

Всероссийская научная конференция с международным участием  
**Параллельные вычислительные технологии (ПаВТ'2023)**  
Санкт-Петербург, 28–30 марта 2023 г.

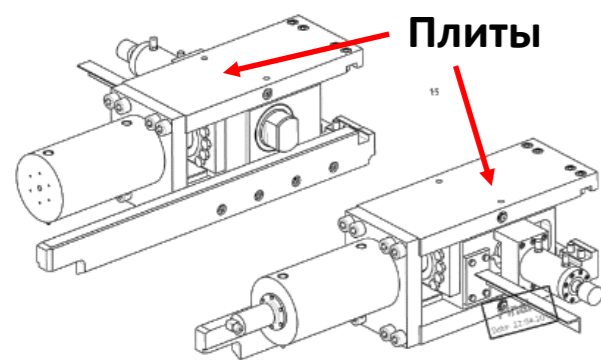
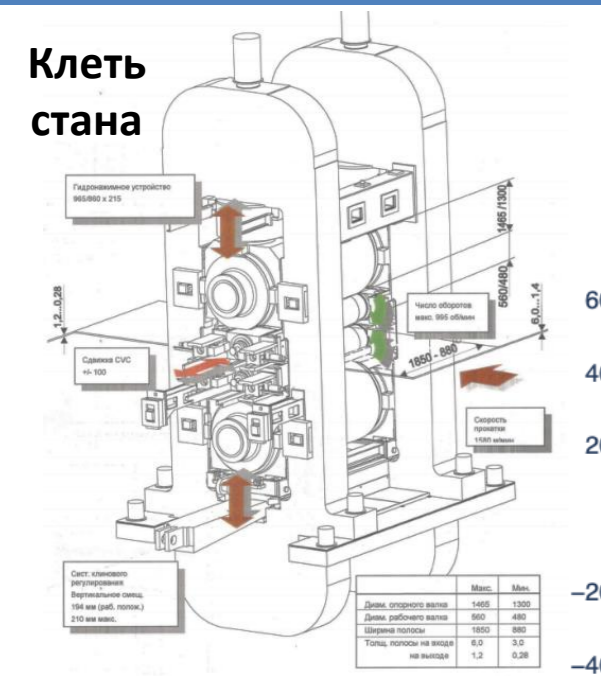
---

# Автоматизированный поиск аномалий временных рядов на графическом процессоре

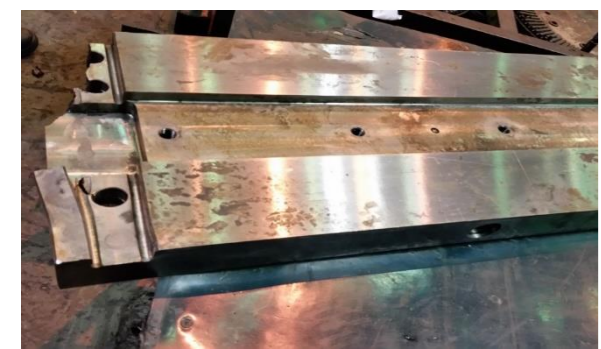
Я.А. Краева, М.Л. Цымблер

Южно-Уральский государственный университет (Челябинск)

# Поиск аномалий во временных рядах из цифровой индустрии

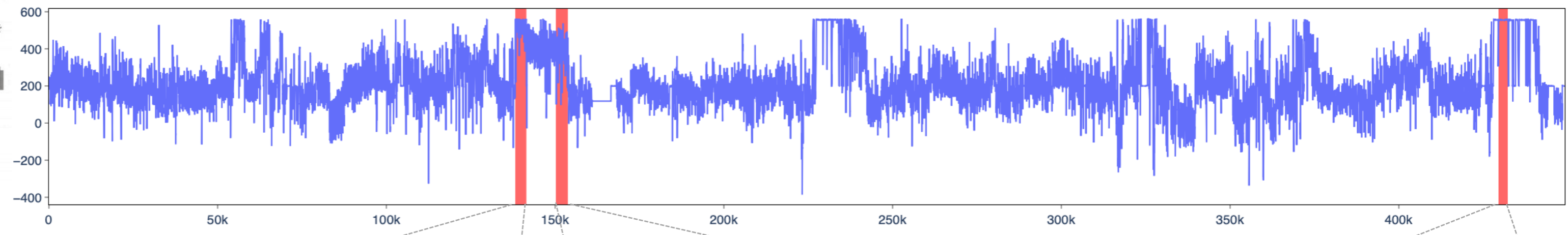


**Система валков с непрерывно изменяемой кривизной**

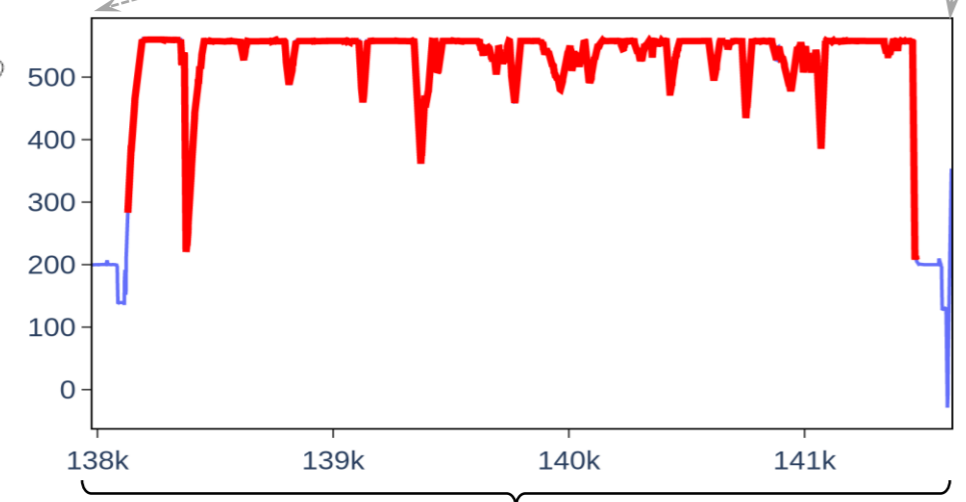


**Разрушение плит**

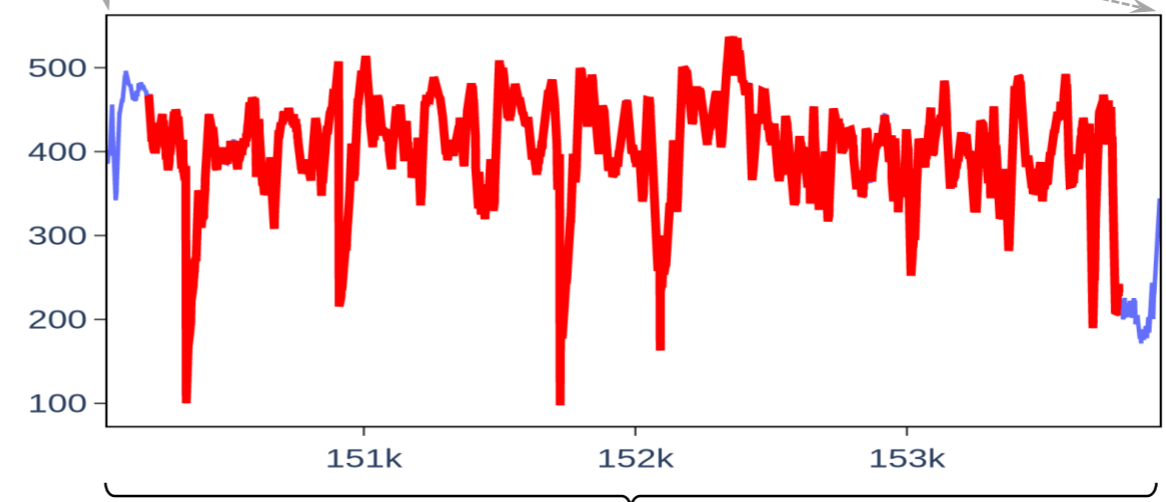
## Фактическое изгибающее усилие прокатки



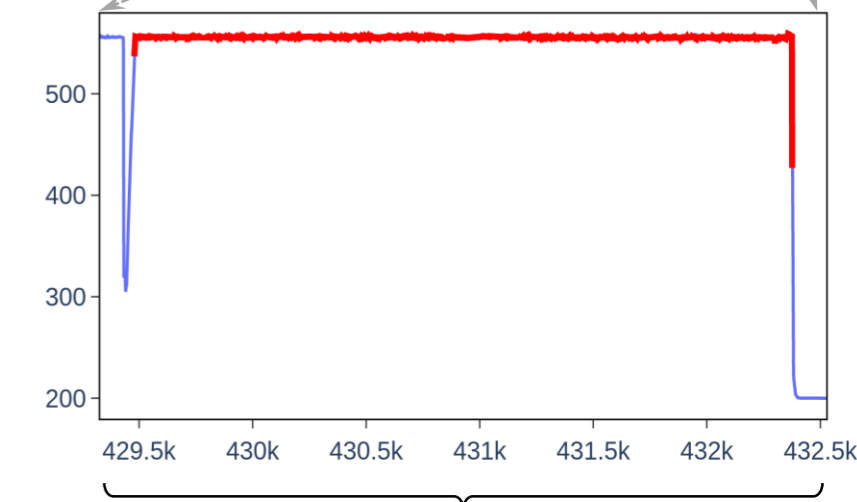
## Аномальные напряжения плит



**55 минут**



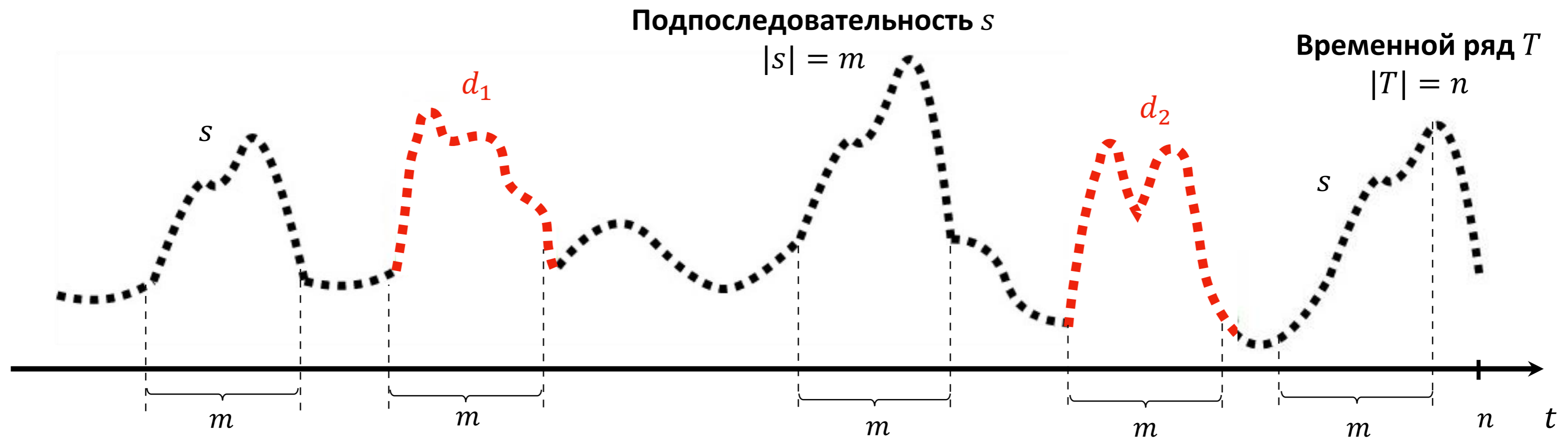
**1 час**



**48 минут**

# Постановка задачи

- **Диссонанс**<sup>1)</sup> – подпоследовательность ряда, расстояние от которой до ближайшего соседа не ниже порога  $r$
- **Дано:** временной ряд  $T$ , длина диссонанса  $m$ , порог  $r$
- **Найти:**  $D = \{d_1, d_2, \dots\}$ ,  $d_i \in D \Leftrightarrow \forall s \in T \min_{s \cap d_i = \emptyset} \text{dist}(d_i, s) \geq r$



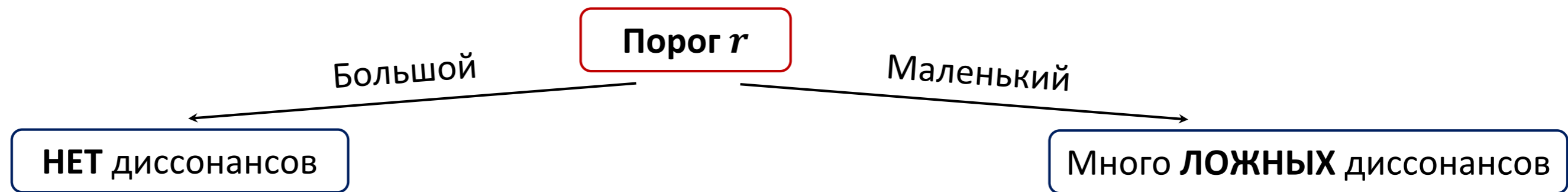
Количество подпоследовательностей:  $N = n - m + 1$

<sup>1)</sup> Yankov D., Keogh E.J., Rebbapragada U. Disk aware discord discovery: finding unusual time series in terabyte sized datasets. Knowl. Inf. Syst. 17(2): 241-262. 2008.

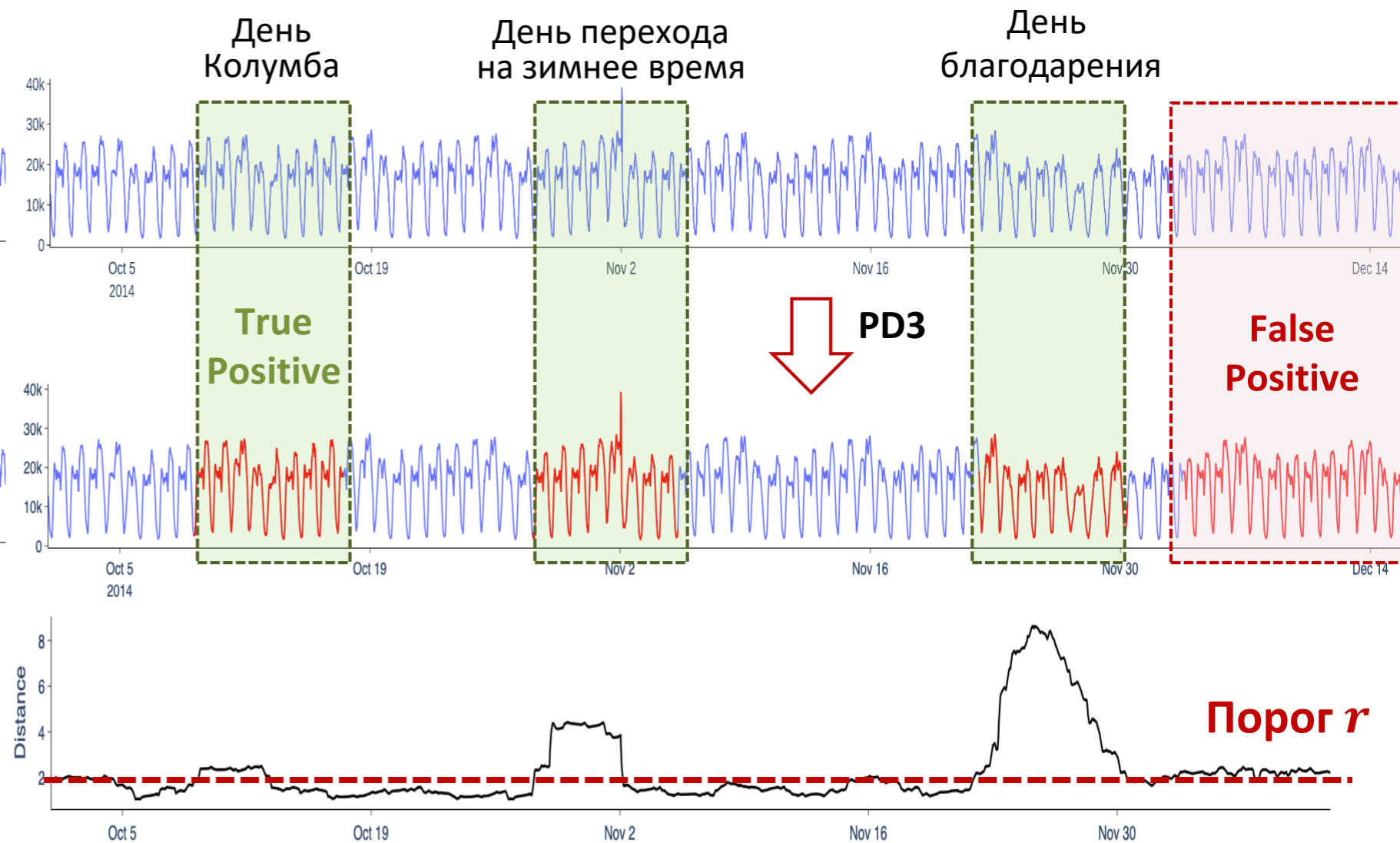
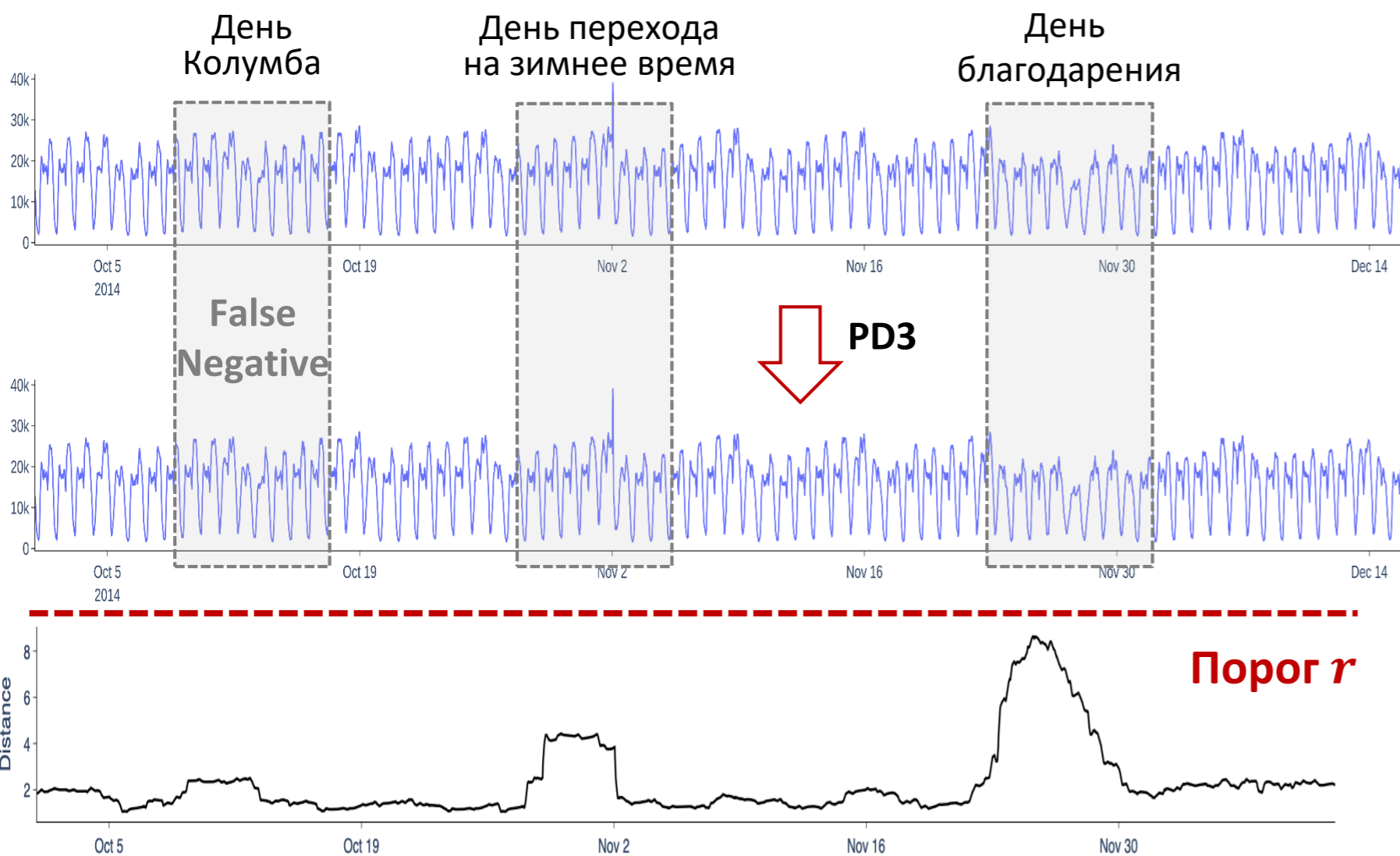
# Основные работы по теме исследования

Алгоритм	Платформа	Критика
<b>Последовательный алгоритм</b>		
Nakamura T., <i>et al.</i> <b>MERLIN</b> : parameter-free discovery of arbitrary length anomalies in massive time series archives. IEEE ICDM 2020. pp. 1190-1195.	CPU	Квадратичная сложность от длины ряда
<b>Параллельные алгоритмы</b>		
<b>DRAG</b> : Yankov D., <i>et al.</i> Disk aware discord discovery: finding unusual time series in terabyte sized datasets. Knowl. Inf. Syst. 17(2): 241-262. 2008.	CPU	Симуляция MapReduce
<b>PDD</b> : Huang T., <i>et al.</i> Parallel discord discovery. PAKDD 2016. LNCS 9652. Springer, 2016. pp. 233-244.	Spark	Низкая производительность ввиду большого количества обменов между узлами
<b>PhiDD</b> : Zymbler M., <i>et al.</i> A Parallel Approach to Discords Discovery in Massive Time Series Data. Computers, Materials & Continua 66(2): 1867-1876. 2021.	Кластер Intel Xeon Phi	Квадратичная пространственная сложность от длины ряда
<b>KBF_GPU</b> : Thuy T.T.H., <i>et al.</i> A new discord definition and an efficient time series discord detection method using GPUs. ICSED 2021. pp. 63-70.	GPU	Полный перебор подпоследовательностей ряда
<b>Zhu B., et al.</b> A GPU Acceleration framework for motif and discord based pattern mining. IEEE Trans. on Parallel and Distr. Systems 32(8): 1987-2004. 2021.	GPU	Поиск одного (самого важного) диссонанса ряда
<b>PD3</b> : Zymbler M., Kraeva Ya. Parallel algorithm for time series discord discovery on a graphics processor. Pattern Recognition and Image Analysis 33(2). 2023.	GPU	Ручной подбор длины диссонанса и порога

# PD3 (Parallel DRAG-based Discord Discovery): Ручной подбор $r$

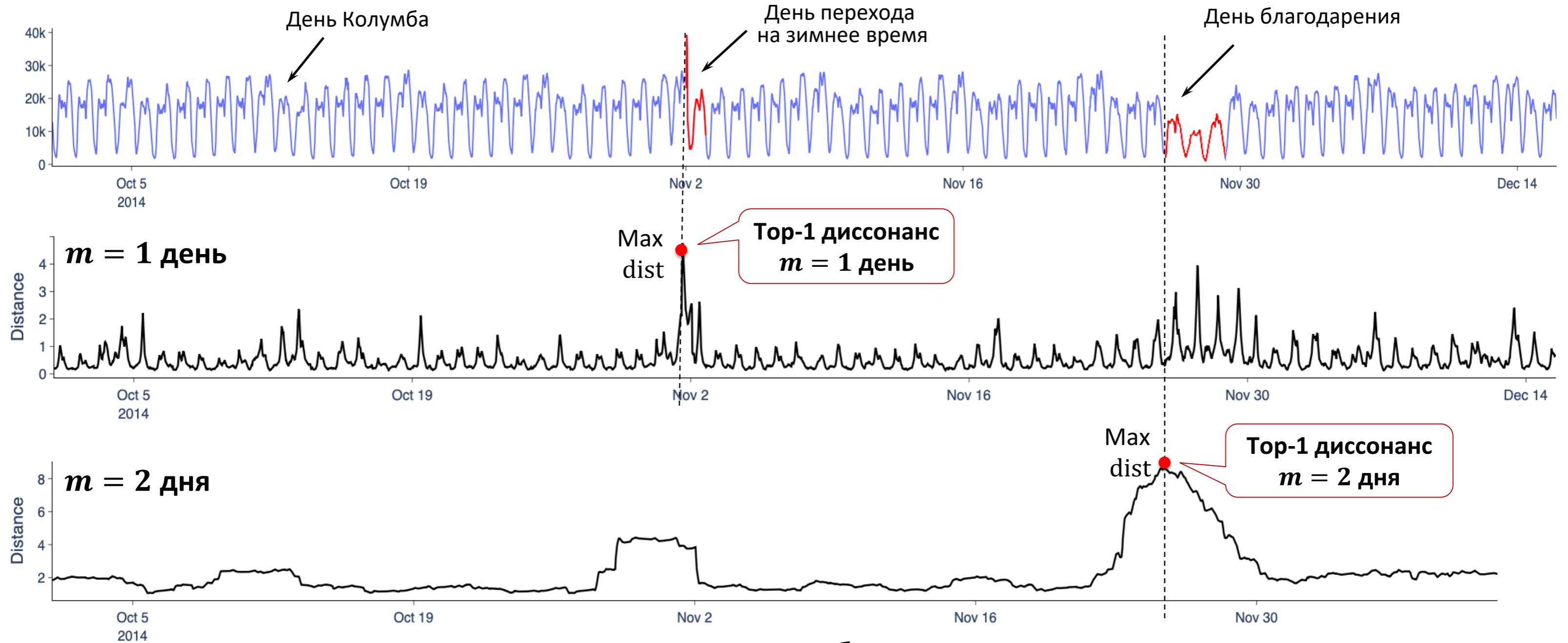


Среднее число пассажиров NY такси (осень 2014 г., каждые полчаса)



# PD3 (Parallel DRAG-based Discord Discovery): Ручной подбор $m$

Среднее число пассажиров NY такси  
(осень 2014 г., каждые полчаса)

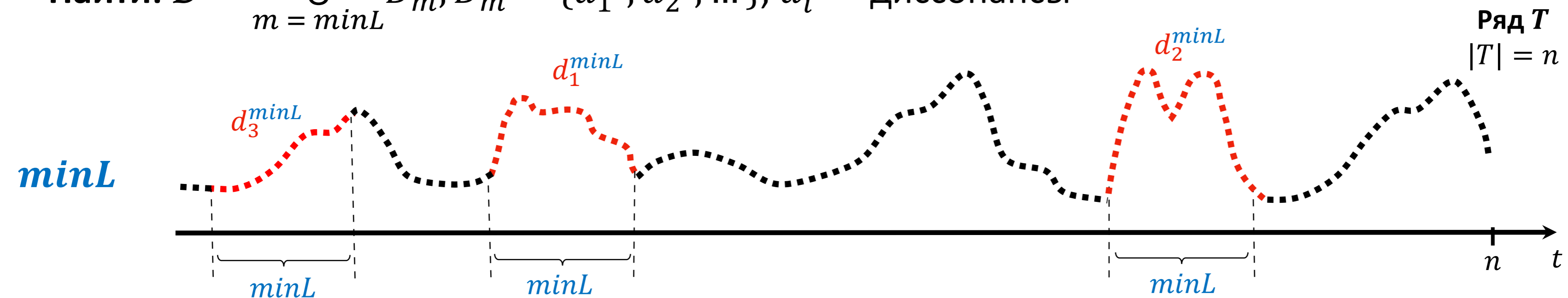


Накладные расходы на подбор параметров  $m$  и  $r$

# Постановка задачи автоматизированного поиска диссонансов<sup>1)</sup>

• **Дано:** временной ряд  $T$ , диапазон длин диссонансов  $minL, \dots, maxL$

• **Найти:**  $\mathcal{D} = \bigcup_{m = minL}^{maxL} D_m, D_m = \{d_1^m, d_2^m, \dots\}, d_i^m$  – диссонансы

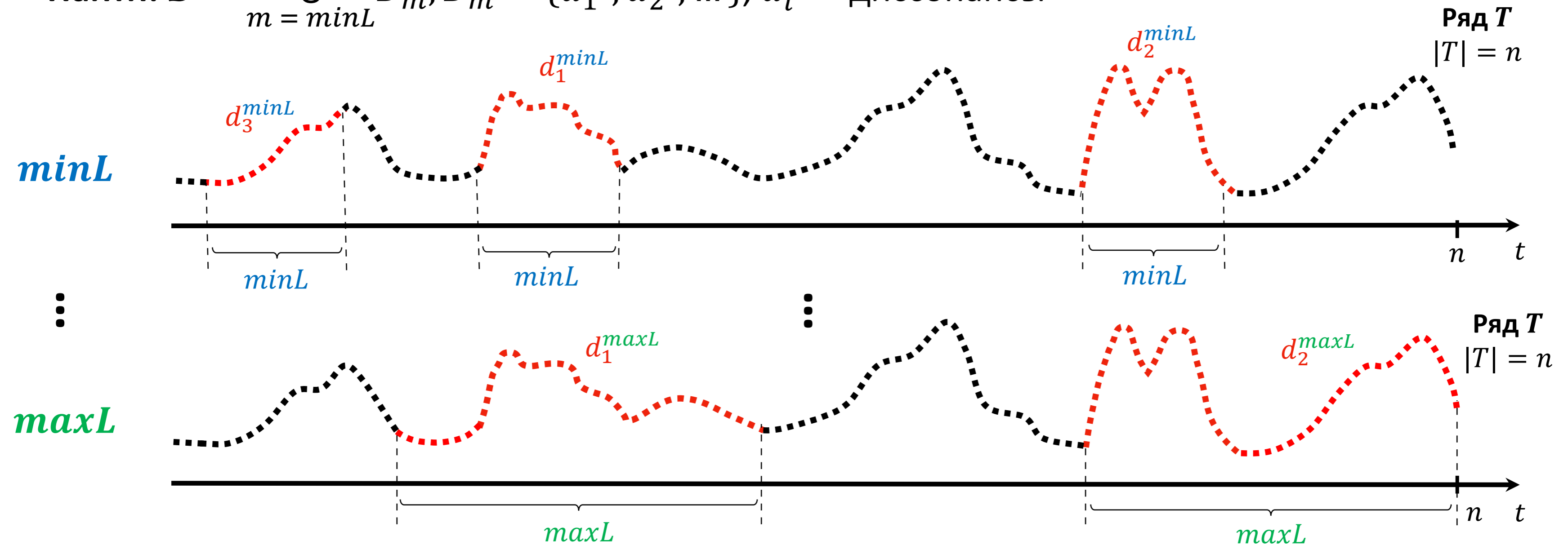


<sup>1)</sup> Nakamura T., et al. MERLIN: parameter-free discovery of arbitrary length anomalies in massive time series archives. IEEE ICDM 2020. pp. 1190-1195.

# Постановка задачи автоматизированного поиска диссонансов<sup>1)</sup>

• **Дано:** временной ряд  $T$ , диапазон длин диссонансов  $minL, \dots, maxL$

• **Найти:**  $\mathcal{D} = \bigcup_{m=minL}^{maxL} D_m, D_m = \{d_1^m, d_2^m, \dots\}, d_i^m$  – диссонансы



<sup>1)</sup> Nakamura T., et al. MERLIN: parameter-free discovery of arbitrary length anomalies in massive time series archives. IEEE ICDM 2020. pp. 1190-1195.



# PALMAD: Parallel Arbitrary Length MERLIN-based Anomaly Discovery

1. Применение  $ED_{\text{norm}}^2$  в качестве функции расстояния<sup>1)</sup>

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

2. Сокращение избыточных вычислений  $\mu$  и  $\sigma$  при вычислении  $ED_{\text{norm}}^2$

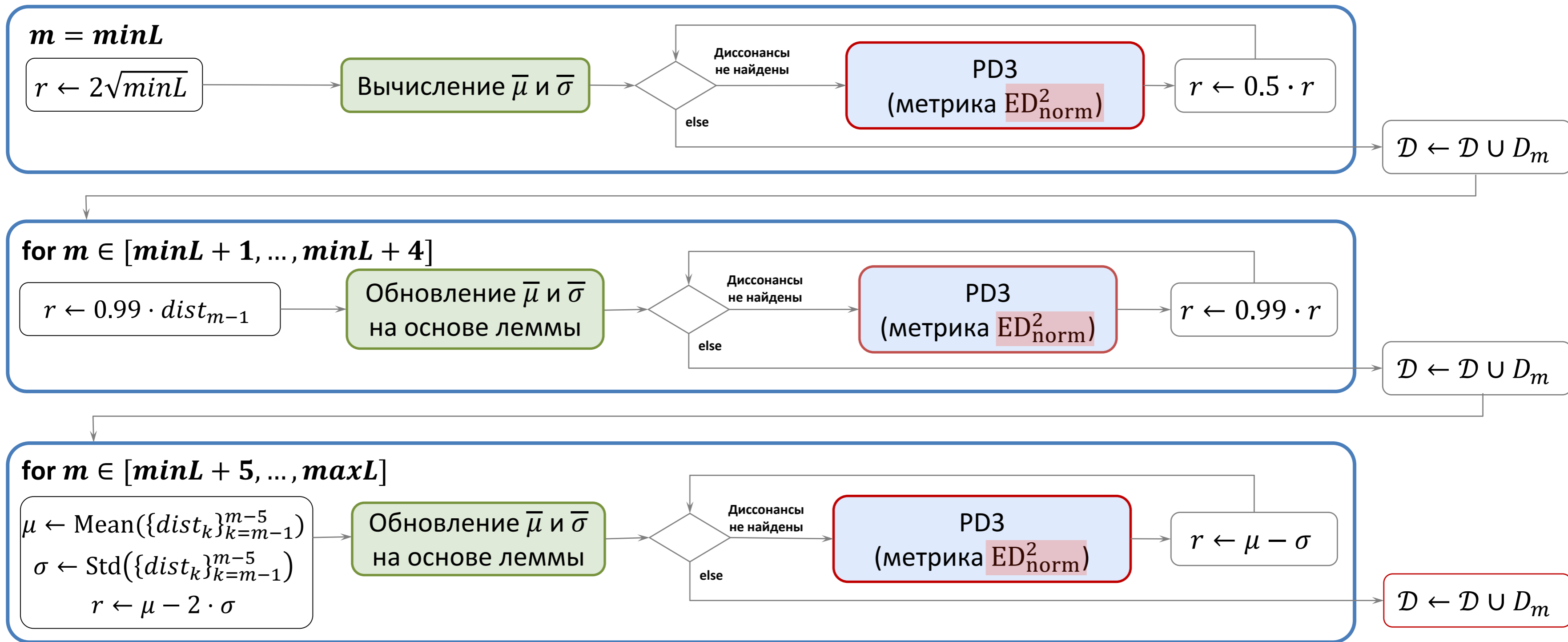
**Лемма.** Пусть даны ряд  $T$ ,  $|T| = n$  и подпоследовательности  $T_{i,m}$  и  $T_{i,m+1}$ . Тогда

$$\mu_{T_{i,m+1}} = \frac{1}{m+1} (m\mu_{T_{i,m}} + t_{i+m}), \quad \sigma_{T_{i,m+1}}^2 = \frac{m}{m+1} \left( \sigma_{T_{i,m}}^2 + \frac{1}{m+1} (\mu_{T_{i,m}} - t_{i+m})^2 \right).$$

3. Автоматизированный подбор порога  $r$
4. Тепловая карта диссонансов

<sup>1)</sup> Mueen A. et al. Fast approximate correlation for massive time-series data. SIGMOD 2010. pp. 171-182. ACM (2010). <https://doi.org/10.1145/1807167.1807188>

# Схема PALMAD



Подбор порога  $r$

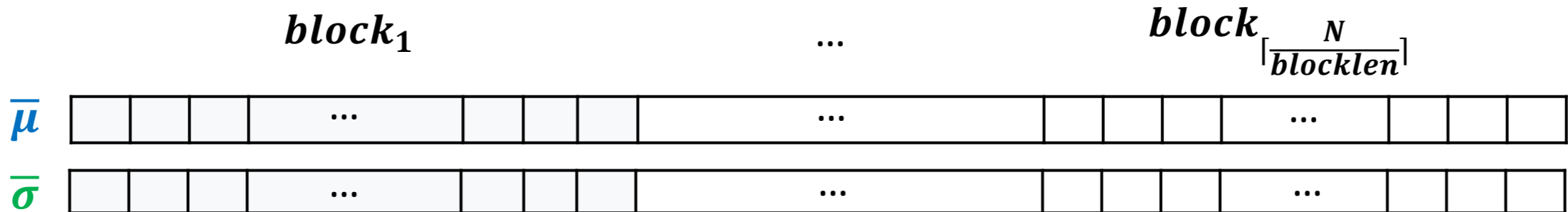
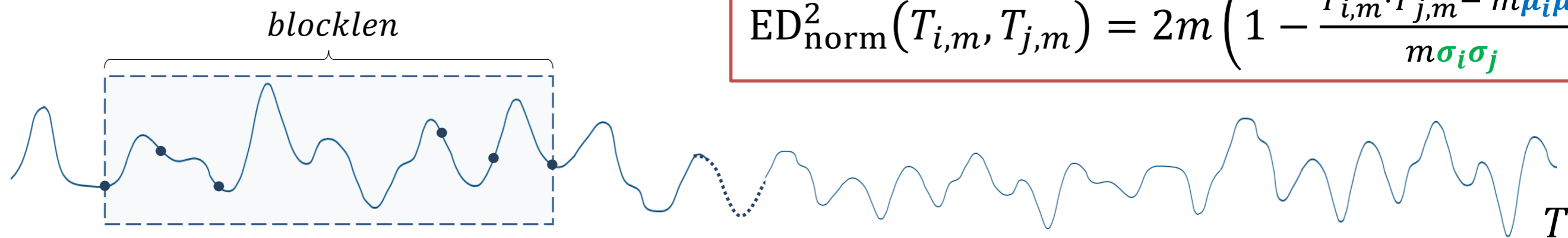
Предобработка  
(сокращение вычислений)

Поиск диссонансов  
(использование PD3 и  $ED_{\text{norm}}^2$ )

# Сокращение избыточных вычислений $\bar{\mu}$ и $\bar{\sigma}$

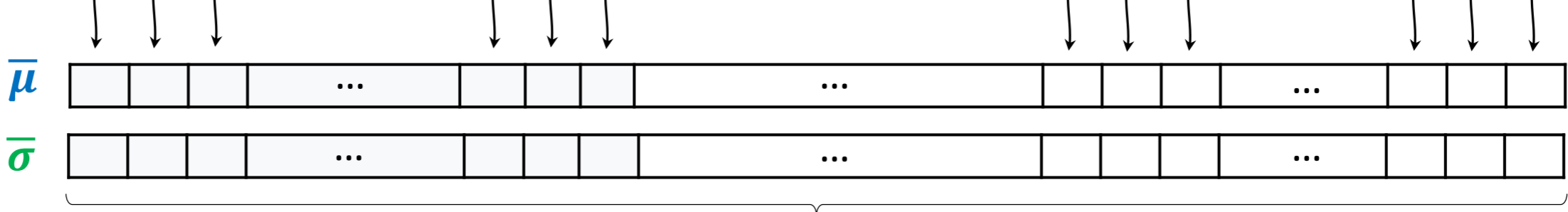
Вычисленные  $\bar{\mu}$  и  $\bar{\sigma}$  для под-ей ряда длины, меньшей на 1

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



**Лемма**

$$\mu_{T_{i,m+1}} = \frac{1}{m+1} (m\mu_{T_{i,m}} + t_{i+m}), \quad \sigma_{T_{i,m+1}}^2 = \frac{m}{m+1} \left( \sigma_{T_{i,m}}^2 + \frac{1}{m+1} (\mu_{T_{i,m}} - t_{i+m})^2 \right)$$

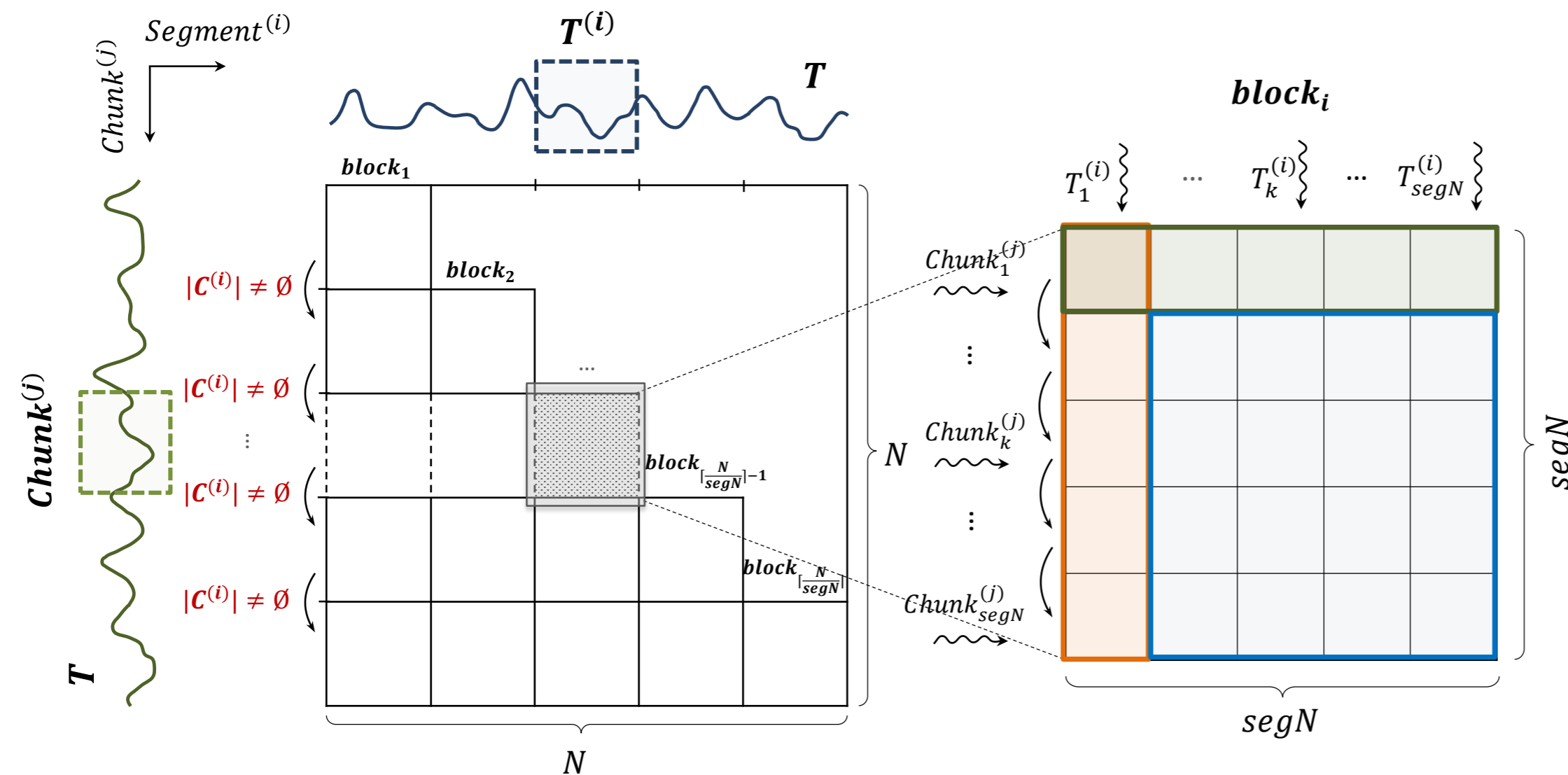


$$N = n - minL + 1$$

# PD3: Parallel DRAG-based Discord Discovery

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

1. PD3 включает 2 фазы: отбор кандидатов в диссонансы и очистка кандидатов
2. Использование концепции параллелизма по данным
3. Эффективное вычисление скалярных произведений  $T_{i,m} \cdot T_{j,m}$

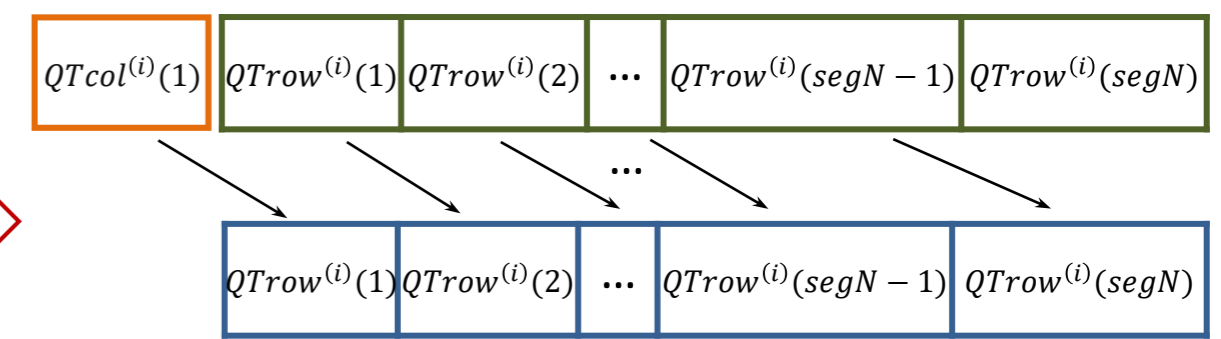


$$\text{QTrow}^{(i)}(tid) = \sum_{k=1}^m T_{tid}^{(i)}(k) \cdot Chunk_1^{(j)}(k)$$

$$\text{QTcol}^{(i)}(tid) = \sum_{k=1}^m T_1^{(i)}(k) \cdot Chunk_{tid}^{(j)}(k)$$

Выч. сложность  $O(1)$  вместо  $O(m)$ !

$$\text{QTrow}^{(i)}(tid) = \text{QTrow}^{(i)}(tid - 1) - T_{tid-1}^{(i)}(1) \cdot Chunk_{tid-1}^{(j)}(1) + T_{tid}^{(i)}(m) \cdot Chunk_{tid}^{(j)}(m)$$



# Эксперименты

- **Конкуренты** (поиск top-1 диссонанса на GPU)

- **KBF\_GPU**: Thuy T.T.H. et al. A new discord definition and an efficient time series discord detection method using GPUs. ICSED 2021. pp. 63–70. <https://doi.org/10.1145/3507473.3507483>.
- **Zhu 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. <https://doi.org/10.1109/TPDS.2021.3055765>.

- **Данные**<sup>1,2)</sup>

Временной ряд	Длина ряда, $n$	Длина диссонанса, $minL = maxL$
Space shuttle	5 000	150
ECG	45 000	200
ECG2	21 600	400
Koski-ECG	100 000	458
Power demand	33 220	750
Respiration	24 125	250
RandomWalk1M	$1 \cdot 10^7$	512
RandomWalk2M	$2 \cdot 10^7$	512

- **Аппаратные платформы**

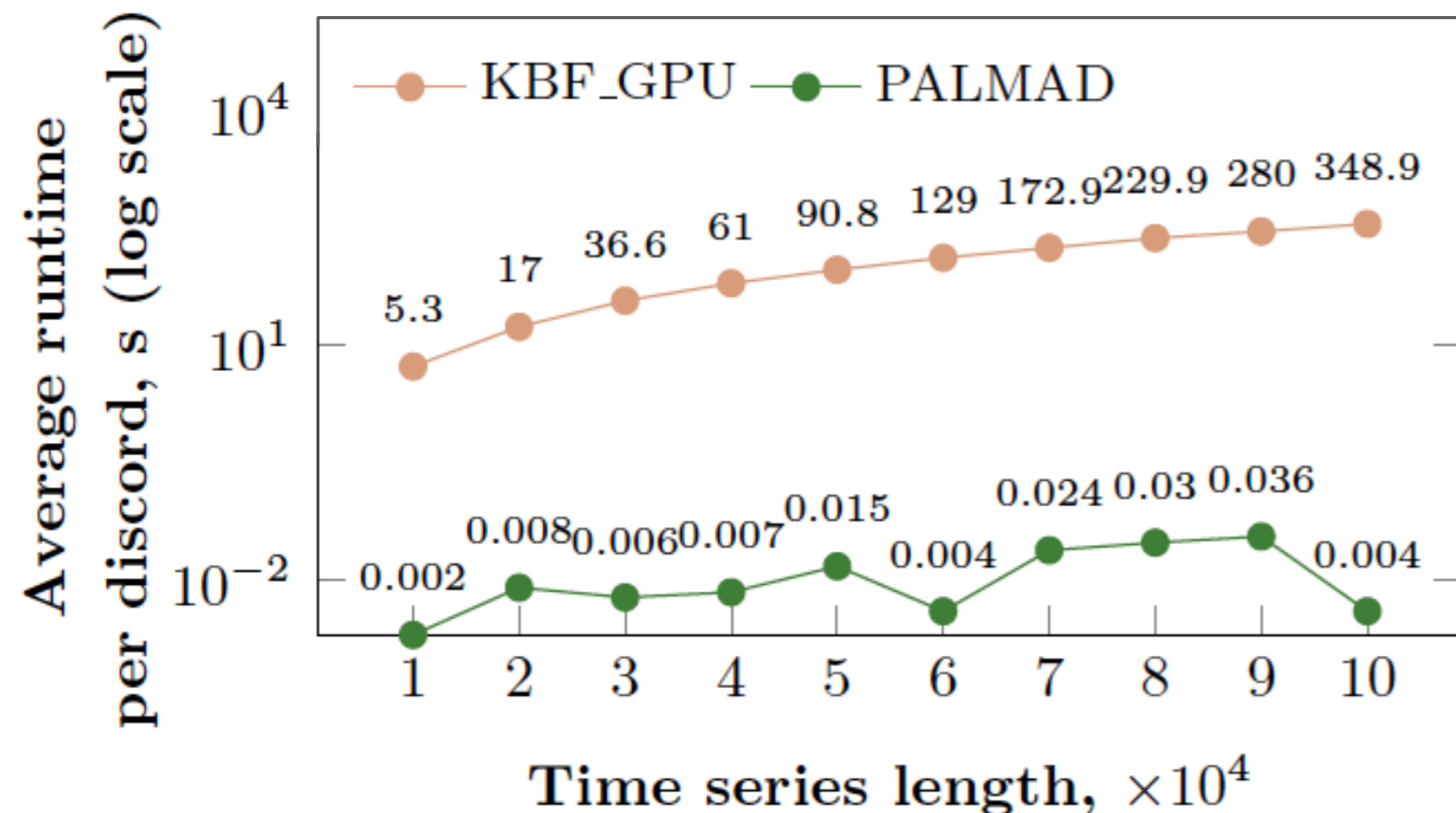
Характеристика	GPU-МГУ	GPU-ЮУрГУ
Производитель, семейство	NVIDIA Tesla	
Модель	P100	V100
# CUDA-ядер	3 584	5 120
Тактовая частота, GHz	1.19	1.3
Оперативная память, Gb	16	32
Пик. пр-ть (double), TFLOPS	4	7

<sup>1)</sup> Keogh E., Lin J., Fu A. HOT SAX: Finding the most unusual time series subsequence: Algorithms and applications. Proc. 5th IEEE Int. Conf. Data Mining 2004: 440–449. URL: <http://www.cs.ucr.edu/~eamonn/discords/>.

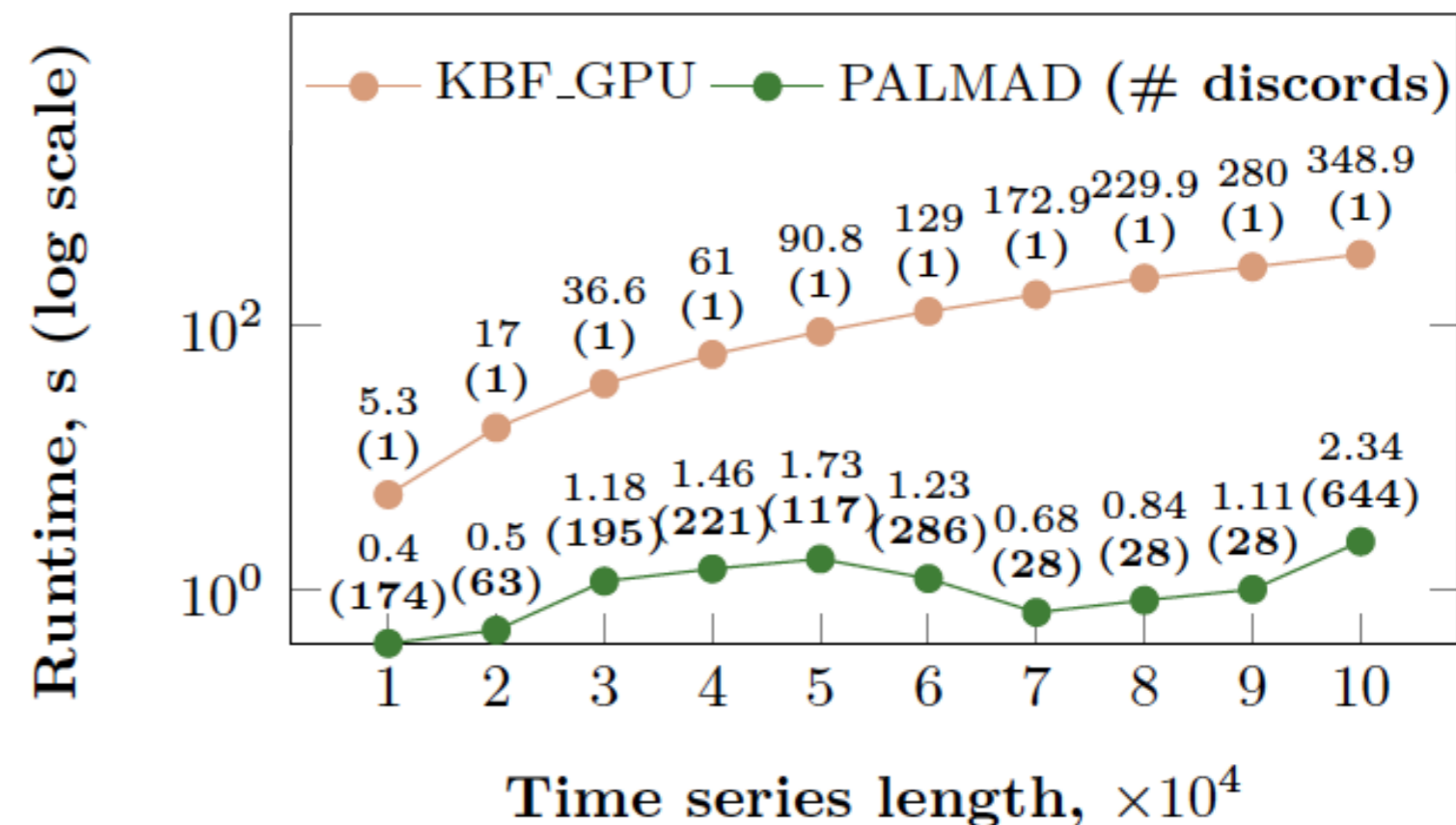
<sup>2)</sup> Pearson K. The problem of the random walk. Nature 72(394). <https://doi.org/10.1038/072342a0>.

# Производительность: сравнение с KBF\_GPU<sup>1)</sup>

Среднее время на поиск **одного** диссонанса



Время на поиск **всех** диссонансов



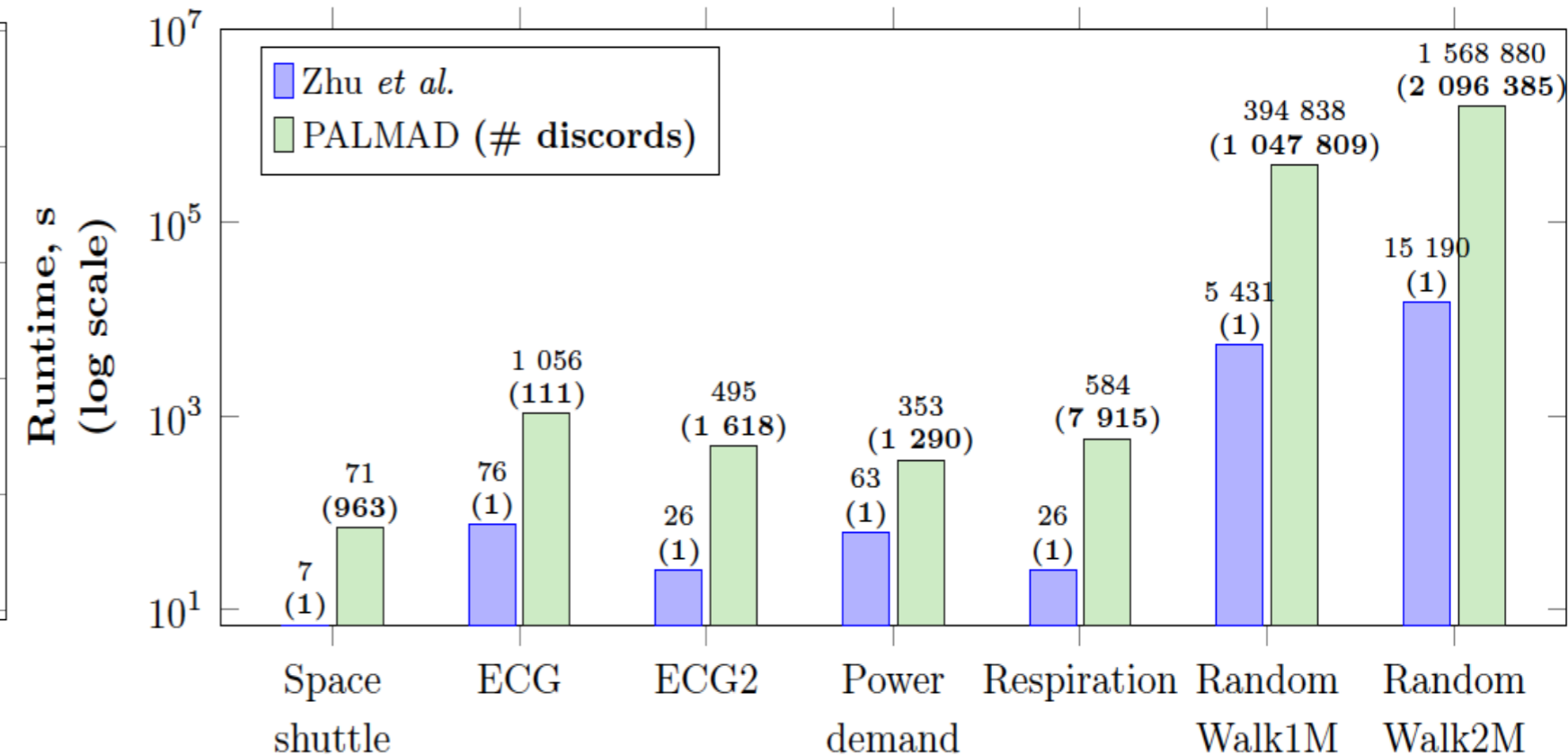
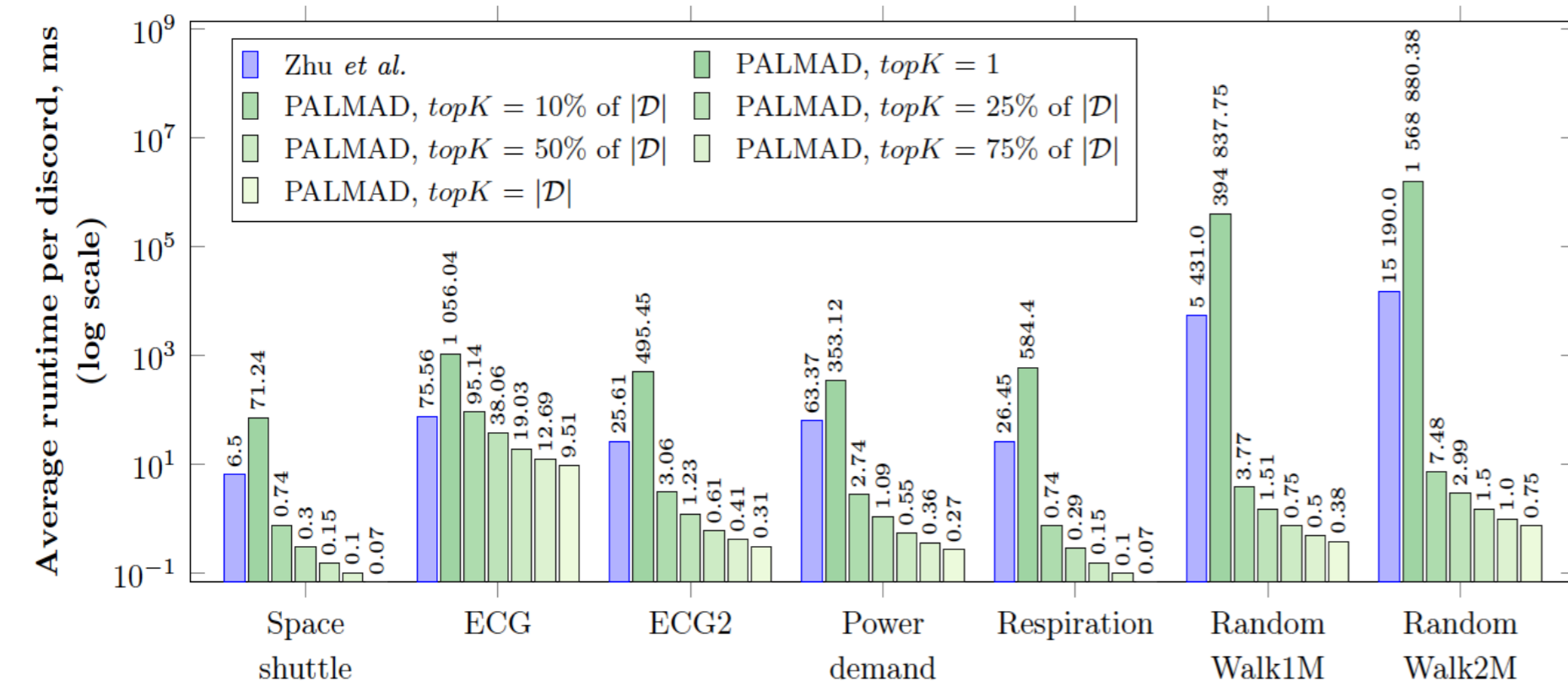
PALMAD опережает KBF\_GPU как по общему времени работы, так и по среднему времени на поиск одного диссонанса

<sup>1)</sup> Thuy T.T.H. et al. A new discord definition and an efficient time series discord detection method using GPUs. ICSED 2021. pp. 63–70. <https://doi.org/10.1145/3507473.3507483>.

# Производительность: сравнение с Zhu et al.<sup>1)</sup>

Среднее время на поиск **одного** диссонанса

Время на поиск **всех** диссонансов



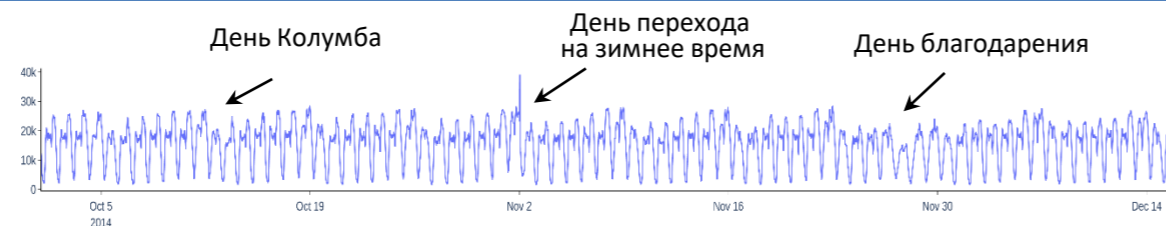
PALMAD значительно опережает алгоритм Zhu et al. по среднему времени на поиск одного диссонанса, начиная с поиска топ-k диссонансов, где k=10% от фактического числа диссонансов

<sup>1)</sup> 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. <https://doi.org/10.1109/TPDS.2021.3055765>.

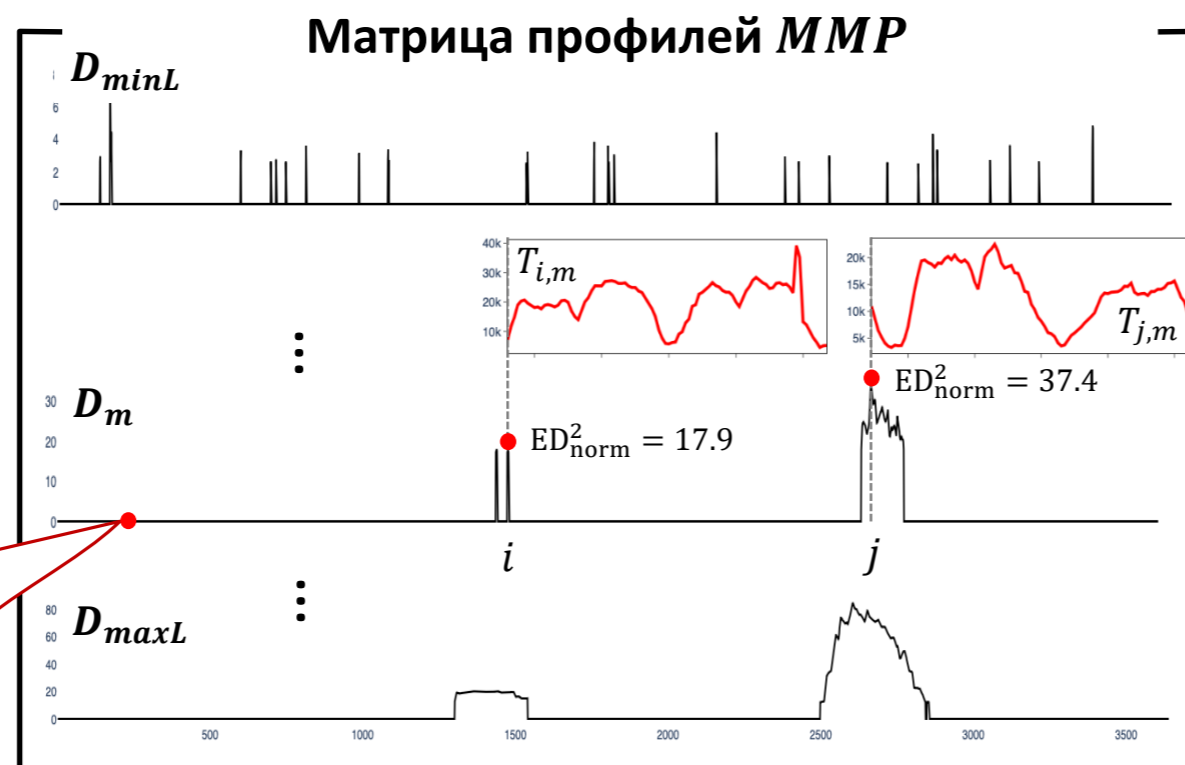
# Тепловая карта диссонансов

**PALMAD**

Исходный ряд



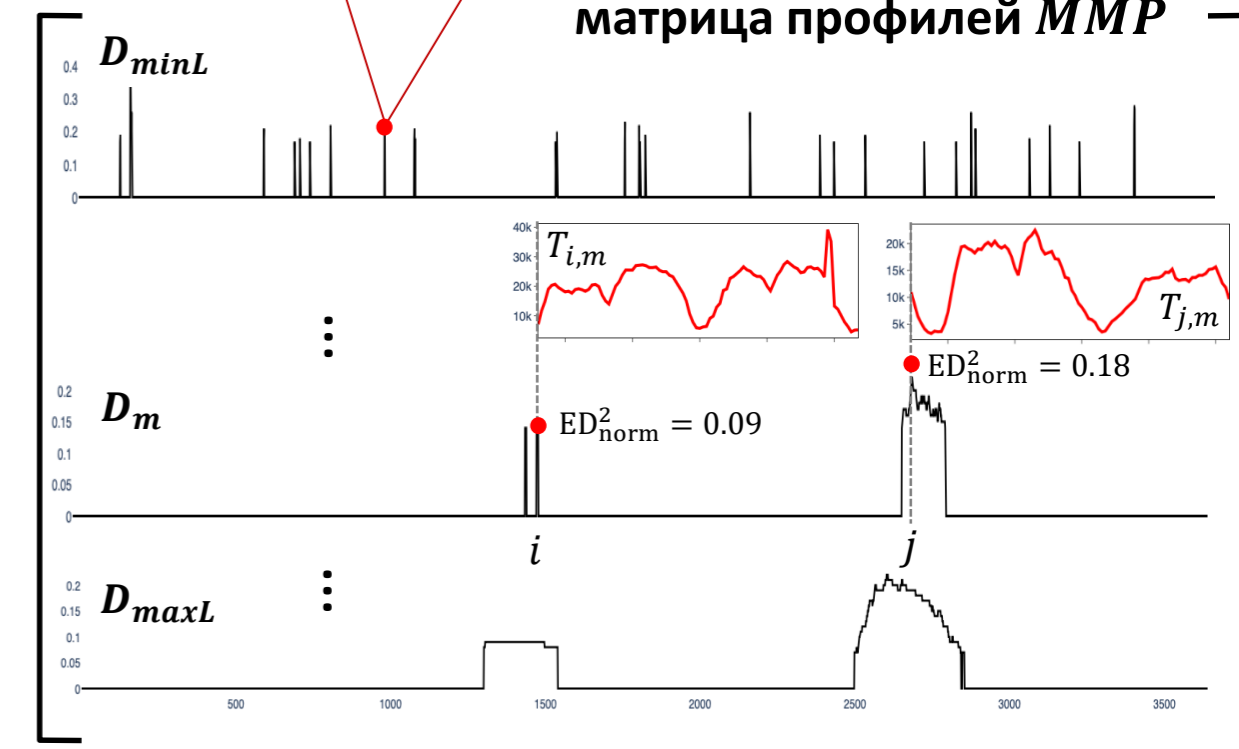
Автомат. поиск диссонансов



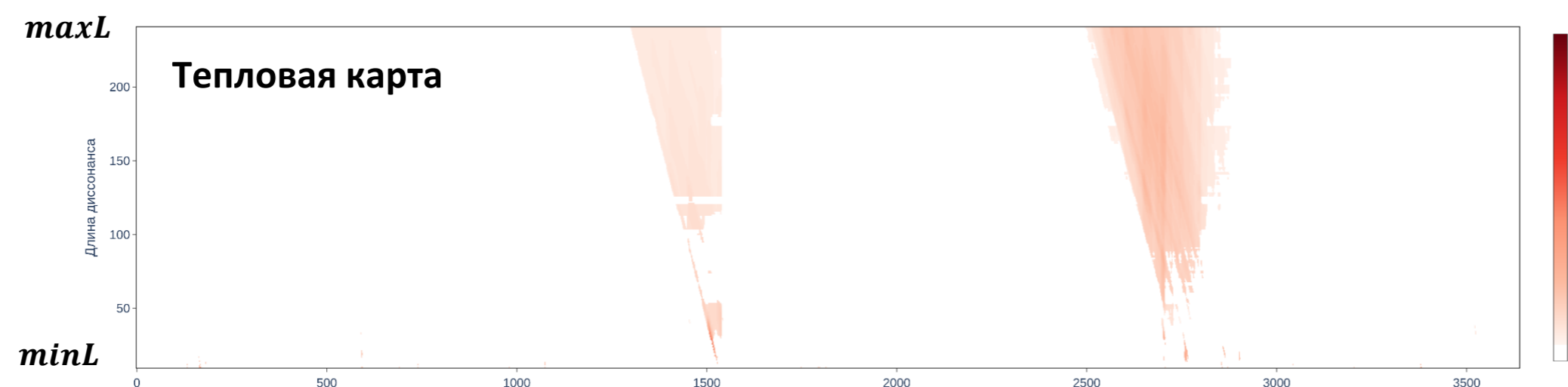
$$MMP(m, i) = \frac{MMP(m, i)}{2m}$$

Нормализованная матрица профилей MMP

Нормализация



$$MMP(m, i) = \begin{cases} ED_{norm}^2(T_{i,m}, NN), & T_{i,m} \in D_m \\ 0, & otherwise \end{cases}$$

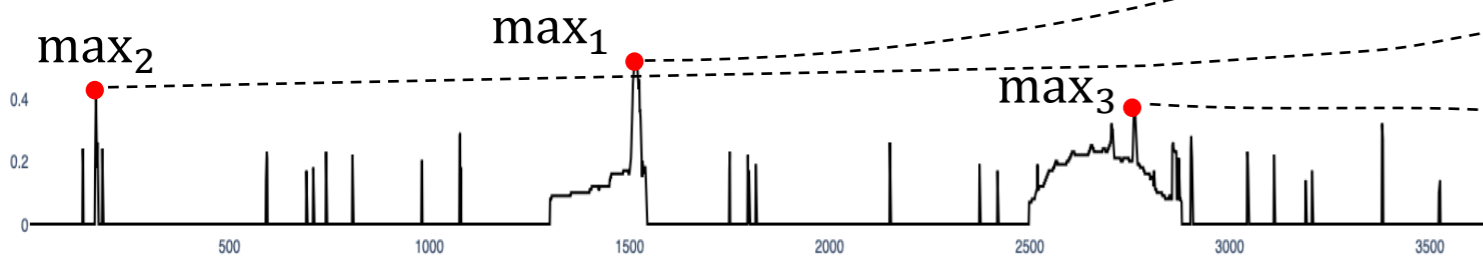
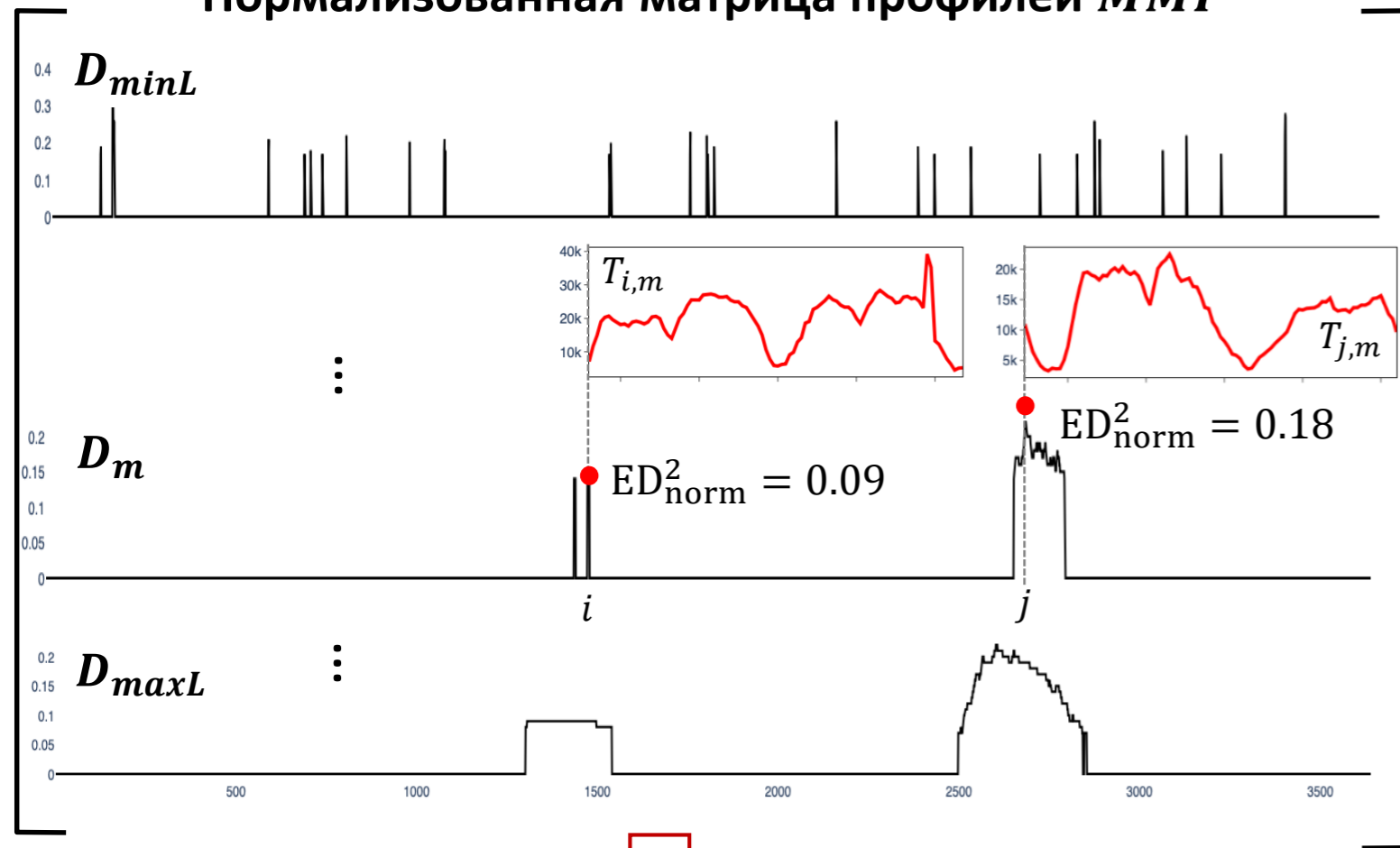


Визуализация (расстояние → яркость)

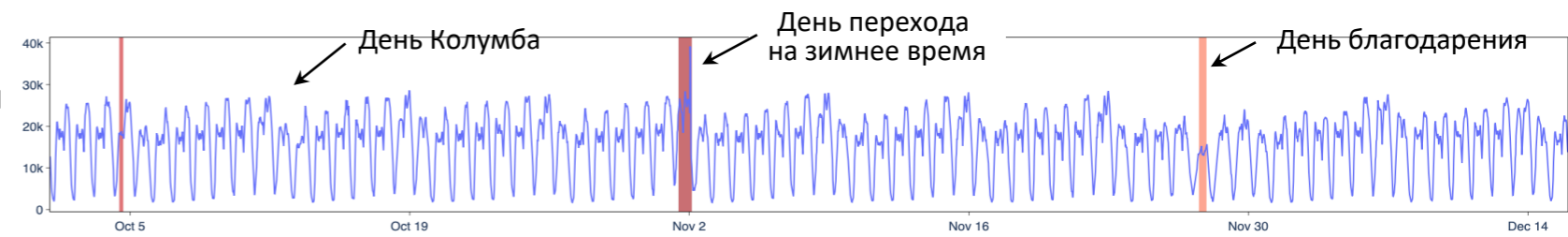


# Ранжирование диссонансов

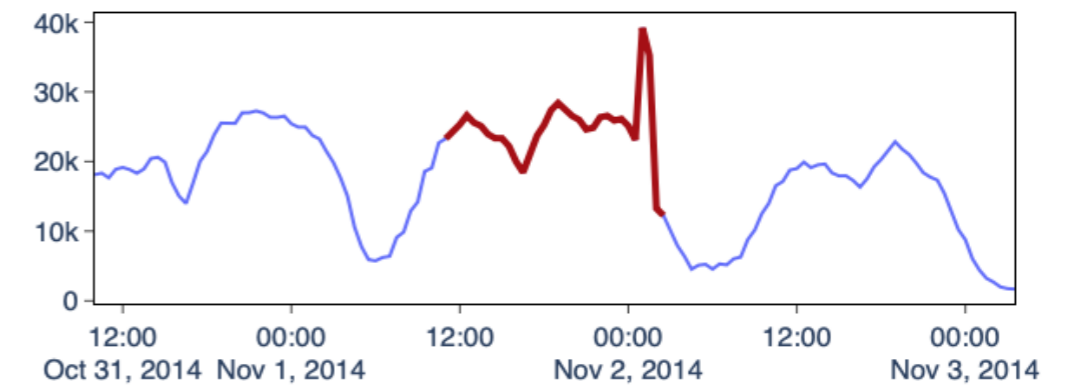
Нормализованная матрица профилей MMP



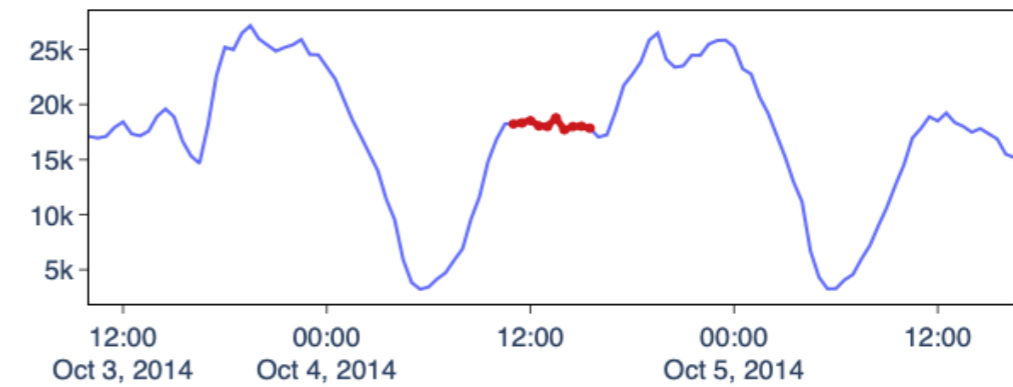
Исходный ряд



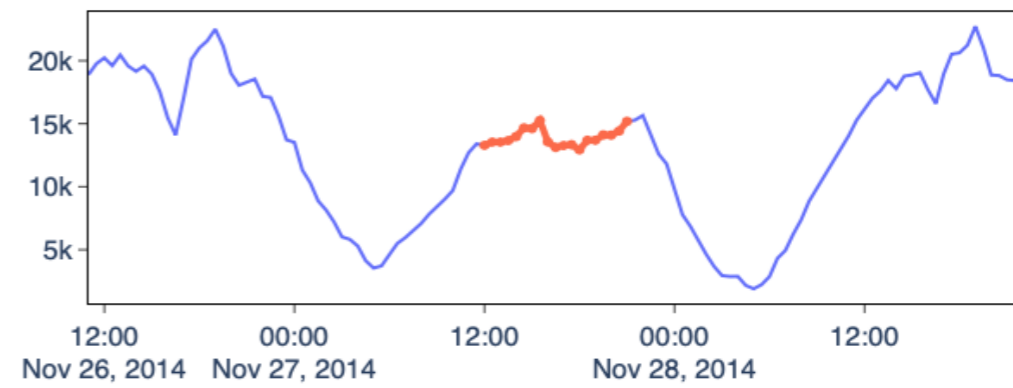
Топ-1 диссонанс,  $m = 32$



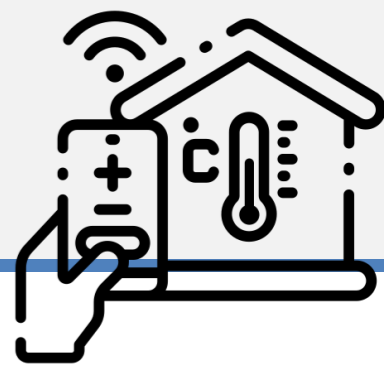
Топ-2 диссонанс,  $m = 10$



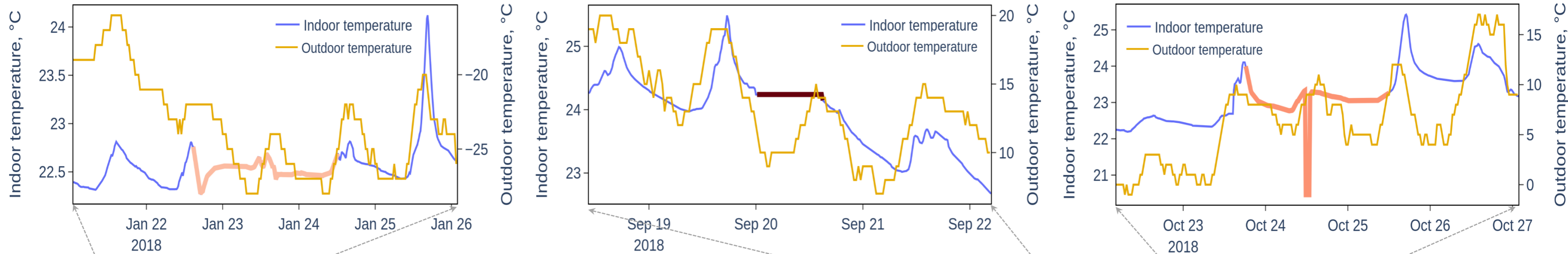
Топ-3 диссонанс,  $m = 19$



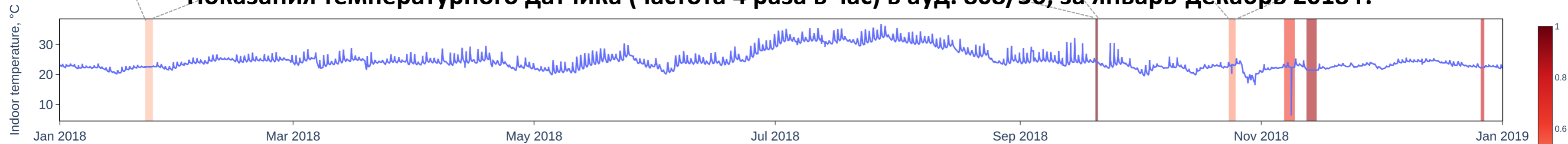
# Выявление аномалий в системе отопления ЮУрГУ



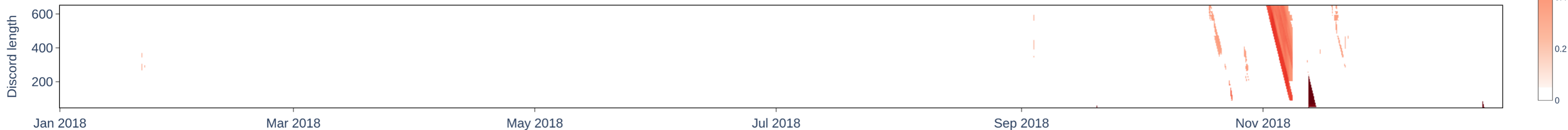
## Примеры найденных аномалий длительностью 0.5-2 суток



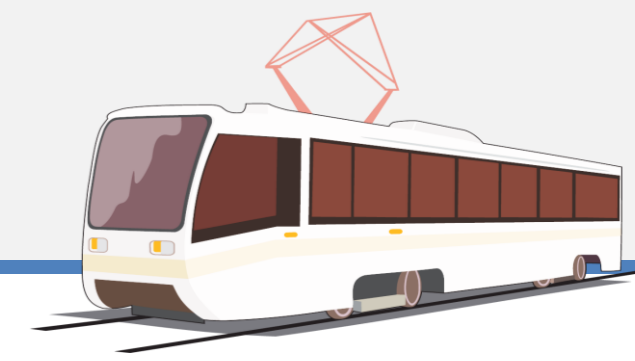
## Показания температурного датчика (частота 4 раза в час) в ауд. 808/36, за январь-декабрь 2018 г.



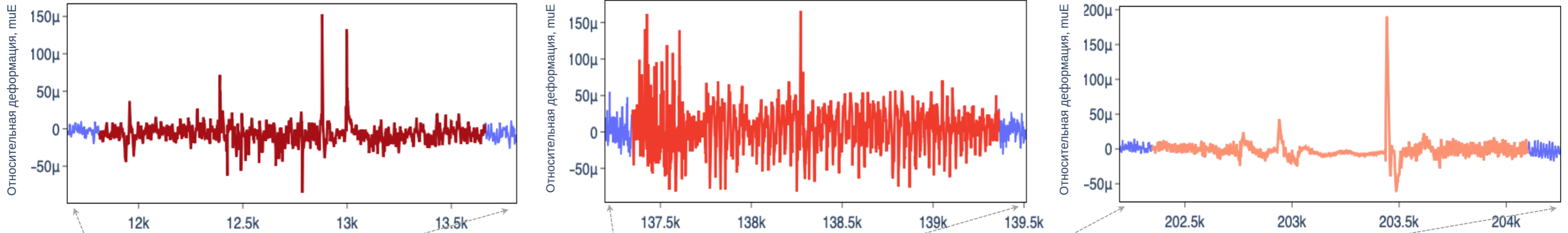
## Тепловая карта найденных аномалий длительностью 0.5-2 суток



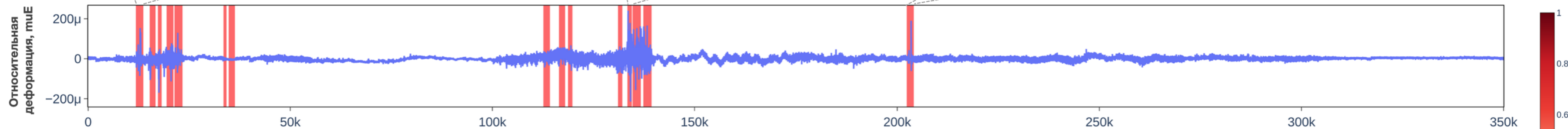
# Выявление аномалий в машиностроении



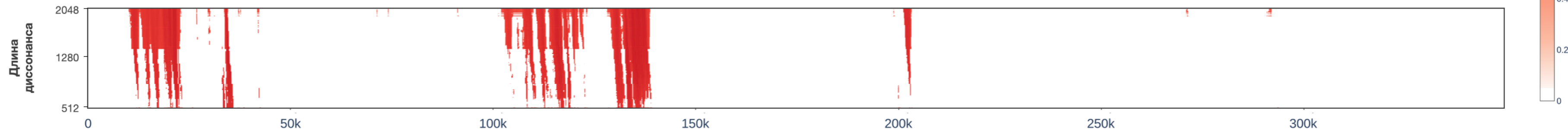
## Примеры найденных аномалий длительностью 0.25-0.5 секунд



## Относительные деформации механизма стыковки вагонов трамвая (частота 4 096 раз в сек) за 1.5 минуты



## Тепловая карта найденных аномалий длительностью 0.25-0.5 секунд



# Заключение

- Предложен новый параллельный алгоритм поиска аномалий временного ряда PALMAD для GPU
- Будущие исследования:
  - Разработка версии PALMAD для кластера с GPU-узлами
  - Применение PALMAD в нейросетевой модели для поиска аномалий временного ряда в режиме реального времени

**Спасибо за внимание! Вопросы?**

Яна Александровна Краева

[kraevaya@susu.ru](mailto:kraevaya@susu.ru)