



HUAWEI System Software Workshop 2022

19 October 2022, Moscow, Russia

Efficient time series analytics

through DBMSs, ANNs, and parallel algorithms

明智是了解事件的人

Wise is the person who understands events.

Chinese proverb



Mikhail Zymbler, Yana Kraeva, Andrey Goglachev, Alexey Yurtin

{[mzym](#), [kraevaya](#), [goglachevai](#), [iurtinaa](#)}@susu.ru

S U S U

Big Data and Machine Learning Lab, South Ural State University, Chelyabinsk, Russia

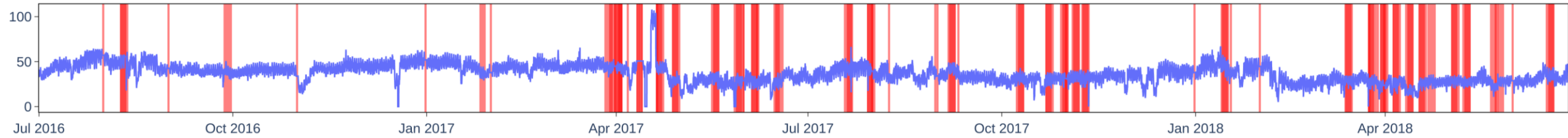
This work was financially supported by the Russian Foundation for Basic Research (grant No. 17-07-00463), and by the Ministry of Science and Higher Education of the Russian Federation (government order FENU-2020-0022)

Outline

- **Introduction**
 - Insights we can understand from time series
 - Offline and online time series analytics
- **Offline time series analytics**
 - Embedding time series analytics into PostgreSQL
- **Online time series analytics**
 - Employing ANNs together with parallel algorithms
 - Parallel time series labeling
 - Online imputation of missing values and anomaly detection
 - Parallel time series anomaly discovery
- **Conclusions**

Insights from time series: Anomalies

2-year power demand (Beijing Guowang Fuda Sci. & Tech. Dev. Co.)¹

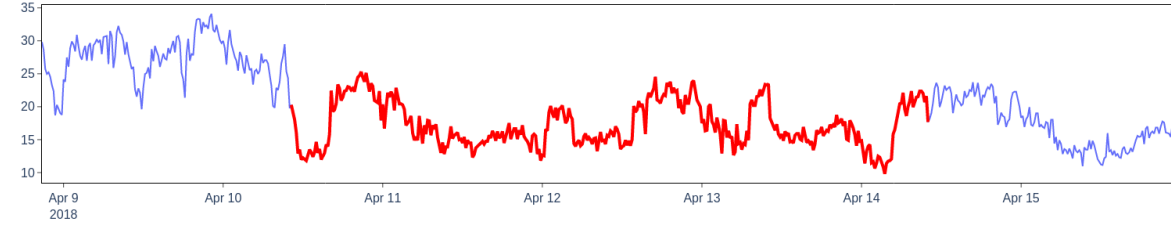
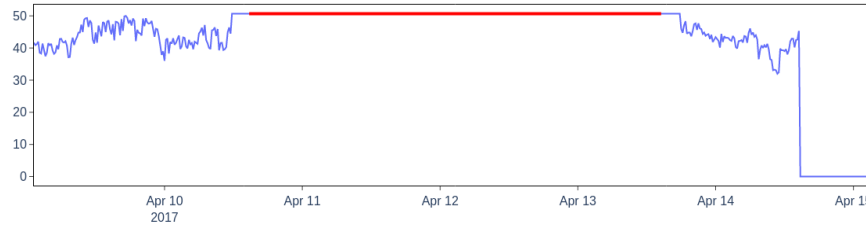
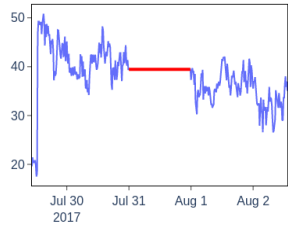


1 day

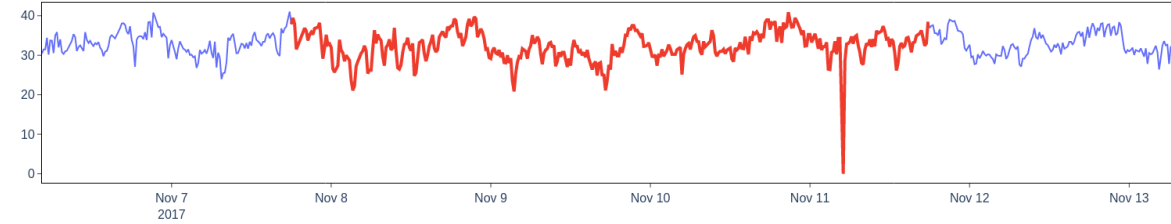
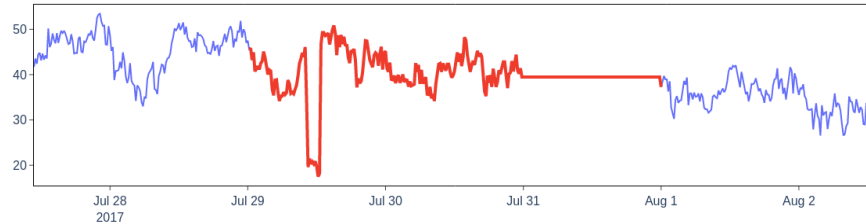
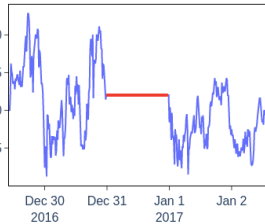
3 days

4 days

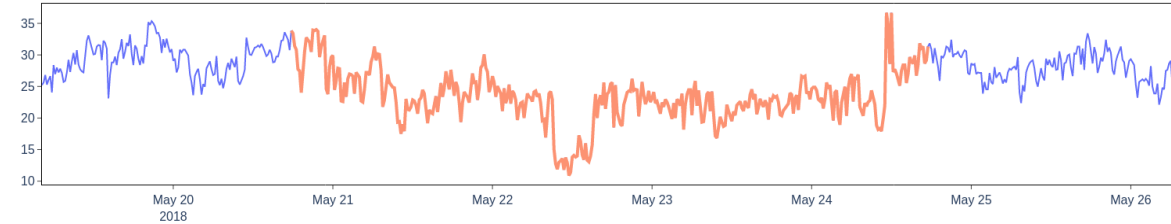
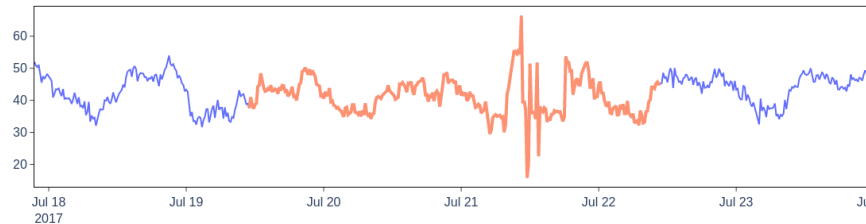
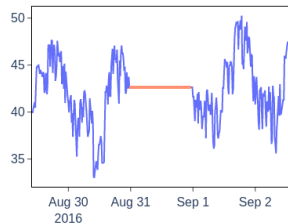
Top-1
anomaly



Top-2
anomaly



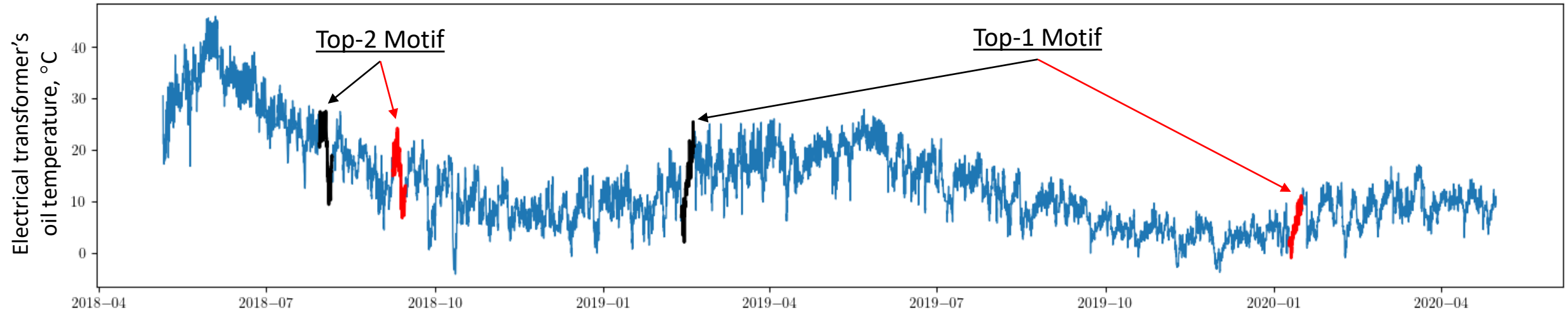
Top-3
anomaly



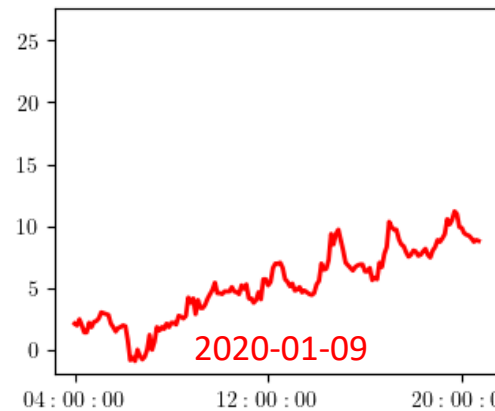
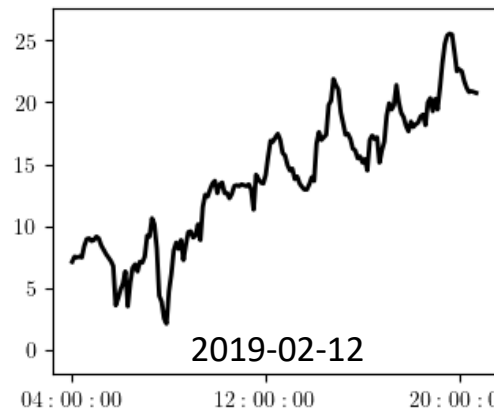
¹Zhou H. *et al.* Informer: beyond efficient transformer for long sequence time-series forecasting. AAAI 2021: 11106-11115. DOI: [10.1609/aaai.v35i12.17325](https://doi.org/10.1609/aaai.v35i12.17325).

Insights from time series: Motifs

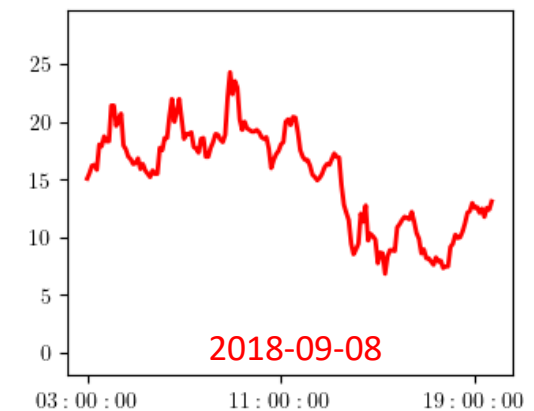
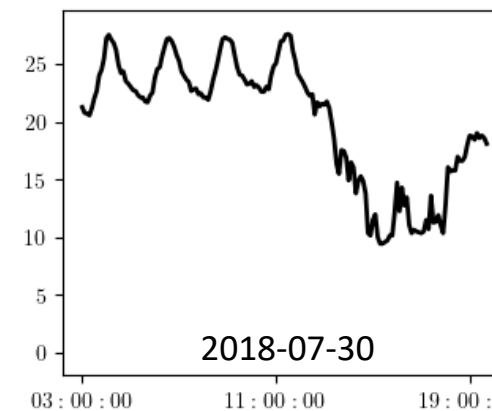
2-year power demand (Beijing Guowang Fuda Sci. & Tech. Dev. Co.)¹



Top-1 Motif



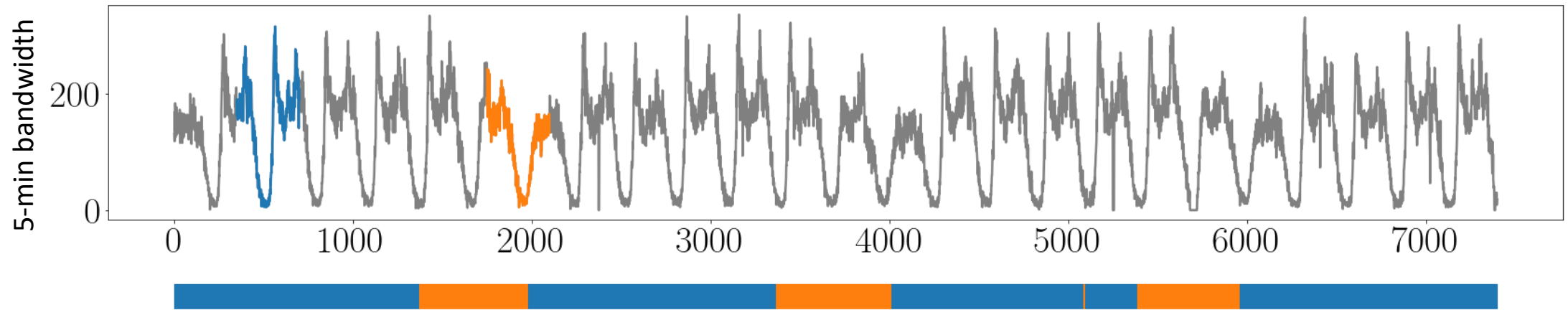
Top-2 Motif



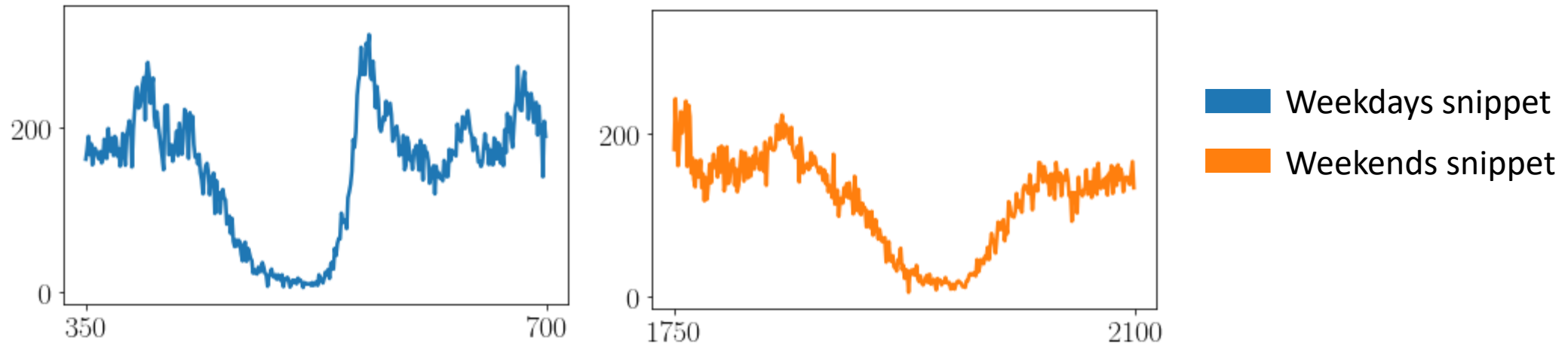
¹Zhou H. *et al.* Informer: beyond efficient transformer for long sequence time-series forecasting. AAAI 2021: 11106-11115. DOI: [10.1609/aaai.v35i12.17325](https://doi.org/10.1609/aaai.v35i12.17325).

Insights from time series: Patterns

One-month urban traffic in Munich (gathered by the Huawei Research Center)¹

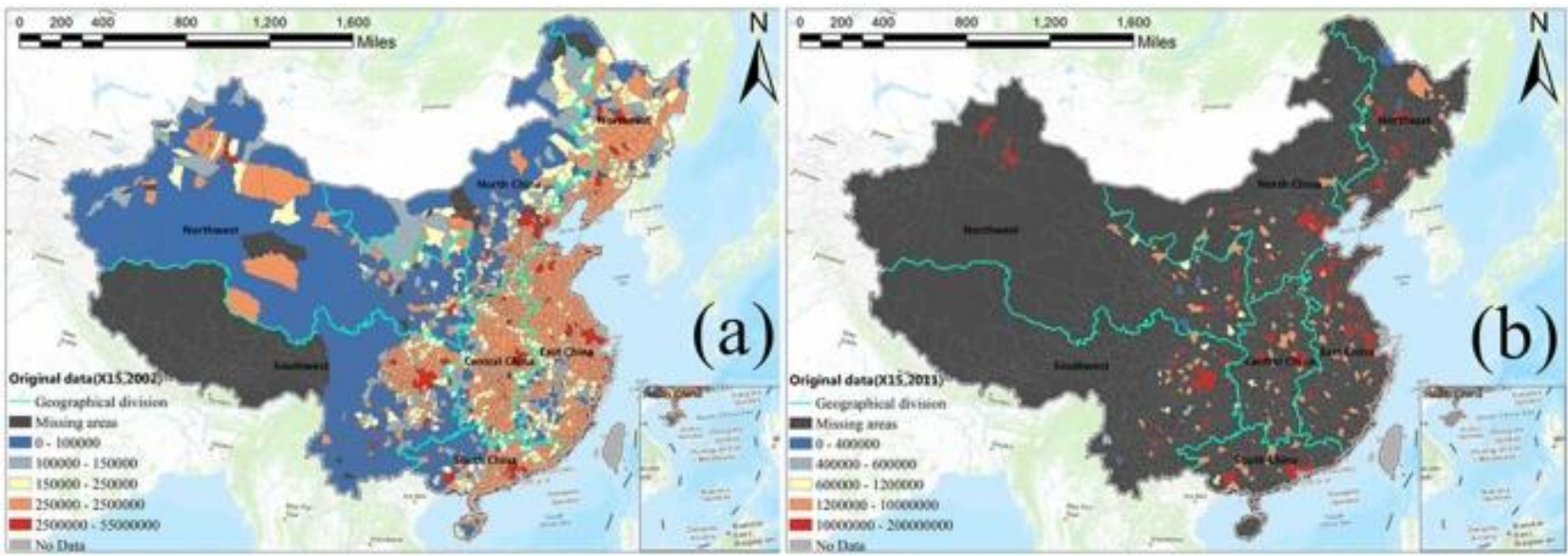


Patterns found

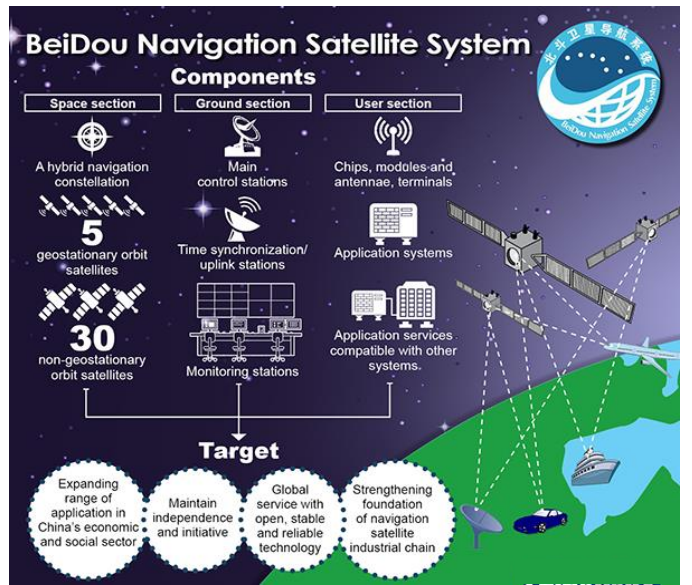


¹ Public (anonymized) road traffic prediction datasets from Huawei Munich Research Center. URL: <https://zenodo.org/record/3653880#.Y0zZi3ZBxPa>

Not an insight, but a headache to get rid of: missing values



China counties with missing official statistical data (one attribute)¹
a) 2002: less than 15% data are missing
b) 2011: more than 85% data are missing

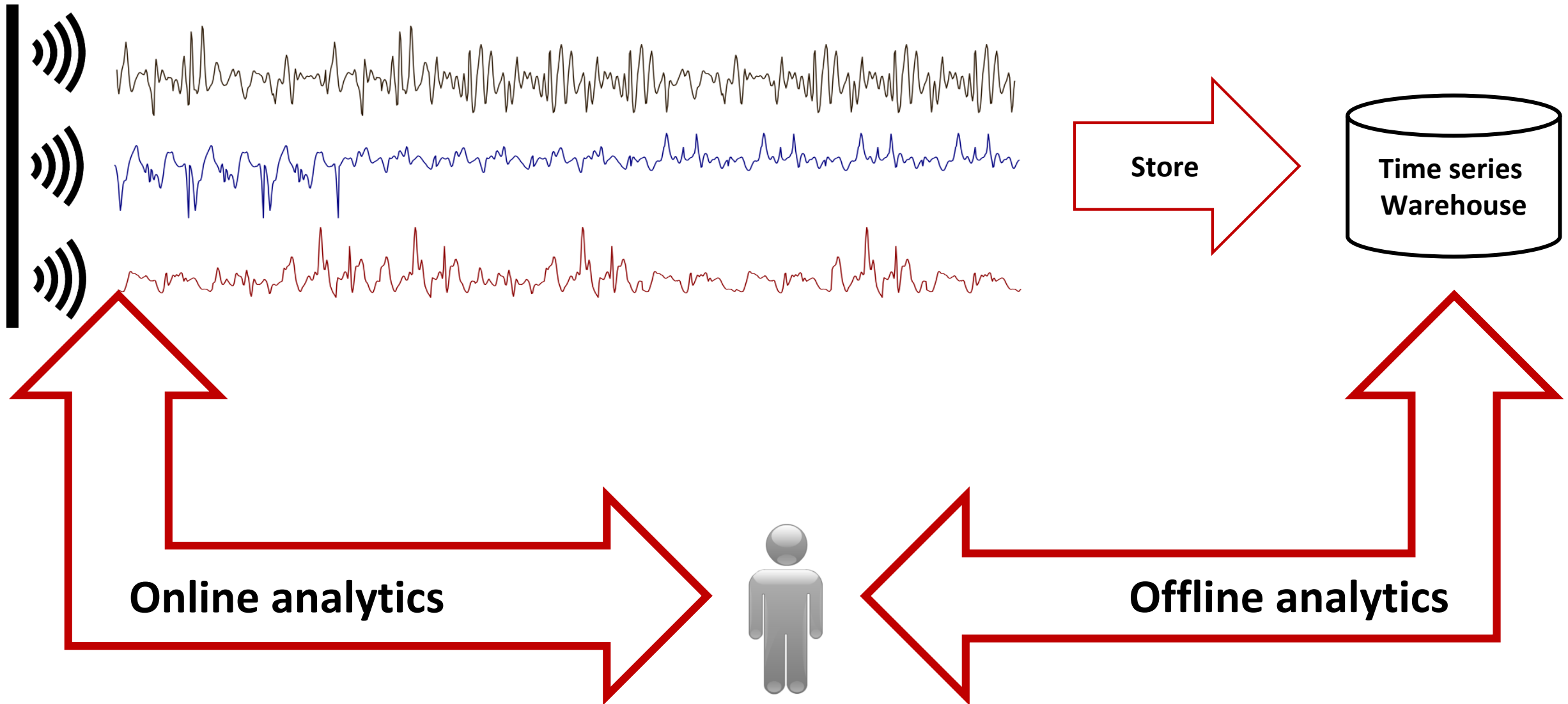


BeiDou satellite transmission link suffers the packet loss²

¹ Song C. *et al.* Estimating missing values in China's official socioeconomic statistics using progressive spatiotemporal Bayesian hierarchical modeling. *Sci. Rep.* 2018. Vol. 8, article 10055. DOI: [10.1038/s41598-018-28322-z](https://doi.org/10.1038/s41598-018-28322-z)

² Liu S. *et al.* A novel BeiDou satellite transmission framework with missing package imputation applied to smart ships. *IEEE Sensors Journal.* 2022. Vol. 22, no. 13. P. 13162-13176. DOI: [10.1109/JSEN.2022.3177167](https://doi.org/10.1109/JSEN.2022.3177167).

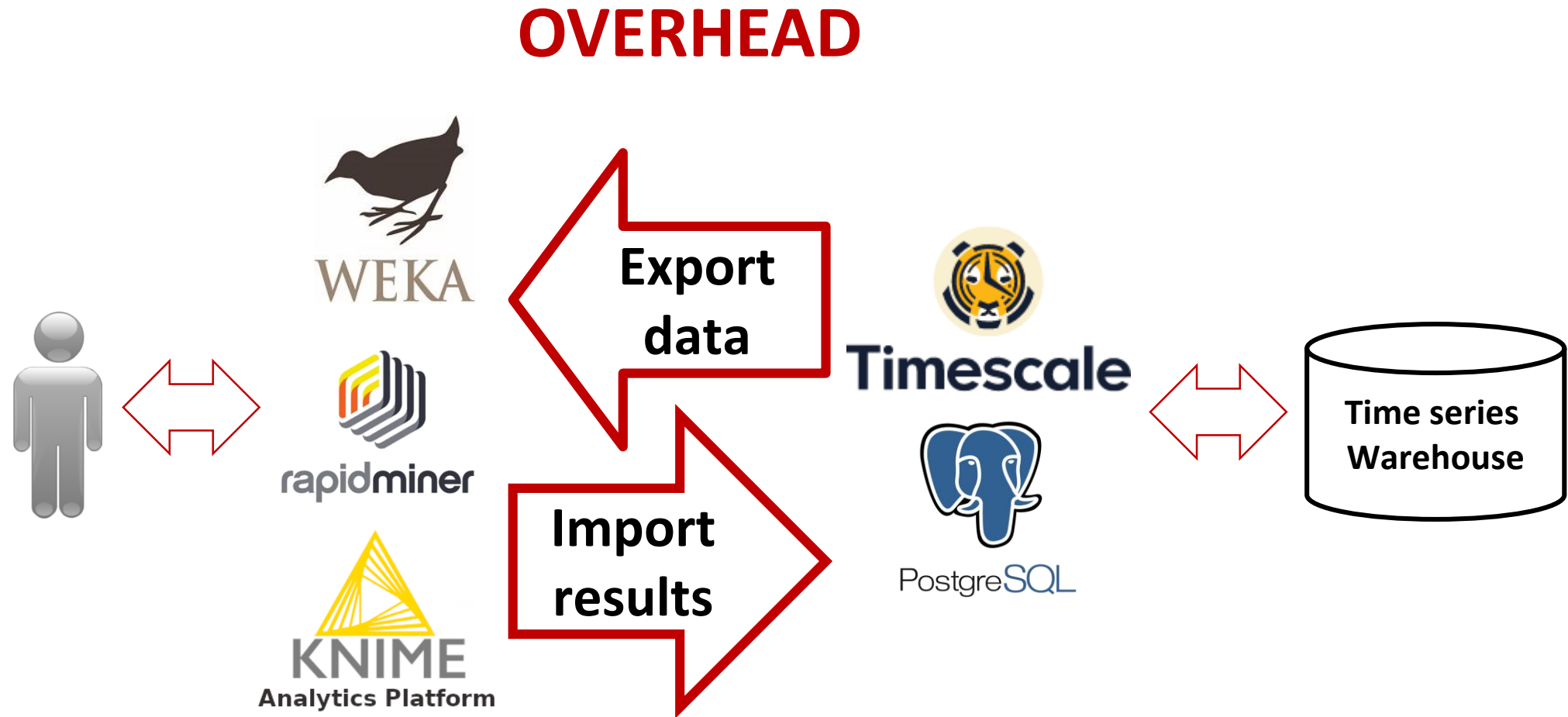
Time series analytics: online vs. offline



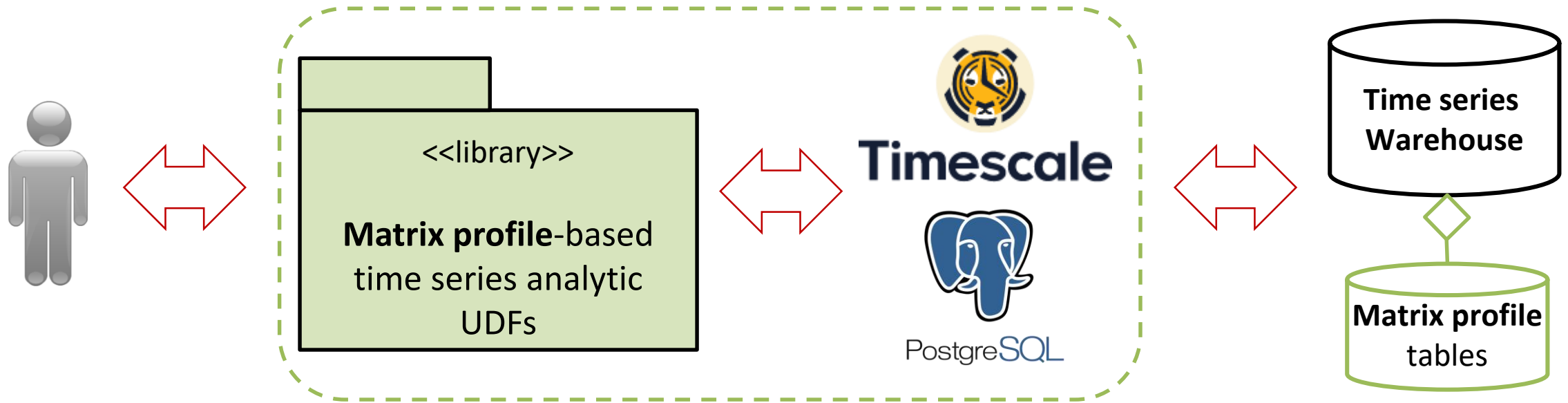
Outline

- Introduction
- **Offline time series analytics**
 - Embedding time series analytics into PostgreSQL
- Online time series analytics
 - Employing ANNs together with parallel algorithms
 - Parallel time series labeling
 - Online imputation of missing values and anomaly detection
 - Parallel time series anomaly discovery
- Conclusions

What's wrong with **offline** time series analytics?



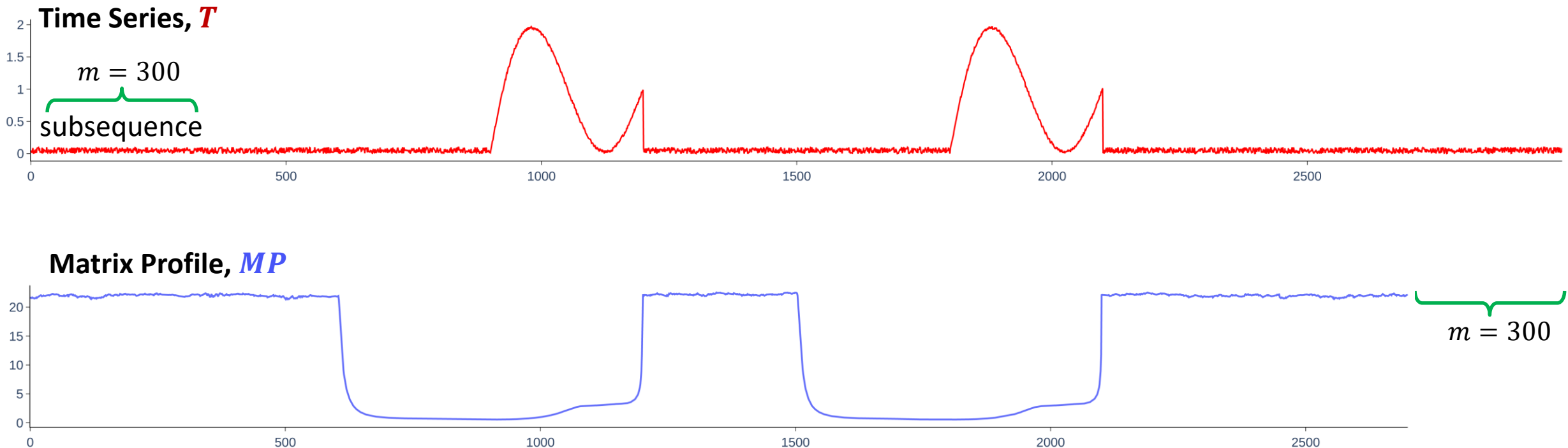
Offline time series analytics inside a DBMS



骑驴找马

Looking for a horse, ride a donkey.
Chinese proverb

Matrix profile^{1,2} of a time series



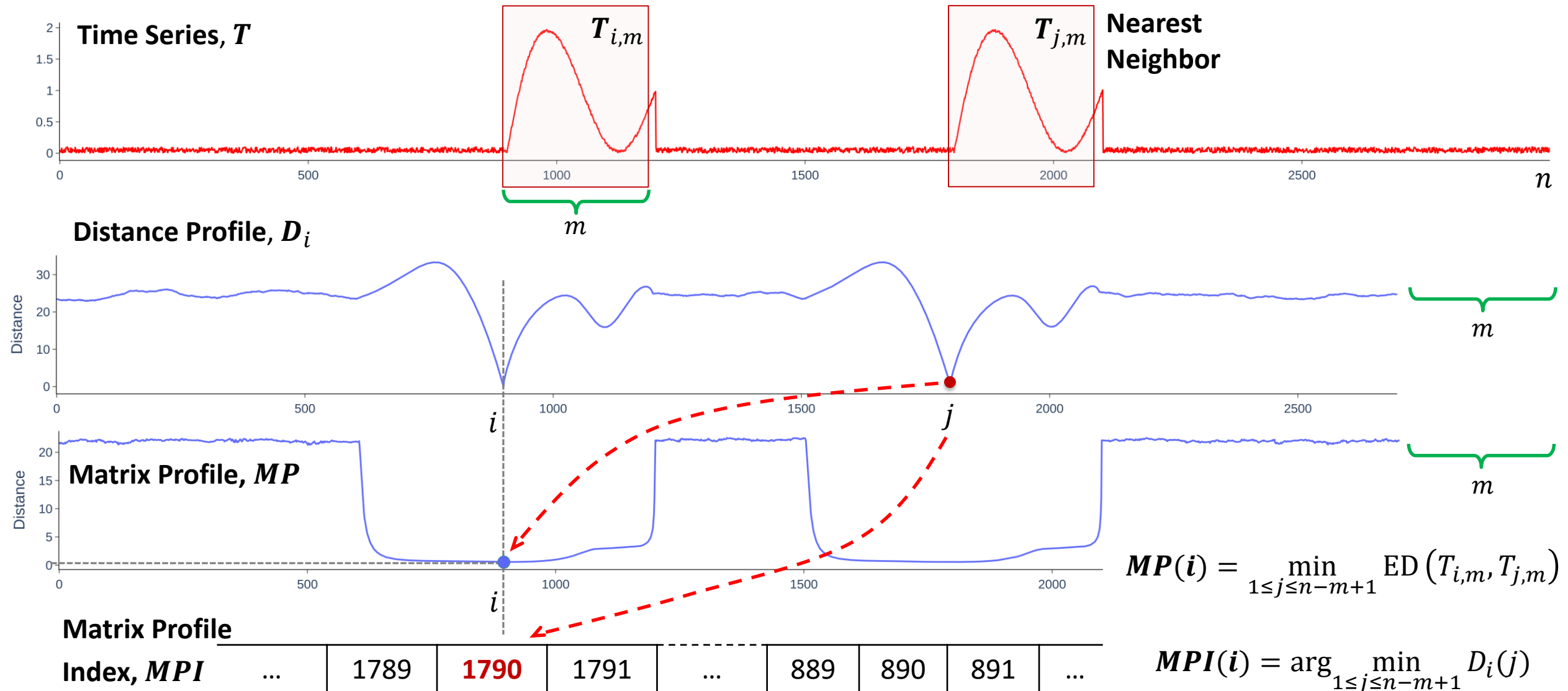
$$MP(i) = \min_{1 \leq j \leq n-m+1} \mathbf{ED}(T_{i,m}, T_{j,m})$$

Matrix profile's i -th element is the **Euclidean distance** between i -th subsequence and its **nearest non-overlapping neighbor**

¹ Yeh C.M. *et al.* Time series joins, motifs, discords and shapelets: A unifying view that exploits the matrix profile. *Data Min. Knowl. Discov.* 32(1), 83-123 (2018). DOI: [10.1007/s10618-017-0519-9](https://doi.org/10.1007/s10618-017-0519-9)

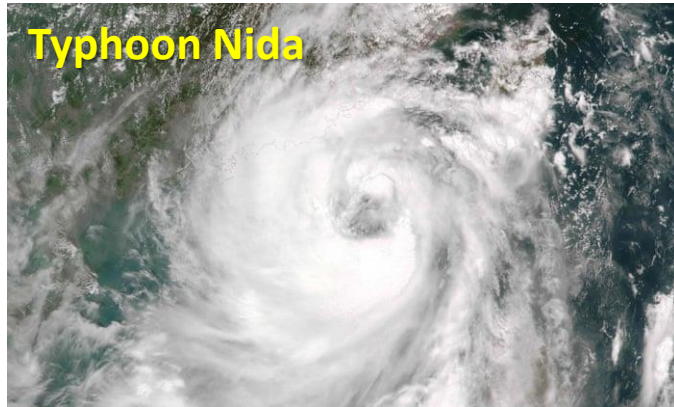
² Zimmerman Z. *et al.* Matrix Profile XIV: Scaling time series motif discovery with GPUs to break a quintillion pairwise comparisons a day and beyond. *SoCC 2019*. DOI: [10.1145/3357223.3362721](https://doi.org/10.1145/3357223.3362721)

Matrix profile of a time series



Insights from the Matrix profile: Discords

Urban traffic speed in Guangzhou¹



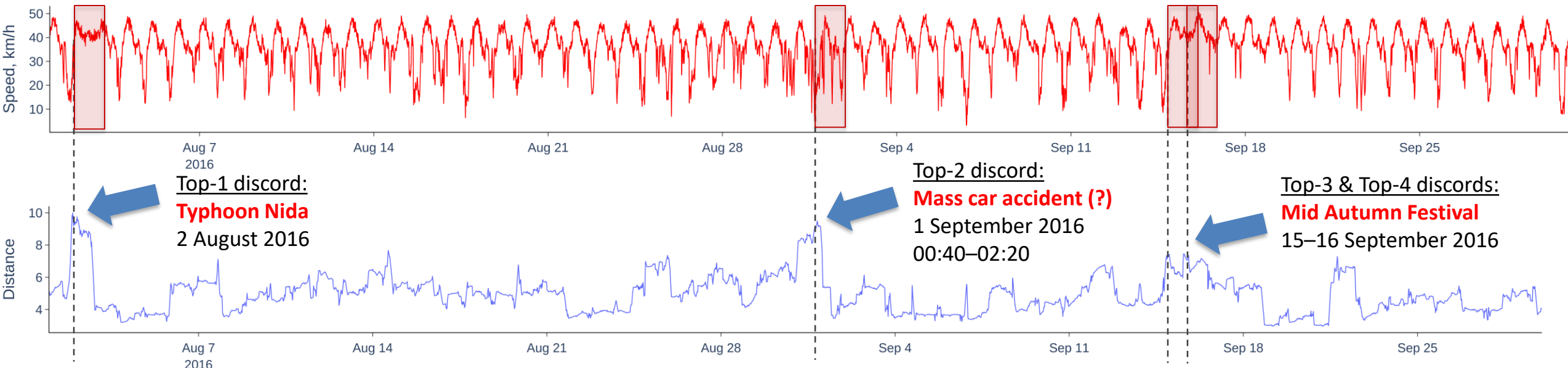
Typhoon Nida



Mass car accident



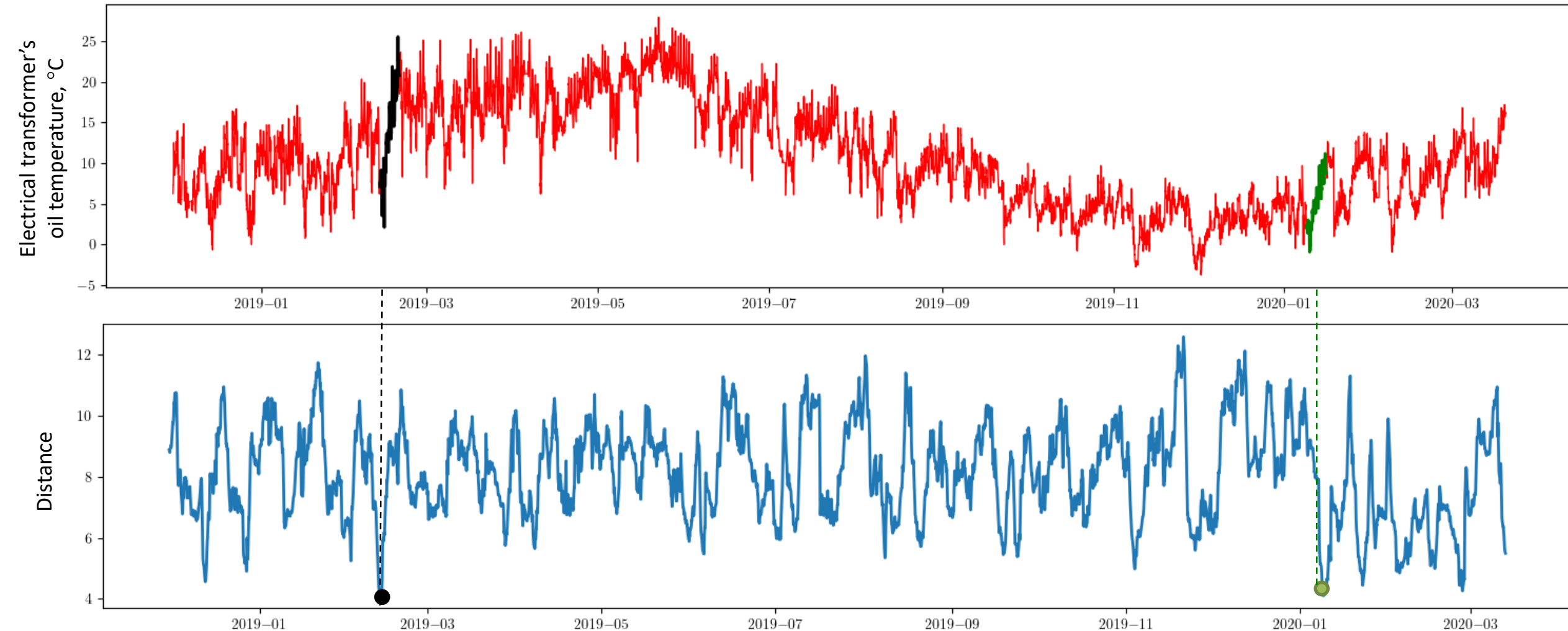
Mid Autumn Festival



¹Chen X, Chen Y, He Z. Urban traffic speed dataset of Guangzhou, China. 2018. DOI: [10.5281/zenodo.1205229](https://doi.org/10.5281/zenodo.1205229).

Insights from the Matrix profile: Motifs

2-year power demand (Beijing Guowang Fuda Sci. & Tech. Dev. Co.)¹



¹Zhou H. *et al.* Informer: beyond efficient transformer for long sequence time-series forecasting. AAAI 2021: 11106-11115. DOI: [10.1609/aaai.v35i12.17325](https://doi.org/10.1609/aaai.v35i12.17325).

Matrix profile-oriented warehouse

TIME SERIES DATA

Time Series Directory		ts_name (a time series)	
<u>ID</u>	BIGSERIAL	<u>num</u>	BIGSERIAL
name	TEXT	stamp	TIMESTAMP
len	BIGINT	val	REAL

MATRIX PROFILE DATA

Matrix Profile Directory		mp_name_subseqLen (a matrix profile)	
<u>ID</u>	BIGSERIAL	<u>num</u>	BIGSERIAL
tsID*	BIGINT	nnDist	REAL
subseqLen	INT	nnIdx	BIGINT

TIME SERIES ANALYTICS API

discoverDiscords
(name TEXT, subseqLen INT)

nnDist	REAL
<u>idx</u>	BIGINT
discord	REAL[subseqLen]

discoverMotifs
(name TEXT, subseqLen INT)

nnDist	REAL
<u>idxLeft</u>	BIGINT
motifLeft	REAL[subseqLen]
<u>idxRight</u>	BIGINT
motifRight	REAL[subseqLen]

discover ...
(name TEXT, subseqLen INT)

...	...
-----	-----

...

matrixProfile
(name TEXT, subseqLen INT)

<u>num</u>	BIGSERIAL
nnDist	REAL
nnIdx	BIGINT

Time series representation

TIME SERIES DATA

Time Series Directory		ts_name (a time series)	
<u>ID</u>	BIGSERIAL	<u>num</u>	BIGSERIAL
name	TEXT	stamp	TIMESTAMP
len	BIGINT	val1	REAL

Time Series Directory

<u>ID</u>	name	len
1	DailyExchangeRate	5000
2	MinutelyPulse	9000

ts_DailyExchangeRate

<u>num</u>	stamp	val
1	10.02.2008	23.35
2	11.02.2008	24.05
...
5000	19.10.2022	63.63

ts_MinutelyPulse

<u>num</u>	stamp	val
1	10.10.2022 00:01	135
2	10.10.2022 00:02	105
...
9000	16.10.2022 06:00	139

Matrix profile representation

MATRIX PROFILE DATA

Matrix Profile Directory		mp_name_subseqLen (Matrix Profile)	
<u>ID</u>	BIGSERIAL	<u>num</u>	BIGSERIAL
tsID*	BIGINT	nnDist	REAL
subseqLen	INT	nnIdx	BIGINT

mp_DailyExchangeRate_7

<u>num</u>	nnDist	nnIdx
1	0.02	2071
2	0.05	2078
...
4994	12.22	184

mp_DailyExchangeRate_30

<u>num</u>	nnDist	nnIdx
1	346.32	3546
2	358.15	4278
...
4971	0.10	1286

Time Series Directory

<u>ID</u>	name	len
1	DailyExchangeRate	5000
2	MinutelyPulse	9000

Matrix Profile Directory

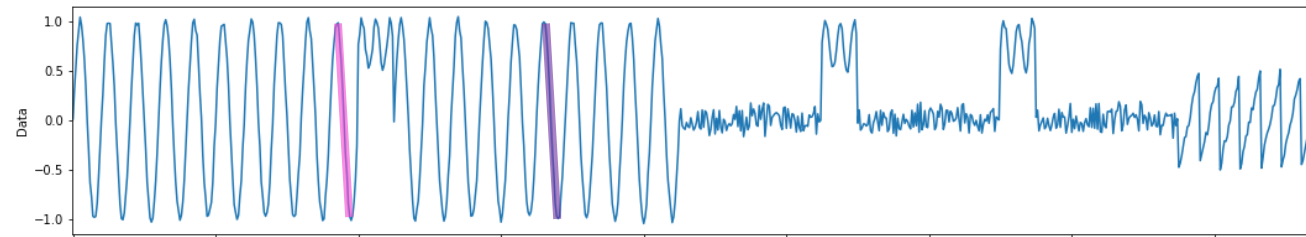
<u>ID</u>	tsID	subseqLen
1	1	7 weekly
2	1	30 monthly
3	2	60 hourly

mp_MinutelyPulse_60

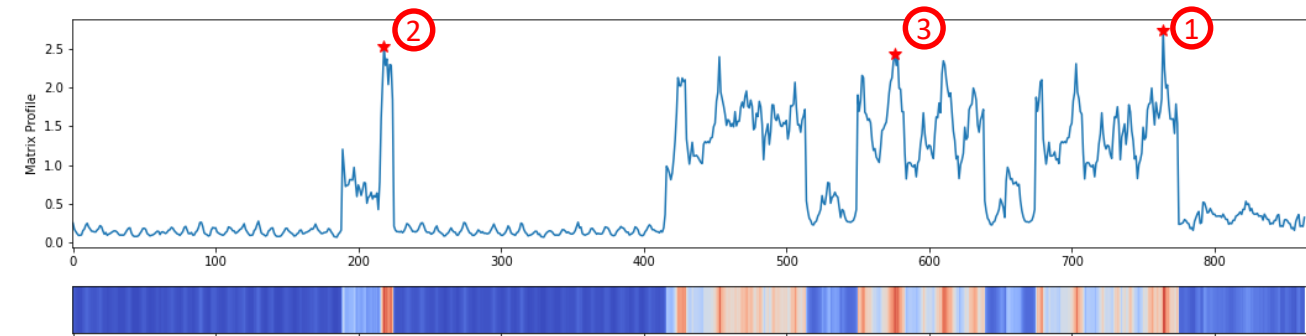
<u>num</u>	nnDist	nnIdx
1	23.04	145
2	39.78	7658
...
8941	11.65	7856

Time series analytics API

SELECT matrixProfile (sensorData, 12);



ts_sensorData



mp_sensorData_12

SELECT discoverDiscords (sensorData, 12, 3);

NNDIST	IDX	DISCORD
2.75 ①	764	[-0.05, 0.02, 0.06, 0.04, 0.04, 0.04, 0.05, -0.18, 0.06, -0.05, 0.01, -0.47]
2.53 ②	218	[0.55, 0.64, 0.74, 0.87, 1.03, 0.98, 0.85, -0.01, 0.35, 0.61, 0.81, 0.96]
2.43 ③	576	[0.05, -0.16, 0.01, -0.08, -0.09, -0.02, 0.05, -0.03, 0.01, -0.01, -0.06, -0.08]

SELECT discoverMotifs (sensorData, 12, 1);




NNDIST	IDX _{LEFT}	MOTIF _{LEFT}	IDX _{RIGHT}	MOTIF _{RIGHT}
0.069	185	[0.97,0.98,0.79,0.59,0.31,-0.01, -0.34,-0.54,-0.85,-0.97,-1.01,-0.96]	330	[0.99,0.97,0.82,0.60,0.32,-0.02, -0.32,-0.56,-0.84,-0.97,-0.99,-0.98]

Outline




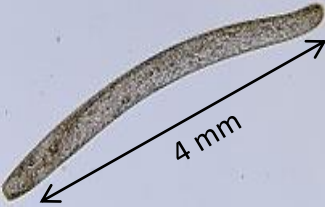
- Introduction
- Offline time series analytics
- **Online time series analytics**
 - Employing ANNs together with parallel algorithms
 - Parallel time series labeling
 - Online anomaly detection and imputation of missing values
 - Parallel time series anomaly discovery
- Conclusions

What's wrong with **online** time series analytics?

- 😊 No data labeling and model learning overheads
- 😞 Fast, but not enough to fit online analytics
- 😞 Accurate, but not like ANNs

Method	Labeling	Learning	Tuning	Mining
Parallel algorithms	子曰、攻乎异端、斯害也己。 Learning the wrong views is harmful Confucius		 Turtle 15-35 km/h, running/swimming	 Cheetah 100 km/h

- 😊 Fast to fit online analytics
- 😊 High accuracy
- 😞 Data labeling and model learning overheads

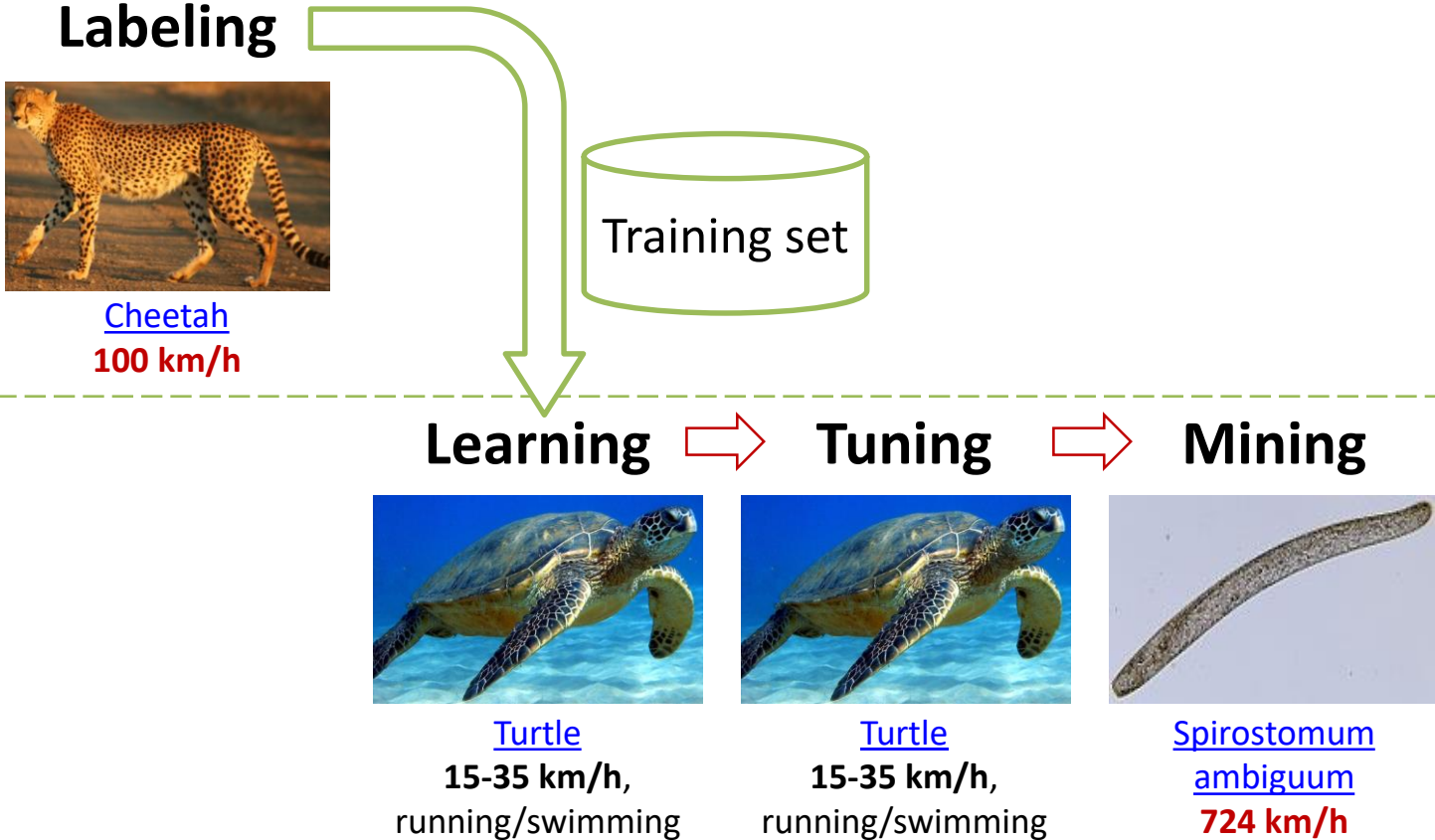
Method	Labeling	Learning	Tuning	Mining
Artificial Neural Networks	 Chén Dìng (racewalker) 15.2 km/h	 Turtle 15-35 km/h, running/swimming	 Turtle 15-35 km/h, running/swimming	 Spirostomum ambiguum 724 km/h

Online time series analytics by parallel algorithms plus ANNs

情义无价
Friendship has no price.
Chinese proverb

Parallel algorithms

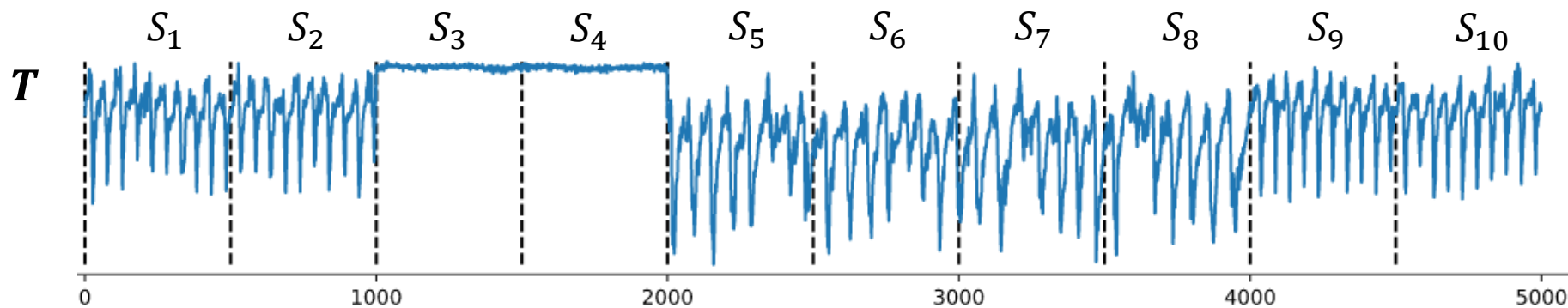
Artificial Neural Networks



Outline

- Introduction
- Offline time series analytics
- **Online time series analytics**
 - Employing ANNs together with parallel algorithms
 - **Parallel time series labeling**
 - Online anomaly detection and imputation of missing values
 - Parallel time series anomaly discovery
- Conclusions

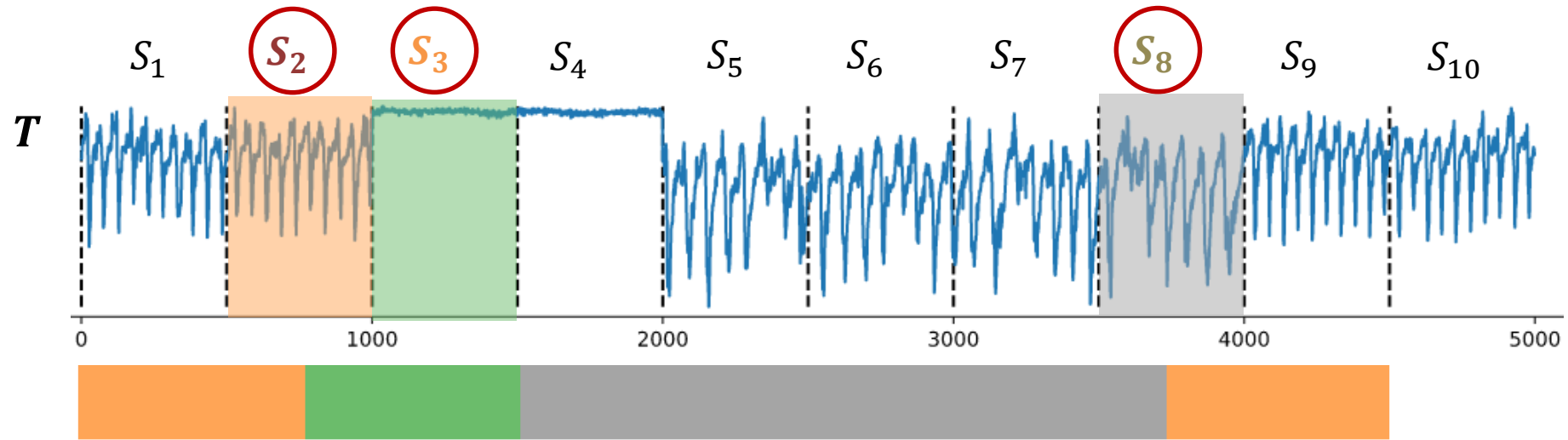
Time series labeling through snippets¹



1. Let us represent a time series as a set of non-overlapped segments of a given length

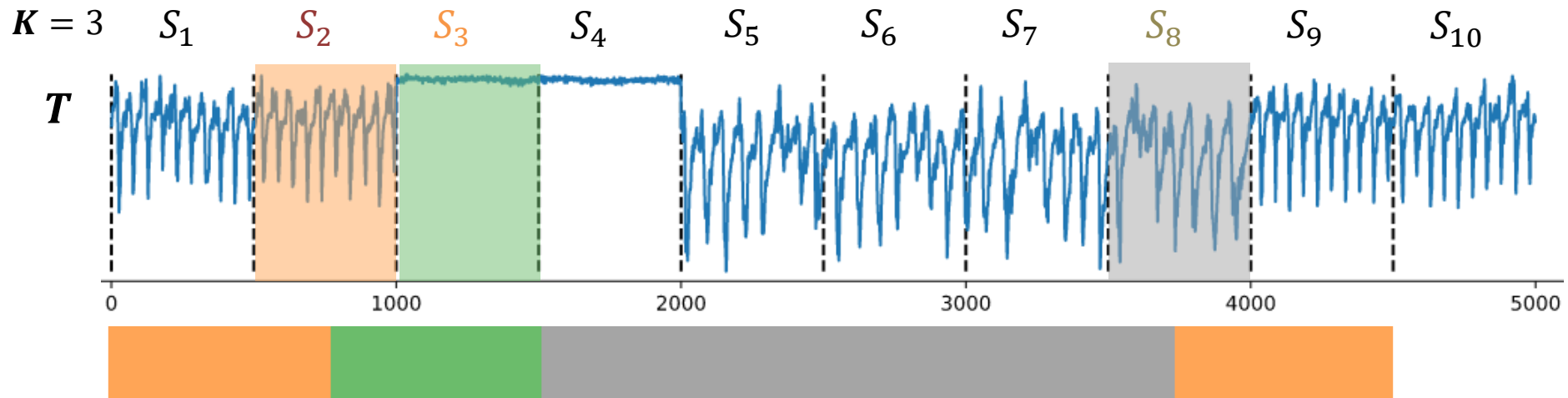
¹ Imani S. *et al.* Introducing time series snippets: a new primitive for summarizing long time series. Data Min. Knowl. Discov. 2020. Vol. 34, no. 6. P. 1713-1743. DOI: [10.1007/s10618-020-00702-y](https://doi.org/10.1007/s10618-020-00702-y)

Time series labeling through snippets

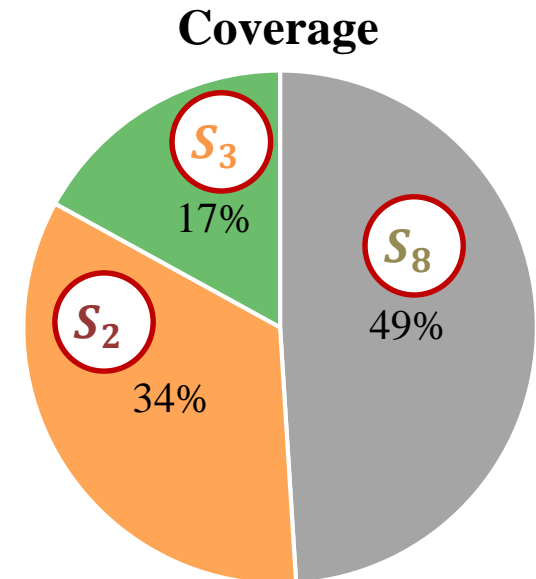


1. Let us represent a time series as a set of non-overlapped segments of a given length
2. For each segment, let us find its **nearest neighbors**

Time series labeling through snippets



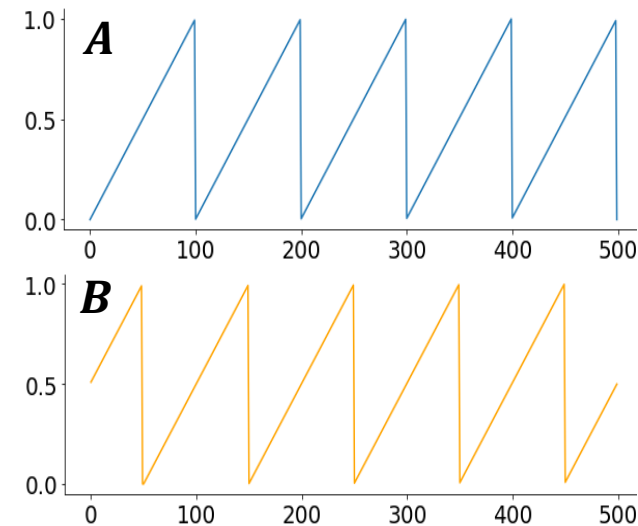
1. Let us represent a time series as a set of non-overlapped segments of a given length
2. For each segment, let us find its nearest neighbors
3. For each segment, let us compute its coverage and take top- K segments



Time series similarity through MPdist¹

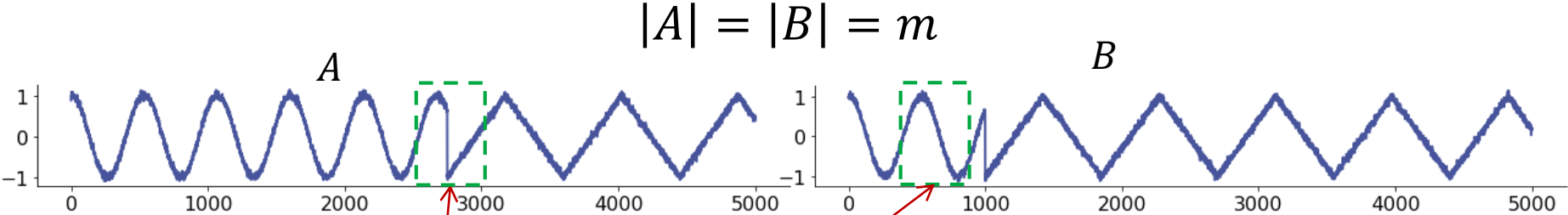
- Two m -length time series are more similar w.r.t. **MPdist**, the more ℓ -length ($3 \leq \ell \leq m$) subsequences close to each other w.r.t. the **normalized Euclidean distance**, are in them
- MPdist is a distance measure (not a metric), i.e. it holds the identity and symmetry axioms but not the triangle inequality
- MPdist is phase-invariant

$ED(A, B)$	11.2
MPdist (A, B)	0

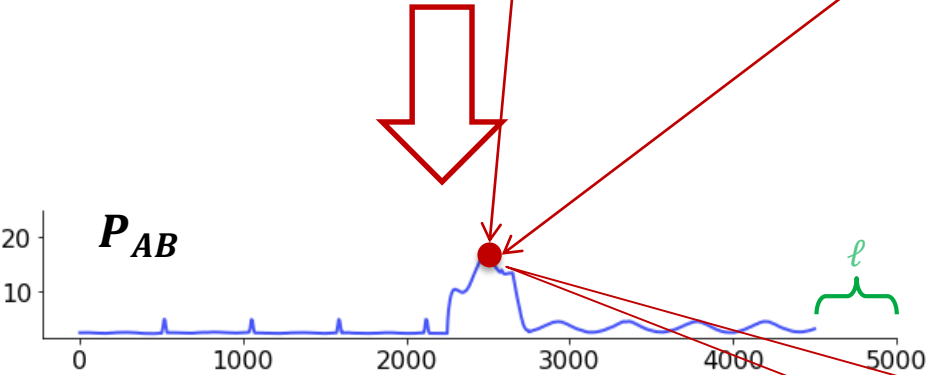


¹ Gharghabi S. *et al.* An ultra-fast time series distance measure to allow data mining in more complex real-world deployments. *Data Min. Knowl. Discov.* 2020. Vol. 34. P. 1104–1135. DOI: [10.1007/s10618-020-00695-8](https://doi.org/10.1007/s10618-020-00695-8)

MPdist: Matrix profile AB

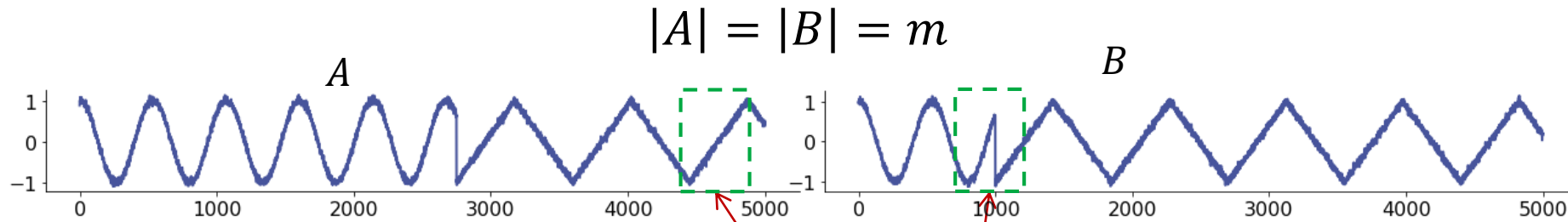


Meaningful subsequence length: $3 \leq \ell \leq m$ (typically, $[0.3m] < \ell \leq [0.8m]$)

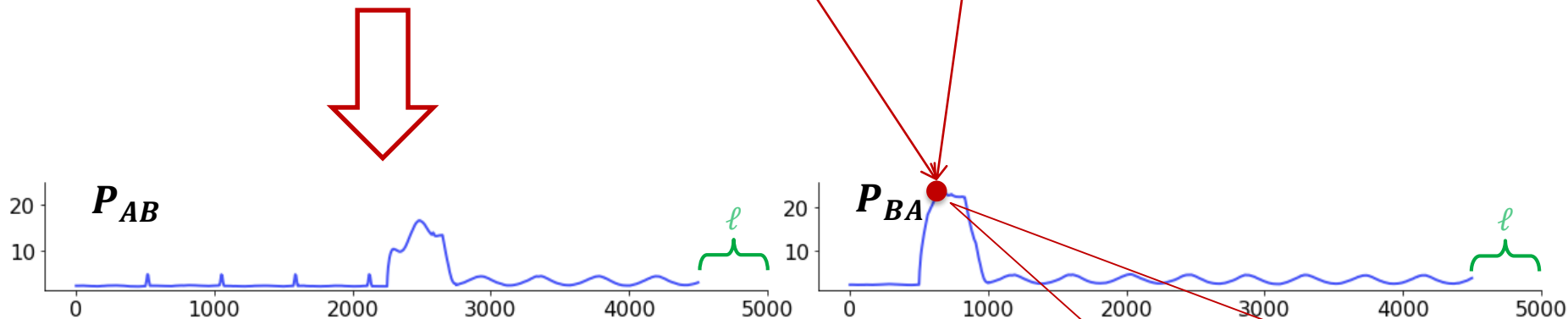


Normalized Euclidean distance between i -th ℓ -length subsequence in A and its nearest ℓ -length subsequence in B

MPdist: Matrix profile BA

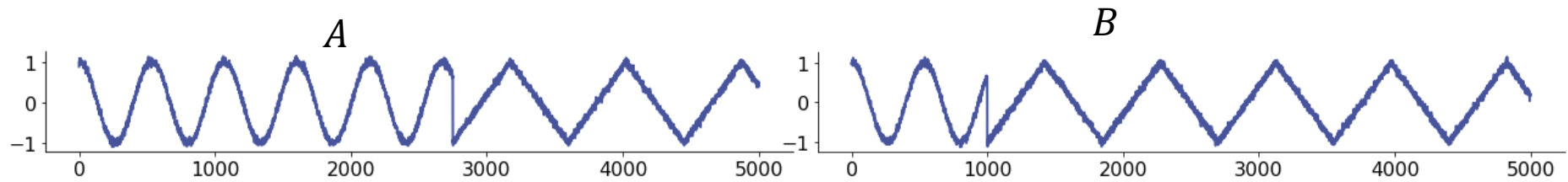


Meaningful subsequence length: $3 \leq \ell \leq m$ (typically, $[0.3m] < \ell \leq [0.8m]$)

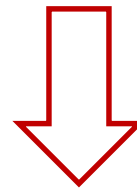
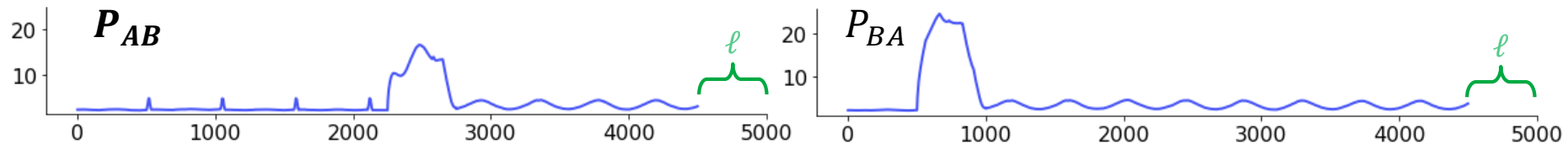


Normalized Euclidean distance between i -th ℓ -length subsequence in A and its nearest ℓ -length subsequence in B

MPdist: Matrix profile ABBA



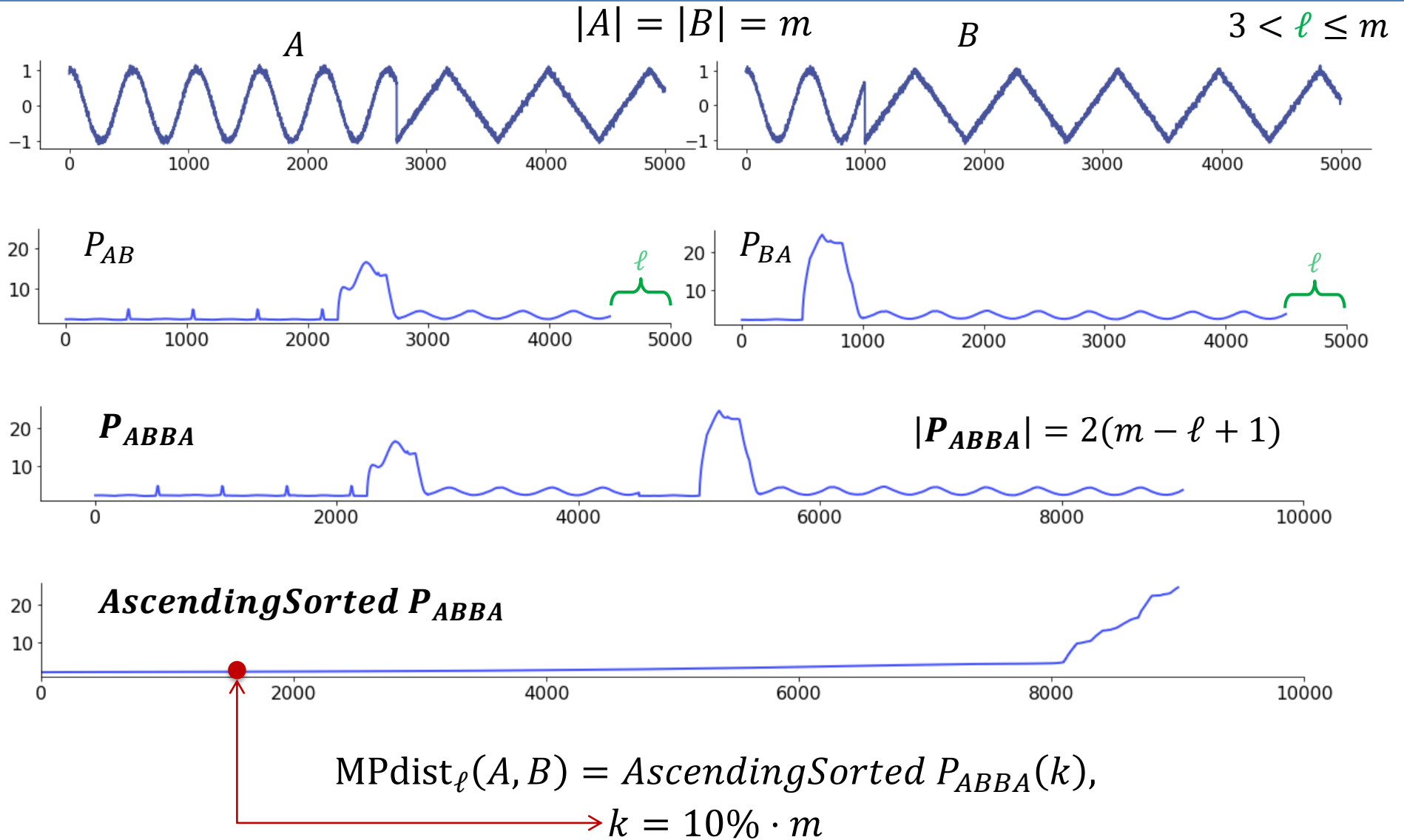
$$3 \leq \ell \leq m$$



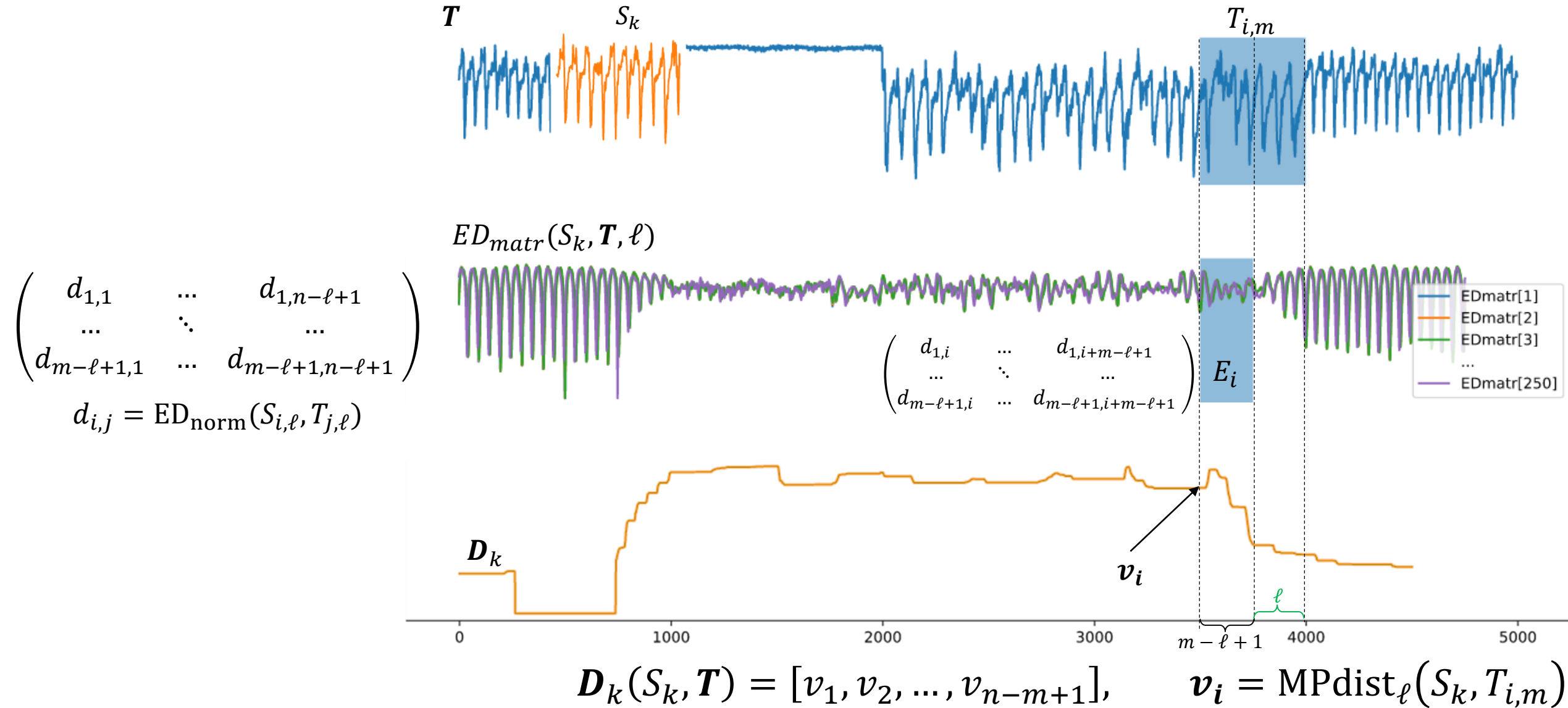
$$P_{ABBA} = P_{AB} \odot P_{BA}$$



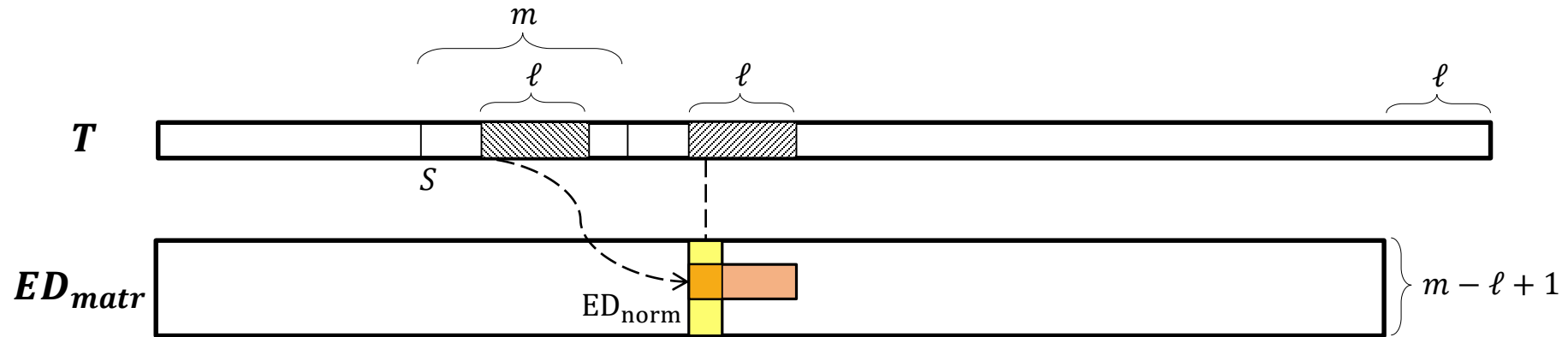
MPdist: Eventual calculation



Parallel labeling: MPdist profile



Parallel snippet discovery: ED_{matr}



$$ED_{norm}(T_{i,m}, T_{j,m}) = \sqrt{2m(1 - P_{i,j})}$$

$$P_{i,j} = \overline{QT}_{i,j} \cdot \frac{1}{\|T_{i,m} - \mu_i\|} \cdot \frac{1}{\|T_{j,m} - \mu_j\|}$$

$$T_{i,m} - \mu_i = (t_i - \mu_i, \dots, t_{i+m-1} - \mu_i),$$

$$\mu_i = \frac{1}{m} \sum_{j=i}^{i+m} t_j,$$

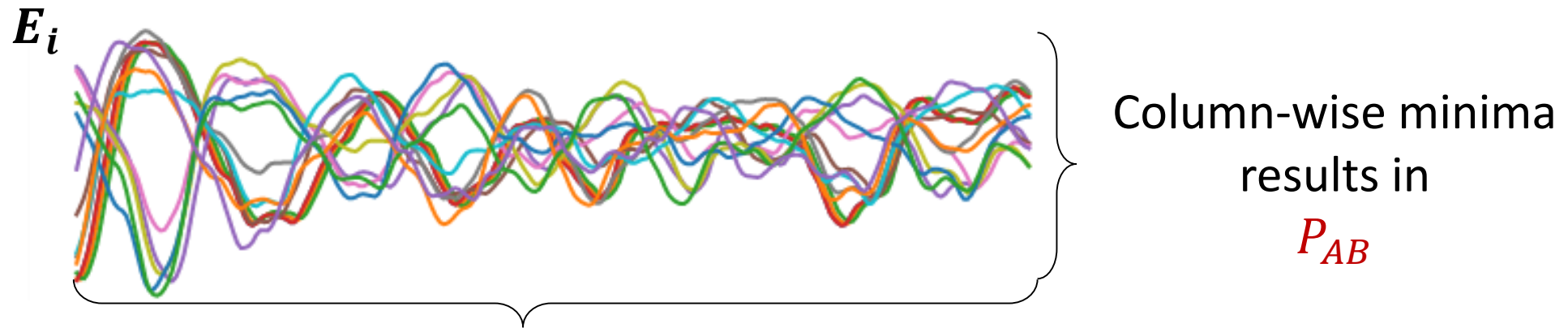
$$dg_0 = 0; dg_i = (t_{i+m-1} - \mu_i) + (t_{i-1} - \mu_{i-1}),$$

$$df_0 = 0; df_i = \frac{t_{i+m-1} - t_{i-1}}{2},$$

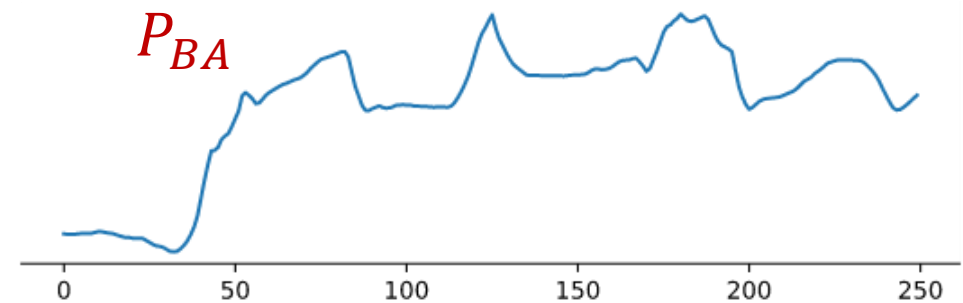
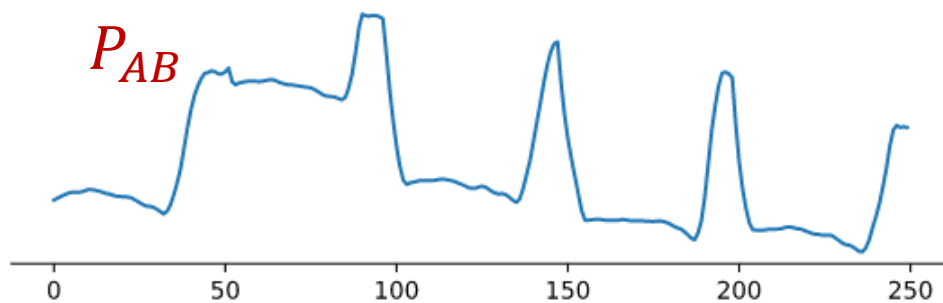
$$\overline{QT}_{i,j} = \overline{QT}_{i-1,j-1} + df_i \cdot dg_j + df_j \cdot dg_i,$$

¹ Zimmerman Z. *et al.* Matrix Profile XIV: Scaling time series motif discovery with GPUs to break a quintillion pairwise comparisons a day and beyond. SoCC 2019. P. 74–86. DOI: [10.1145/3357223.3362721](https://doi.org/10.1145/3357223.3362721)

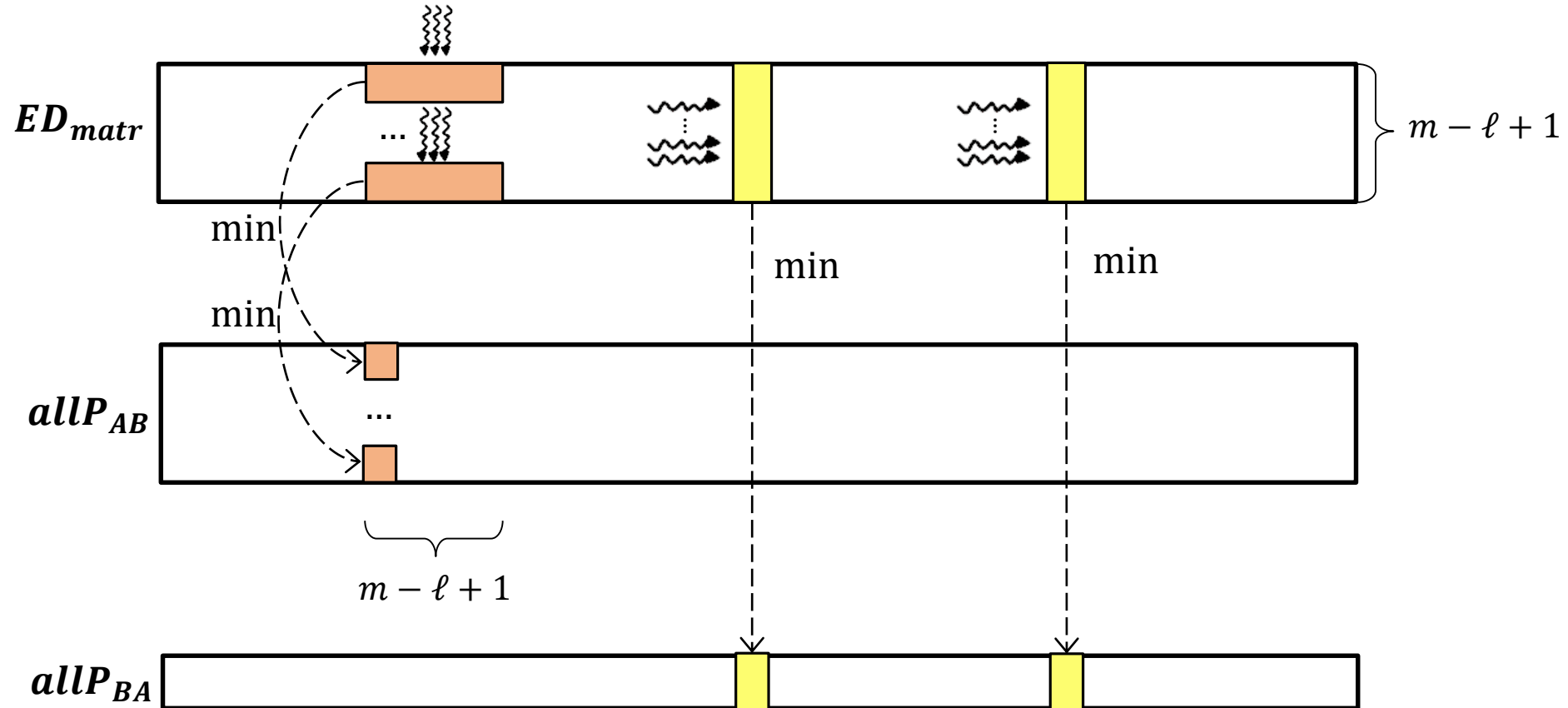
Parallel labeling: MPdist profile



Row-wise minima results in P_{BA}



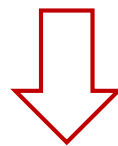
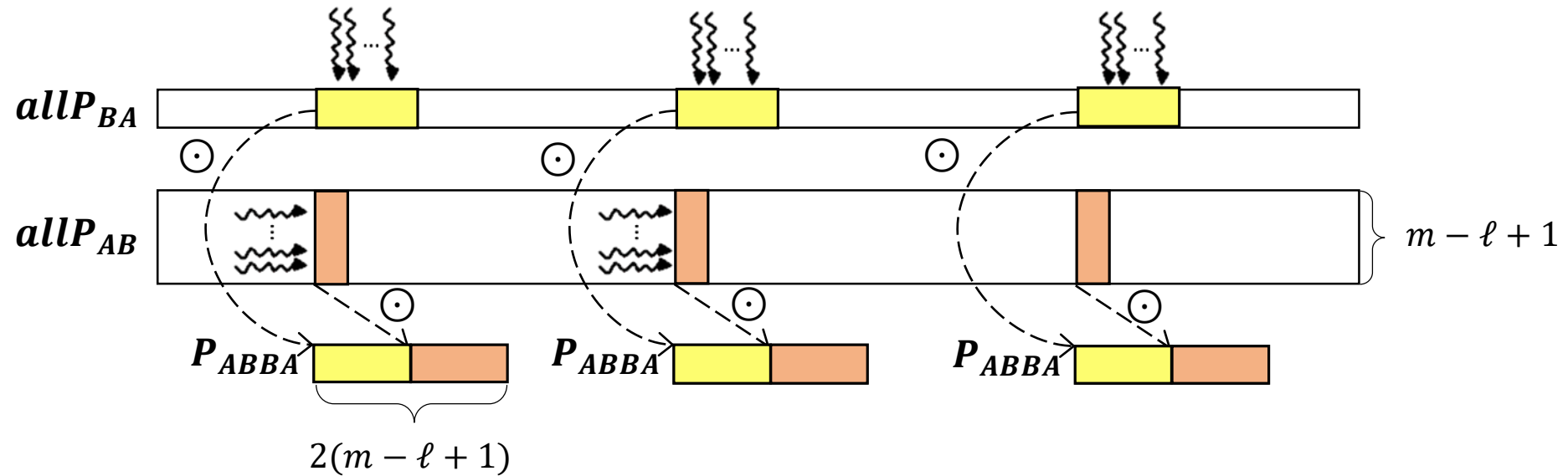
Parallel labeling: $allP_{AB}$ and $allP_{BA}$



$$allP_{AB}(i, j) = \min_{j \leq c \leq j+m-\ell+1} ED_{matr}(i, c)$$

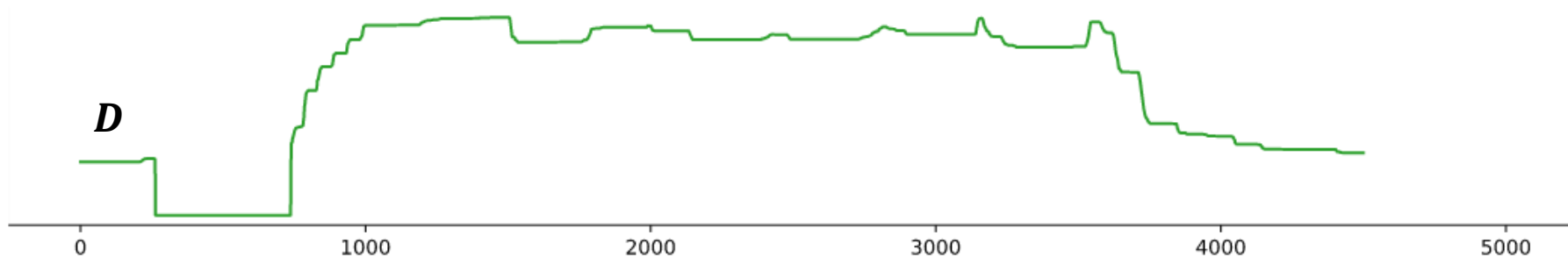
$$allP_{BA}(j) = \min_{1 \leq i \leq m-\ell+1} ED_{matr}(i, j)$$

Parallel labeling: P_{ABBA}

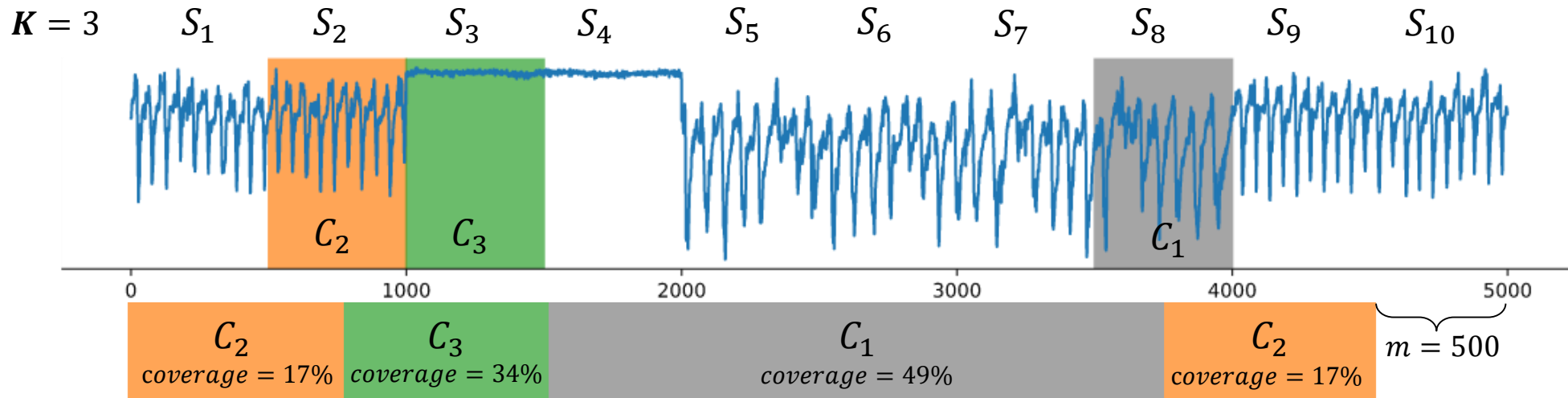


$$MPdist_{\ell}(A, B) = \text{AscendingSorted } P_{ABBA}(k),$$

where $k = 10\% \cdot m$



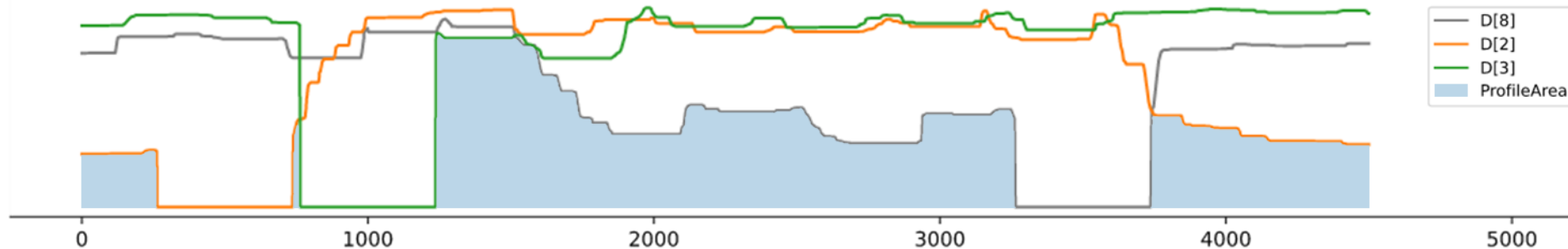
Parallel labeling: Results



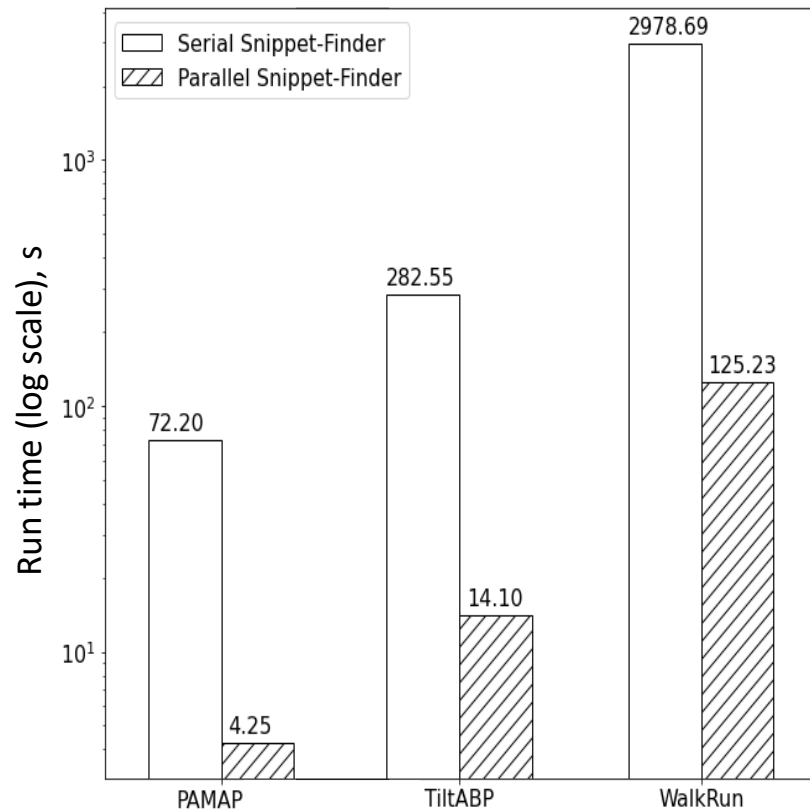
$$i = 1: \quad C_1 = \arg \min_{1 \leq j \leq n/m} ProfileArea(\{D_j\})$$

$$i = 2: \quad C_2 = \arg \min_{1 \leq j \leq n/m} ProfileArea(\{D_{C_1}, D_j\})$$

$$3 \leq i \leq K: \quad C_i = \arg \min_{1 \leq j \leq n/m} ProfileArea(\{D_{C_1}, \dots, D_{C_{i-1}}, D_j\})$$



Parallel labeling: Experiments on performance



Parallel labeling algorithm
is **20 times ahead**
of serial analog

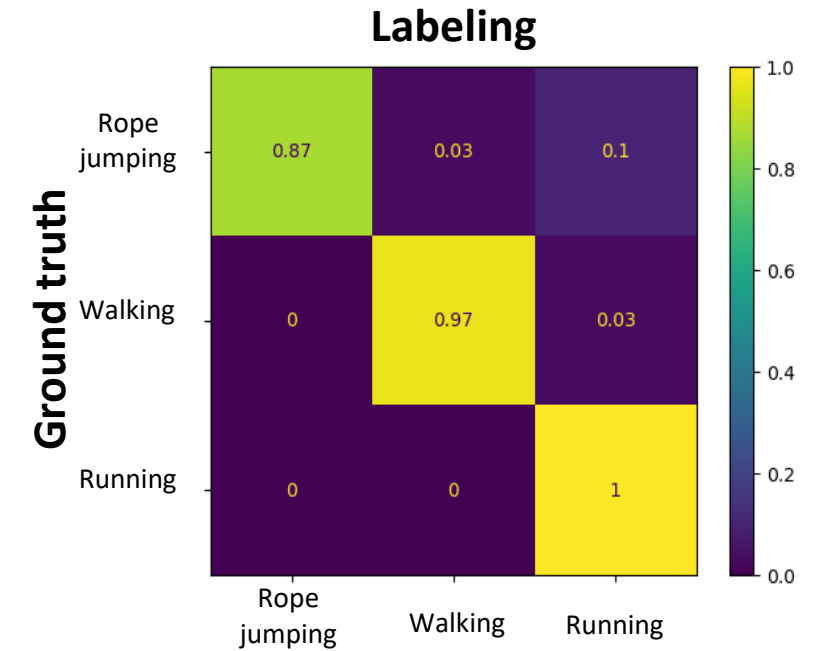
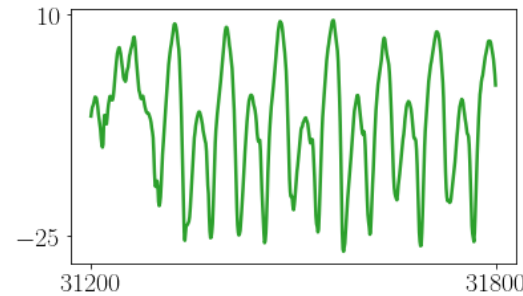
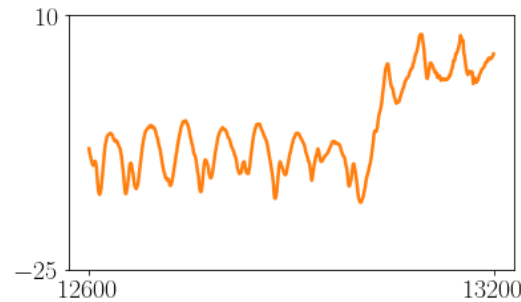
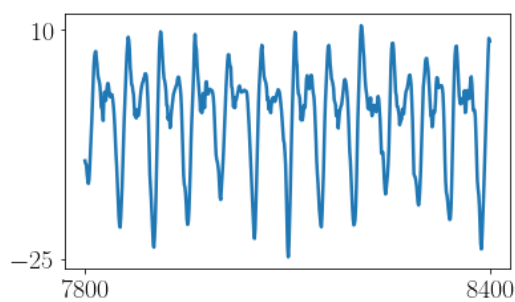
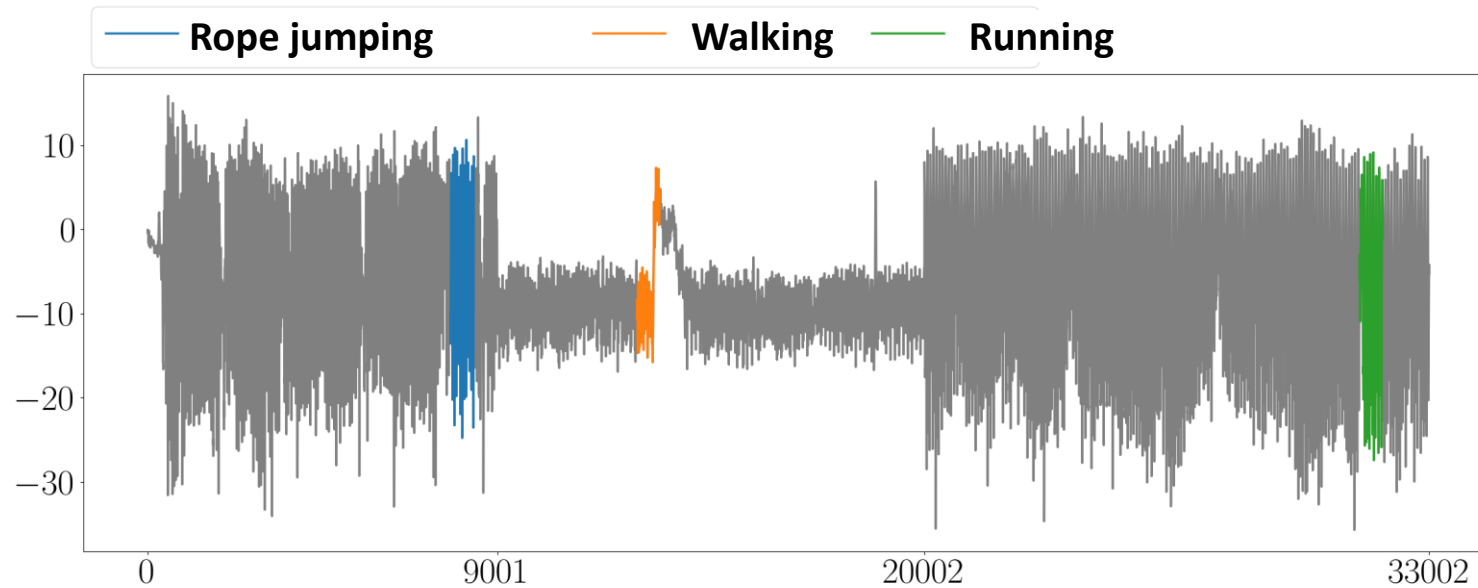
Hardware: NVIDIA Tesla V100 SXM2 (5120 cores @1.3 GHz)

Time series	Length	Snippet length	Domain
TiltABP ¹	40 000	630	Human blood pressure during rapid tilts
PAMAP ²	20 002	600	Wearable accelerometer readings during various types of human physical activity
WalkRun ²	100 000	240	

¹ Imani S. *et al.* Introducing time series snippets: a new primitive for summarizing long time series. *Data Min. Knowl. Discov.* 2020. Vol. 34, no. 6. P. 1713-1743. DOI: [10.1007/s10618-020-00702-y](https://doi.org/10.1007/s10618-020-00702-y)

² Reiss A., Stricker D. Introducing a new benchmarked dataset for activity monitoring. *ISWC 2012*. P. 108–109. DOI: [10.1109/ISWC.2012.13](https://doi.org/10.1109/ISWC.2012.13)

Parallel labeling: Experiments on accuracy



Parallel labeling algorithm demonstrates good accuracy

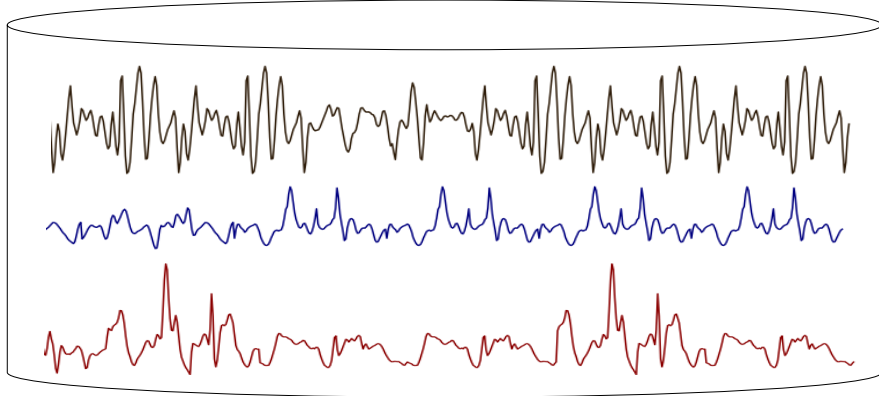
* Reiss A., Stricker D. Introducing a new benchmarked dataset for activity monitoring. ISWC 2012. P. 108–109. DOI: [10.1109/ISWC.2012.13](https://doi.org/10.1109/ISWC.2012.13)

Outline

- Introduction
- Offline time series analytics
- **Online time series analytics**
 - Employing ANNs together with parallel algorithms
 - Parallel time series labeling
 - **Online imputation of missing values and anomaly detection**
 - Parallel time series anomaly discovery
- Conclusions

Online time series imputation

Representative fragment of time series



Parallel Labeling

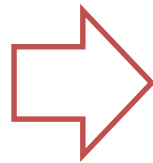
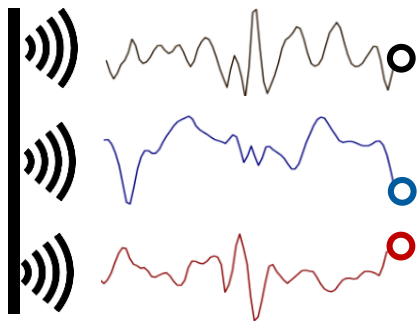


Training sets

recognize

reconstruct

Subsequence with missing values



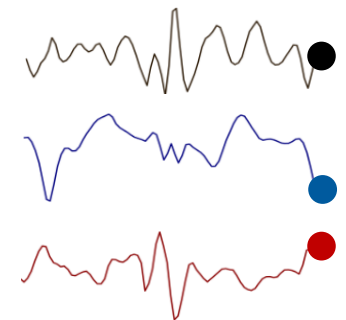
Recognizer ANNs



Reconstructor ANNs

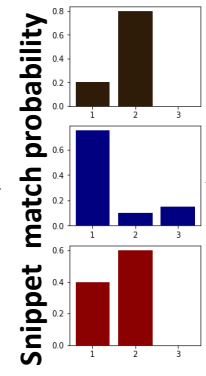
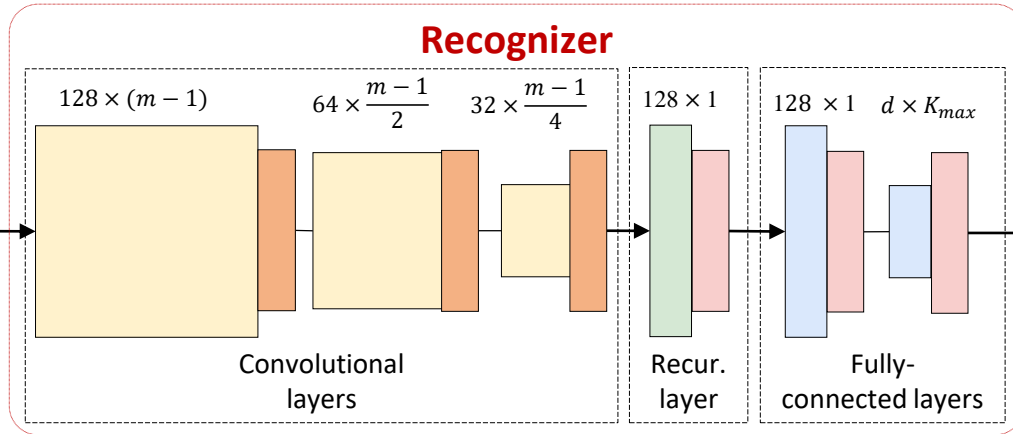
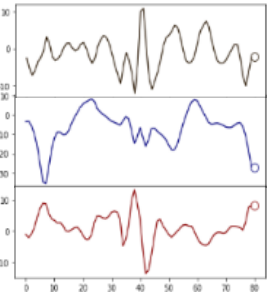


Subsequence with imputed values

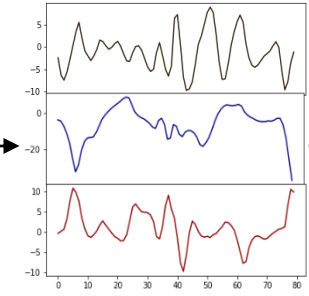


SANNI: Snippet & ANN-based Imputation

Subsequences with missing values

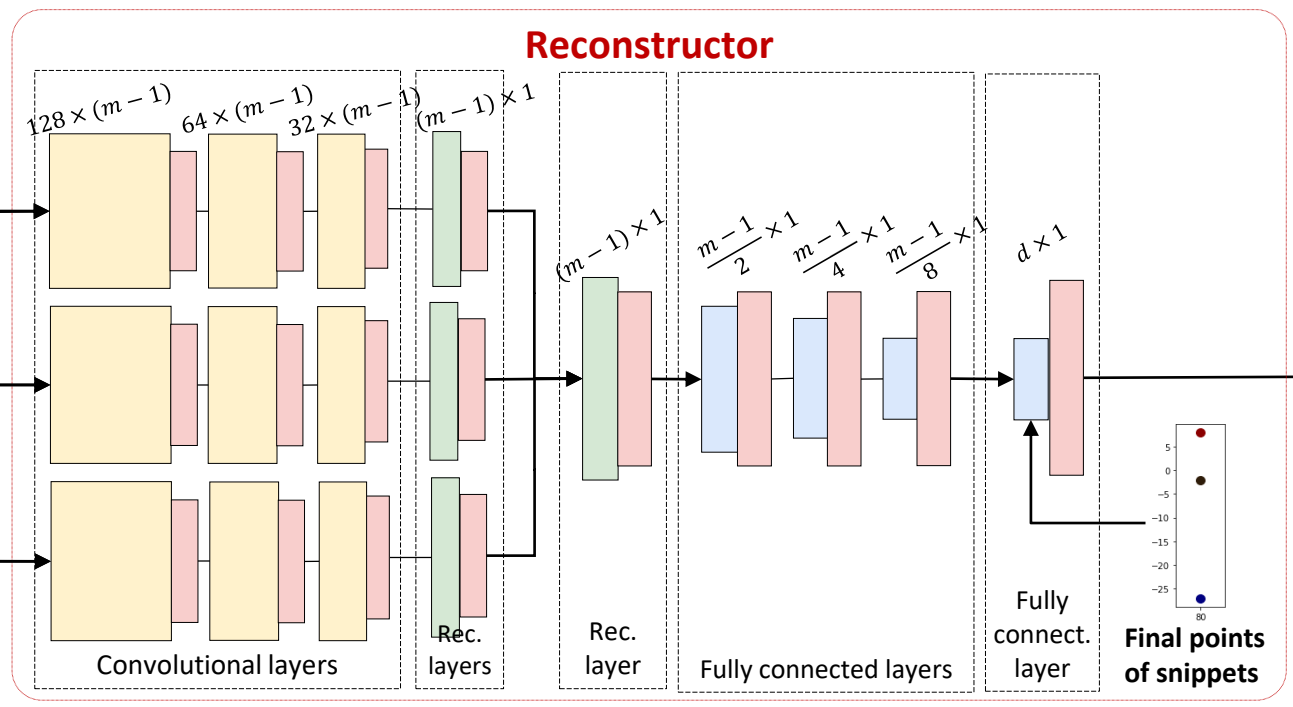
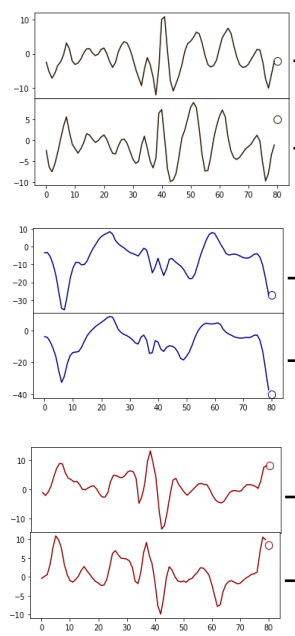


Matching snippets

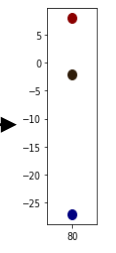


- Conv1D
- MaxPool+ReLU
- Leaky ReLU
- GRU
- Fully connected

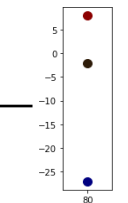
Subsequences and snippets with missing values



Imputed values



Final points of snippets

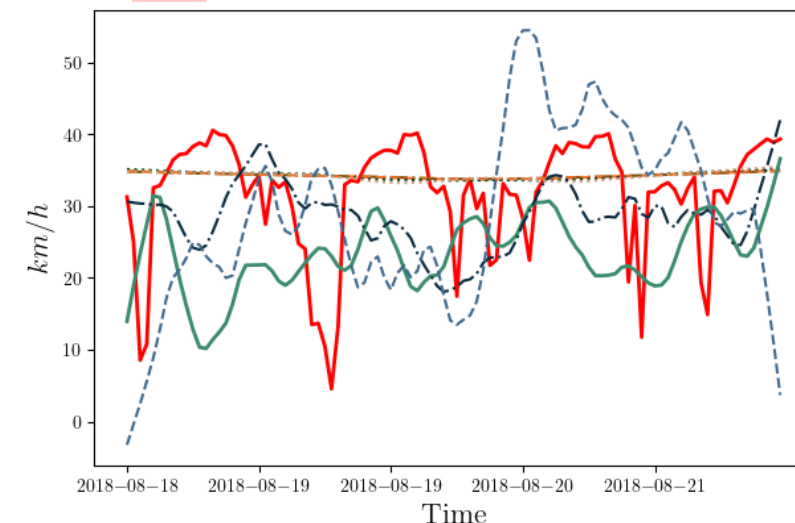
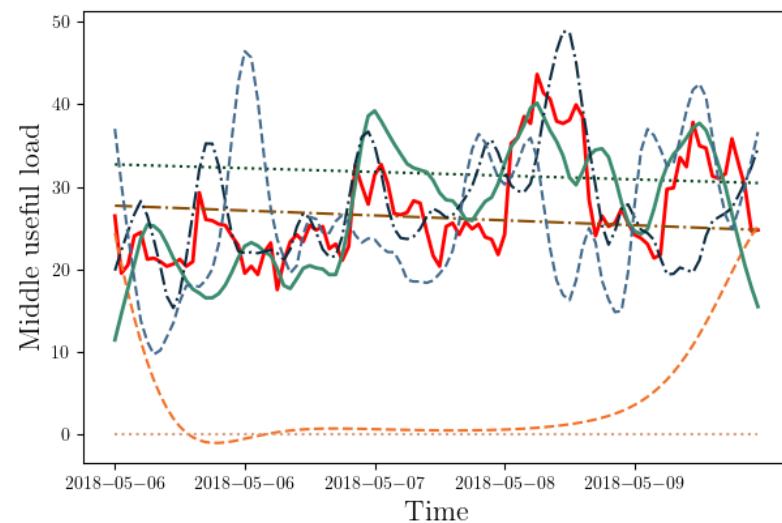
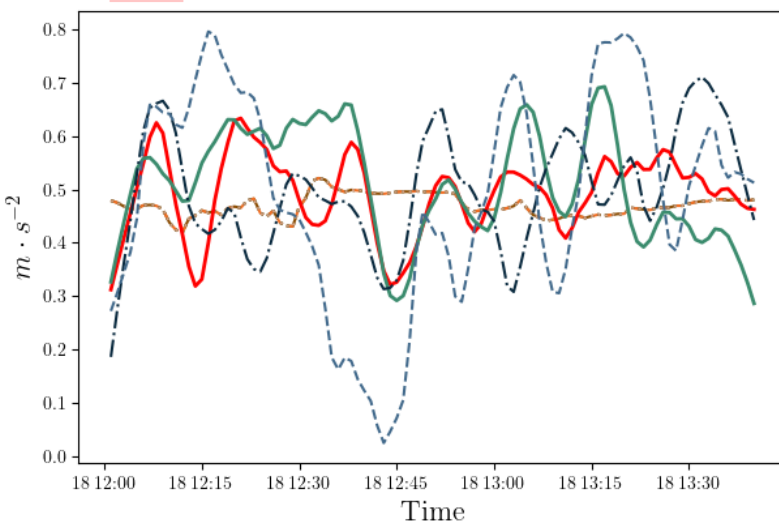
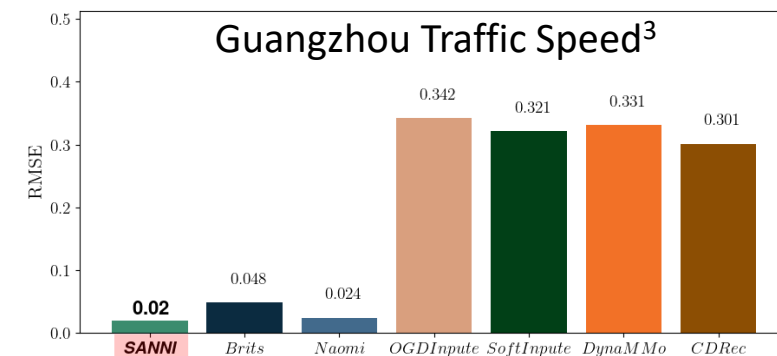
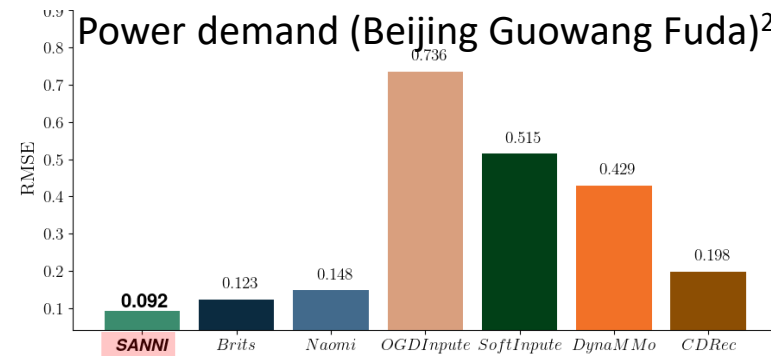
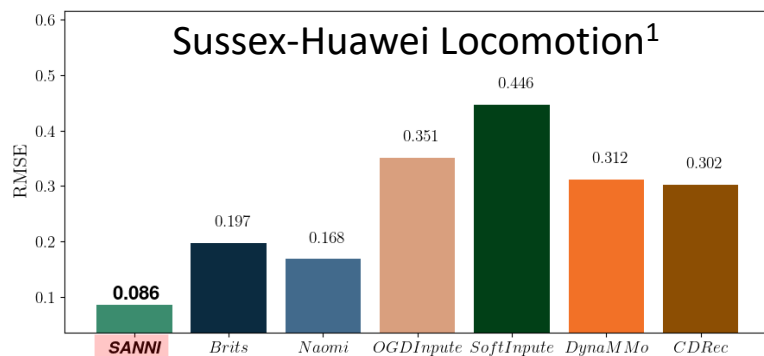


Online imputation: Experiments on accuracy

$$\text{RMSE} = \sqrt{\frac{1}{h} \sum_{i=1}^h (t_i - \tilde{t}_i)^2}$$

Lower is better

■ **SANNI**
 ■ *Naomi*
 ■ *Brits*
 ■ *MRNN*
 ■ *SoftInpute*
 ■ *CDRec*
 ■ *OGDInpute*



— *Data*
 - - - *CDRec*
 - - - *DynaMMo*
 - - - *SoftInpute*
 - - - *OGDInpute*
— **SANNI**
- - - *Brits*
- - - *Naomi*

¹ Rogger D. *et al.* The University of Sussex-Huawei Locomotion (SHL) dataset. URL: <http://www.shl-dataset.org/>

² Zhou H. *et al.* Informer: Beyond Efficient Transformer for Long Sequence Time-Series Forecasting. AAAI 2021. P. 11106-11115. DOI: [10.1609/aaai.v35i12.17325](https://doi.org/10.1609/aaai.v35i12.17325)

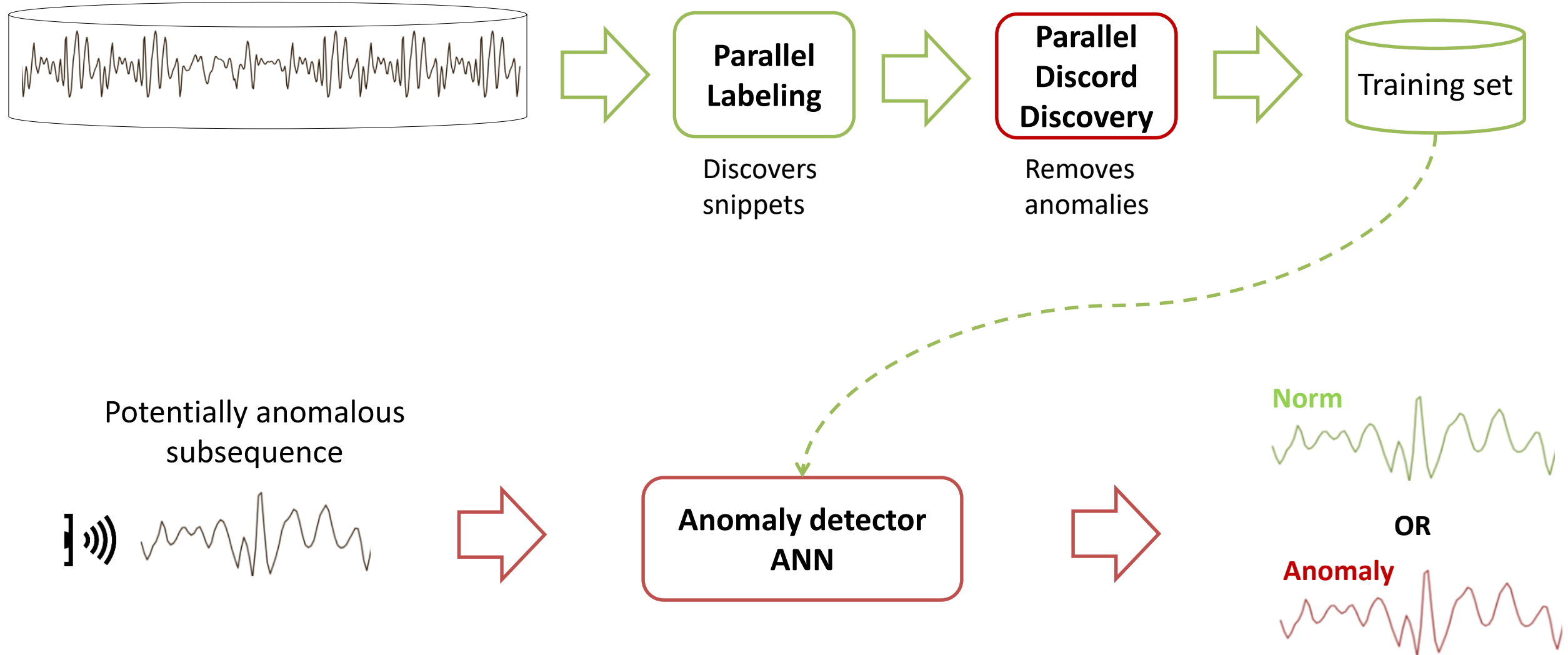
³ Chen X. *et al.* A Bayesian tensor decomposition approach for spatiotemporal traffic data imputation. 2018 Transportation Research Part C: Emerging Technologies. DOI: [10.1016/J.TRC.2018.11.003](https://doi.org/10.1016/J.TRC.2018.11.003)

Outline

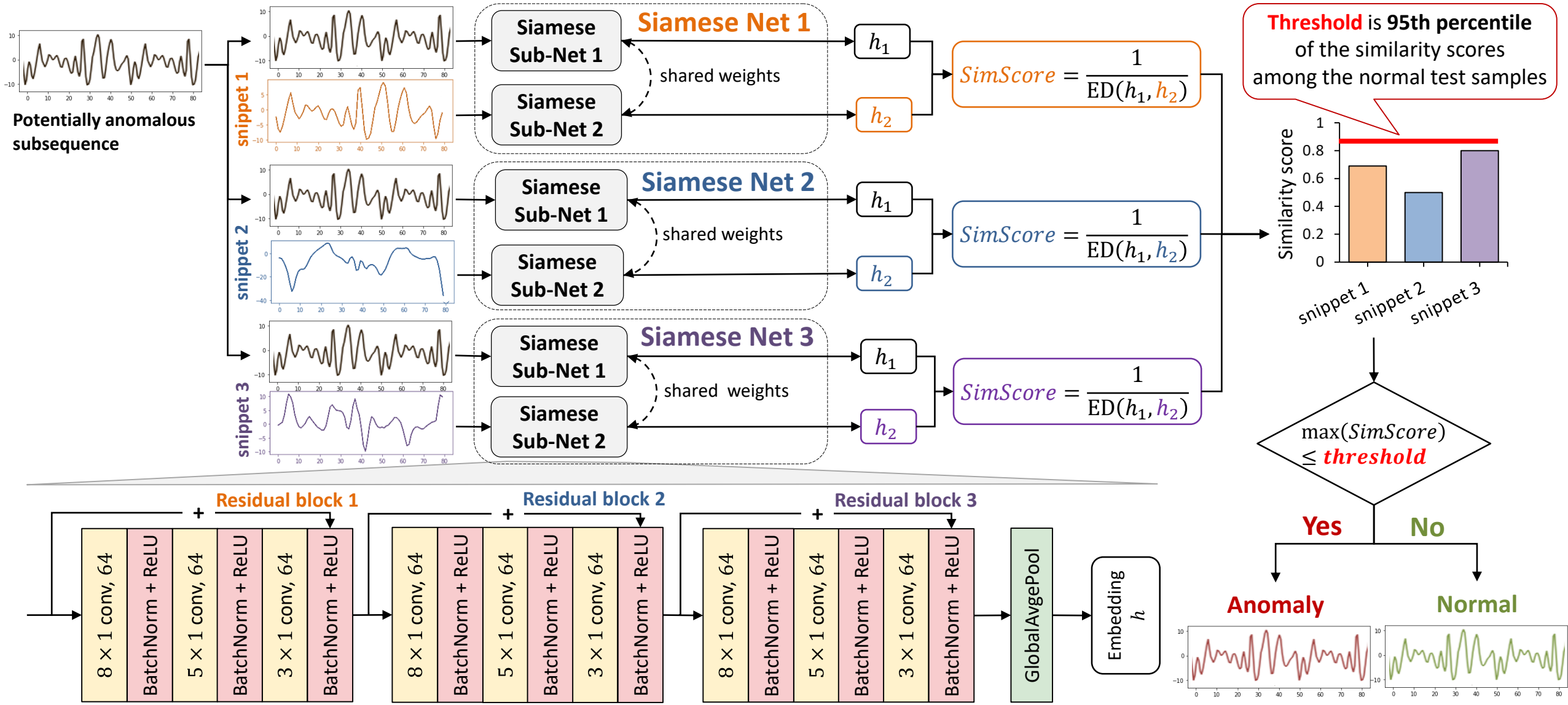
- Introduction
- Offline time series analytics
- **Online time series analytics**
 - Employing ANNs together with parallel algorithms
 - Parallel time series labeling
 - **Online** imputation of missing values and **anomaly detection**
 - Parallel time series anomaly discovery
- Conclusions

Online time series anomaly detection

Representative fragment of time series



DiSSiD: Discord, Snippet, and Siamese Net-based Detector of anomalies

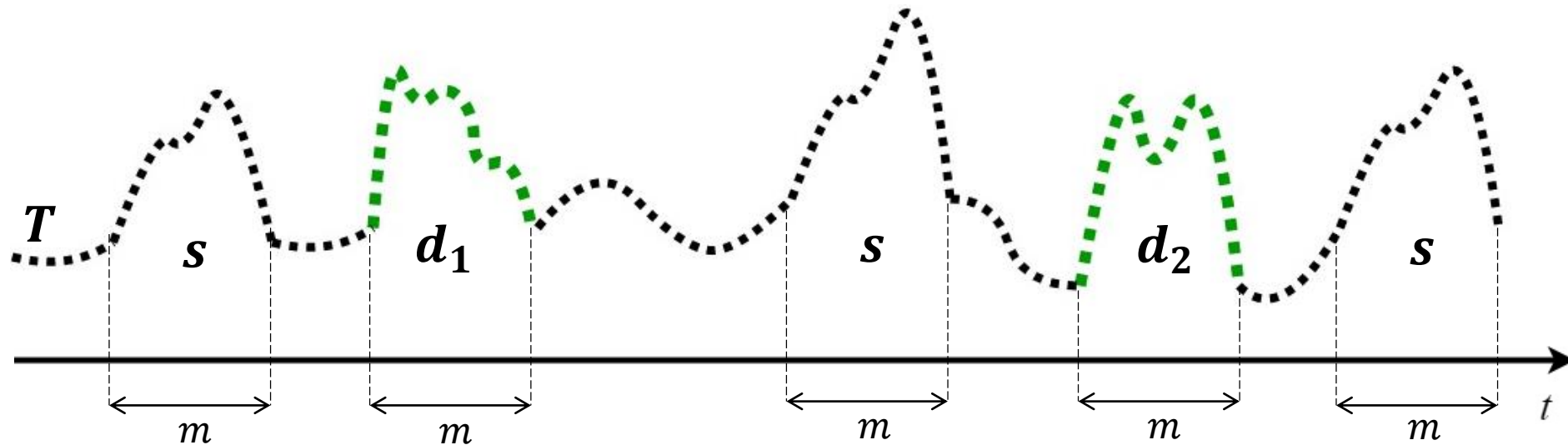


Outline

- Introduction
- Offline time series analytics
- **Online time series analytics**
 - Employing ANNs together with parallel algorithms
 - Parallel time series labeling
 - Online imputation of missing values and anomaly detection
 - **Parallel time series anomaly discovery**
- Conclusions

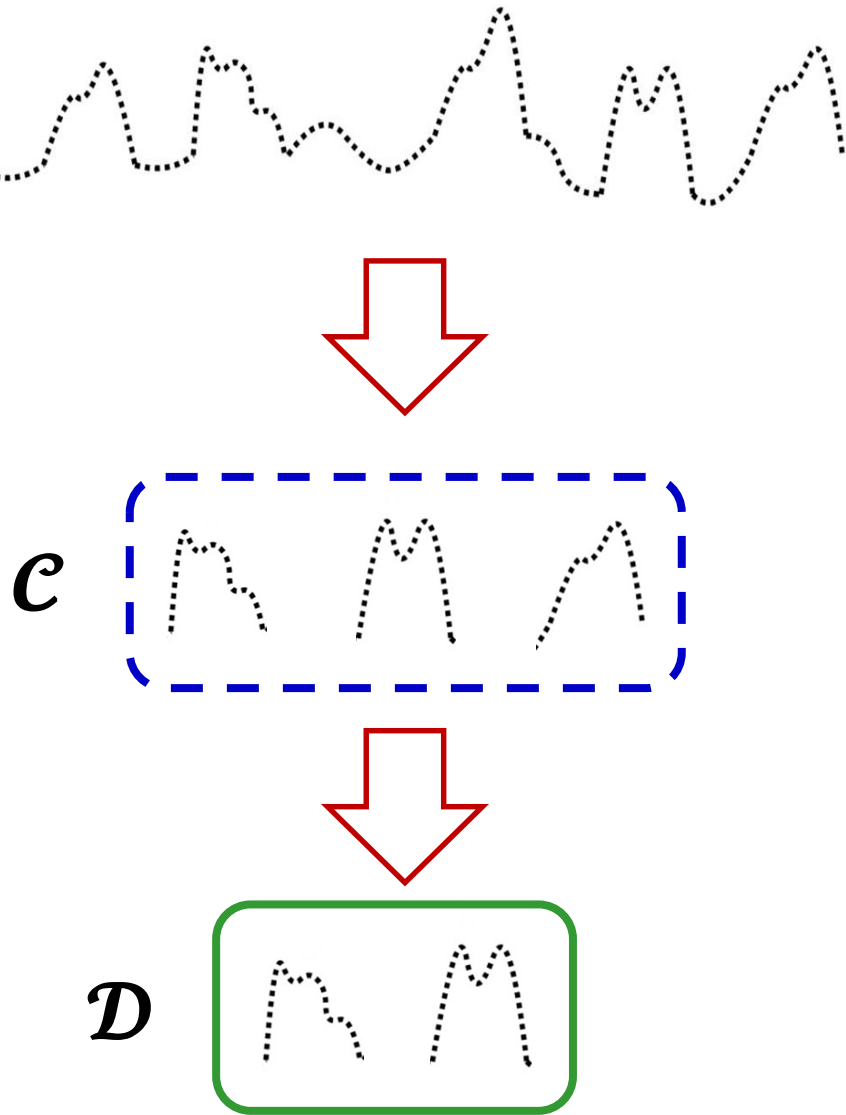
How to formalize anomalies? Discords¹

- Discord is a subsequence whose nearest neighbor is at least at a given threshold far away
- We are given: T , discord length m , threshold r
- We are to find: $\mathcal{D} = \{d_1, d_2, \dots\}$, $d_i \in \mathcal{D} \Leftrightarrow \min_{s \in T, s \cap d_i = \emptyset} \text{ED}(d_i, s) \geq r$



¹Yankov D. *et al.* Disk aware discord discovery: finding unusual time series in terabyte sized datasets. *Knowl. Inf. Syst.* 17(2): 241–262. 2008. DOI: [10.1007/s10115-008-0131-9](https://doi.org/10.1007/s10115-008-0131-9)

Discord discovery



1. Selection

Through one full scan of the time series, create a **set of candidates** to discords

2. Refinement

Through one full scan of the time series, **prune false positives** from the set above

Discord discovery: Selection

while not end of T

get current subsequence s

$isCandidate := \text{TRUE}$

for each $c_i \in \mathcal{C}$ and $s \cap c_i = \emptyset$

if $\text{ED}(s, c_i) < r$ then

$\mathcal{C} := \mathcal{C} \setminus c_i$; $isCandidate := \text{FALSE}$

if $isCandidate = \text{TRUE}$ then $\mathcal{C} := \mathcal{C} \cup s$

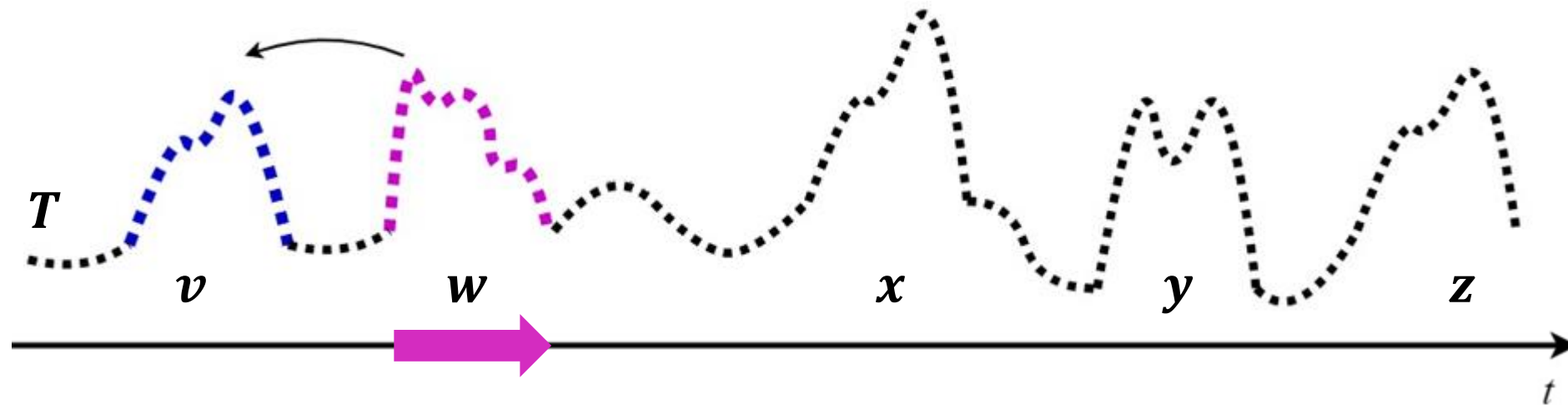
$\mathcal{C} = \{v\}$



$\text{ED}(w, v) \geq r$



$\mathcal{C} = \{v, w\}$



Discord discovery: Selection

```
while not end of  $T$ 
  get current subsequence  $s$ 
   $isCandidate := TRUE$ 
  for each  $c_i \in \mathcal{C}$  and  $s \cap c_i = \emptyset$ 
    if  $ED(s, c_i) < r$  then
       $\mathcal{C} := \mathcal{C} \setminus c_i$ ;  $isCandidate := FALSE$ 
  if  $isCandidate = TRUE$  then  $\mathcal{C} := \mathcal{C} \cup s$ 
```

$$\mathcal{C} = \{v, w\}$$

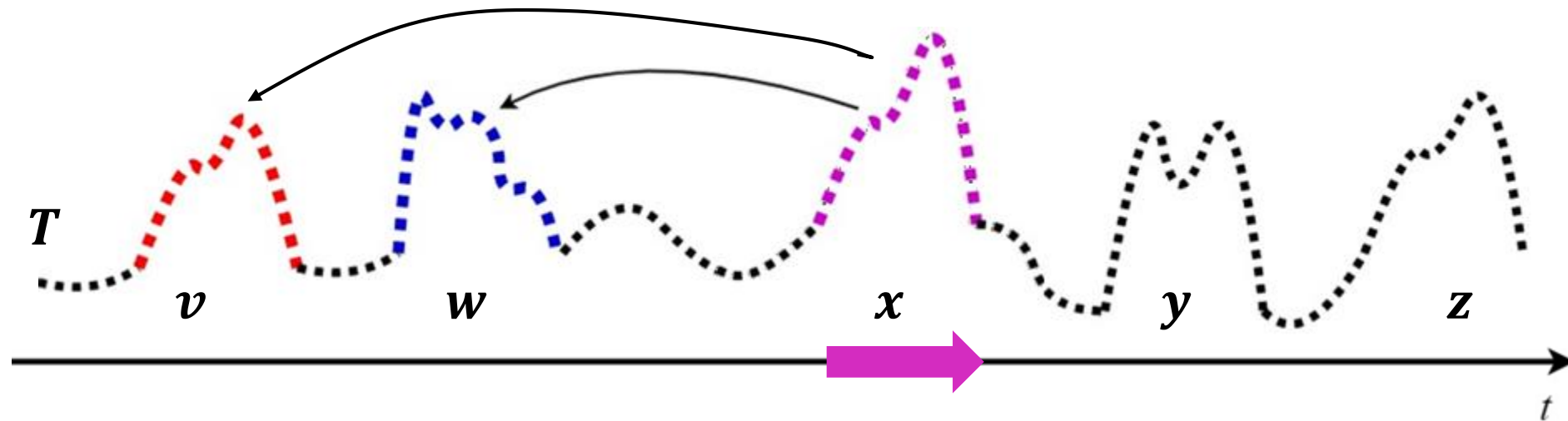


$$ED(x, v) < r$$

$$ED(x, w) \geq r$$



$$\mathcal{C} = \{w\}$$



Discord discovery: Selection

while not end of T

get current subsequence s

$isCandidate := \text{TRUE}$

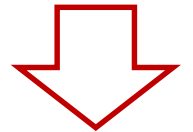
for each $c_i \in \mathcal{C}$ and $s \cap c_i = \emptyset$

if $\text{ED}(s, c_i) < r$ then

$\mathcal{C} := \mathcal{C} \setminus c_i$; $isCandidate := \text{FALSE}$

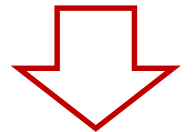
if $isCandidate = \text{TRUE}$ then $\mathcal{C} := \mathcal{C} \cup s$

$$\mathcal{C} = \{w, y\}$$

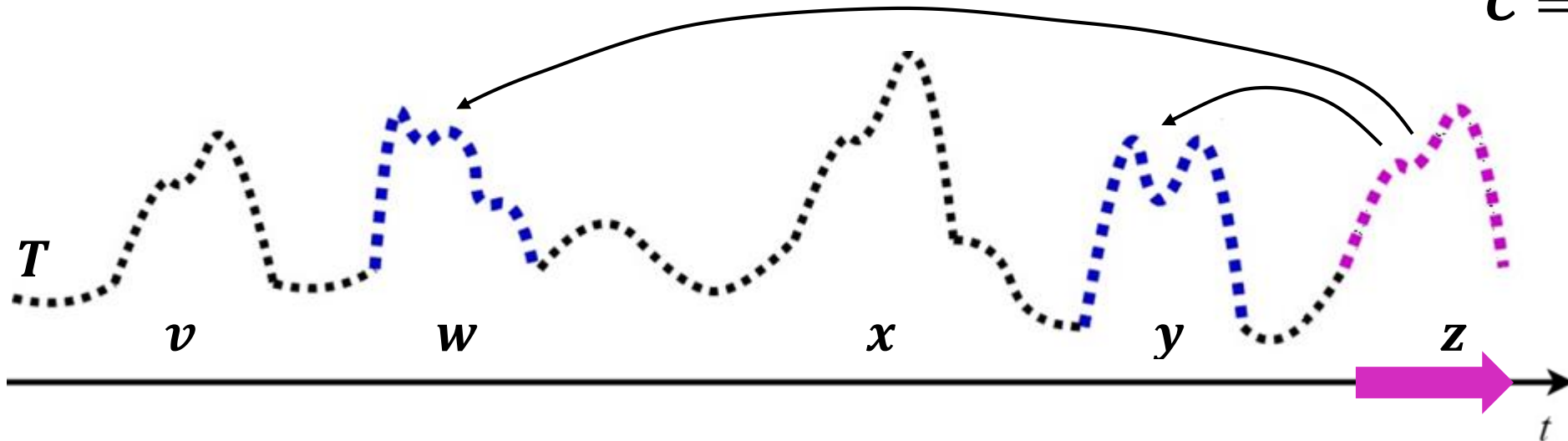


$$\text{ED}(z, w) \geq r$$

$$\text{ED}(z, y) \geq r$$



$$\mathcal{C} = \{w, y, z\}$$



Discord discovery: Refinement

$\mathcal{D} := \mathcal{C}$

while not end of T

get current subsequence s

for each $d_i \in \mathcal{D}$ and $s \cap d_i = \emptyset$

if $ED(s, d_i) < r$ then

$\mathcal{D} := \mathcal{D} \setminus d_i$

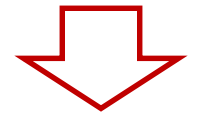
$\mathcal{D} = \{w, y, z\}$



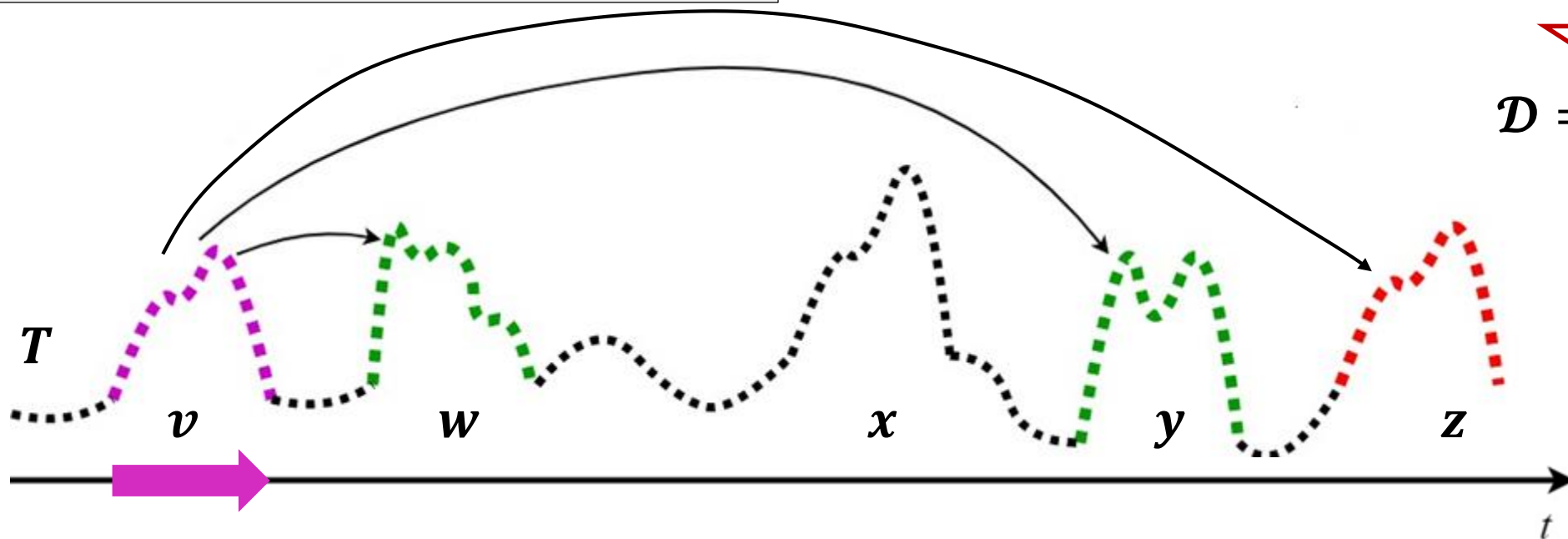
$ED(v, w) \geq r$

$ED(v, y) \geq r$

$ED(v, z) < r$



$\mathcal{D} = \{w, y\}$



Discord discovery: Refinement

$\mathcal{D} := \mathcal{C}$

while not end of T

get current subsequence s

for each $d_i \in \mathcal{D}$ and $s \cap d_i = \emptyset$

if $ED(s, d_i) < r$ then

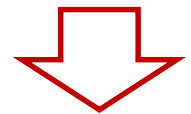
$\mathcal{D} := \mathcal{D} \setminus d_i$

$\mathcal{D} = \{w, y\}$

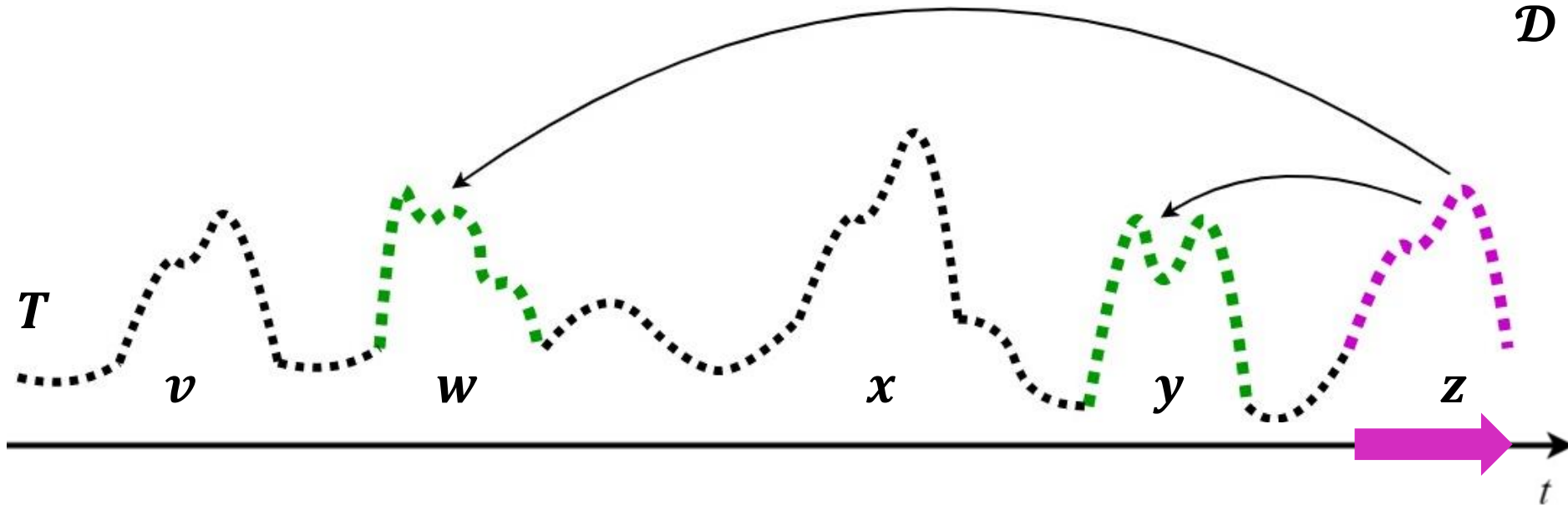


$ED(z, w) \geq r$

$ED(z, y) \geq r$



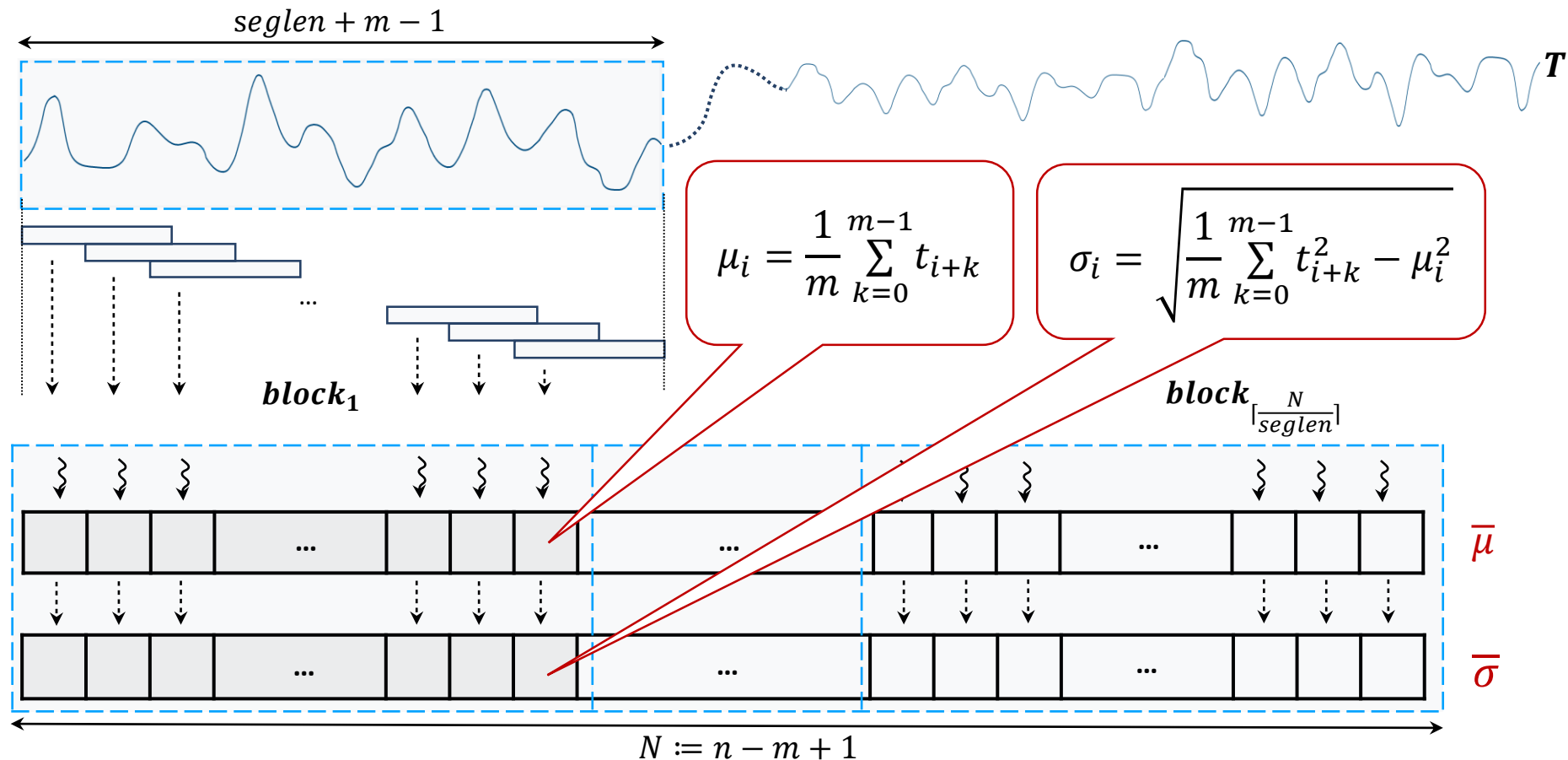
$\mathcal{D} = \{w, y\}$



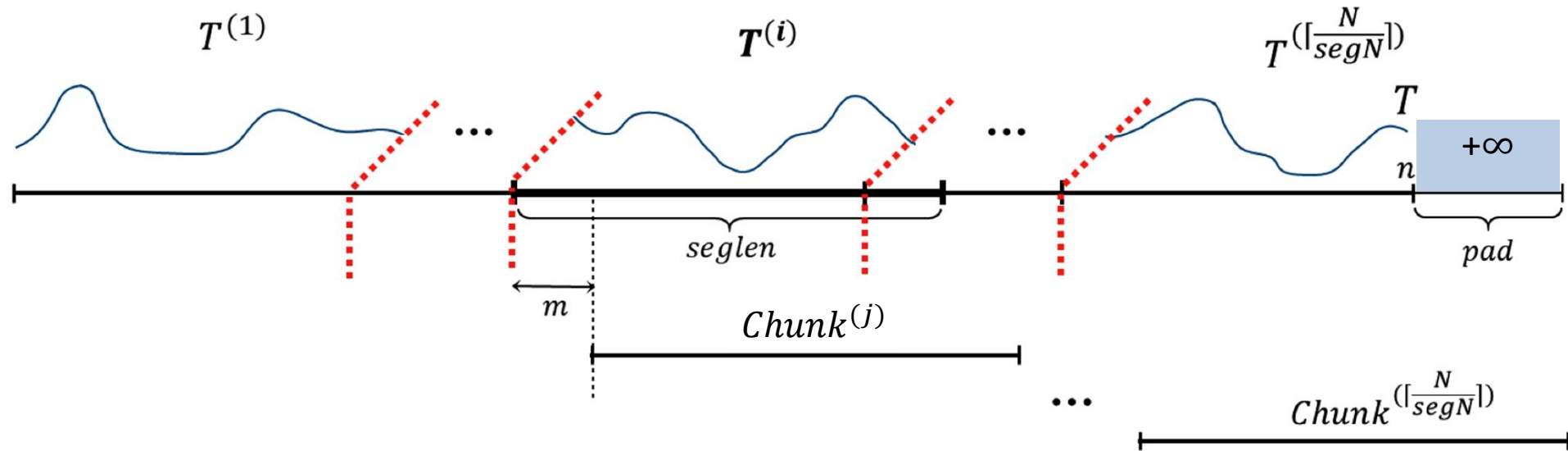
Parallel discord discovery: Preprocessing

For highest possible performance, let us use **quadratic ED**

$$ED_{\text{norm}}^2(T_{i,m}, T_{j,m}) = 2m \left(1 - \frac{\langle T_{i,m}, T_{j,m} \rangle - m\mu_i\mu_j}{m\sigma_i\sigma_j} \right)$$

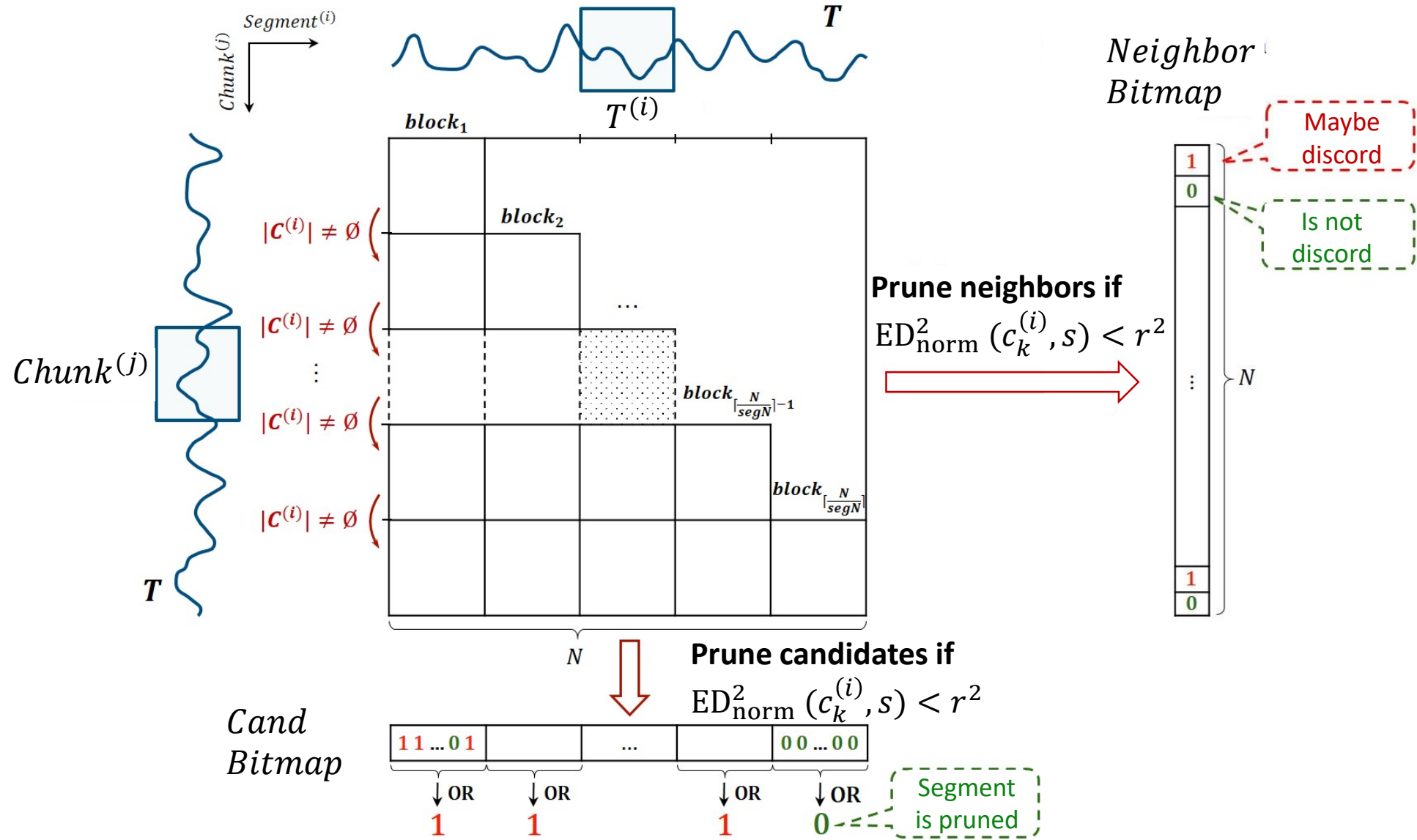


Parallel discord discovery: Segmentation



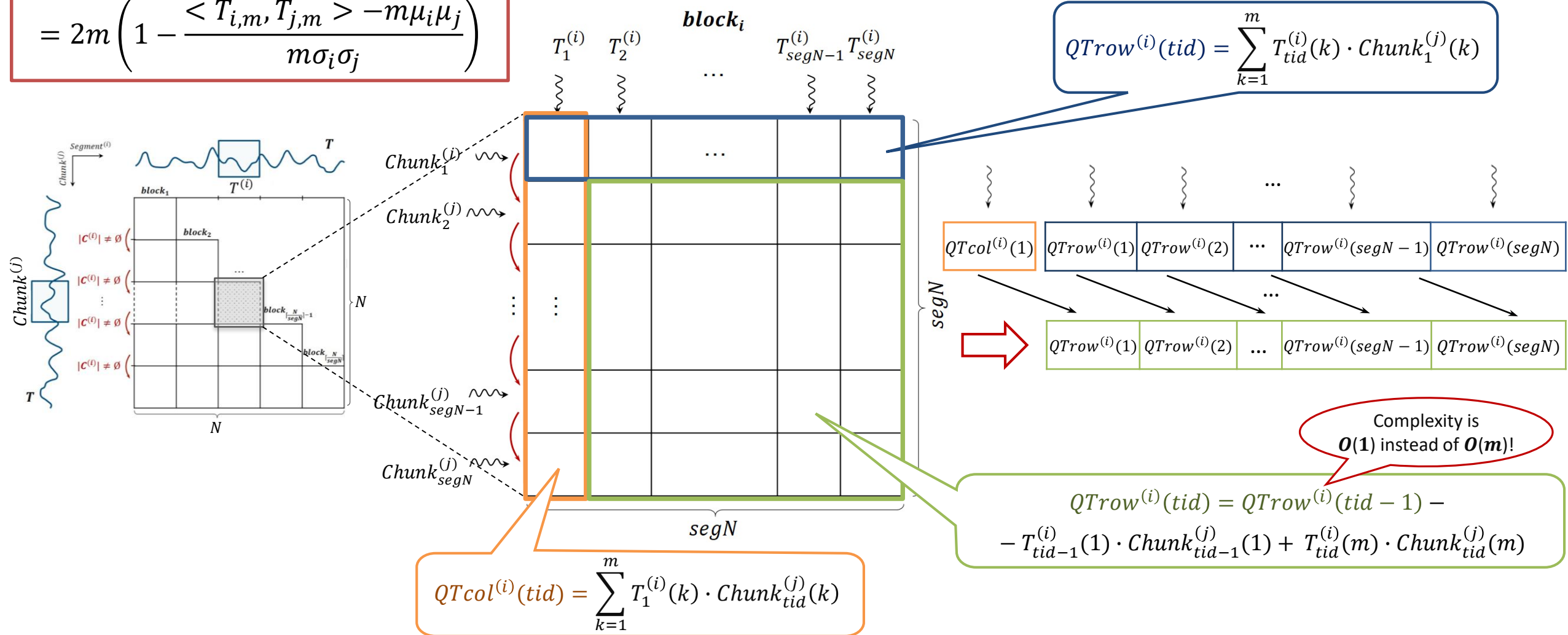
Parameter	Semantic
$T^{(i)}$	Segment to select (prune) candidates
$seglen = segN + m - 1$	Segment length
$segN = k \cdot warpsize$	# candidates in a segment
$warpsize = 32$	# threads in a group within a thread block
$Chunk^{(j)}$	Interval to test its subsequences against a segment candidates
pad	# dummy elements

Parallel discord discovery: Selection (blocks)



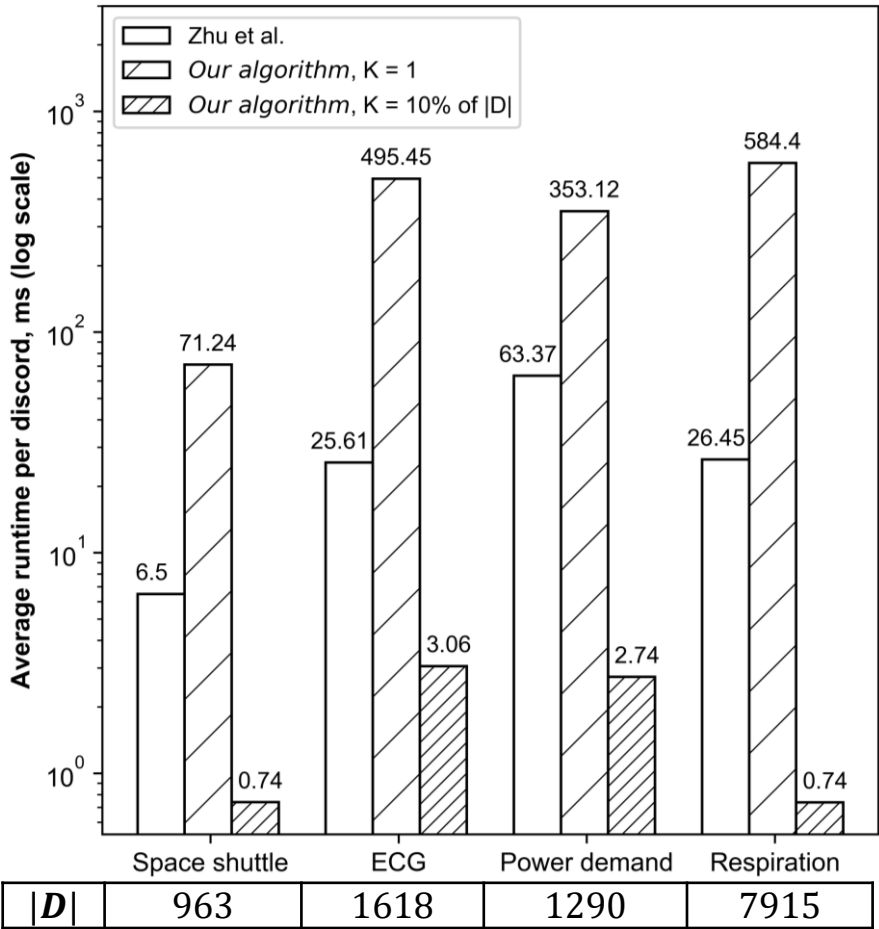
Parallel discord discovery: Selection (threads)

$$d_{i,j} = ED_{\text{norm}}^2(T_{i,m}, T_{j,m}) = 2m \left(1 - \frac{\langle T_{i,m}, T_{j,m} \rangle - m\mu_i\mu_j}{m\sigma_i\sigma_j} \right)$$



Parallel discords discovery: Experiments on performance

Parallel algorithm
is up to **30 times ahead** of rival
w.r.t. average running time to discover one discord



Hardware: NVIDIA Tesla P100 (3584 cores @1.19 GHz)

Time series	Length	Discord	Domain
Space shuttle	5 000	150	Measurements of a sensor on the NASA spacecraft ¹
ECG	45 000	200	ECG of an adult patient ²
Power demand	33 220	750	Annual energy consumption of an office ³
Respiration	24 125	250	Human breathing by chest expansion ⁴

D	963	1618	1290	7915
---	-----	------	------	------

¹ Ferrell B., et al. NASA shuttle valve data 2005. URL: <http://www.cs.fit.edu/~pkc/nasa/data/>

² Goldberger A., et al. PhysioBank, PhysioToolkit, and PhysioNet: components of a new research resource for complex physiologic signals. Circulation 101(23): 215–220. DOI: [10.1161/01.CIR.101.23.e215](https://doi.org/10.1161/01.CIR.101.23.e215)

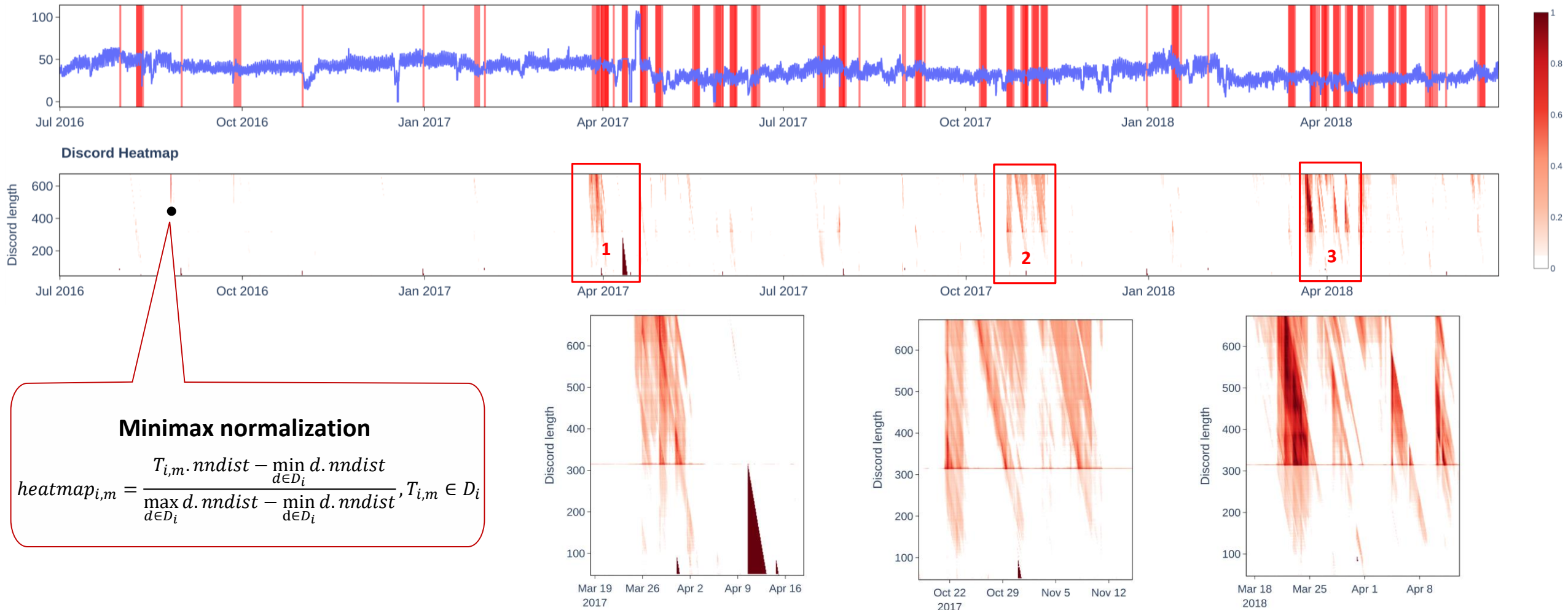
³ van Wijk J.J., et al. Cluster and calendar based visualization of time series data. INFOVIS'99: 4–9. DOI: [10.1109/INFVIS.1999.801851](https://doi.org/10.1109/INFVIS.1999.801851)

⁴ Keogh E., et al. HOT SAX: Finding the most unusual time series subsequence: Algorithms and applications. ICDM 2004: 440–449. URL: <http://www.cs.ucr.edu/~eamonn/discords/>

⁵ Zhu B. et al. A GPU Acceleration framework for motif and discord based pattern mining. IEEE Transactions on Parallel and Distributed Systems 32(8): 1987–2004. 2021. DOI: [10.1109/TPDS.2021.3055765](https://doi.org/10.1109/TPDS.2021.3055765)

Parallel discords discovery: Experiments on accuracy

2-year power demand (Beijing Guowang Fuda Sci. & Tech. Dev. Co.)¹



¹Zhou H. *et al.* Informer: beyond efficient transformer for long sequence time-series forecasting. AAAI 2021: 11106-11115. DOI: [10.1609/aaai.v35i12.17325](https://doi.org/10.1609/aaai.v35i12.17325).

Conclusions: We are ready for further collaboration

- Applying time series analytics to Smart DBMS
 - Anomaly detection to monitor database activities
 - Hardware lifecycle prediction
 - Workload prediction based on resource usage patterns
- Embedding time series analytics into DBMS
 - In-DBMS matrix profile support
 - Table data imputation: on the fly and/or in background
- Online time series analytics of mobile users
 - Activity recognition
 - Anomaly detection
 - Data imputation



Big Data
and Machine Learning
Laboratory



**Mikhail
Zymbler**
Dr.Sci.
Assist. Prof.



**Yana
Kraeva**
MSc,
2yr PhD student



**Andrey
Goglachev**
MSc,
1yr PhD student



**Alexey
Yurtin**
MSc,
1yr PhD student