

Scientific seminar
on information technologies

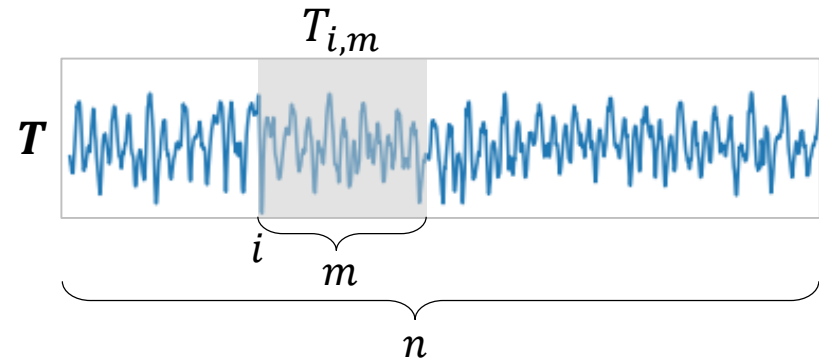
Parallel Algorithm for Discovery Typical Subsequences of a Time Series on Graphical Processor

Andrey Goglachev, Mikhail Zymbler

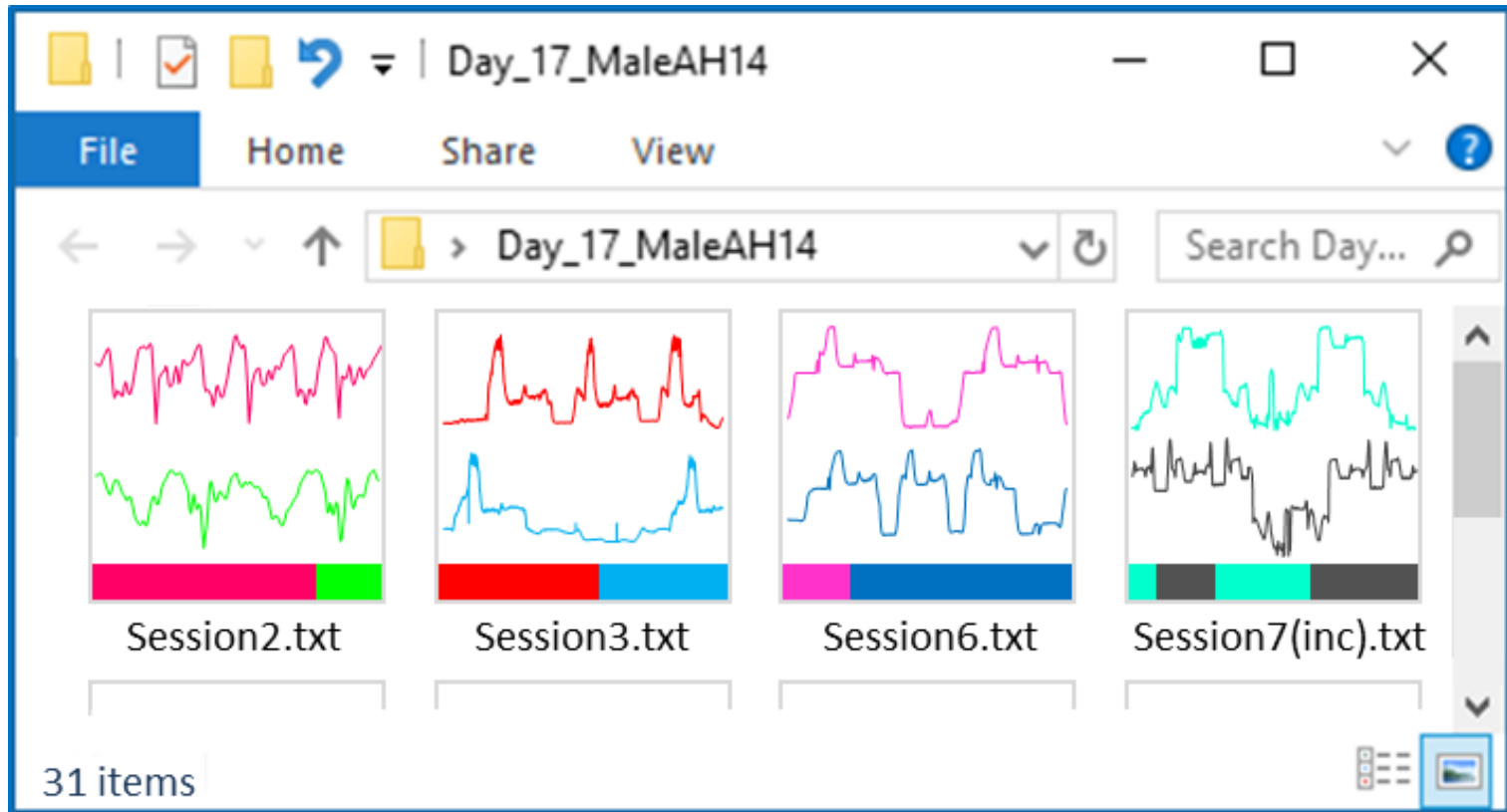
Chelyabinsk–2021

Informal Problem Statement

- We are given:
 - an n -length time series T
 - a subsequence length m
- We must find:
 - a set of subsequences that reflects the respective process/activity
- Application: annotating and visualization of long time series
 - monitoring of human functional diagnostics indicators;
 - monitoring the technical conditions of complex machines and mechanisms;
 - etc.



Examples





Summarizing the patient's motor activity according to the indications of the hip accelerometer

Examples

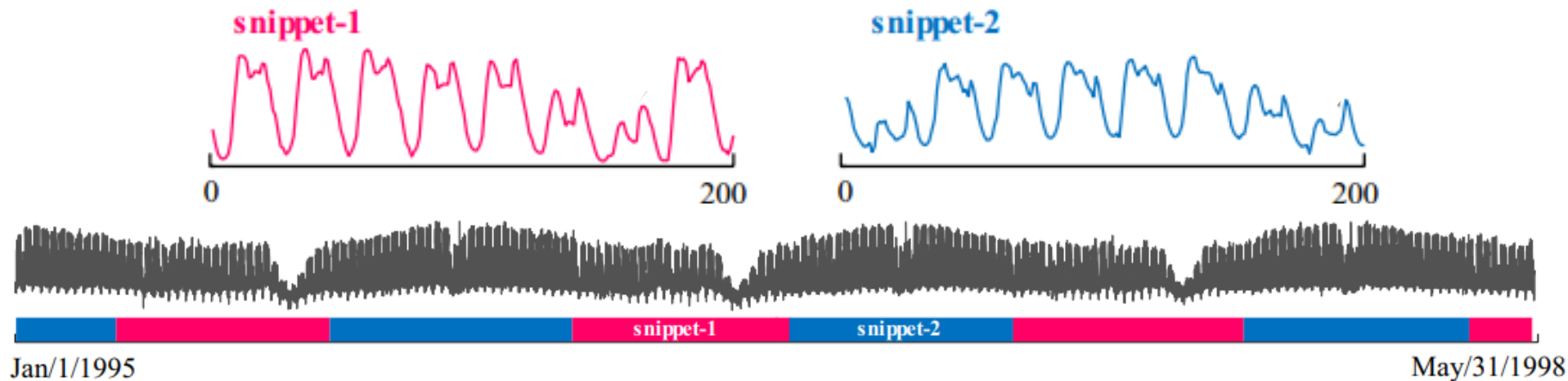


Summarizing the patient's motor activity
according to the indications of the chest accelerometer

Patient Smith slept for 7.2 hours. This ten-second snippet () accounts for 78% of his respiration, and this () ten-second snippet accounts for 17% of his respiration. His maximum temperature was 98.7°...

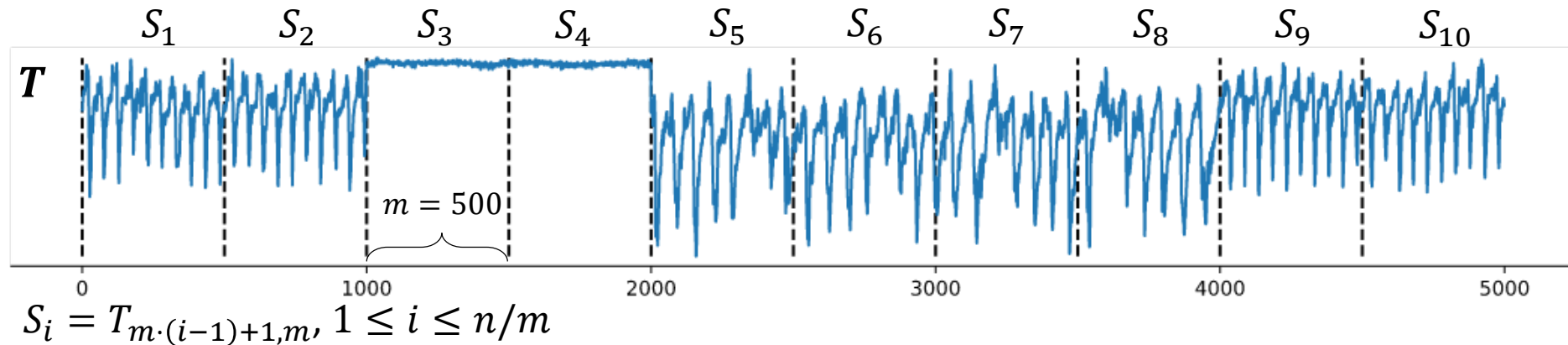
Summarizing the patient's respiratory activity
in studies of apnea syndrome

Examples



Summary of hourly energy consumption in Italy for 3 years.
Typical subsequences are weekly intervals in warm and cold seasons

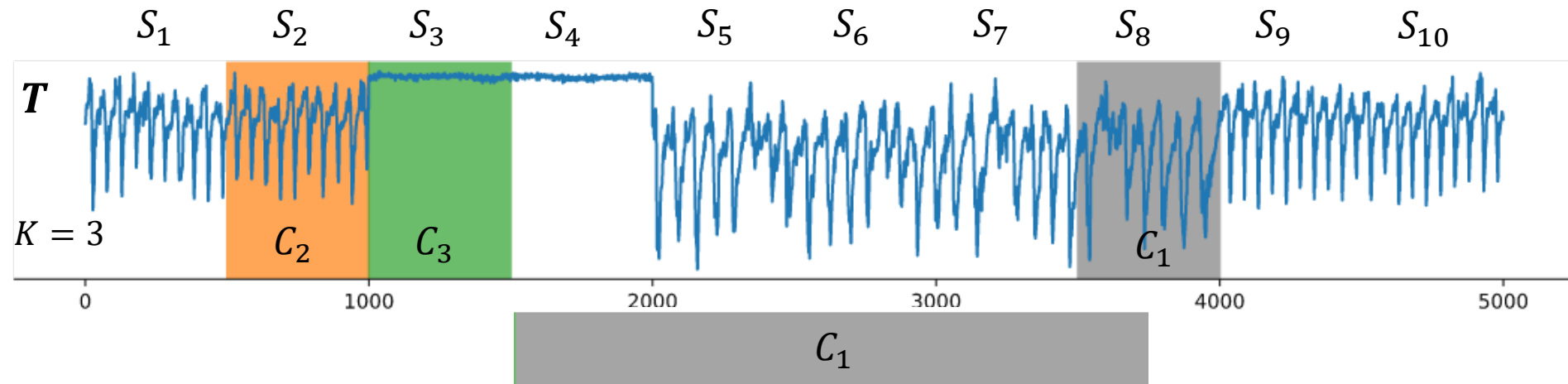
Formalization: The Snippet Concept*



1. Let us represent a time series as a set of n/m -length non-overlapped segments
 - if n is not a multiple of m , then pad the time series right by zeroes

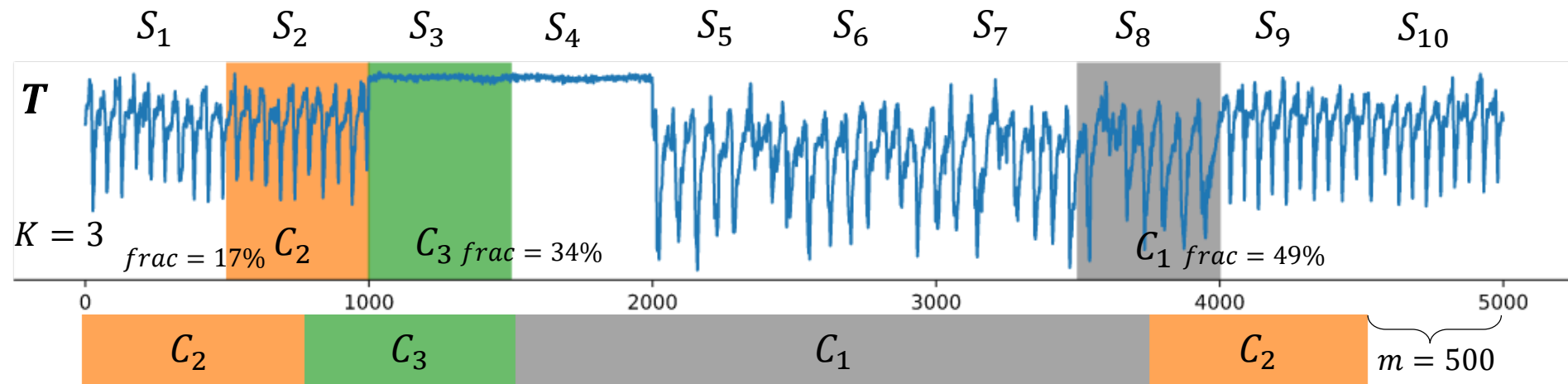
* Imani S., Madrid F., Ding W., Crouter S.E., Keogh E.J. Introducing time series snippets: a new primitive for summarizing long time series. Data Min. Knowl. Discov. 34(6): 1713-1743 (2020). doi: [10.1007/s10618-020-00702-y](https://doi.org/10.1007/s10618-020-00702-y)

Formalization: The Snippet Concept



1. Let us represent a time series as a set of n/m -length non-overlapped segments
2. For each segment, let us find the most similar subsequences (**nearest neighbors**)

Formalization: The Snippet Concept



1. Let us represent a time series as a set of n/m -length non-overlapped segments
2. For each segment, let us find the most similar subsequences (**nearest neighbors**)
3. Let us identify the segment (**snippet**) by its nearest neighbors
4. Let us take the top- K snippets in descending order of the number of their nearest neighbors (**coverage**)

MPdist*: A Subsequence Similarity Measure

Two ***m*-length time series** are the more similar by the **MPdist measure**,

the more ***l*-length** ($3 \leq l \leq m$) **normalized subsequences** close to each other by the **Euclidean metric**, are in them

Metric

- Measure
1. Identity of indiscernibles: $d(x, y) = 0 \Leftrightarrow x = y$
 2. Symmetry: $d(x, y) = d(y, x)$
 3. Triangle inequality: $d(x, z) \leq d(x, y) + d(y, z)$

* Gharghabi S., Imani S., Bagnall A.J., Darvishzadeh A., Keogh E.J.: Matrix Profile XII: MPdist: A Novel Time Series Distance Measure to Allow Data Mining in More Challenging Scenarios. ICDM 2018: 965-970. DOI: [10.1109/ICDM.2018.00119](https://doi.org/10.1109/ICDM.2018.00119)

Normalization

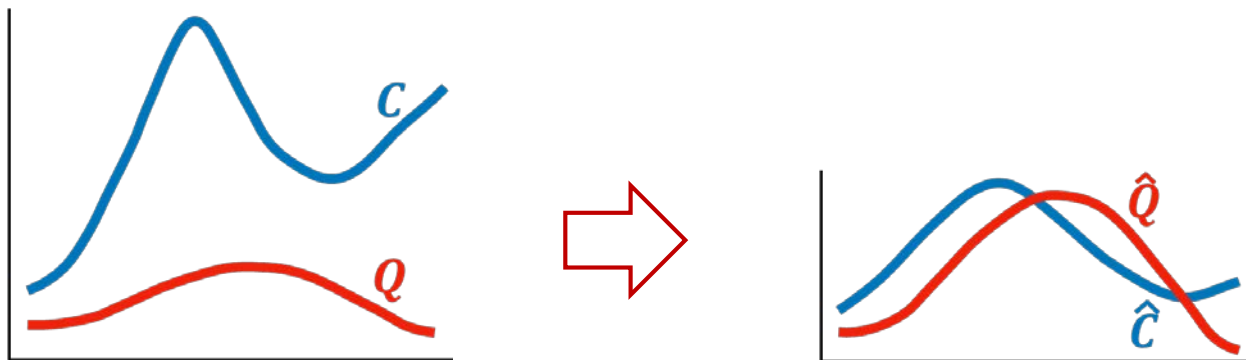
- Provides a correct comparison of subsequences with different amplitudes

$$\hat{T} = (\hat{t}_1, \dots, \hat{t}_m)$$

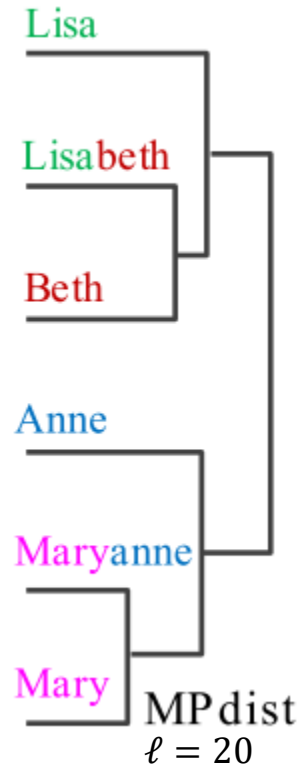
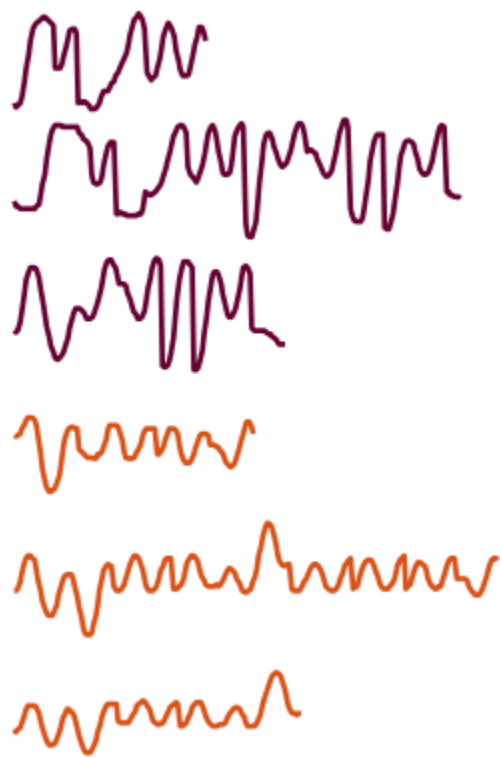
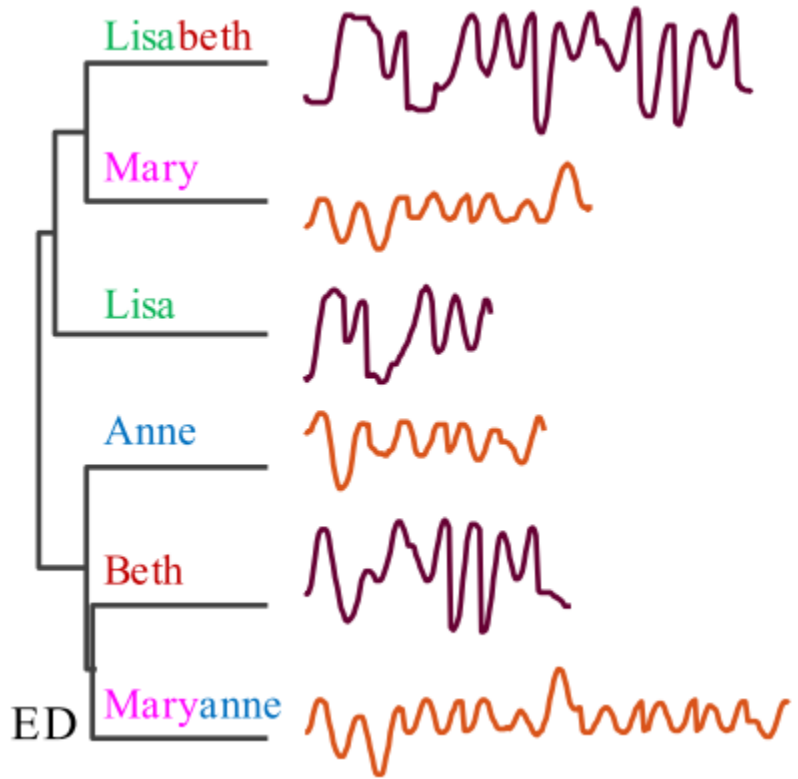
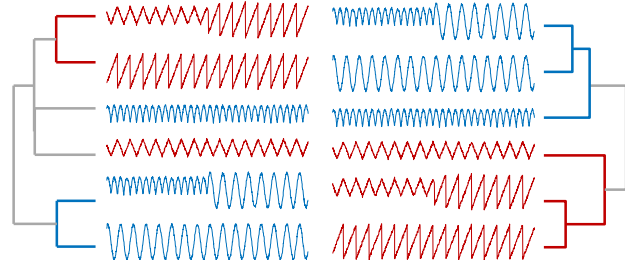
$$\hat{t}_i = \frac{t_i - \mu}{\sigma}$$

$$\mu = \frac{1}{n} \sum_{i=1}^m t_i$$

$$\sigma = \sqrt{\frac{1}{m} \sum_{i=1}^m t_i^2 - \mu^2}$$

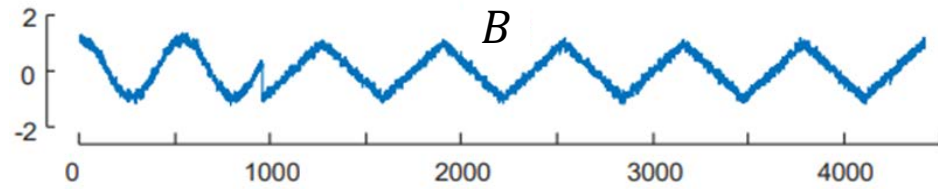
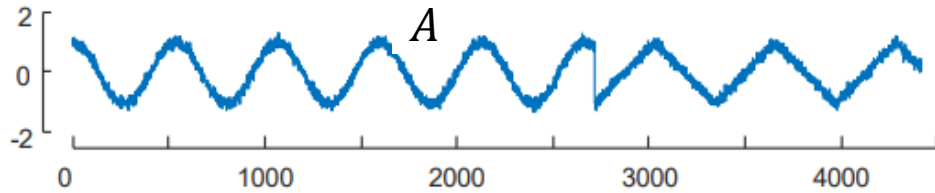


MPdist vs. Euclid



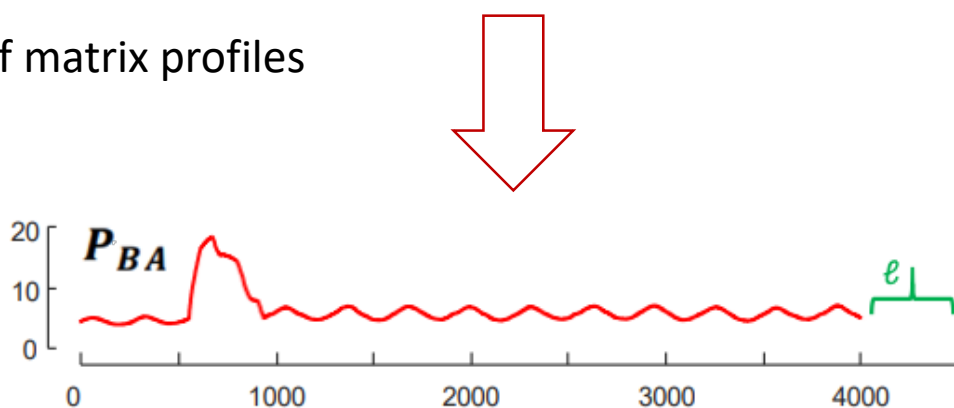
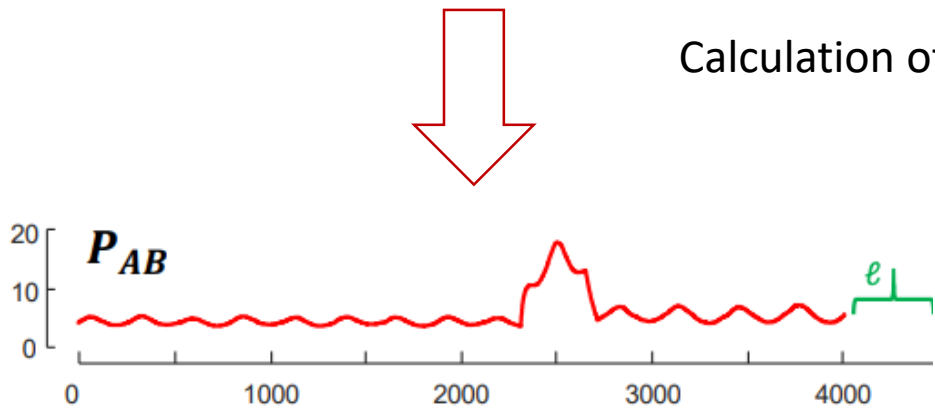
The MPdist Measure

$$|A| = |B| = m$$



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

Calculation of matrix profiles



$$\{P_{AB}(i) = \text{ED}_{\text{norm}}(A_{i,\ell}, B_{j,\ell})\}_{i=1}^{m-\ell+1},$$

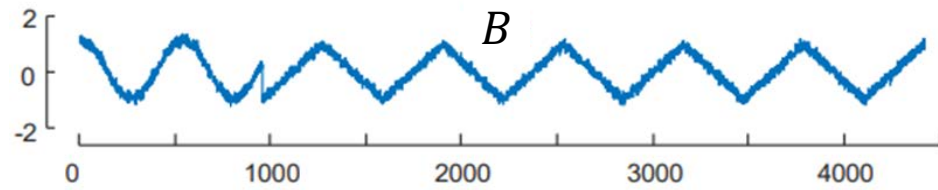
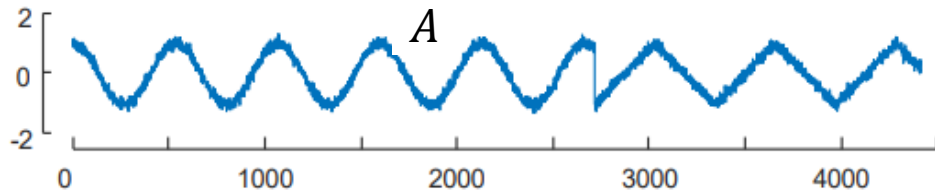
$$B_{j,\ell} = \arg \min_{1 \leq q \leq m-\ell+1} \text{ED}_{\text{norm}}(A_{i,\ell}, B_{q,\ell})$$

$$\{P_{BA}(i) = \text{ED}_{\text{norm}}(B_{i,\ell}, A_{j,\ell})\}_{i=1}^{m-\ell+1},$$

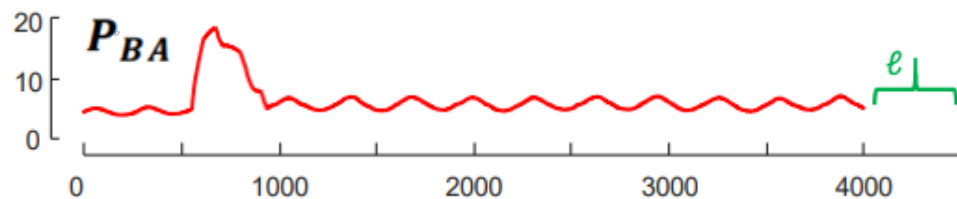
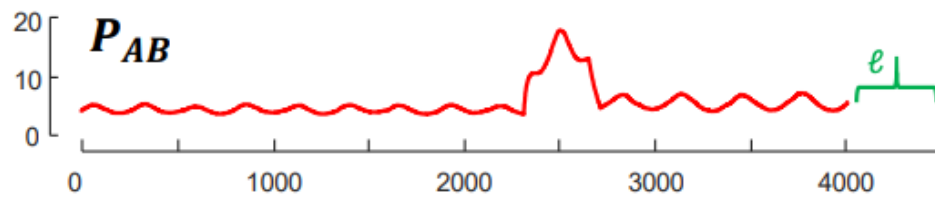
$$A_{j,\ell} = \arg \min_{1 \leq q \leq m-\ell+1} \text{ED}_{\text{norm}}(B_{i,\ell}, A_{q,\ell})$$

The MPdist Measure

$$|A| = |B| = m$$

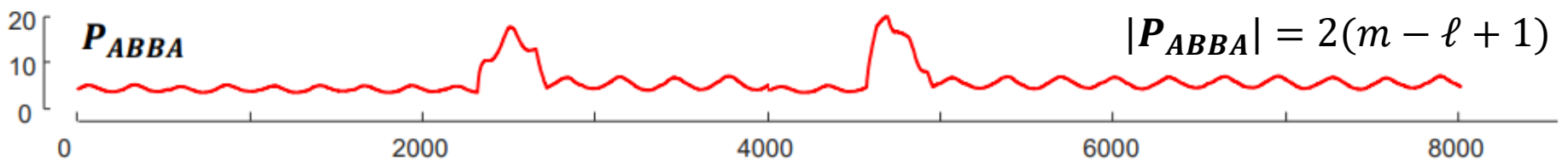


$$3 \leq \ell \leq m$$

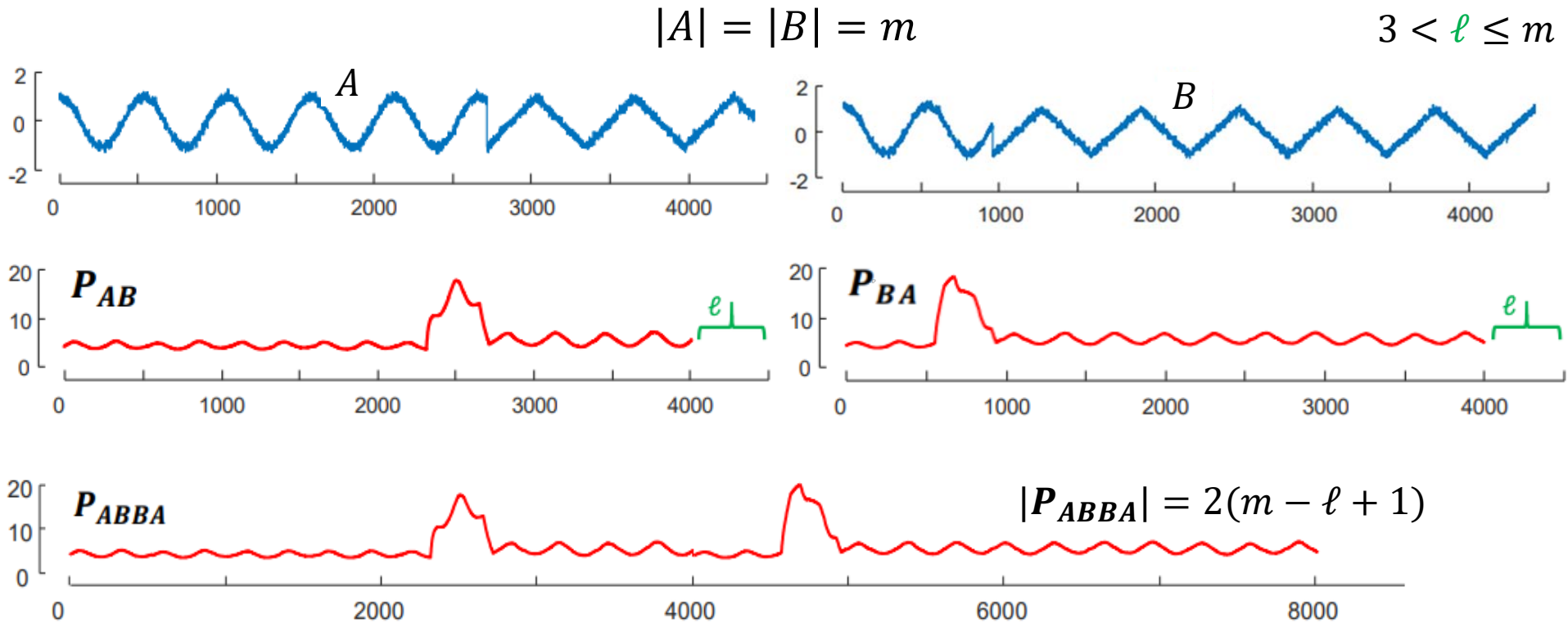


Concatenation of matrix profiles

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



The MPdist Measure



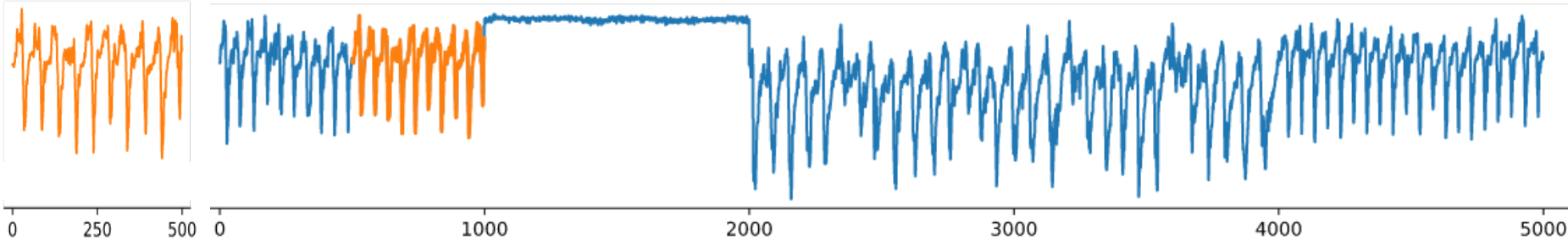
$$\text{MPdist}(A, B, \ell) = \begin{cases} \text{Sorted}P_{ABBA}(k), & |P_{ABBA}| > k \\ \text{Sorted}P_{ABBA}(2(m - \ell + 1)), & |P_{ABBA}| \leq k \end{cases}$$

где $k = \lceil 0.05 \cdot 2m \rceil = \lceil 0.1m \rceil$.

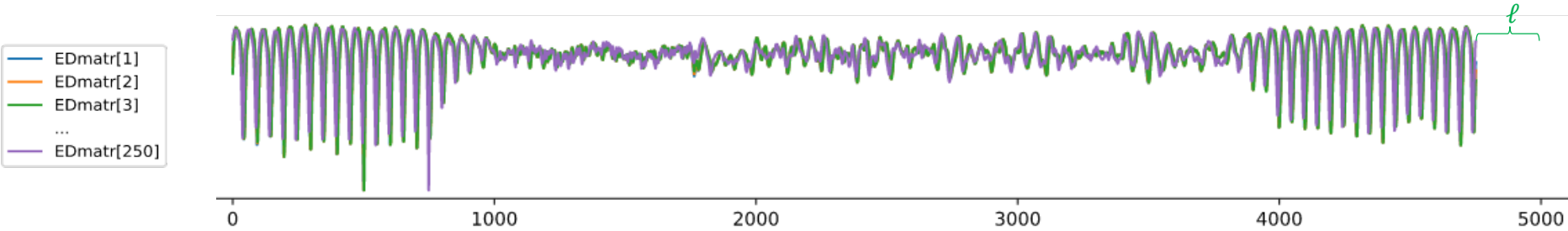
MPdist Profile of a Segment

Segment S
 $|S| = m$

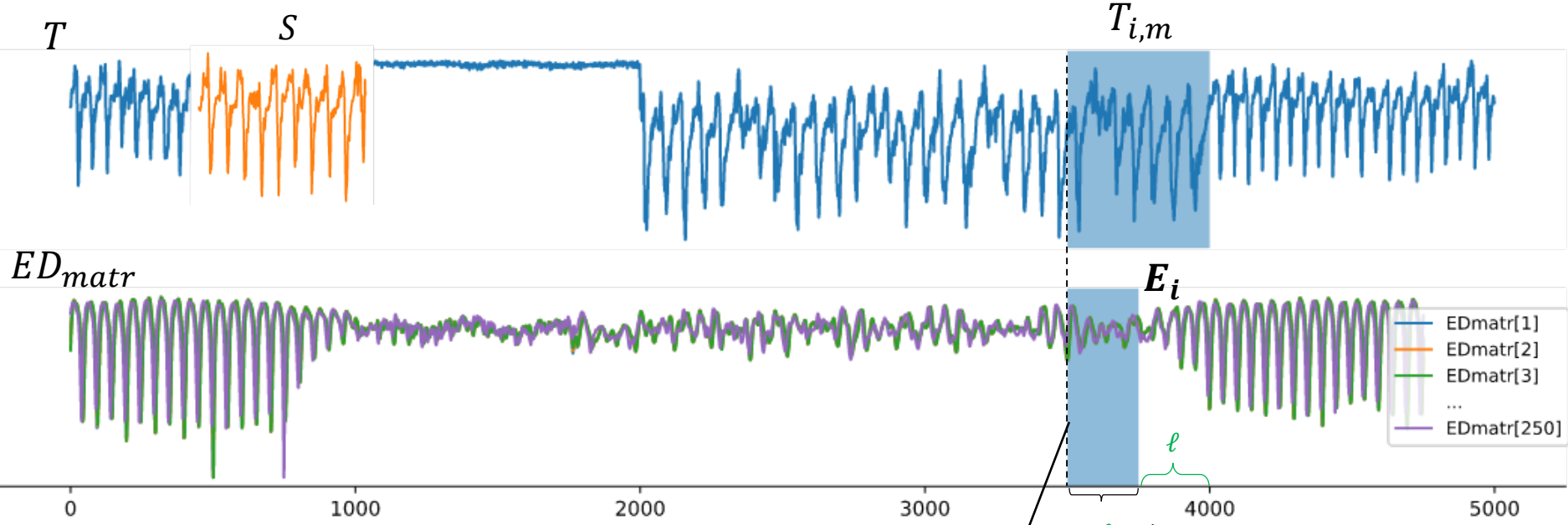
Time series T
 $|T| = n \gg m$



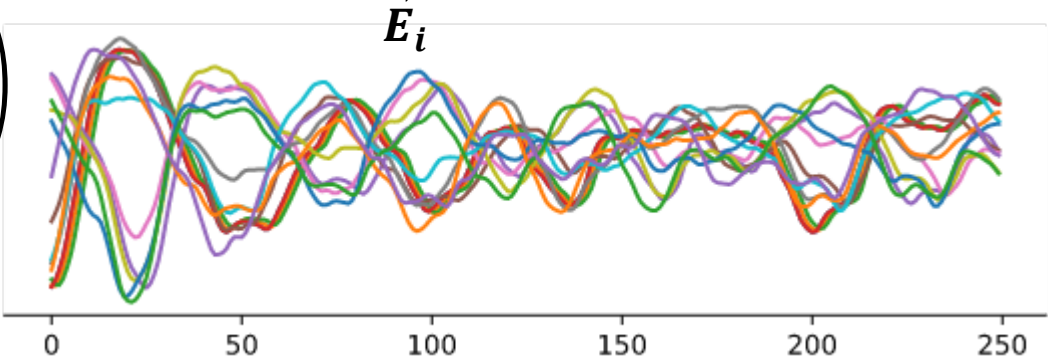
$$ED_{matr}(S, T, \ell) = \begin{pmatrix} d_{1,1} & \dots & d_{1,n-\ell+1} \\ \dots & \ddots & \dots \\ d_{m-\ell+1,1} & \dots & d_{m-\ell+1,n-\ell+1} \end{pmatrix}, \quad d_{i,j} = ED_{norm}(S_{i,\ell}, T_{j,\ell})$$



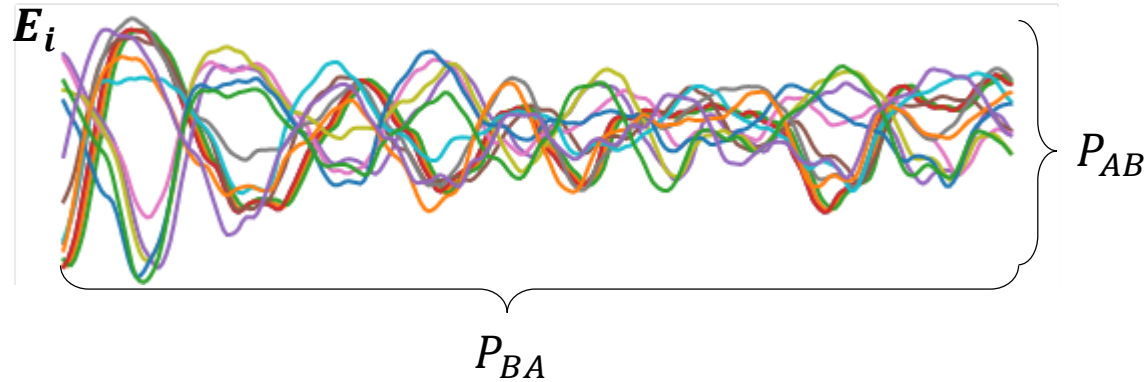
MPdist Profile of a Segment



$$E_i = \begin{pmatrix} d_{1,i} & \dots & d_{1,i+m-\ell+1} \\ \dots & \ddots & \dots \\ d_{m-\ell+1,i} & \dots & d_{m-\ell+1,i+m-\ell+1} \end{pmatrix}$$



MPdist Profile of a Segment

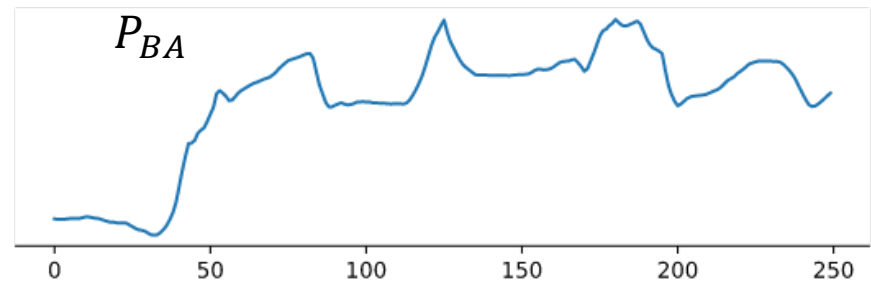
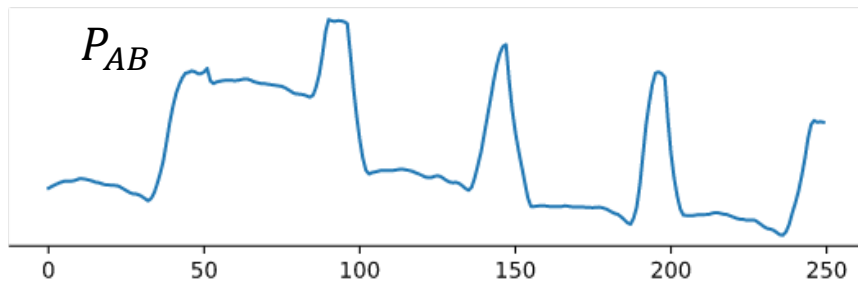


$$P_{AB}(i) = \min_{1 \leq j \leq m - \ell + 1} E(i, j),$$

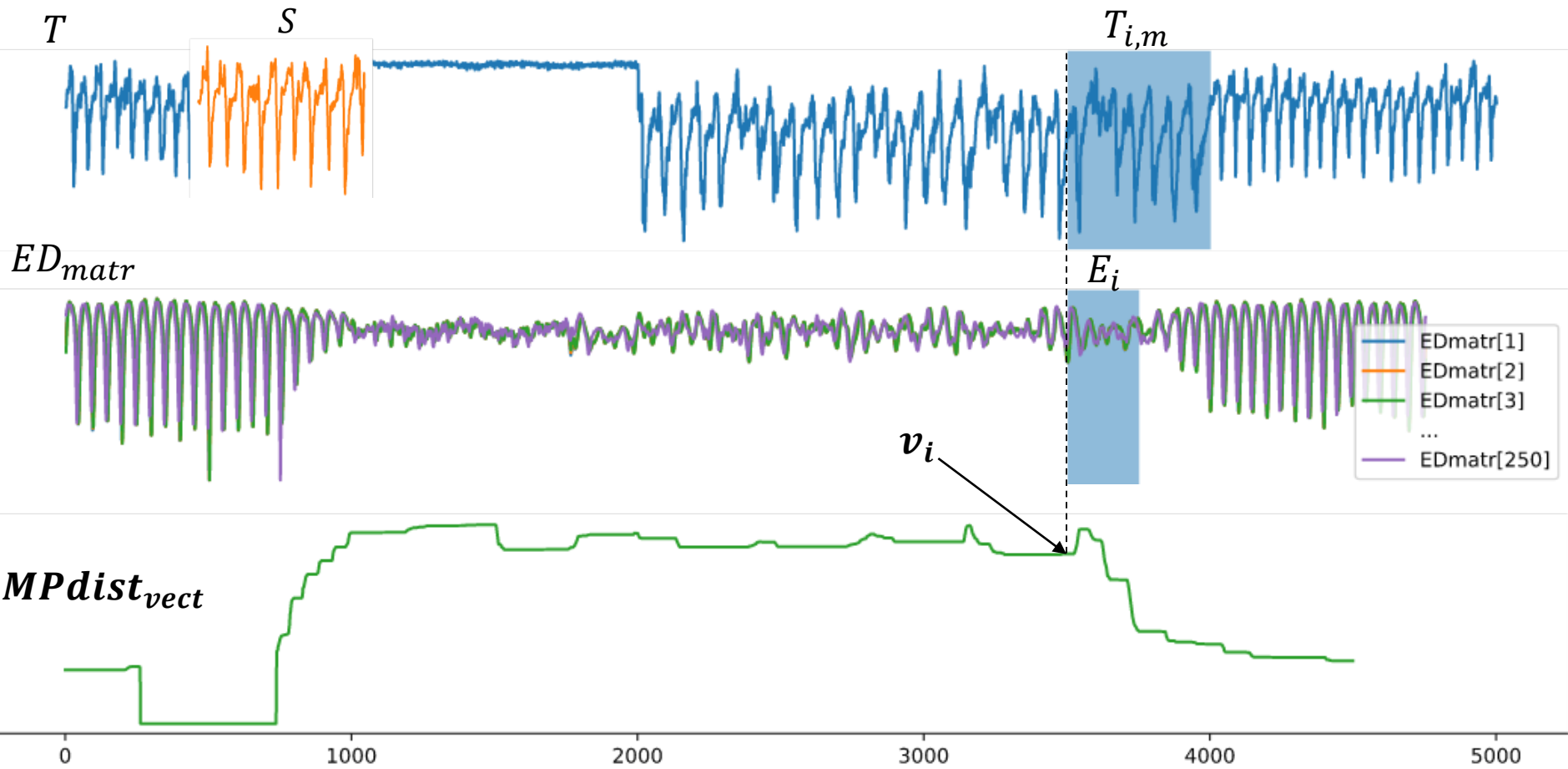
$$1 \leq i \leq m - \ell + 1$$

$$P_{BA}(j) = \min_{1 \leq i \leq m - \ell + 1} E(i, j),$$

$$1 \leq j \leq m - \ell + 1$$



MPdist-profile of a Segment



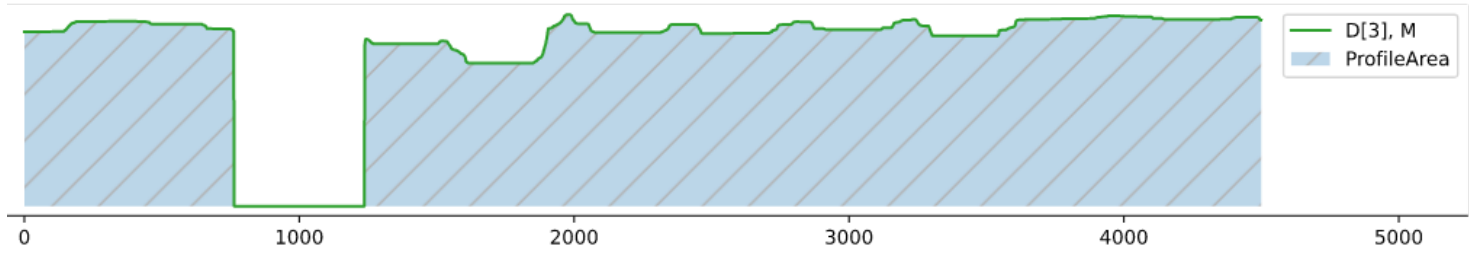
$$MPdist_{vect}(Q, T, \ell) = [v_1, v_2, \dots, v_{n-m+1}], v_i = MPdist(Q, T_{i,m}, \ell)$$

Discovery of the Top-1 Snippet

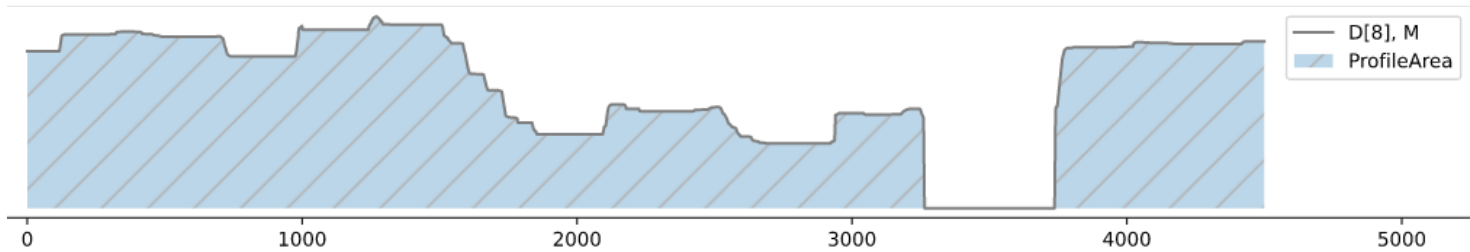
Поиск C_1

i	$ProfileArea$
1	60813
2	60371
3	74451
4	75141
5	56766
6	57729
7	58713
8	53769
9	62127
10	61286

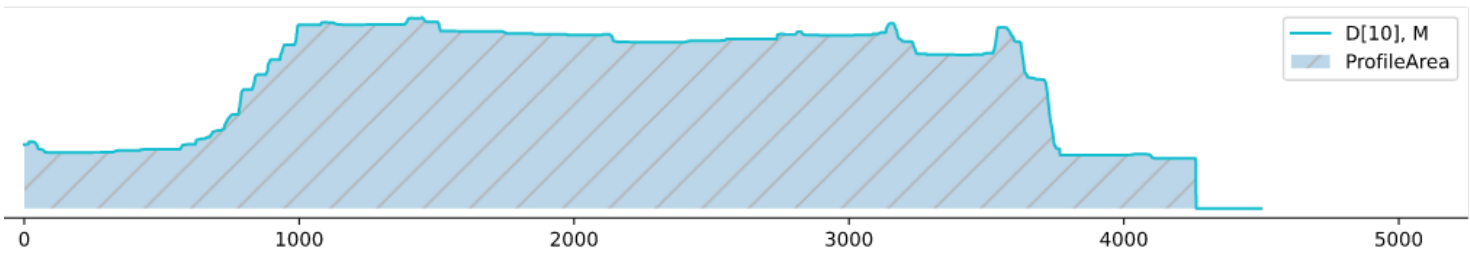
$C_1.index = 8$



$ProfileArea(\{D_3\}) = 74451$

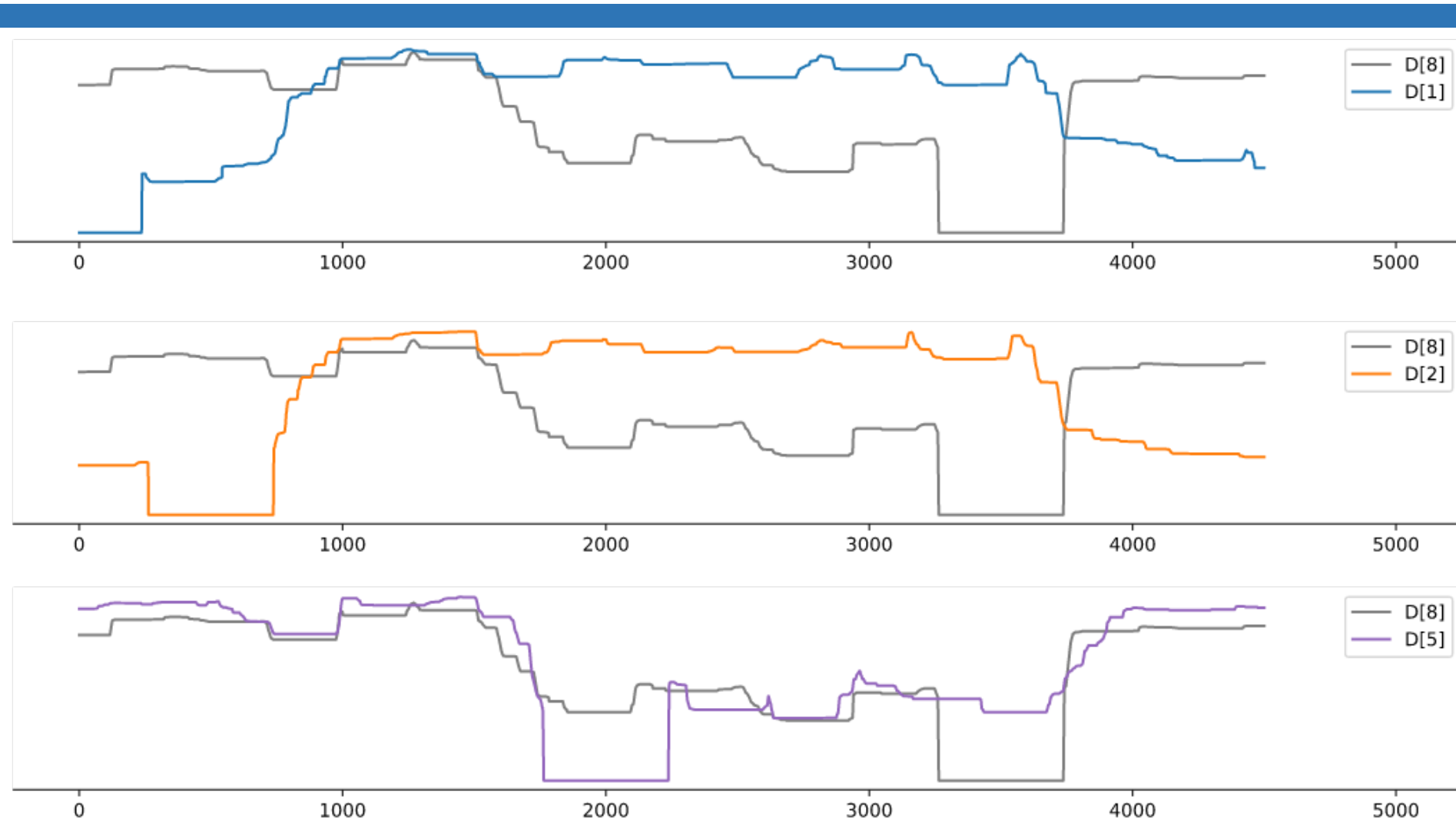


$ProfileArea(\{D_8\}) = 53769$



$ProfileArea(\{D_{10}\}) = 61286$

Discovery of the Top-2 Snippet

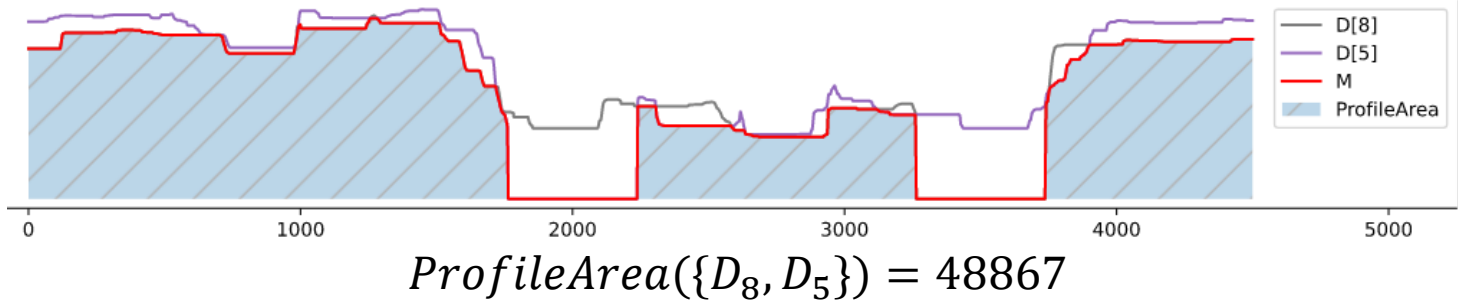
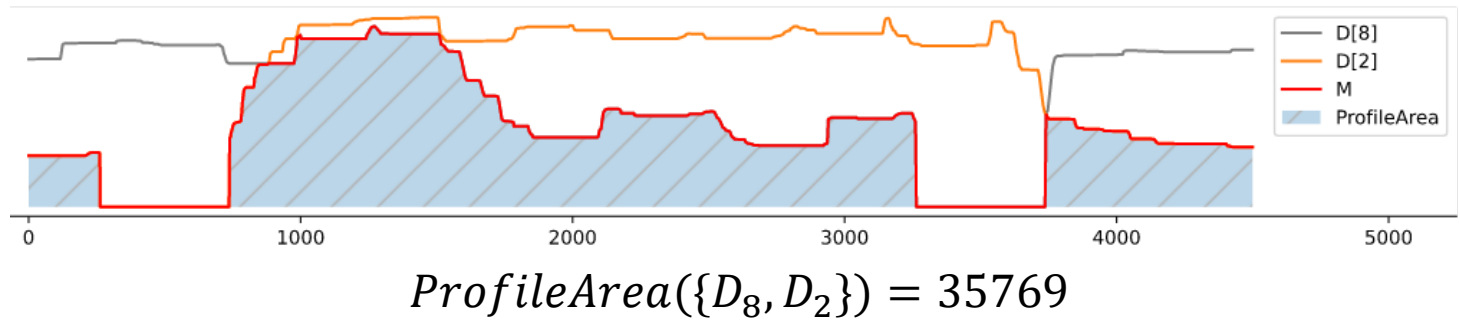
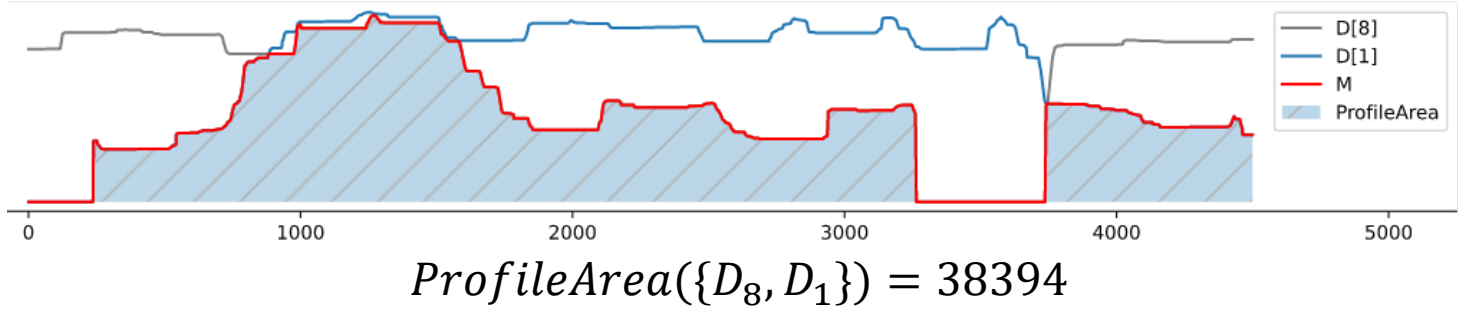


Discovery of the Top-2 Snippet

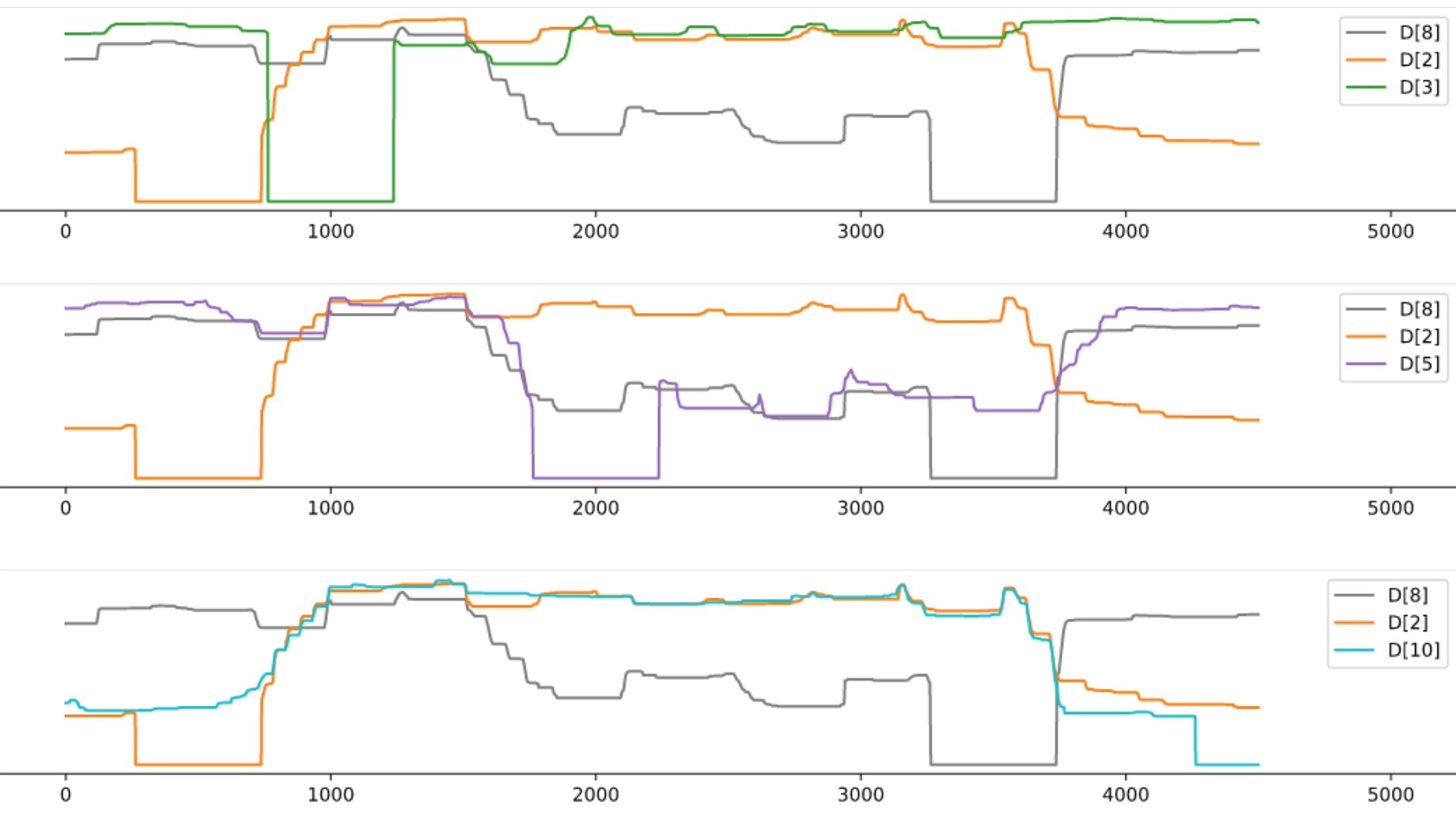
C_2

i	ProfileArea
1	38394
2	35769
3	45629
4	45908
5	48857
6	49264
7	48975
9	36684
10	36482

$C_2.index = 2$



Discovery of the Top-3 Snippet

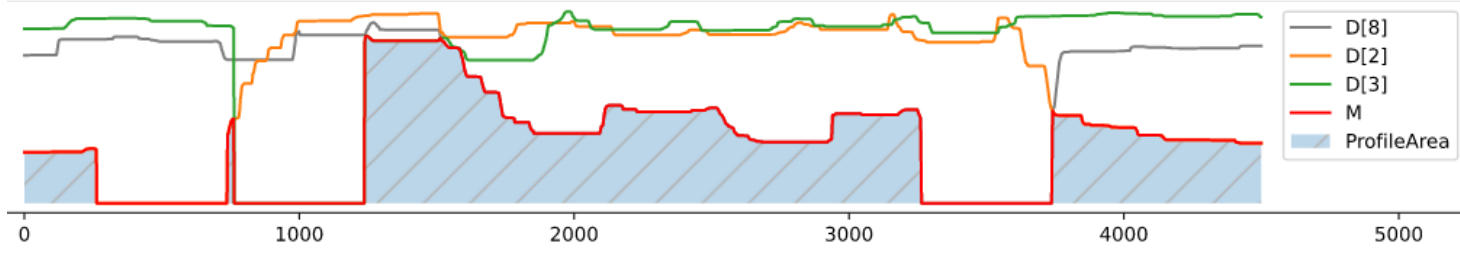


Discovery of the Top-3 Snippet

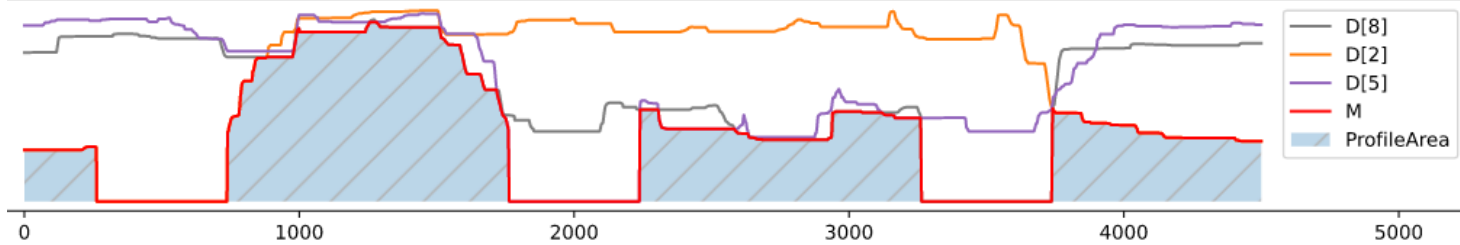
C_3

i	ProfileArea
1	34475
3	27899
4	27908
5	31168
6	31532
7	31672
9	31654
10	33044

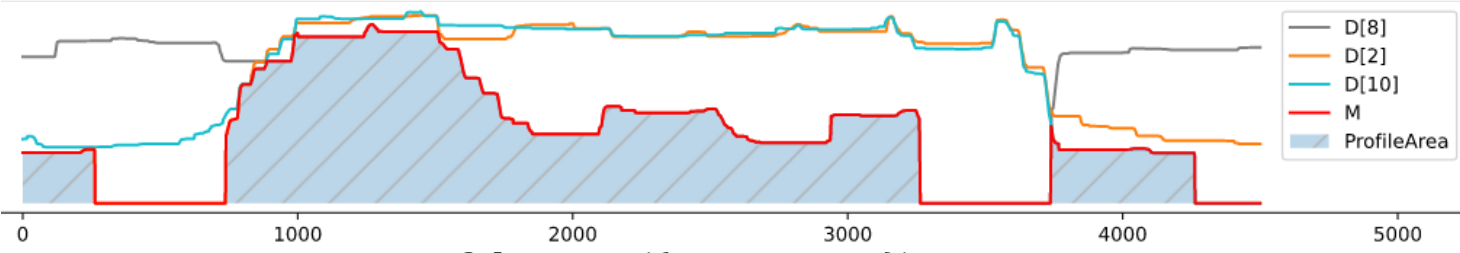
$C_3.index = 3$



$ProfileArea(\{D_8, D_2, D_3\}) = 27899$

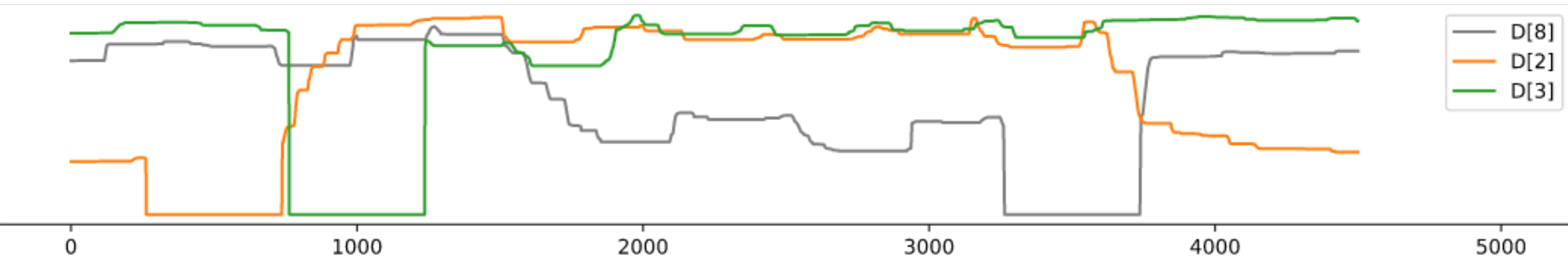
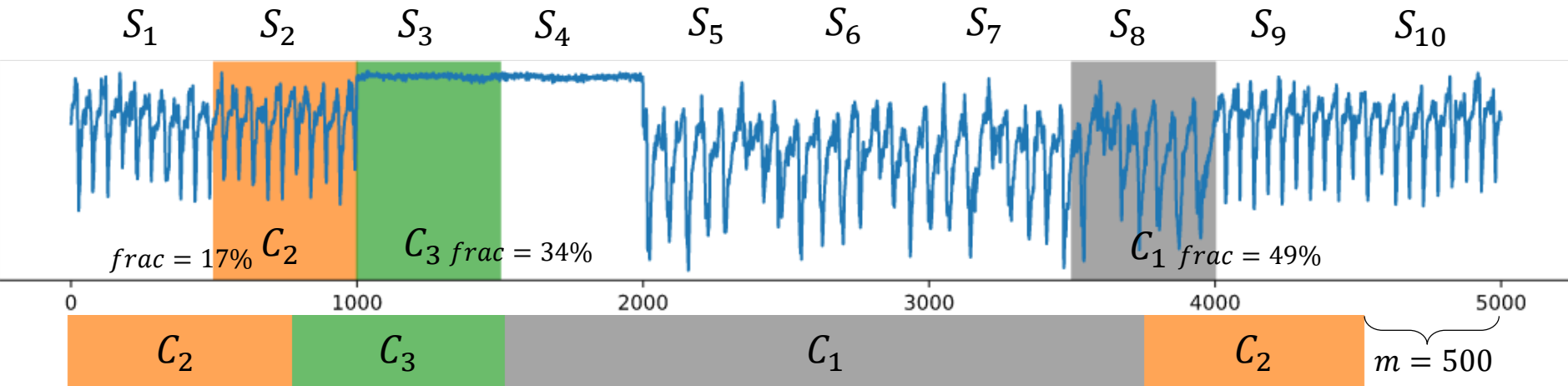


$ProfileArea(\{D_8, D_2, D_5\}) = 31168$

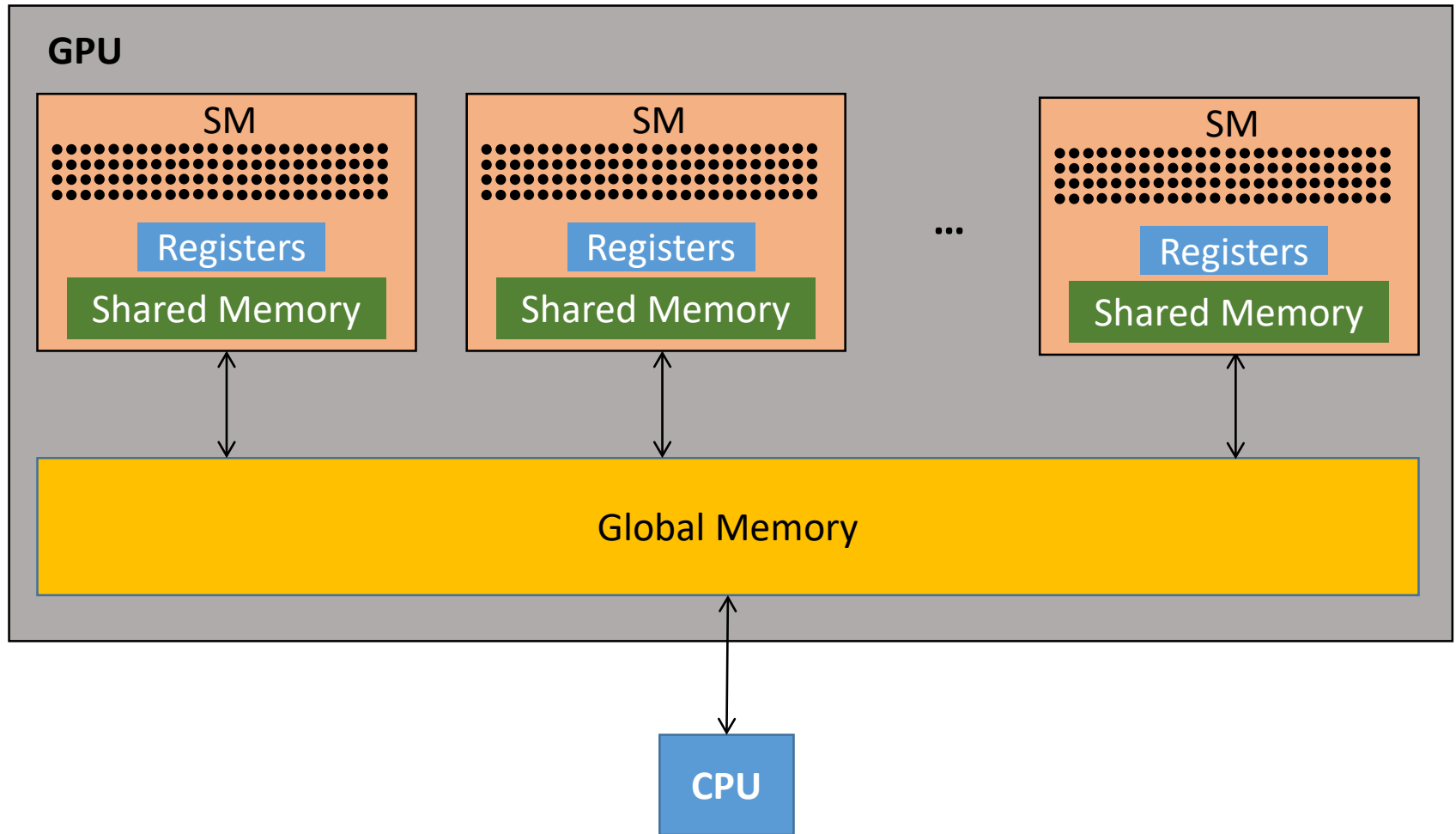


$ProfileArea(\{D_8, D_2, D_{10}\}) = 33044$

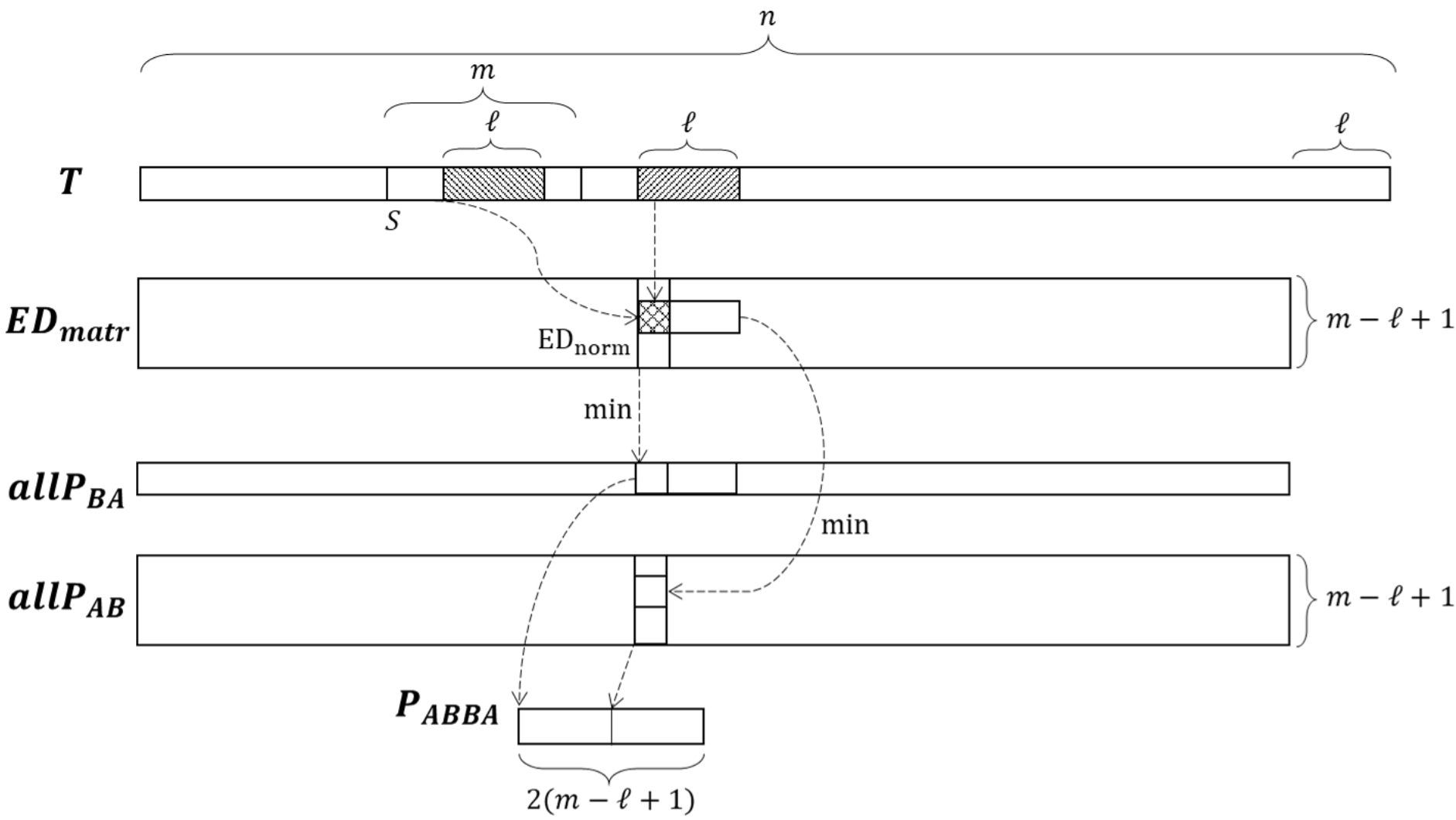
Resulting Snippets



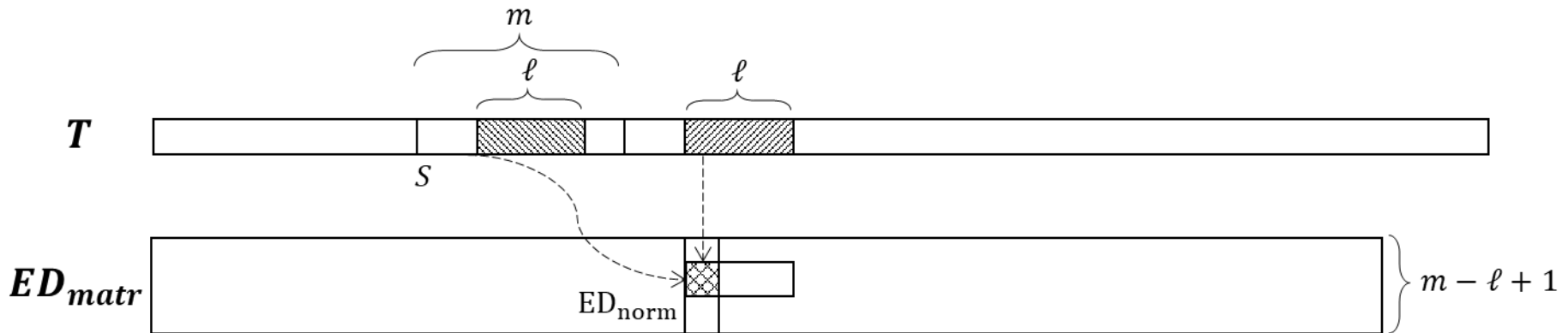
Hardware Architecture



Parallelizing: Data Structures



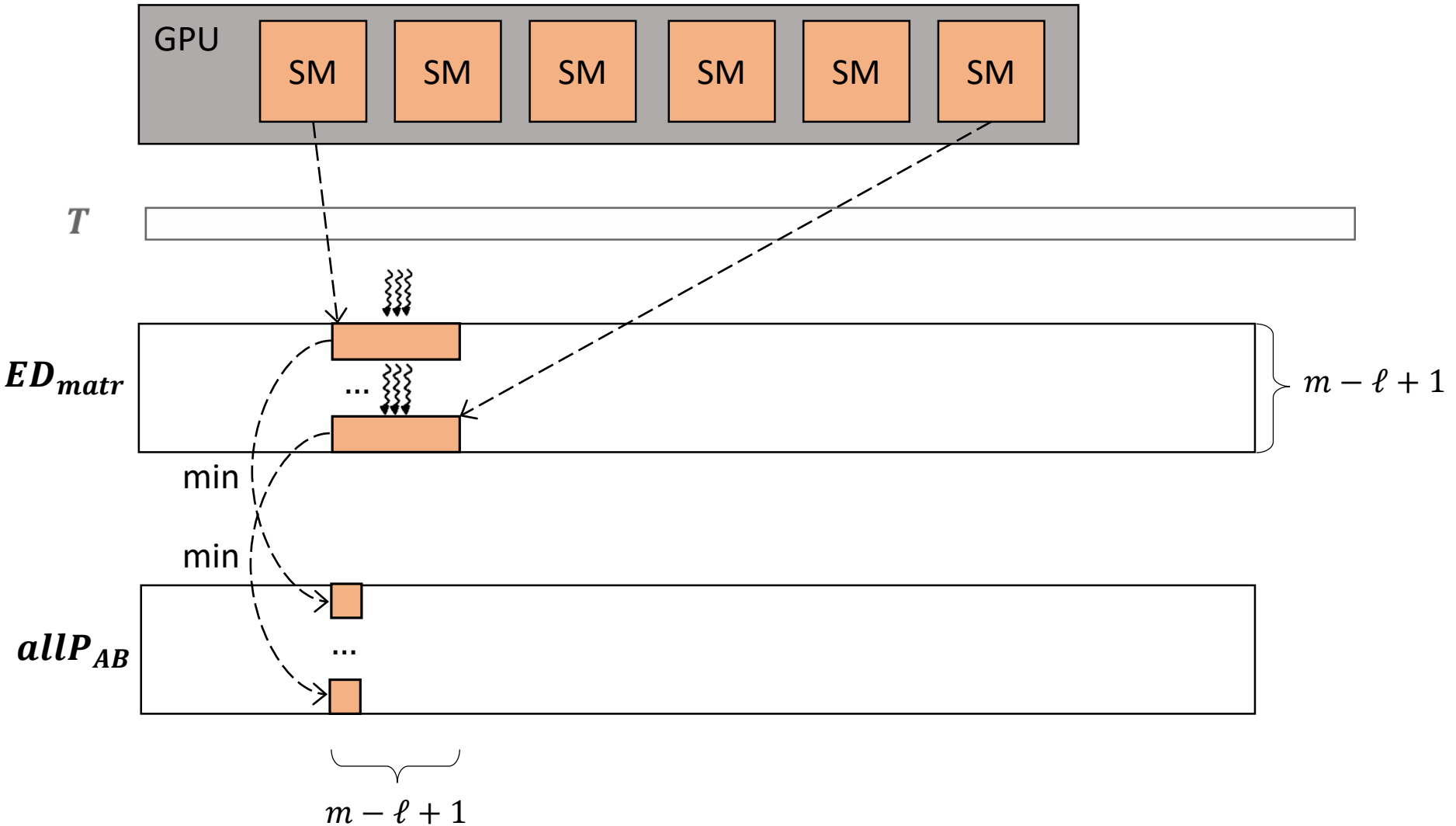
Calculation of ED_{matr}



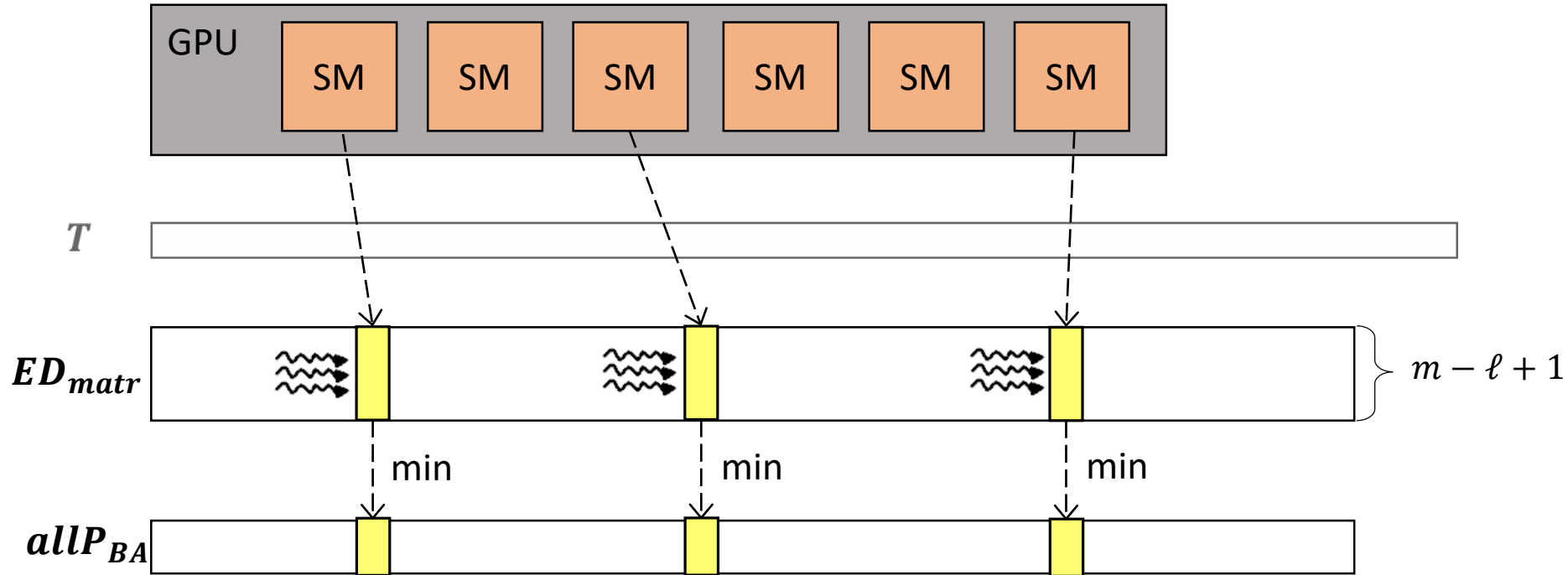
$$\begin{aligned} \overline{QT}_{i,j} &= \overline{QT}_{i-1,j-1} + df_i \cdot dg_j + df_j \cdot dg_i, \\ df_0 &= 0; df_i = \frac{t_{i+m-1} - t_{i-1}}{2}, \\ dg_0 &= 0; dg_i = (t_{i+m-1} - \mu_i) + (t_{i-1} - \mu_{i-1}), \\ \mu_i &= \frac{1}{m} \sum_{j=i}^{i+m} t_j, \\ T_{i,m} - \mu_i &= (t_i - \mu_i, \dots, t_{i+m-1} - \mu_i), \\ P_{i,j} &= \overline{QT}_{i,j} \cdot \frac{1}{\|T_{i,m} - \mu_i\|} \cdot \frac{1}{\|T_{j,m} - \mu_j\|}, \\ ED_{norm}(T_{i,m}, T_{j,m}) &= \sqrt{2m(1 - P_{i,j})} \end{aligned}$$

* Zimmerman Z., Kamgar K., Senobari N.S. et al. Matrix Profile XIV: Scaling Time Series Motif Discovery with GPUs to Break a Quintillion Pairwise Comparisons a Day and Beyond. ACM SoCC'2019. DOI: [10.1145/3357223.3-362721](https://doi.org/10.1145/3357223.3-362721).

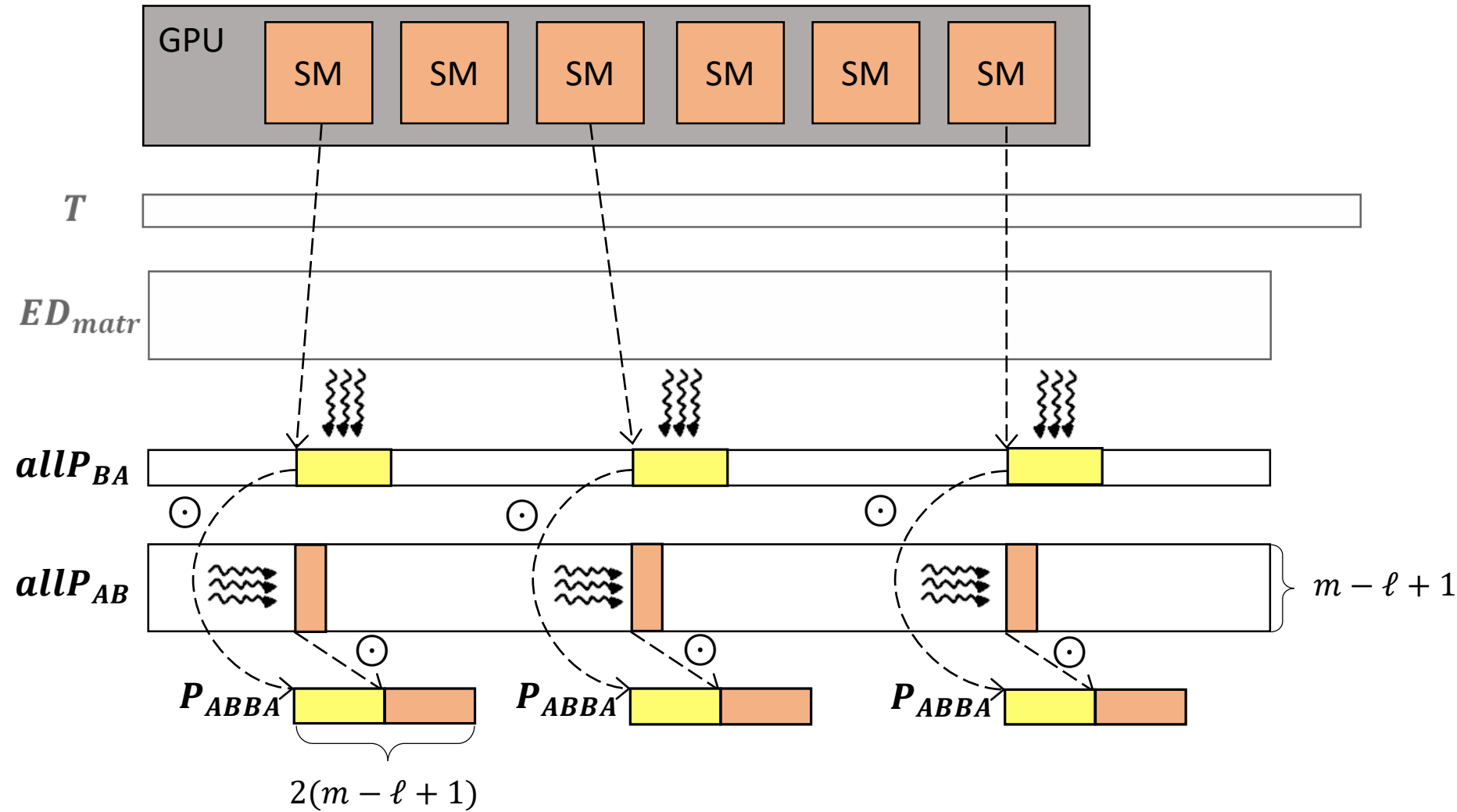
Calculation of $allP_{AB}$



Calculation of $allP_{BA}$



Calculation of P_{ABBA}



Experiments: Hardware

- CPU:
 - Intel Xeon Gold 6254@4 GHz
 - Cores: 18 (but only one was employed)
 - RAM: 64 Gb
 - Peak performance: 1.2 TFLOPS
- GPU:
 - NVIDIA Tesla V100 SXM2
 - Cores: 5120 @1.312 GHz (84 streaming multiprocessors)
 - RAM: 32 Gb
 - Peak performance: 15.7 TFLOPS

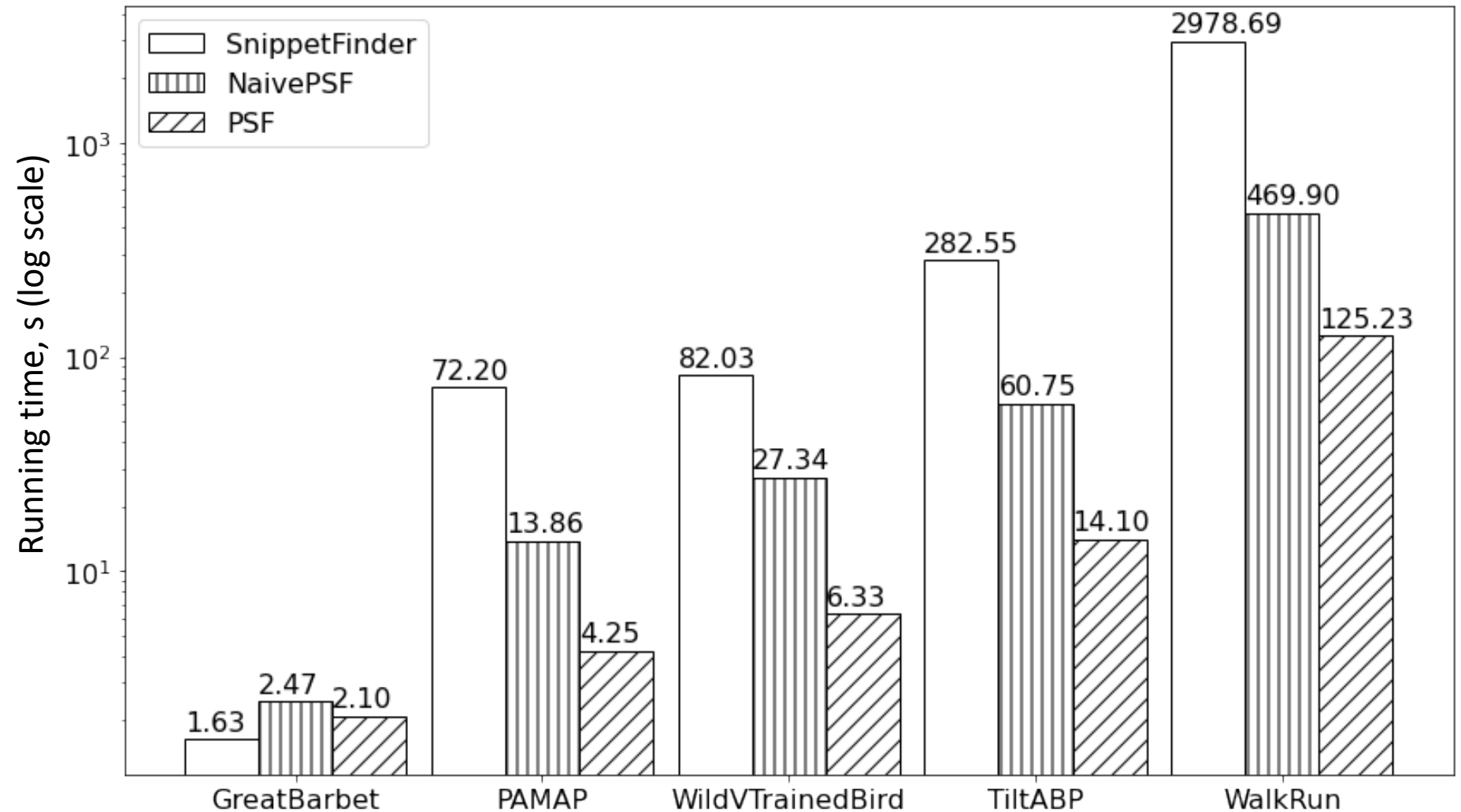
Experiments: Data

Time series	Length n	Segment length m	Description
GreatBarbet ⁽¹⁾	2 801	150	Physiological indicators of bird vital activity
WildVTrainedBird ⁽¹⁾	20 002	900	
PAMAP ⁽²⁾	20 002	600	Wearable accelerometer readings during various types of human physical activity
WalkRun ⁽²⁾	100 000	240	
TiltABP ⁽¹⁾	40 000	630	Human blood pressure readings during rapid tilts

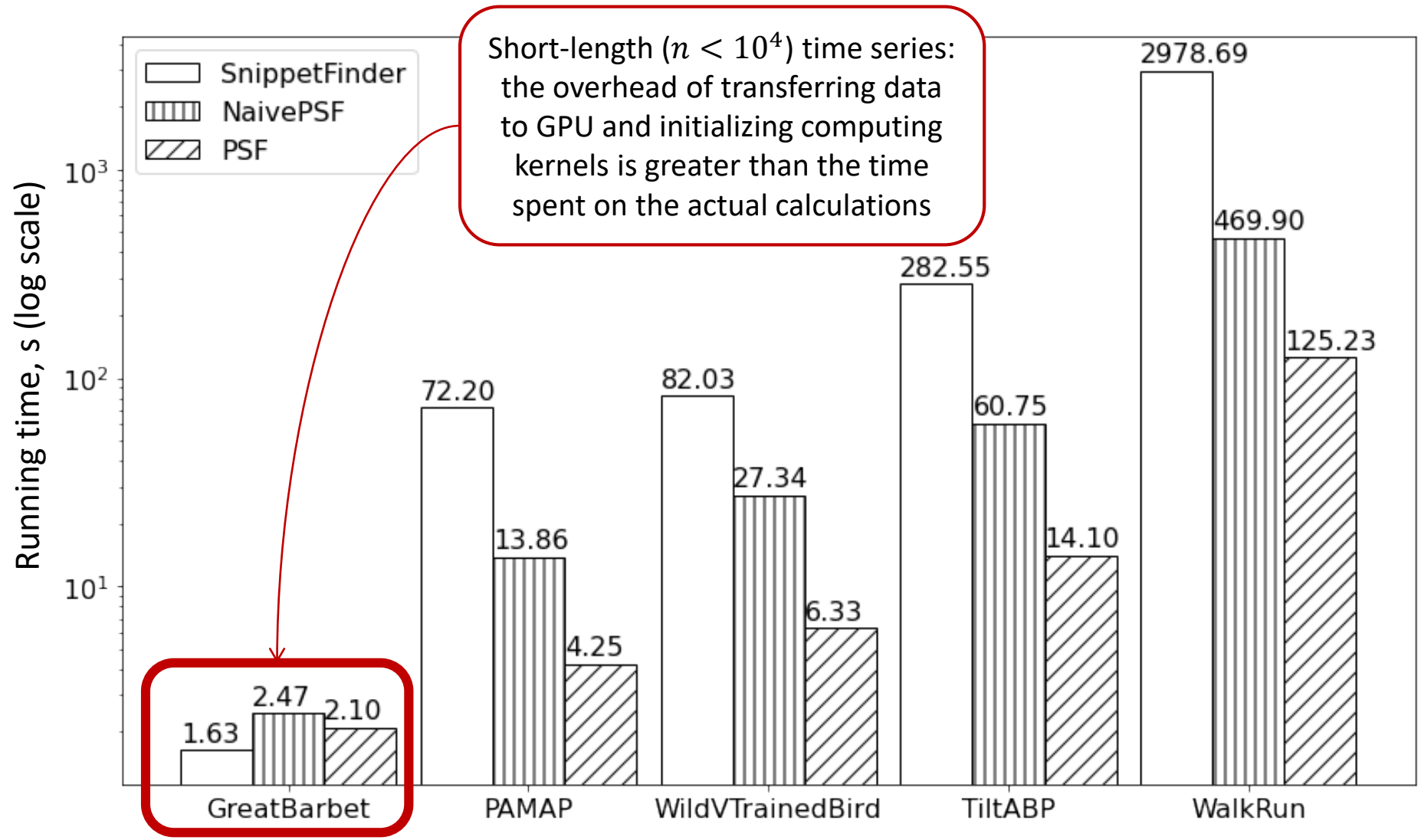
⁽¹⁾ Imani S., Madrid F., Ding W., Crouter S.E., Keogh E.J. Introducing time series snippets: a new primitive for summarizing long time series. *Data Min. Knowl. Discov.* 34(6): 1713-1743 (2020). doi: [10.1007/s10618-020-00702-y](https://doi.org/10.1007/s10618-020-00702-y)

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

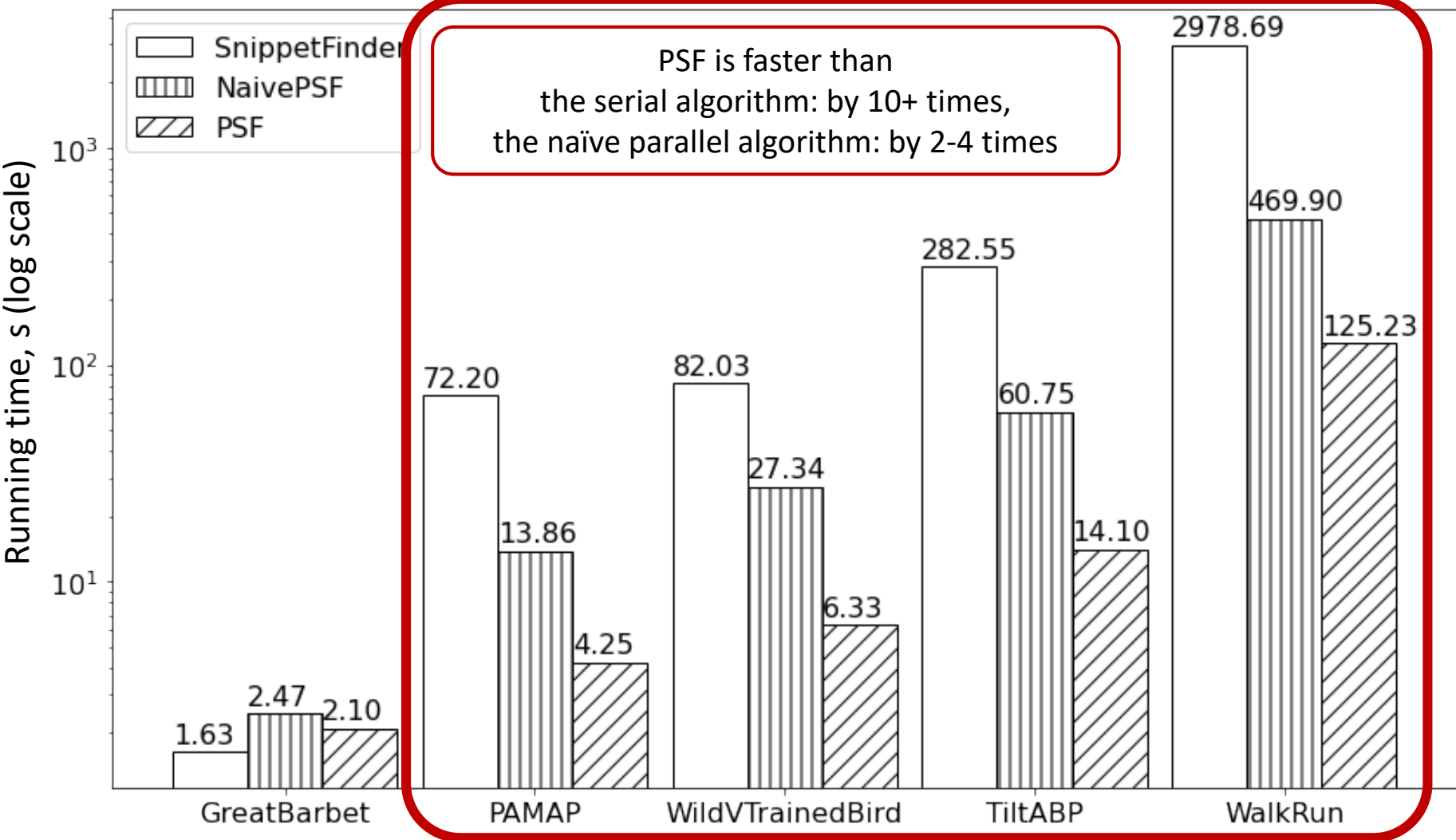
Experiments: Performance



Experiments: Performance

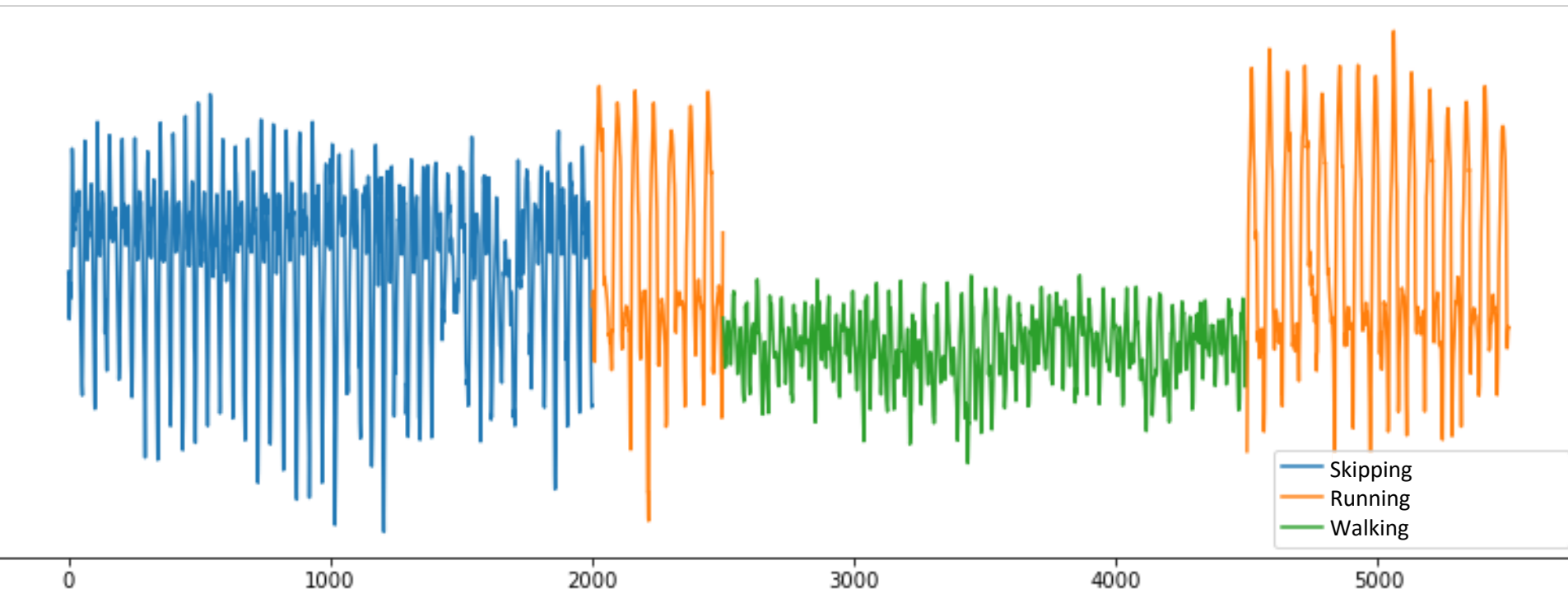


Experiments: Performance



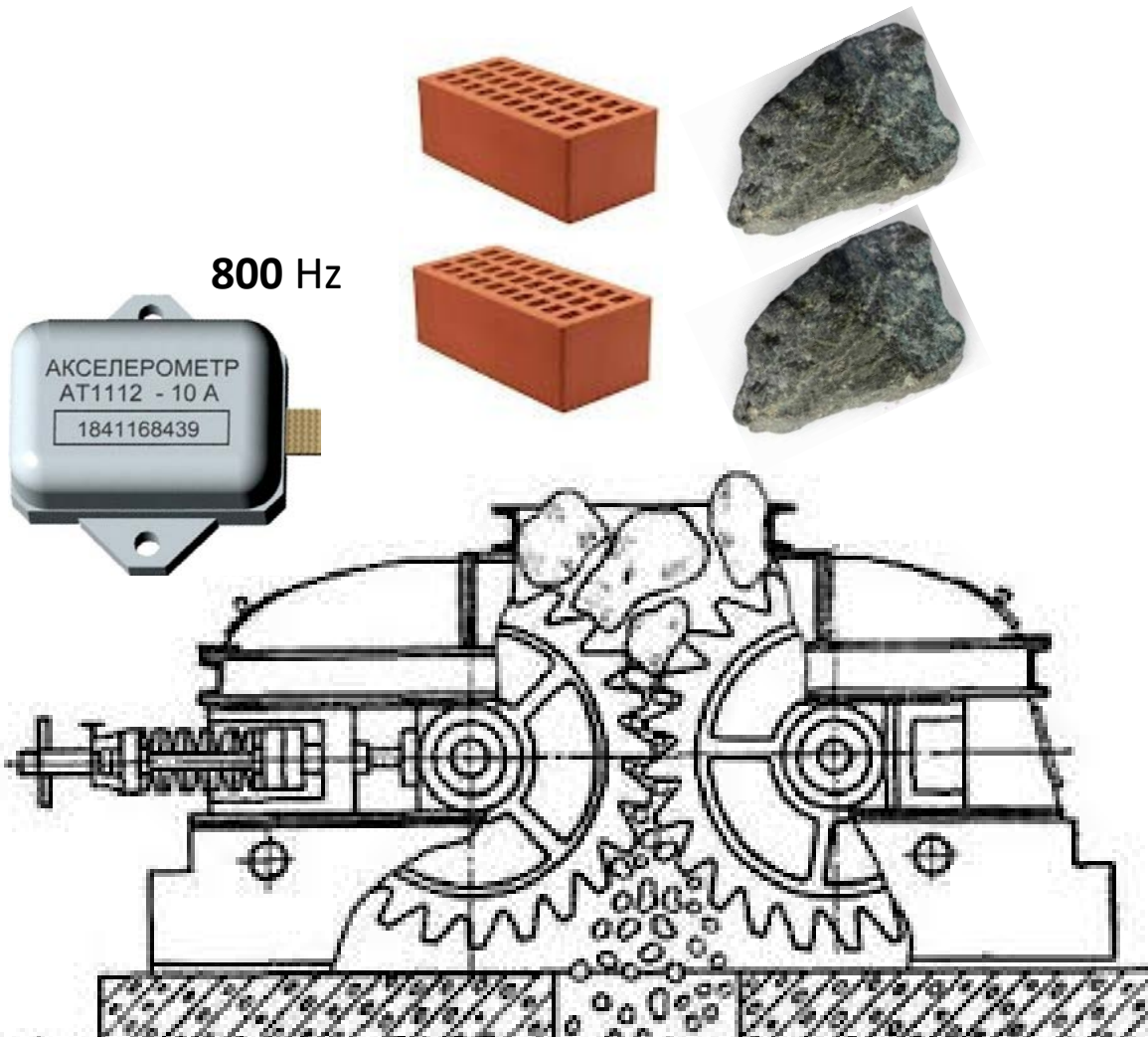
Experiments: Visualization

Time series: PAMAP *



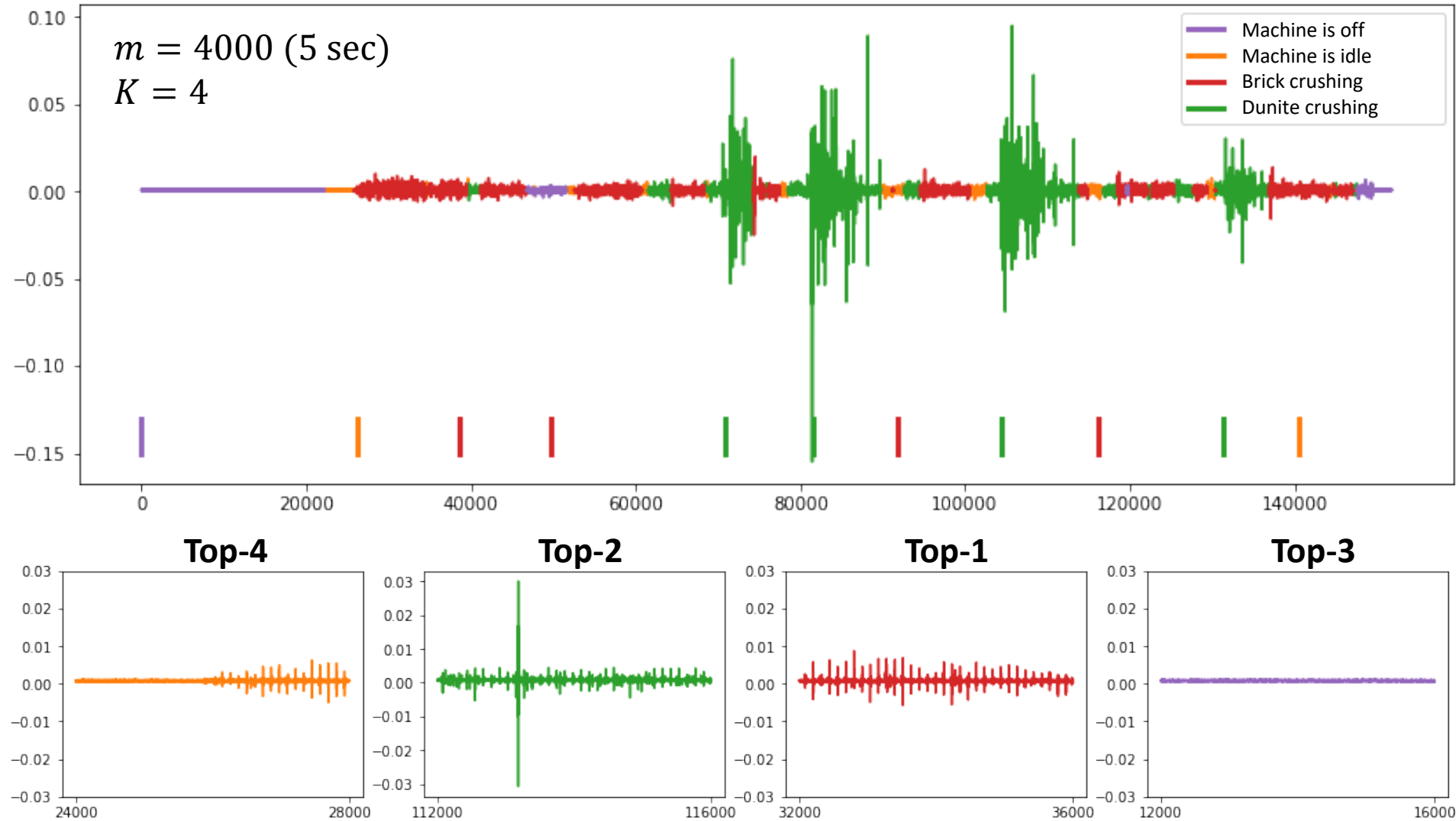
* Reiss A., Stricker D. Introducing a New Benchmarked Dataset for Activity Monitoring, ISWC'2012. DOI: [10.1109/iswc.2012.13](https://doi.org/10.1109/iswc.2012.13)

Case Studies: Small-sized Crushing Machine



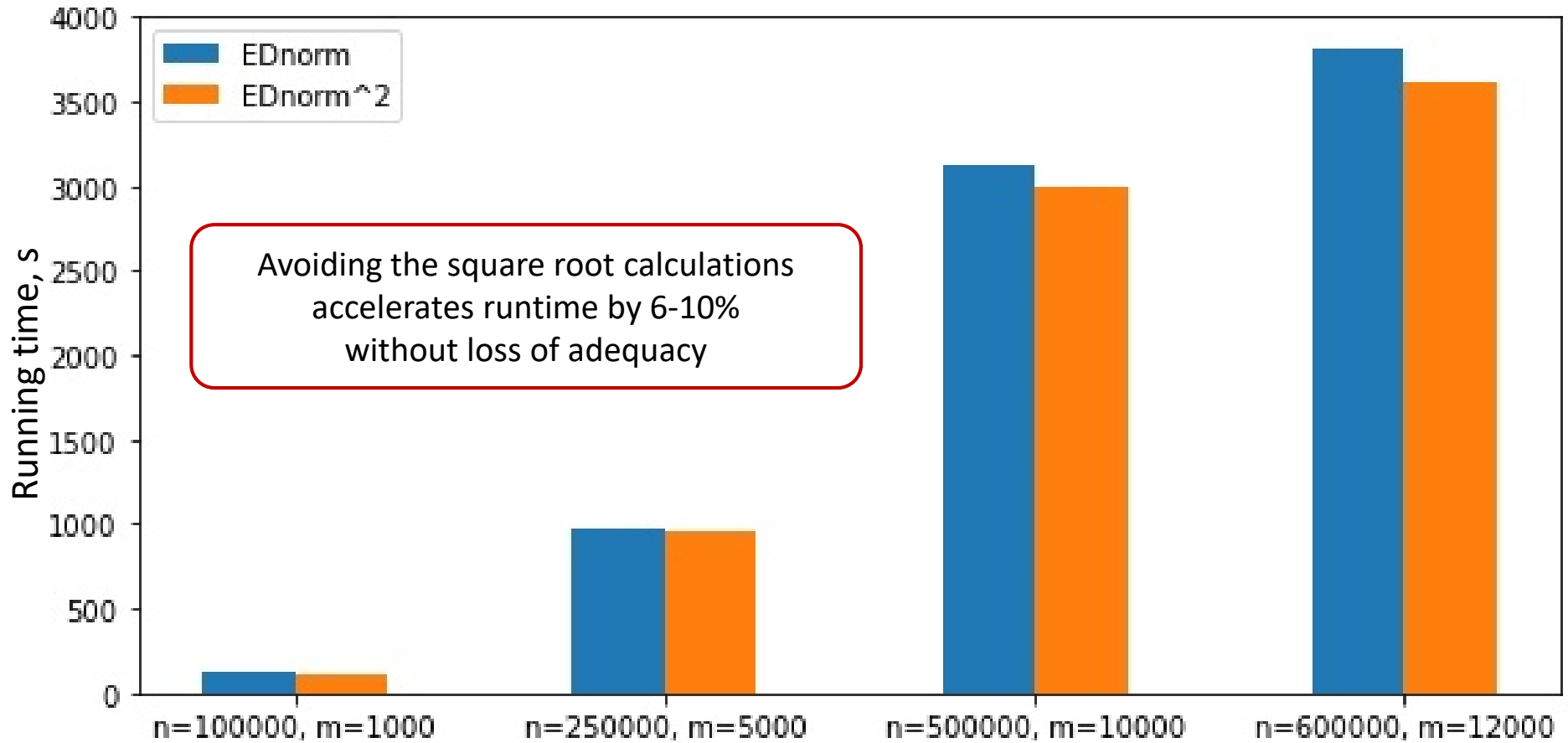
Time	Event
00:00:00	Machine is off
00:32.60	Turning the machine on, machine is idle
00:48.34	Loading with bricks
01:02.09	Loading with bricks
01:28.05	Loading with dunite
01:41.86	Loading with dunite
01:54.68	Loading with bricks
02:10.49	Loading with dunite
02:25.32	Loading with bricks
02:44.20	Loading with dunite
02:55.49	Finishing crushing, machine is idle
03:07.06	Turning the machine off

Case Studies: Small-sized Crushing Machine



Experiments: ED_{norm}^2 vs. ED_{norm}

Time series: Random Walk *



* Pearson K. The problem of the random walk. Nature. 72(1865), 294 (1905). DOI: [10.1038/072342A0](https://doi.org/10.1038/072342A0)

Conclusions

- PSF (Parallel Snippet Finder) is a novel parallel algorithm to discover snippets of a time series on GPU
- PSF showed high performance in the experiments
- Further study: PSF for HPC-cluster

- Thank you for paying attention! Any questions?
 - Andrey Goglachev, goglachevai@susu.ru