# A Wasserstein Subsequence Kernel for Time Series

**Christian Bock, Matteo Togninalli**, Elisabetta Ghisu, Thomas Gumbsch, Bastian Rieck, Karsten Borgwardt

Department of Biosystems Science and Engineering
Machine Learning and Computational Biology Group

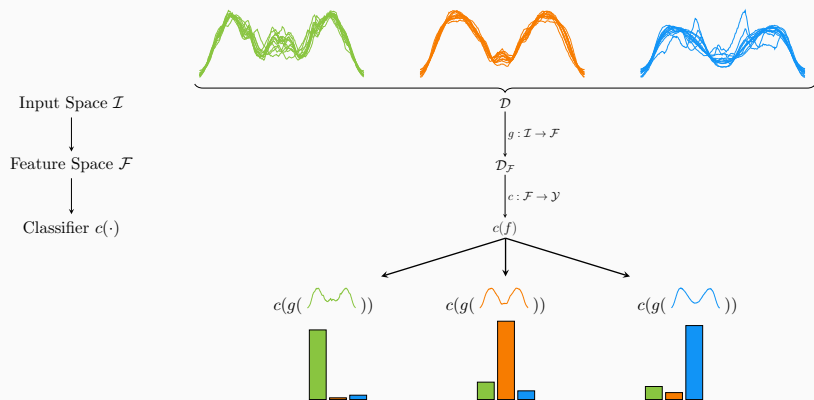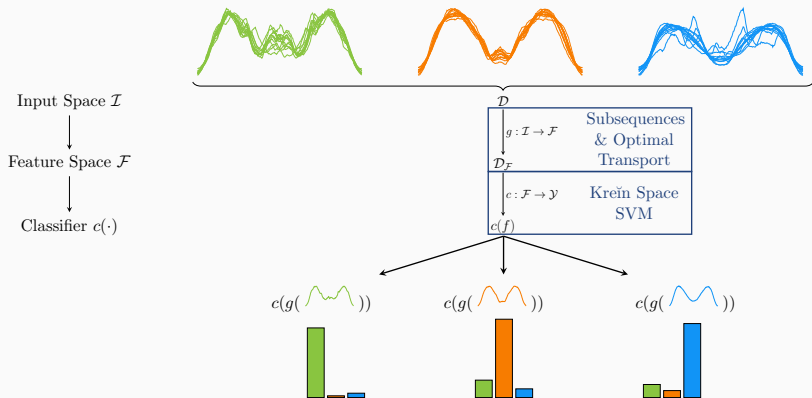ETHa ETHb ETHc ETHd ETHe ETHf ETHg ETHh ETHi

## Table of contents (internal)

# Motivation

Input Space $\mathcal{I}$

Feature Space $\mathcal{F}$

Classifier $c(\cdot)$

$\mathcal{D}$

$g : \mathcal{I} \to \mathcal{F}$ — Subsequences & Optimal Transport

$\mathcal{D}_{\mathcal{F}}$

$c : \mathcal{F} \to \mathcal{Y}$ — Kreĭn Space SVM

$c(f)$

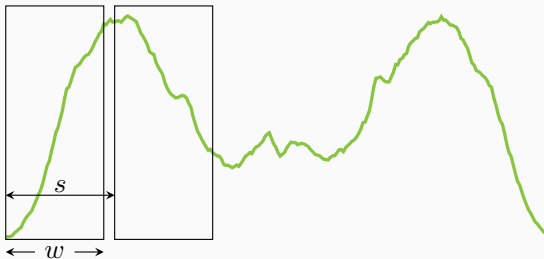$c(g(\quad))$ $\quad$ $c(g(\quad))$ $\quad$ $c(g(\quad))$

# Feature Space: Subsequences and Optimal Transport

# Subsequence Extraction
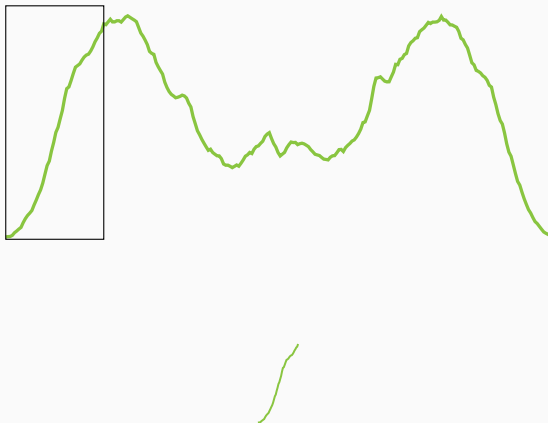
## The Sliding Window Approach

- Window Size $w$
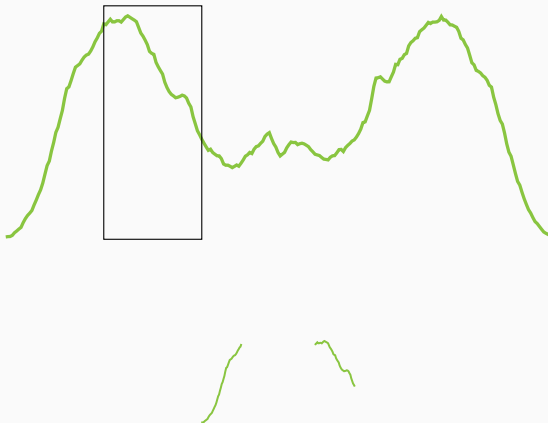- Stride $s$

# Subsequence Extraction

## The Sliding Window Approach

- Window Size $w$
- Stride $s$

# Subsequence Extraction

## The Sliding Window Approach
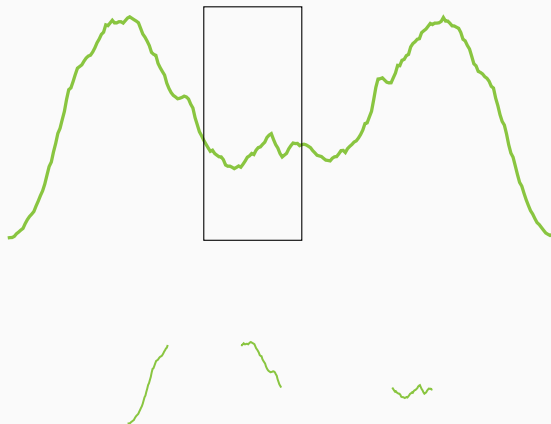
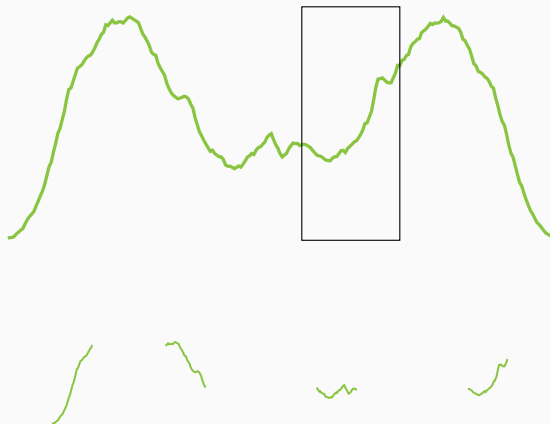- Window Size $w$
- Stride $s$

## The Sliding Window Approach

- Window Size $w$
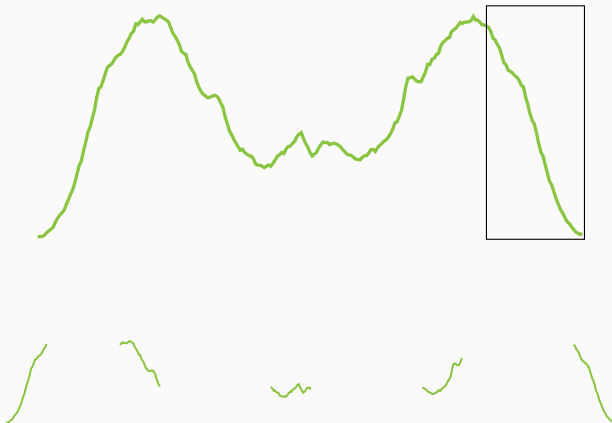- Stride $s$

# Subsequence Extraction

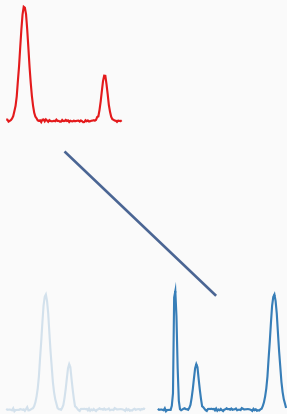## The Sliding Window Approach

- Window Size $w$
- Stride $s$

## The Sliding Window Approach

- Window Size $w$
- Stride $s$

# Optimal Transport and the Wasserstein Distance
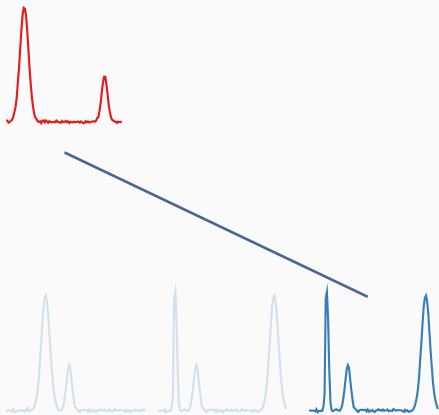
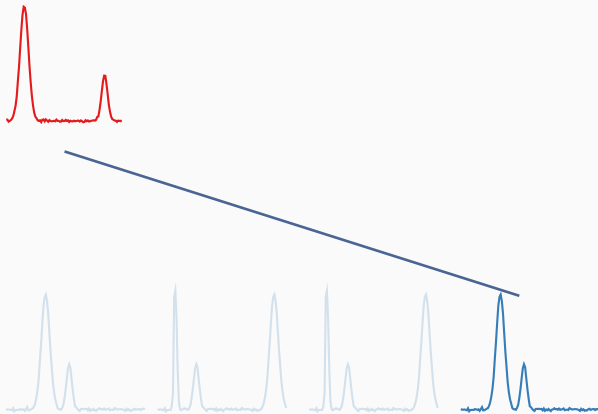# A Wasserstein Kernel for Time Series

## From Distance to Kernels

**Definition (Wasserstein time series kernel)**
Let $T_i$ and $T_j$ be two time series, and $s_{i1}, \ldots, s_{iU}$ as well as $s_{j1}, \ldots, s_{jV}$ be their respective subsequences. Moreover, let $D$ be a $U \times V$ matrix that contains the pairwise distances of all subsequences, such that $D_{uv} := \text{dist}(s_{iu}, s_{jv})$, where $\text{dist}(\cdot, \cdot)$ denotes the usual Euclidean distance. The optimisation problem

$$W_1(T_i, T_j) := \min_{P \in \Gamma(T_i, T_j)} \langle D, P \rangle_F, \tag{1}$$

yields the optimal transport cost to transform $T_i$ into $T_j$ by means of their subsequences. Then, given $\lambda \in_{>0}$, we can define

$$\text{WTK}(T_i, T_j) := \exp(-\lambda W_1(T_i, T_j)), \tag{2}$$

which we refer to as our *Wasserstein-based subsequence kernel*;

## Wasserstein Time Series Kernel

---

**Algorithm 1** Wasserstein Time Series Kernel

---

**Input:** Time series for training and testing $\mathcal{T}_{\text{train}}$, $\mathcal{T}_{\text{test}}$; subsequence length $w$; kernel weight factor $\lambda$

**Output:** $\mathcal{K}^{\text{train}}, \mathcal{K}^{\text{test}}$

1: $\mathcal{S}^{\text{train}} \leftarrow \textsc{Subsequences}(\mathcal{T}_{\text{train}}, w)$      // Extract subsequences
2: $\mathcal{S}^{\text{test}} \leftarrow \textsc{Subsequences}(\mathcal{T}_{\text{test}}, w)$      // Extract subsequences
3: **for** $T_i \in \mathcal{T}_{\text{train}}$ **do**
4:      **for** $T_j \in \mathcal{T}_{\text{train}}$ **do**
5:          $\mathcal{D}_{ij}^{\text{train}} \leftarrow W_1\left(\mathcal{S}_i^{\text{train}}, \mathcal{S}_j^{\text{train}}\right)$      // Wasserstein distance calculation (train)
6:      **end for**
7:      **for** $T_k \in \mathcal{T}_{\text{test}}$ **do**
8:          $\mathcal{D}_{ik}^{\text{test}} \leftarrow W_1\left(\mathcal{S}_i^{\text{train}}, \mathcal{S}_k^{\text{test}}\right)$      // Wasserstein distance calculation (test)
9:      **end for**
10: **end for**
11: $\mathcal{K}^{\text{train}} \leftarrow \exp\left(-\lambda \mathcal{D}^{\text{train}}\right)$      // Kernel matrix calculation
12: $\mathcal{K}^{\text{test}} \leftarrow \exp\left(-\lambda \mathcal{D}^{\text{test}}\right)$      // Kernel matrix calculation
13: **return** $\mathcal{K}^{\text{train}}, \mathcal{K}^{\text{test}}$

---

# Results

## Experiments

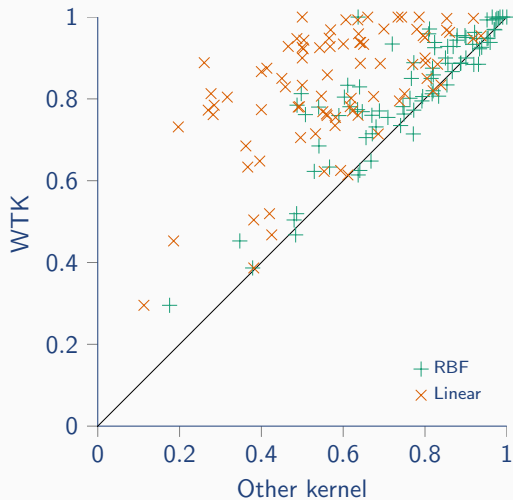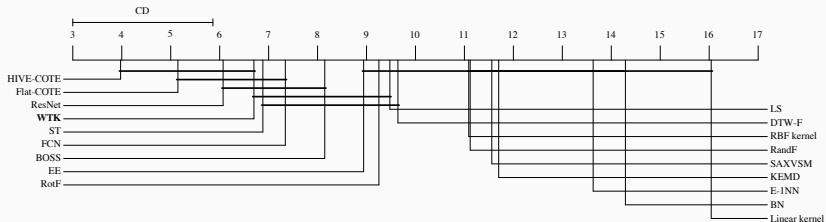|  |  |
|--:|:--|
| **Datasets** | *UCR Time Series Archive* |
|  | 85 datasets |
|  | predetermined train/test splits |
| **Hyperparameters** | Selected via a 5-fold cross-validation on the training set |
| **Evaluation metric** | Classification accuracy |
| **Comparison partners** | – Other kernels |
|  | – DTW-1NN |
|  | – State-of-the-art methods |

# Comparison with Other Kernels

## Critical Difference Plot

The classification performances of methods sharing horizontal bars are not significantly different.

# Take aways

**Questions?**

**References**