40 ± 0.02** |
| Crop | Image | 0.29 | **0.28 ± 0.01** |
| DiatomSizeReduction | Image | 0.07 | **0.06 ± 0.03** |
| DistalPhalanxOutlineAgeGroup | Image | 0.37 | **0.23 ± 0.03** |
| DistalPhalanxOutlineCorrect | Image | 0.28 | **0.23 ± 0.01** |
| DistalPhalanxTW | Image | 0.37 | **0.29 ± 0.04** |
| DodgerLoopDay | Sensor | 0.41 | 0.85 ± 0.03 |
| DodgerLoopGame | Sensor | 0.07 | 0.43 ± 0.09 |
| DodgerLoopWeekend | Sensor | 0.02 | 0.35 ± 0.09 |
| ECG200 | Medical Data | 0.12 | 0.16 ± 0.04 |
| ECG5000 | Medical Data | 0.07 | **0.08 ± 0.01** |
| ECGFiveDays | Medical Data | 0.20 | **0.19 ± 0.03** |
| EOGHorizontalSignal | Medical Data | 0.52 | **0.38 ± 0.02** |
| EOGVerticalSignal | Medical Data | 0.52 | **0.49 ± 0.01** |
| Earthquakes | Sensor | 0.27 | 0.32 ± 0.03 |
| ElectricDevices | Device | 0.38 | **0.25 ± 0.00** |
| EthanolLevel | Spectro | 0.72 | **0.71 ± 0.01** |
| FaceAll | Image | 0.19 | **0.05 ± 0.01** |
| FaceFour | Image | 0.11 | 0.17 ± 0.05 |
| FacesUCR | Image | 0.09 | 0.11 ± 0.01 |
| FiftyWords | Image | 0.24 | 0.31 ± 0.01 |
| Fish | Image | 0.15 | 0.21 ± 0.02 |
| FordA | Sensor | 0.31 | **0.27 ± 0.01** |
| FordB | Sensor | 0.39 | **0.30 ± 0.01** |
| FreezerRegularTrain | Sensor | 0.09 | **0.08 ± 0.02** |
| FreezerSmallTrain | Sensor | 0.32 | **0.24 ± 0.02** |

**Table 5**. Test error of 1-NN classification using Learned DTW or Learned OT (continued).

| Name | Type | DTW error | OT error |
|---|---|---|---|
| Fungi | Medical Data | 0.18 | $0.41 \pm 0.05$ |
| GesturePebbleZ1 | Sensor | 0.17 | $0.84 \pm 0.01$ |
| GesturePebbleZ2 | Sensor | 0.22 | $0.83 \pm 0.02$ |
| Ham | Spectro | 0.40 | $\mathbf{0.28 \pm 0.03}$ |
| HandOutlines | Image | 0.14 | $\mathbf{0.15 \pm 0.01}$ |
| Herring | Image | 0.47 | $\mathbf{0.49 \pm 0.04}$ |
| HouseTwenty | Device | 0.06 | $0.25 \pm 0.04$ |
| InsectEPGRegularTrain | Medical Data | 0.17 | $\mathbf{0.00 \pm 0.00}$ |
| InsectEPGSmallTrain | Medical Data | 0.31 | $\mathbf{0.00 \pm 0.00}$ |
| InsectWingbeatSound | Sensor | 0.42 | $0.45 \pm 0.01$ |
| ItalyPowerDemand | Sensor | 0.04 | $0.07 \pm 0.02$ |
| LargeKitchenAppliances | Device | 0.21 | $0.41 \pm 0.01$ |
| Lightning2 | Sensor | 0.13 | $0.22 \pm 0.05$ |
| Lightning7 | Sensor | 0.29 | $0.35 \pm 0.04$ |
| Meat | Spectro | 0.07 | $\mathbf{0.01 \pm 0.01}$ |
| MedicalImages | Image | 0.25 | $0.29 \pm 0.02$ |
| MelbournePedestrian | Traffic | 0.18 | $\mathbf{0.05 \pm 0.00}$ |
| MiddlePhalanxOutlineAgeGroup | Image | 0.48 | $\mathbf{0.31 \pm 0.03}$ |
| MiddlePhalanxOutlineCorrect | Image | 0.23 | $\mathbf{0.22 \pm 0.01}$ |
| MiddlePhalanxTW | Image | 0.49 | $\mathbf{0.45 \pm 0.03}$ |
| MixedShapesSmallTrain | Image | 0.17 | $\mathbf{0.19 \pm 0.02}$ |
| MoteStrain | Sensor | 0.13 | $\mathbf{0.15 \pm 0.03}$ |
| NonInvasiveFetalECGThorax1 | Medical Data | 0.19 | $\mathbf{0.18 \pm 0.01}$ |
| OSULeaf | Image | 0.39 | $\mathbf{0.42 \pm 0.04}$ |
| OliveOil | Spectro | 0.13 | $\mathbf{0.14 \pm 0.03}$ |
| PLAID | Device | 0.17 | $0.88 \pm 0.04$ |
| PhalangesOutlinesCorrect | Image | 0.24 | $\mathbf{0.23 \pm 0.01}$ |
| Phoneme | Sensor | 0.77 | $0.88 \pm 0.01$ |
| PickupGestureWiimoteZ | Sensor | 0.34 | $0.90 \pm 0.02$ |
| PigAirwayPressure | Medical Data | 0.90 | $\mathbf{0.85 \pm 0.01}$ |
| PigArtPressure | Medical Data | 0.80 | $\mathbf{0.78 \pm 0.02}$ |
| PigCVP | Medical Data | 0.84 | $0.86 \pm 0.02$ |
| Plane | Sensor | 0.00 | $0.05 \pm 0.01$ |
| PowerCons | Power | 0.08 | $\mathbf{0.04 \pm 0.01}$ |
| ProximalPhalanxOutlineAgeGroup | Image | 0.21 | $0.24 \pm 0.02$ |
| ProximalPhalanxOutlineCorrect | Image | 0.21 | $\mathbf{0.21 \pm 0.02}$ |
| ProximalPhalanxTW | Image | 0.24 | $0.27 \pm 0.02$ |
| RefrigerationDevices | Device | 0.56 | $\mathbf{0.49 \pm 0.02}$ |
| ScreenType | Device | 0.59 | $\mathbf{0.57 \pm 0.01}$ |
| ShakeGestureWiimoteZ | Sensor | 0.16 | $0.90 \pm 0.01$ |
| ShapesAll | Image | 0.20 | $0.27 \pm 0.02$ |
| SmallKitchenAppliances | Device | 0.33 | $0.37 \pm 0.02$ |
| SonyAIBORobotSurface1 | Sensor | 0.30 | $\mathbf{0.14 \pm 0.03}$ |
| SonyAIBORobotSurface2 | Sensor | 0.14 | $\mathbf{0.14 \pm 0.02}$ |
| Strawberry | Spectro | 0.05 | $\mathbf{0.04 \pm 0.01}$ |
| SwedishLeaf | Image | 0.15 | $0.23 \pm 0.01$ |
| Symbols | Image | 0.06 | $0.14 \pm 0.03$ |
| Trace | Sensor | 0.01 | $0.21 \pm 0.03$ |
| TwoLeadECG | Medical Data | 0.13 | $0.31 \pm 0.06$ |
| Wafer | Sensor | 0.00 | $\mathbf{0.01 \pm 0.00}$ |
| Wine | Spectro | 0.39 | $\mathbf{0.13 \pm 0.07}$ |
| WordSynonyms | Image | 0.26 | $0.36 \pm 0.02$ |
| Yoga | Image | 0.16 | $0.17 \pm 0.01$ |

# B. HIERARCHICAL CLUSTERING COMPARISON

**Table 6**. RI and time taken for DTW- and OT-based Hierarchical clustering methods on datasets in the UCR time-series benchmark collection. Noticeable improvements in boldface.

| Name | Type | DTW RI | DTW time | OT RI | OT time |
|------|------|--------|----------|-------|---------|
| Adiac | Image | 0.79 | 220.82 | 0.61 | **32.83** |
| ArrowHead | Image | 0.35 | 63.86 | **0.36** | **14.33** |
| Beef | Spectro | 0.42 | 14.85 | **0.60** | **2.56** |
| BeetleFly | Image | 0.59 | 12.97 | **0.62** | **1.68** |
| BirdChicken | Image | 0.50 | 13.07 | **0.51** | **1.93** |
| Car | Sensor | 0.48 | 94.91 | **0.66** | **11.69** |
| Chinatown | Traffic | 0.59 | 42.02 | **0.61** | **5.93** |
| ChlorineConcentration | Sensor | 0.42 | 200.71 | **0.50** | **32.12** |
| Coffee | Spectro | 0.50 | 5.59 | **0.57** | **1.64** |
| Computers | Device | 0.50 | 2,495.25 | 0.50 | **94.26** |
| Crop | Image | 0.87 | 90.27 | **0.92** | **17.05** |
| DiatomSizeReduction | Image | 0.30 | 250.90 | **0.31** | **31.80** |
| DistalPhalanxOutlineAgeGroup | Image | 0.71 | 103.39 | 0.71 | **22.53** |
| DistalPhalanxOutlineCorrect | Image | 0.53 | 103.47 | 0.51 | **22.43** |
| DistalPhalanxTW | Image | 0.79 | 102.95 | **0.87** | **22.32** |
| ECG200 | Medical Data | 0.54 | 18.45 | **0.61** | **5.63** |
| ECG5000 | Medical Data | 0.87 | 161.50 | **0.89** | **27.41** |
| ECGFiveDays | Medical Data | 0.50 | 157.94 | **0.51** | **26.80** |
| Earthquakes | Sensor | 0.51 | 1,594.13 | **0.68** | **59.17** |
| ElectricDevices | Device | 0.56 | 114.64 | 0.53 | **23.78** |
| FaceAll | Image | 0.57 | 147.43 | **0.70** | **27.11** |
| FaceFour | Image | 0.54 | 30.19 | **0.75** | **6.03** |
| FacesUCR | Image | 0.57 | 149.10 | **0.78** | **27.17** |
| FiftyWords | Image | 0.92 | 400.92 | **0.95** | **41.76** |
| Fish | Image | 0.17 | 486.88 | **0.64** | **39.02** |
| FordA | Sensor | 0.50 | 1,220.33 | 0.50 | **66.80** |
| FordB | Sensor | 0.51 | 1,222.25 | 0.51 | **66.06** |
| FreezerRegularTrain | Sensor | 0.51 | 474.91 | **0.56** | **46.98** |
| FreezerSmallTrain | Sensor | 0.50 | 480.85 | **0.66** | **48.18** |
| Fungi | Medical Data | 0.96 | 43.36 | **0.98** | **12.53** |
| Ham | Spectro | 0.50 | 161.65 | 0.50 | **19.27** |
| Herring | Image | 0.51 | 121.49 | **0.52** | **11.87** |
| InsectEPGRegularTrain | Medical Data | 1.00 | 669.58 | 1.00 | **37.38** |
| InsectEPGSmallTrain | Medical Data | 1.00 | 625.76 | 1.00 | **32.23** |
| InsectWingbeatSound | Sensor | 0.68 | 426.09 | **0.85** | **40.46** |
| ItalyPowerDemand | Sensor | 0.50 | 78.79 | 0.50 | **12.79** |
| LargeKitchenAppliances | Device | 0.34 | 2,817.32 | **0.57** | **91.34** |
| Lightning2 | Sensor | 0.50 | 114.12 | **0.61** | **12.99** |
| Lightning7 | Sensor | 0.62 | 43.17 | **0.82** | **8.52** |
| Meat | Spectro | 0.77 | 62.18 | **0.77** | **10.92** |
| MedicalImages | Image | 0.60 | 119.65 | 0.57 | **24.08** |
| MelbournePedestrian | Traffic | 0.64 | 79.33 | **0.83** | **12.70** |
| MiddlePhalanxOutlineAgeGroup | Image | 0.70 | 103.33 | 0.70 | **22.24** |
| MiddlePhalanxOutlineCorrect | Image | 0.52 | 103.52 | 0.52 | **22.27** |
| MiddlePhalanxTW | Image | 0.80 | 102.36 | **0.81** | **22.51** |
| MoteStrain | Sensor | 0.50 | 105.35 | **0.68** | **22.53** |
| NonInvasiveFetalECGThorax1 | Medical Data | 0.83 | 3,068.55 | **0.93** | **95.15** |
| NonInvasiveFetalECGThorax2 | Medical Data | 0.92 | 3,112.45 | **0.94** | **94.65** |
| OSULeaf | Image | 0.58 | 700.76 | **0.69** | **49.36** |
| OliveOil | Spectro | 0.74 | 27.54 | **0.78** | **2.93** |
| PhalangesOutlinesCorrect | Image | 0.53 | 108.40 | 0.53 | **22.23** |
| Plane | Sensor | 0.96 | 30.50 | 0.91 | **8.11** |
| PowerCons | Power | 0.50 | 89.67 | **0.52** | **19.95** |
| ProximalPhalanxOutlineAgeGroup | Image | 0.77 | 104.00 | **0.78** | **22.52** |
| ProximalPhalanxOutlineCorrect | Image | 0.53 | 107.66 | **0.56** | **22.33** |
| ProximalPhalanxTW | Image | 0.85 | 107.30 | **0.87** | **22.04** |
| RefrigerationDevices | Device | 0.34 | 2,757.08 | **0.44** | **91.89** |
| ScreenType | Device | 0.34 | 2,845.35 | **0.52** | **91.20** |
| ShapesAll | Image | 0.83 | 1,979.40 | **0.94** | **67.31** |
| SmallKitchenAppliances | Device | 0.34 | 2,646.25 | **0.37** | **91.76** |
| SonyAIBORobotSurface1 | Sensor | 0.50 | 96.16 | **0.75** | **21.52** |
| SonyAIBORobotSurface2 | Sensor | 0.52 | 93.90 | **0.67** | **21.03** |
| Strawberry | Spectro | 0.52 | 327.31 | **0.54** | **38.26** |
| SwedishLeaf | Image | 0.37 | 167.14 | **0.52** | **26.56** |
| Symbols | Image | 0.87 | 778.76 | 0.87 | **60.72** |
| Trace | Sensor | 0.87 | 65.96 | 0.76 | **14.12** |
| TwoLeadECG | Medical Data | 0.50 | 104.65 | **0.51** | **22.13** |
| Wafer | Sensor | 0.53 | 183.66 | **0.80** | **29.94** |
| Wine | Spectro | 0.50 | 16.15 | 0.50 | **3.76** |
| WordSynonyms | Image | 0.85 | 407.55 | **0.89** | **41.29** |
| Yoga | Image | 0.50 | 865.61 | **0.51** | **58.34** |
| Total | - | 0.60 | 521.74 | **0.67** | **32.15** |

## C. DEFERRED PROOFS

### C.1. Proof of theorem 1

First we will prove the following intermediate result:

**Lemma 3.** *Let $a, b$ be two positive sequences. Without loss of generality assume $\sum_{j=1}^{n} a_j \leq \sum_{j=1}^{n} b_j$ and let $A(i) = \sum_{j=1}^{i} a_j, B(i) = \sum_{j=1}^{i} b_j$ be the partial sums of the sequences. It holds that:*

$$\widehat{W}_m(a, b) \leq \begin{cases} \min_c W(a, c) + m|A(n) - B(n)| \\ \text{subject to } 0 \leq c \leq b, \sum_{j=1}^{n} c_j = \sum_{j=1}^{n} a_j \end{cases} \tag{12}$$

*Where $W(a, c)$ is the original (eq. (1)) Optimal Transport distance (note that $a, b$ sum up to the same value).*

*Proof.* Let $c$ such that $0 \leq c \leq b$ and $\sum_{j=1}^{n} c_j = \sum_{j=1}^{n} a_j$, Let $T^\star$ be the optimal transport map between $a$ and $c$. That means $T^\star \geq 0, T^\star \mathbb{1} = a, \mathbb{1}^\top T^\star = c^\top$. We now build a feasible transport map for the unbalanced problem eq. (5) based on $T^\star$:

$$\widehat{T} = \begin{bmatrix} T^\star & 0 \\ b^\top - c^\top & \sum_{j=1}^{n} a_i \end{bmatrix} \tag{13}$$

then it is easy to see that $\widehat{T}\mathbb{1} = \hat{a}$ and $\mathbb{1}^\top \widehat{T} = \hat{b}$, so $\widehat{T}$ is a feasible transport map for eq. (5). Computing the objective value we obtain:

$$\begin{aligned} \widehat{W}_m(a, b) &\leq \left\langle \widehat{T}, D(m) \right\rangle \\ &= \langle T^\star, D \rangle + m \left( \sum_{j=1}^{n} b_j - \sum_{j=1}^{n} c_j \right) \\ &= \langle T^\star, D \rangle + m \left( \sum_{j=1}^{n} b_j - \sum_{j=1}^{n} a_j \right) \\ &= W(a, c) + m|A(n) - B(n)| \end{aligned} \tag{14}$$

$\square$

minimizing the right hand side we obtain:

$$\widehat{W}_m(a, b) \leq \begin{cases} \min_c W(a, c) + m|A(n) - B(n)| \\ \text{subject to } 0 \leq c \leq b, \sum_{j=1}^{n} c_j = \sum_{j=1}^{n} a_j \end{cases} \tag{15}$$

Now we will define a nonnegative sequence $c$ as follows: let $i^\star \in [n]$ be such that

$$\sum_{j=1}^{i^\star} b_j \geq \sum_{j=1}^{n} a_j, \qquad \sum_{j=1}^{i^\star - 1} b_j < \sum_{j=1}^{n} a_j \tag{16}$$

that is, $i^\star$ is exactly the index where the cumulative mass of $b$ exceeds the total mass of $a$. This index exists because without loss of generality we can assume $\sum_{j=1}^{n} a_j \leq \sum_{j=1}^{n} b_j$. Now we define

$$c_i = \begin{cases} b_i & : \text{if } i = 1, \dots, i^\star - 1 \\ \sum_{j=1}^{n} - \sum_{j=1}^{i^\star} b_j & : \text{if } i = i^\star \\ 0 & : \text{if } i > i^\star \end{cases} \tag{17}$$

then it is easy to see that $\sum_{j=1}^{n} c_i = \sum_{j=1}^{n} a_i$. Then for this choice of $c$ it holds that

$$\widehat{W}_m(a, b) \leq W(a, c) + m|A(n) - B(n)| \tag{18}$$

Now we have a formula for $W(a, c)$ according to [33]:

$$W(a, c) = \sum_{j=1}^{n} |A(j) - C(j)|, \qquad C(j) = \sum_{i=1}^{j} c_i \tag{19}$$

Now we have

$$\sum_{j=1}^{n} |A(j) - C(j)| = \sum_{j=1}^{n-1} |A(j) - C(j)| \qquad \text{(because } A(n) = C(n)\text{)}$$

$$= \sum_{j=1}^{i^\star-1} |A(j) - C(j)| + \sum_{j=i^\star}^{n-1} C(j) - A(j) \qquad \text{(because } C(j) \geq A(j) \text{ for } j \geq i^\star)$$

$$\leq \sum_{j=1}^{i^\star-1} |A(j) - C(j)| + \sum_{j=i^\star}^{n-1} B(j) - A(j) \qquad \text{(because } C(j) \leq B(j) \text{ for } j \geq i^\star) \tag{20}$$

$$= \sum_{j=1}^{i^\star-1} |A(j) - B(j)| + \sum_{j=i^\star}^{n-1} B(j) - A(j) \qquad \text{(because } C(j) = B(j) \text{ for } j < i^\star)$$

$$= \sum_{j=1}^{n-1} |A(j) - B(j)|$$

Hence, by lemma 3 we have

$$\widehat{W}_m(a, b) \leq \sum_{j=1}^{n-1} |A(j) - B(j)| + m|A(n) - B(n)| \tag{21}$$

## C.2. Proof of lemma 1

Without loss of generality we again assume that $\sum_{i=1}^{n} a_i \leq \sum_{i=1}^{n} b_i$. We will show that

$$|\widehat{W}_m(a, b') - \widehat{W}_m(a, b)| \leq t \left( \sum_{i=1}^{n} a_i \right) \tag{22}$$

To this end we will first show that

$$\widehat{W}(a, b') \leq \widehat{W}(a, b) + t \left( \sum_{i=1}^{n} a_i \right) \tag{23}$$

then by symmetry (we can swap the roles of $b$ and $b'$ and follow the same argument) we can conclude:

$$\widehat{W}(a, b) \leq \widehat{W}(a, b') + t \left( \sum_{i=1}^{n} a_i \right) \tag{24}$$

so that

$$|\widehat{W}_m(a, b') - \widehat{W}_m(a, b)| = \max \left( \widehat{W}_m(a, b') - \widehat{W}_m(a, b), \widehat{W}_m(a, b) - \widehat{W}_m(a, b') \right)$$

$$\leq t \left( \sum_{i=1}^{n} a_i \right) \tag{25}$$

Leading to the result. Going back to prove eq. (23), let $T^\star$ be an unbalanced optimal transport map between $a$ and $b$. From this map, we will construct a transport map between $a, b'$ which has cost equal to the right hand side of eq. (23). Recall that unbalanced transport maps (eq. (5)) can be split in the following block form

$$T^\star = \begin{bmatrix} T' & \tilde{a} \\ \tilde{b}^\top & c \end{bmatrix} \tag{26}$$

where $T$ has dimensions $n \times n$, $\tilde{a}$ is a column vector of length $n$, $\tilde{b}$ is a vector of length $n$ and $c$ is a scalar. This follows the structure of the matrix

$$D(m) = \begin{bmatrix} D & m\mathbb{1} \\ m\mathbb{1}^\top & 0 \end{bmatrix} \tag{27}$$

The fact that the column-sums of the transport map should be equal to $b$, (this is a constraint in the definition eq. (5)) which by assumption is such that $b_{n-t+1} = \ldots = b_n = 0$ implies that the left-most block in eq. (26) can be written as:

$$\begin{bmatrix} T' \\ \tilde{b}^\top \end{bmatrix} = \begin{bmatrix} T'' & \mathbf{0}_1 \\ \tilde{b}'^\top & \mathbf{0}_2^\top \end{bmatrix} \tag{28}$$

where $T''$ has dimensions $n \times (n-t)$, $\tilde{b}'$ is a vector of length $n-t$, $\mathbf{0}_1$ is a matrix of all-zeros with dimensions $n \times t$ and $\mathbf{0}_2$ is a zero vector of length $t$. The new transport map between $a$ and $b'$ is defined as

$$\widehat{T} = \begin{bmatrix} T''' & \tilde{a} \\ \tilde{b}''^\top & c \end{bmatrix}, \qquad \begin{bmatrix} T''' \\ \tilde{b}''^\top \end{bmatrix} = \begin{bmatrix} \mathbf{0}_1 & T'' \\ \mathbf{0}_2 & \tilde{b}'^\top \end{bmatrix} \tag{29}$$

that is, the new transport map is obtained by shifting to the right the left block in the decomposition eq. (26) by $t$ units. The following observations hold: (1) after this transformation there is no change in the row-sums and (2) because $b'$ is obtained precisely by shifting $b$ to the right by $t$ units, the new unbalanced transport map satisfies

$$\mathbb{1}^\top \widehat{T} = \begin{bmatrix} b' | \sum_{i=1}^n a_i \end{bmatrix} = \hat{b}', \qquad \widehat{T}\mathbb{1} = \begin{bmatrix} a | \sum_{i=1}^n b_i \end{bmatrix} = \begin{bmatrix} a | \sum_{i=1}^n b_i' \end{bmatrix} = \hat{a}, \tag{30}$$

Hence it is indeed the case that $\widehat{T}$ is an unbalanced transport map between $a$ and $b'$. We now compute the objective value:

$$\widehat{W}(a, b') \leq \left\langle \widehat{T}, D(m) \right\rangle = \left\langle \begin{bmatrix} T''' & \tilde{a} \\ \tilde{b}''^\top & c \end{bmatrix}, \begin{bmatrix} D & m\mathbb{1} \\ m\mathbb{1}^\top & 0 \end{bmatrix} \right\rangle$$

$$= \langle T''', D \rangle + m \sum_{i=1}^n \tilde{a}_i + m \sum_{i=1}^n \tilde{b}''_i \tag{31}$$

$$= \langle T''', D \rangle + m \sum_{i=1}^n \tilde{a}_i + m \sum_{i=1}^n \tilde{b}_i \qquad (\tilde{b}'' \text{ is a shifted version of } \tilde{b})$$

We will now show the following:

$$\langle T''', D \rangle \leq \langle T', D \rangle + t \left( \sum_{i=1}^n a_i \right) \tag{32}$$

$$\langle T''', D\rangle = \sum_{j=1}^{n}\sum_{i=1}^{n} T'''_{i,j} D_{i,j} = \sum_{j=t+1}^{n}\sum_{i=1}^{n} T''_{i,j-t}|i-j|$$

$$= \sum_{j=t+1}^{n}\sum_{i=1}^{n} T''_{i,j-t}|i-(j-t)+t|$$

$$\leq \sum_{j=t+1}^{n}\sum_{i=1}^{n} T''_{i,j-t}|i-(j-t)| + t\sum_{j=t+1}^{n}\sum_{i=1}^{n} T''_{i,j-t}$$

$$= \sum_{j=1}^{n-t}\sum_{i=1}^{n} T''_{i,j}|i-j| + t\sum_{j=t+1}^{n}\sum_{i=1}^{n} T''_{i,j-t} \tag{33}$$

$$= \sum_{j=1}^{n}\sum_{i=1}^{n} T'_{i,j}|i-j| + t\sum_{j=t+1}^{n}\sum_{i=1}^{n} T''_{i,j-t}$$

$$\leq \langle T', D\rangle + t\mathbb{1}^{\top} T'\mathbb{1}$$

$$\leq \langle T', D\rangle + t\mathbb{1}^{\top}(T'\mathbb{1}+\tilde{a}) \qquad \text{(because } \tilde{a} \text{ is positive)}$$

$$= \langle T', D\rangle + t\mathbb{1}^{\top} a \qquad \text{(row-sums constraint)}$$

$$= \langle T', D\rangle + t\left(\sum_{i=1}^{n} a_i\right)$$

which shows that eq. (32) holds. Now notice

$$\widehat{W}_m(a,b) = \langle T^{\star}, D(m)\rangle = \left\langle \begin{bmatrix} T' & \tilde{a} \\ \tilde{b}^{\top} & c \end{bmatrix}, \begin{bmatrix} D & m\mathbb{1} \\ m\mathbb{1}^{\top} & 0 \end{bmatrix} \right\rangle$$
$$= \langle T', D\rangle + m\sum_{i=1}^{n}\tilde{a}_i + m\sum_{i=1}^{n}\tilde{b}_i \tag{34}$$

so that combining eqs. (31), (32) and (34) we have:

$$\widehat{W}(a,b') \leq \langle T''', D\rangle + m\sum_{i=1}^{n}\tilde{a}_i + m\sum_{i=1}^{n}\tilde{b}_i$$

$$\leq \langle T', D\rangle + t\left(\sum_{i=1}^{n} a_i\right) + m\sum_{i=1}^{n}\tilde{a}_i + m\sum_{i=1}^{n}\tilde{b}_i \tag{35}$$

$$= \widehat{W}_m(a,b) + t\left(\sum_{i=1}^{n} a_i\right)$$

we have proven that eq. (23) holds, and the result follows as explained in eqs. (24) and (25)

### C.3. Proof of lemma 2

We assume $m=1$. When $s=1$ we have

$$A_1(i) := \sum_{j=1}^{i} a_j - \sum_{j=1}^{i-1} a_j = a_i, \qquad B_s(i) := \sum_{j=1}^{i} b_j - \sum_{j=1}^{i-1} b_j = b_i \tag{36}$$

hence

$$\text{OTW}_{1,1}(a,b) := |A_1(n) - B_1(n)| + \sum_{i=1}^{n-1}|A_1(i) - B_1(i)|$$

$$= |a_n - b_n| + \sum_{i=1}^{n-1}|a_i - b_i| = \|a - b\|_1 \tag{37}$$

Now when $s = n$ we have

$$A_n(i) := \sum_{j=1}^{i} a_j - \sum_{j=1}^{i-n} a_j = \sum_{j=1}^{i} a_j = A(i), \qquad B_n(i) := \sum_{j=1}^{i} b_j - \sum_{j=1}^{i-n} b_j = \sum_{j=1}^{i} b_j = B(i) \tag{38}$$

where this is due by convention: $i \leq n$ so that $i - n \leq 0$, and the sums from $j = 1$ to $j = i - n \leq 0$ is an empty sum (so it has zero value). Then, we get

$$\mathrm{OTW}_{1,n}(a, b) := |A_n(n) - B_n(n)| + \sum_{i=1}^{n-1} |A_n(i) - B_n(i)|$$

$$= |A(n) - B(n)| + \sum_{i=1}^{n-1} |A(i) - B(i)| = \mathrm{OTW}_m(a, b) \tag{39}$$