

Международная научная конференция
Параллельные вычислительные технологии (ПаВТ'2022)
Дубна, 29–31 марта 2022 г.

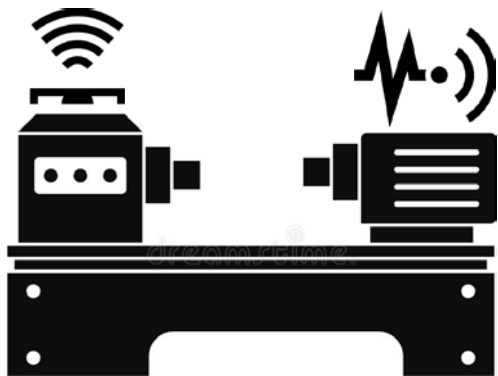
Параллельный алгоритм восстановления пропущенных значений потокового временного ряда в режиме реального времени

М.Л. Цымблер¹, А.Н. Полуянов²,

¹Южно-Уральский государственный университет (Челябинск),

²Институт математики им. С.Л. Соболева СО РАН (Омск)

Восстановление потоковых данных онлайн



Предиктивное ТО
промышленных агрегатов



Интернет
вещей



Персональная
медицина



Мониторинг погоды
и моделирование климата

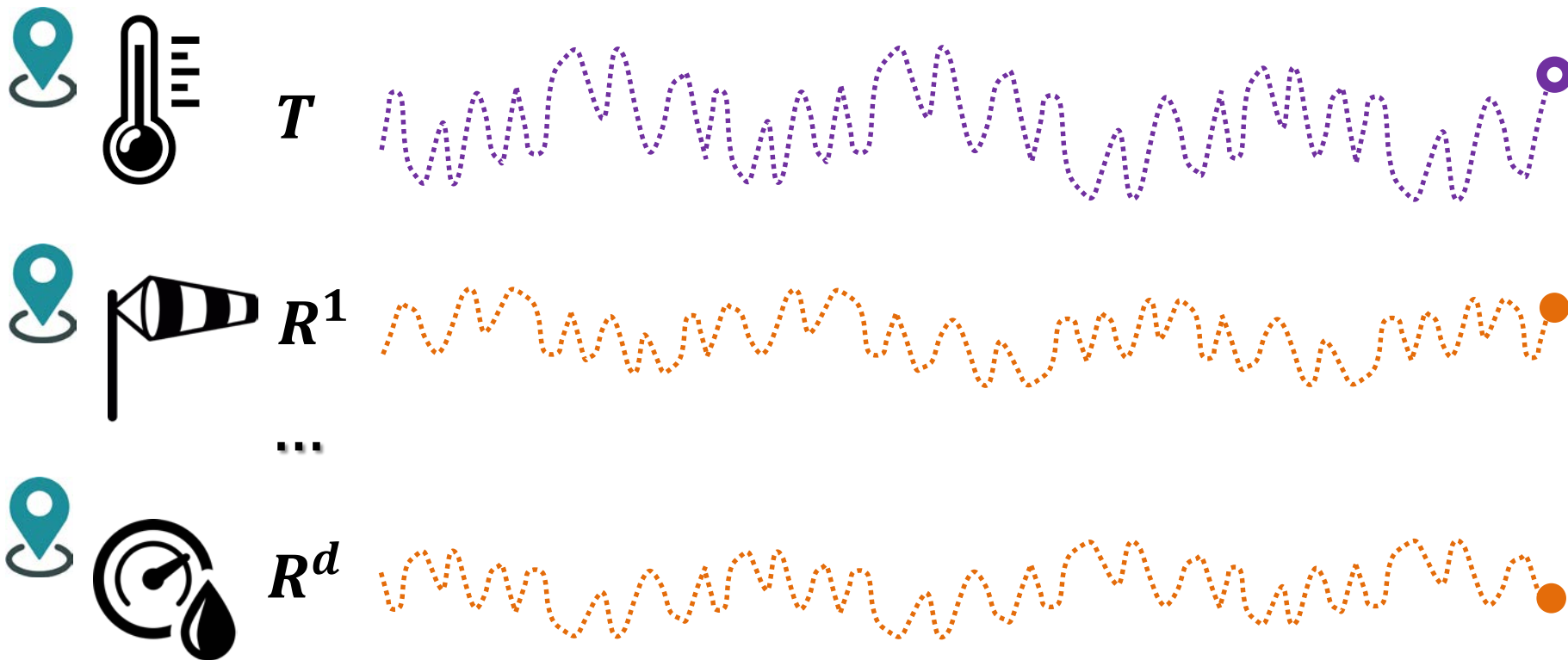


T



$t_n = \text{NULL}$

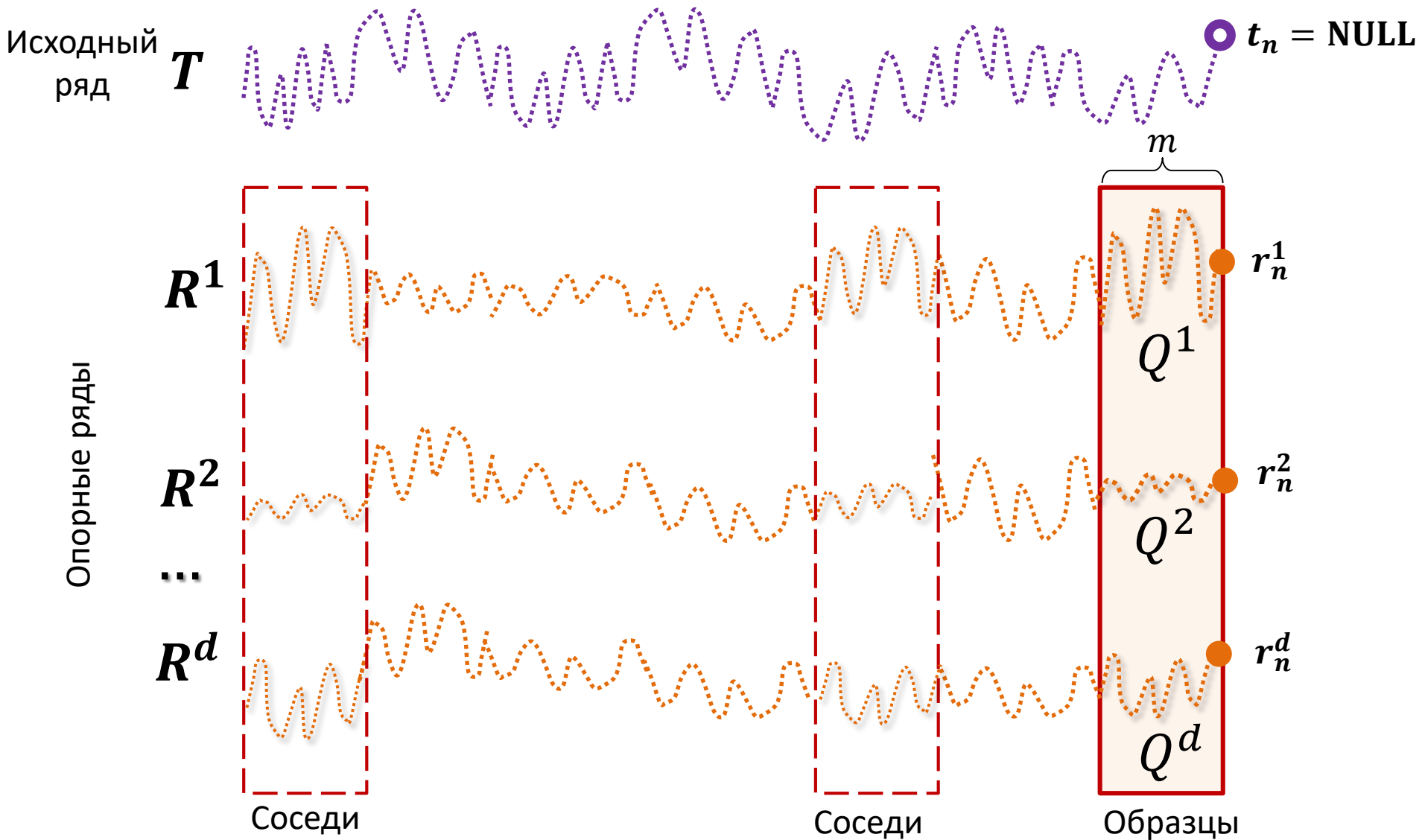
Восстановление по опорным рядам



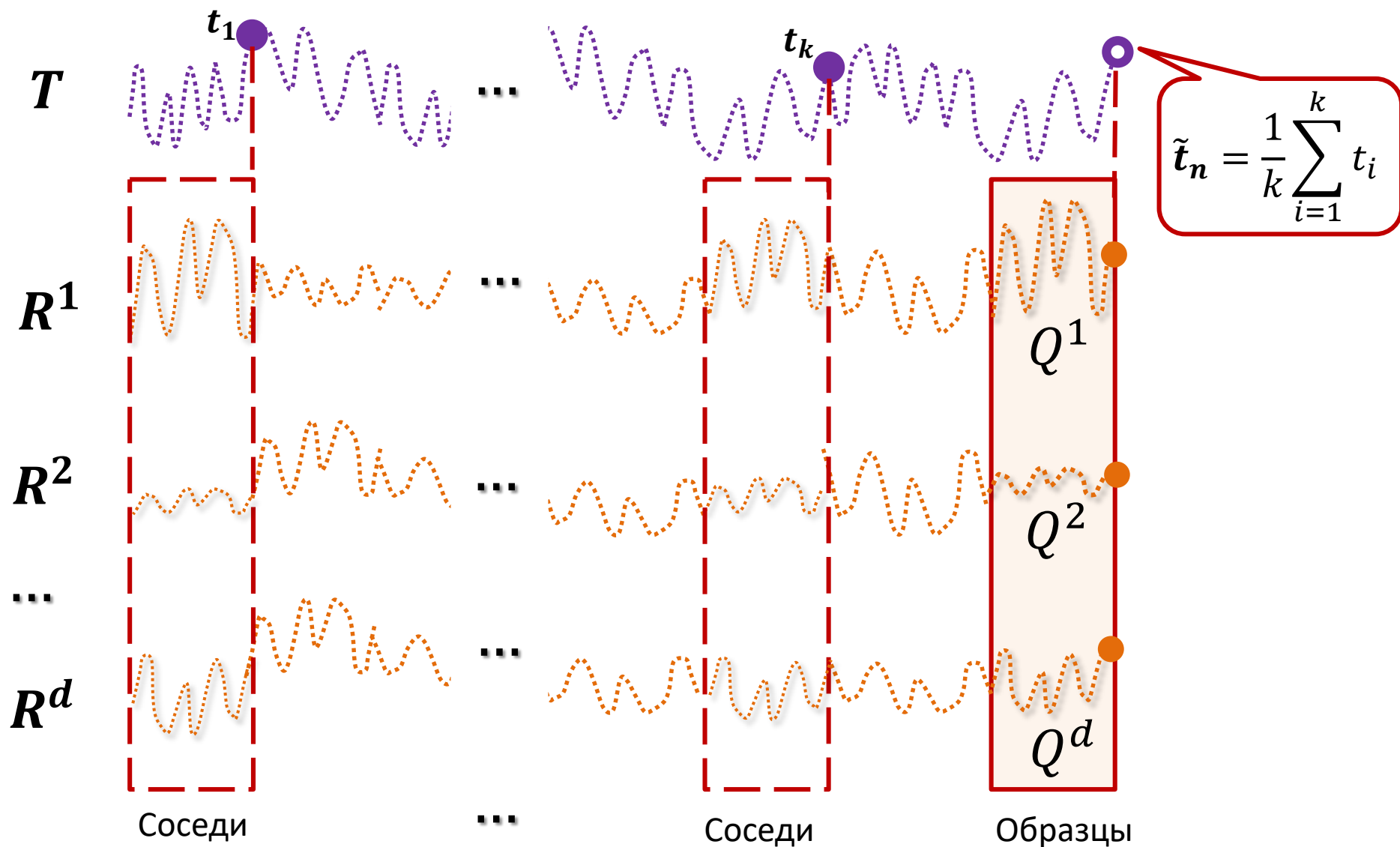
Эвристика

похожие промежутки в исходном ряде возникают в тех же интервалах времени, что и в опорных рядах

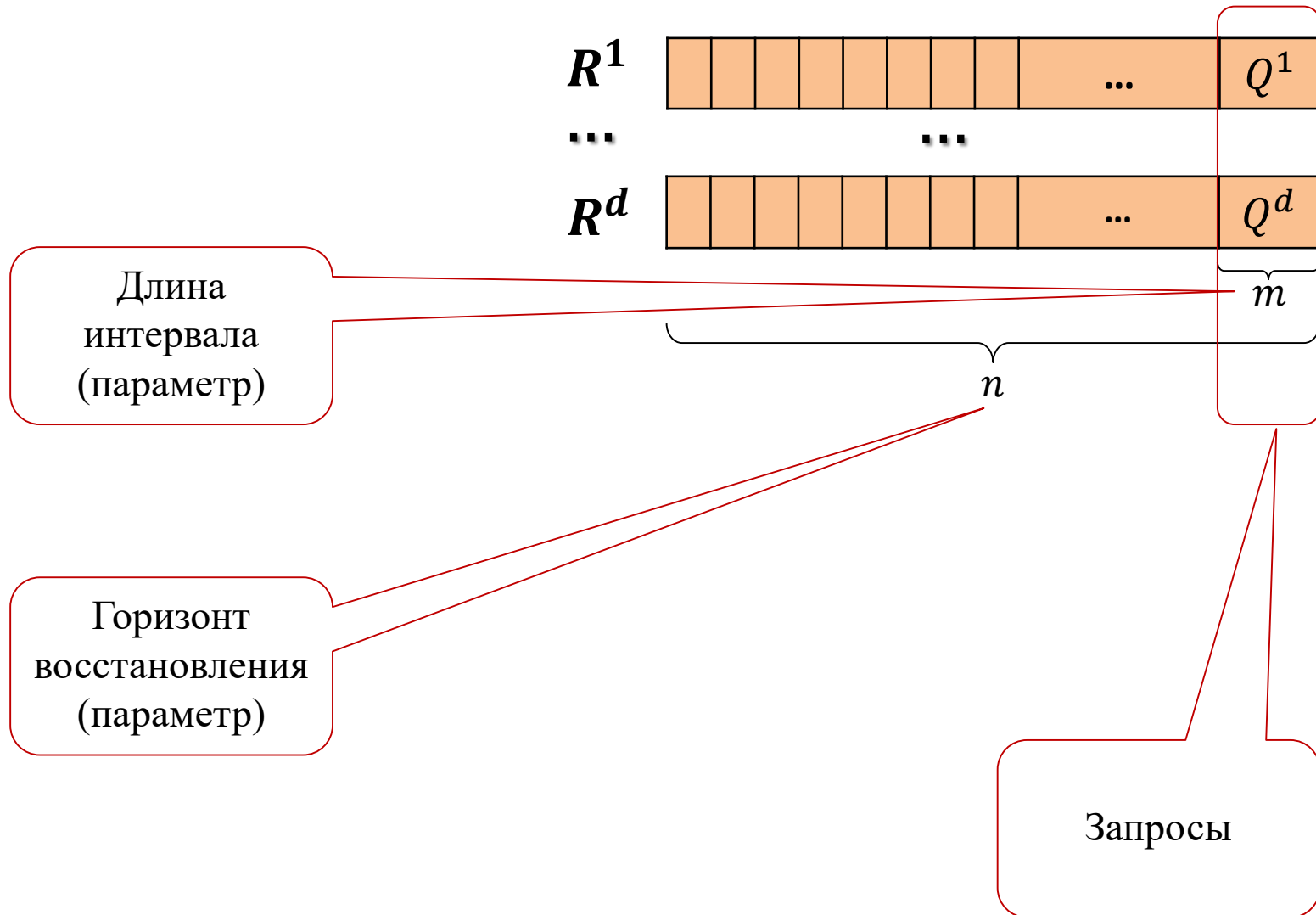
Поиск k ближайших соседей



Реконструкция пропущенного значения



Алгоритм IDK (Imputation by DTW & kNN)



Алгоритм IDK: поиск соседей по образцу

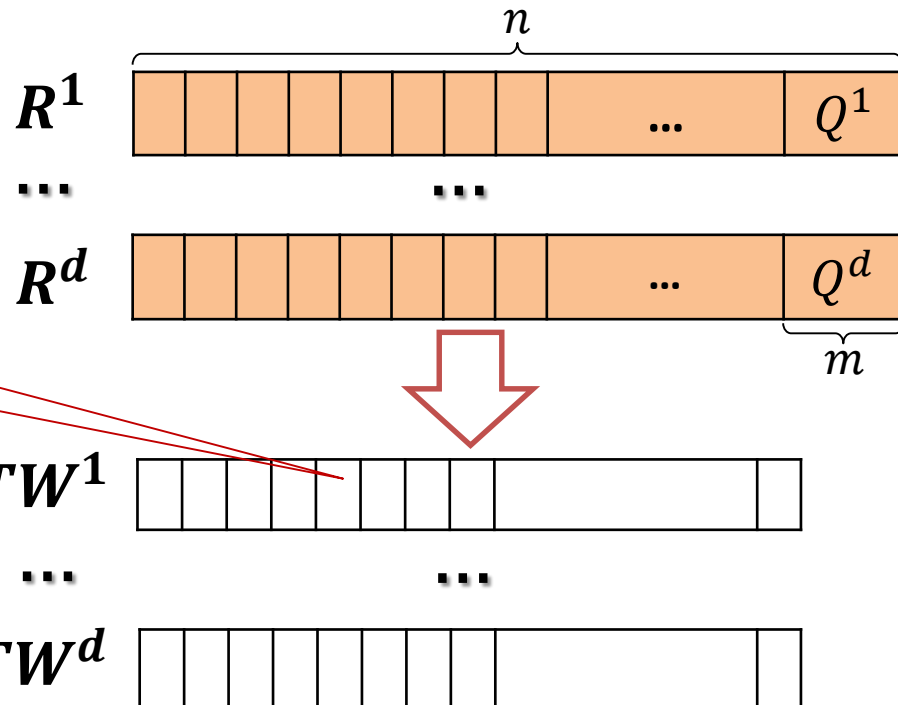
Расстояние между интервалом $R^j[i:m]$ и образцом Q^j в смысле меры DTW (Dynamic Time Warping)

1. Поиск соседей по образцу

for $j := 1$ to d do

for $i := 1$ to $n - m + 1$ do

$DTW^j(R^j[i:m], Q^j)$



Алгоритм IDK: скоринг интервалов

1. Поиск соседей по образцу

for $j := 1$ to d do

for $i := 1$ to $n - m + 1$ do

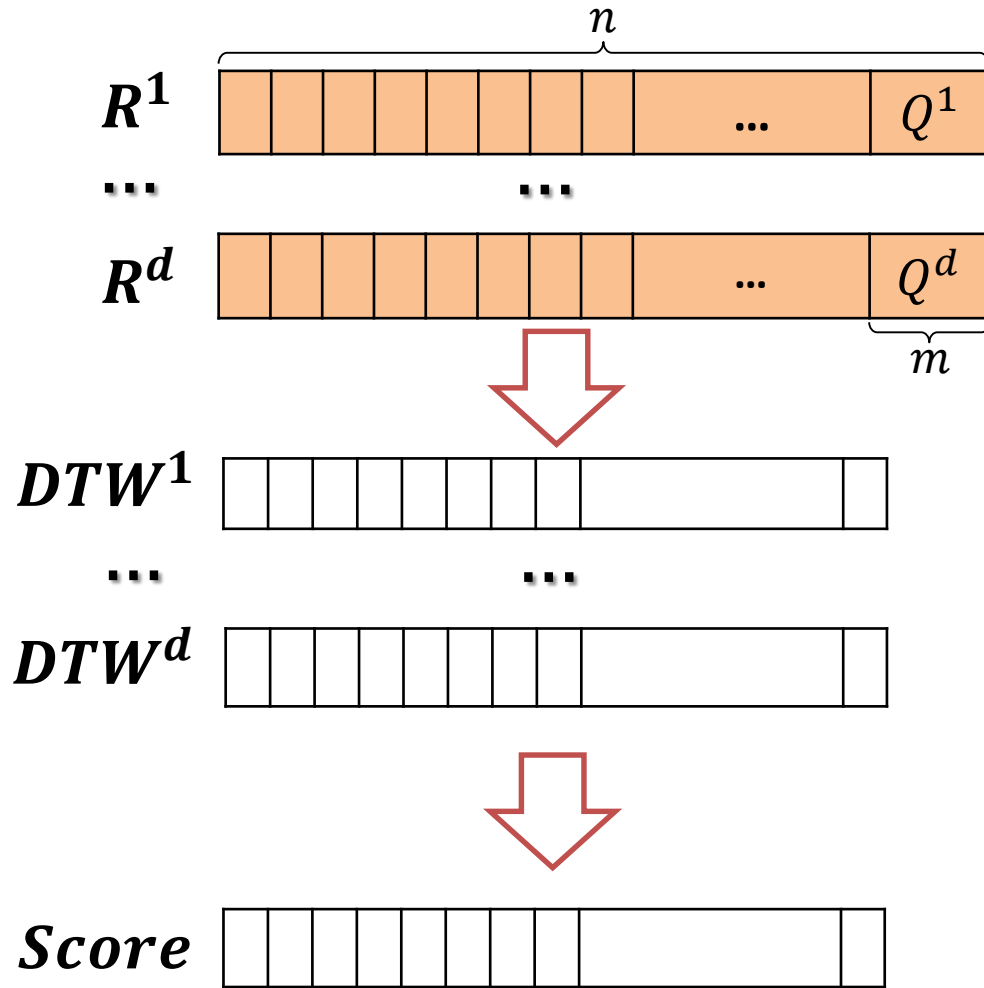
$$DTW^j(R^j[i:m], Q^j)$$

2. Скоринг интервалов

for $i := 1$ to $n - m + 1$ do

for $j := 1$ to d do

$$Score(i) := Score(i) + \frac{1}{DTW^j(i) + \varepsilon}$$



Алгоритм IDK: отбор интервалов

1. Поиск соседей по образцу

for $j := 1$ to d do

for $i := 1$ to $n - m + 1$ do

$$DTW^j(R^j[i:m], Q^j)$$

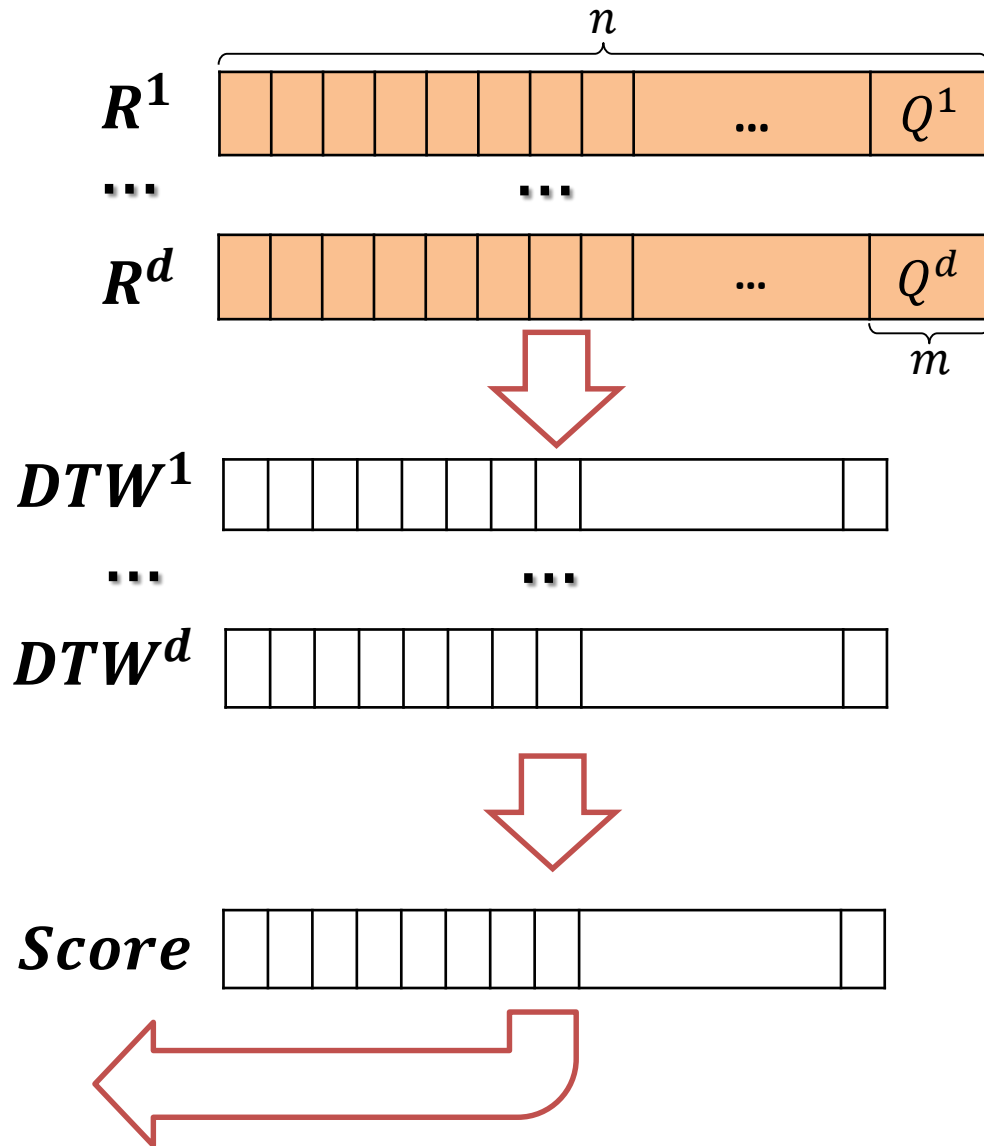
2. Скоринг интервалов

for $i := 1$ to $n - m + 1$ do

for $j := 1$ to d do

$$Score(i) := Score(i) + \frac{1}{DTW^j(i)+\epsilon}$$

3. Отбор TOP- k интервалов, реконструкция

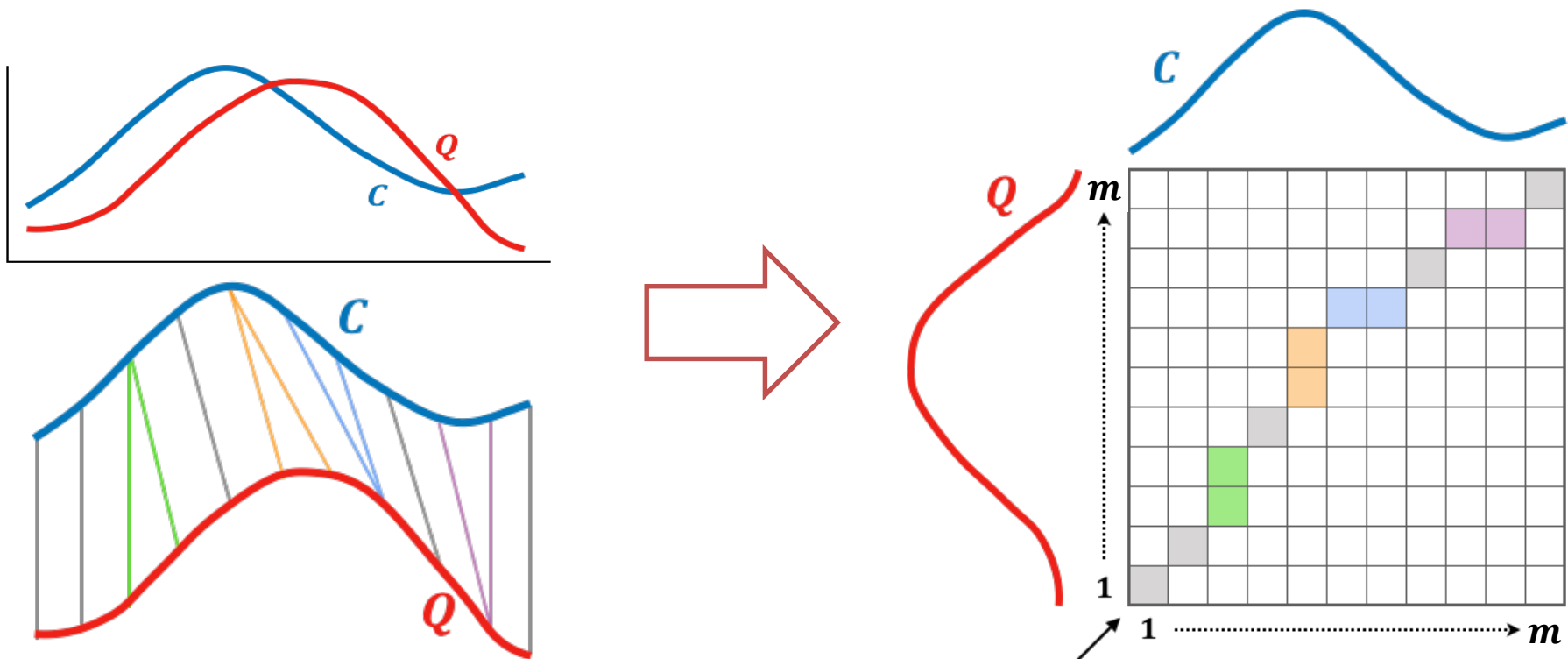


Расстояние DTW (Dynamic Time Warping)

$$\mathbf{DTW}(Q, C) = d(m, m)$$

$$d(i, j) = (q_i - c_j)^2 + \min \begin{cases} d(i-1, j) \\ d(i, j-1) \\ d(i-1, j-1) \end{cases}$$

$$d(0, 0) = 0, d(i, 0) = d(0, j) = +\infty; 1 \leq i, j \leq m$$



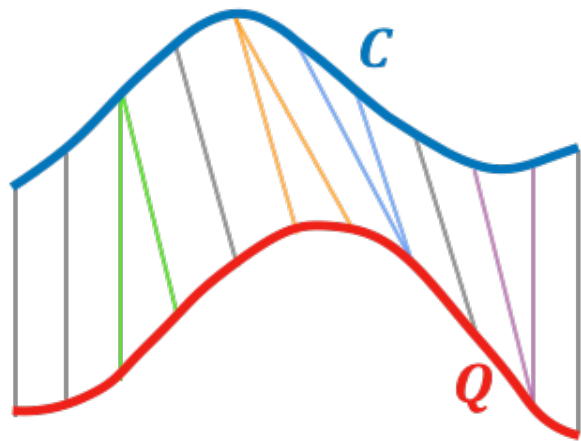
Расстояние DTW: точность vs. сложность

Сложность
 $O(m^2)$

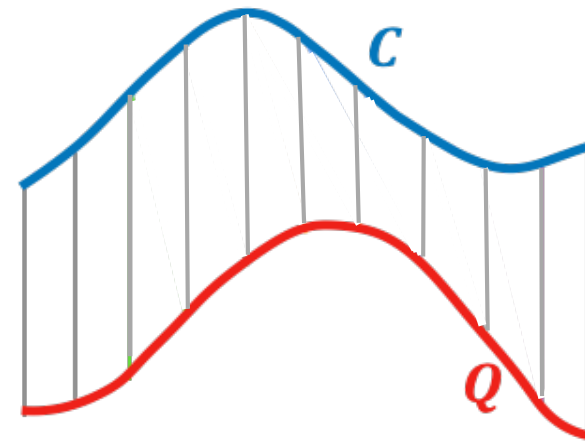
$$\text{DTW}(Q, C) = d(m, m)$$

$$d(i, j) = (q_i - c_j)^2 + \min \begin{cases} d(i-1, j) \\ d(i, j-1) \\ d(i-1, j-1) \end{cases}$$

$$d(0, 0) = 0, d(i, 0) = d(0, j) = +\infty; 1 \leq i, j \leq m$$



DTW лучше определяет схожесть по форме, чем Евклидова метрика



Ускорение подсчета DTW: полоса Сако-Чибя

Сложность
 $O(mr)$

$$\text{DTW}(Q, C) = d(m, m)$$

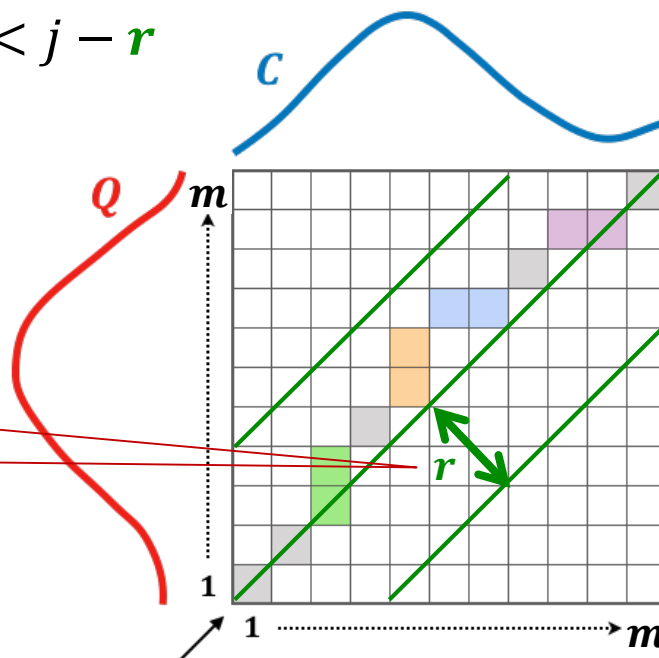
$$d(i, j) = (q_i - c_j)^2 + \min \begin{cases} d(i-1, j) \\ d(i, j-1) \\ d(i-1, j-1) \end{cases}$$

$$d(0, 0) = 0, d(i, 0) = d(0, j) = +\infty; 1 \leq i, j \leq m;$$

$$0 \leq r \leq m - 1, j - r \leq i \leq j + r$$

$$d(i, j) = +\infty, j + r < i < j - r$$

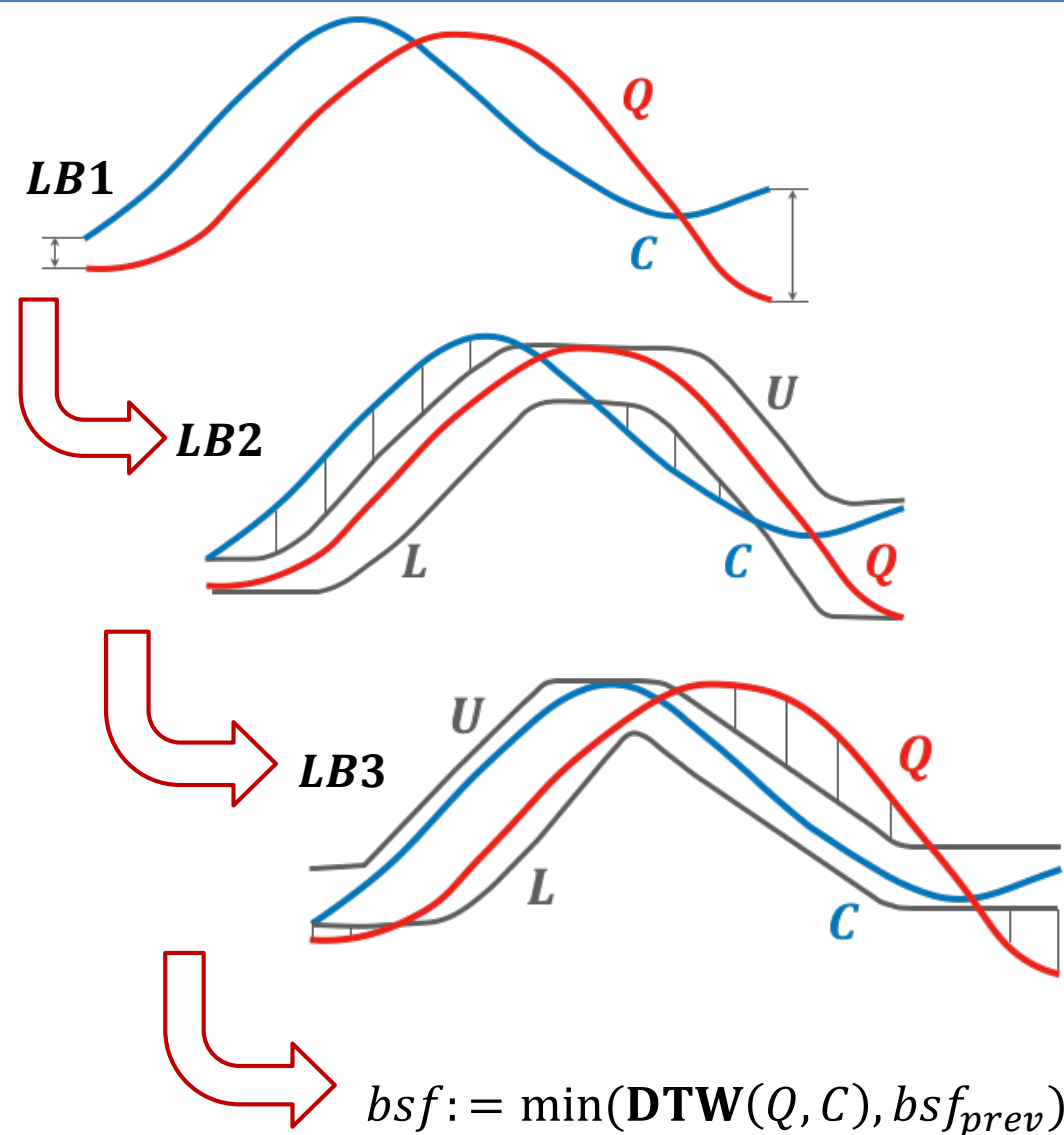
Сокращает сложность
за счет огрубления
схожести



Ускорение подсчета DTW: отбрасывание

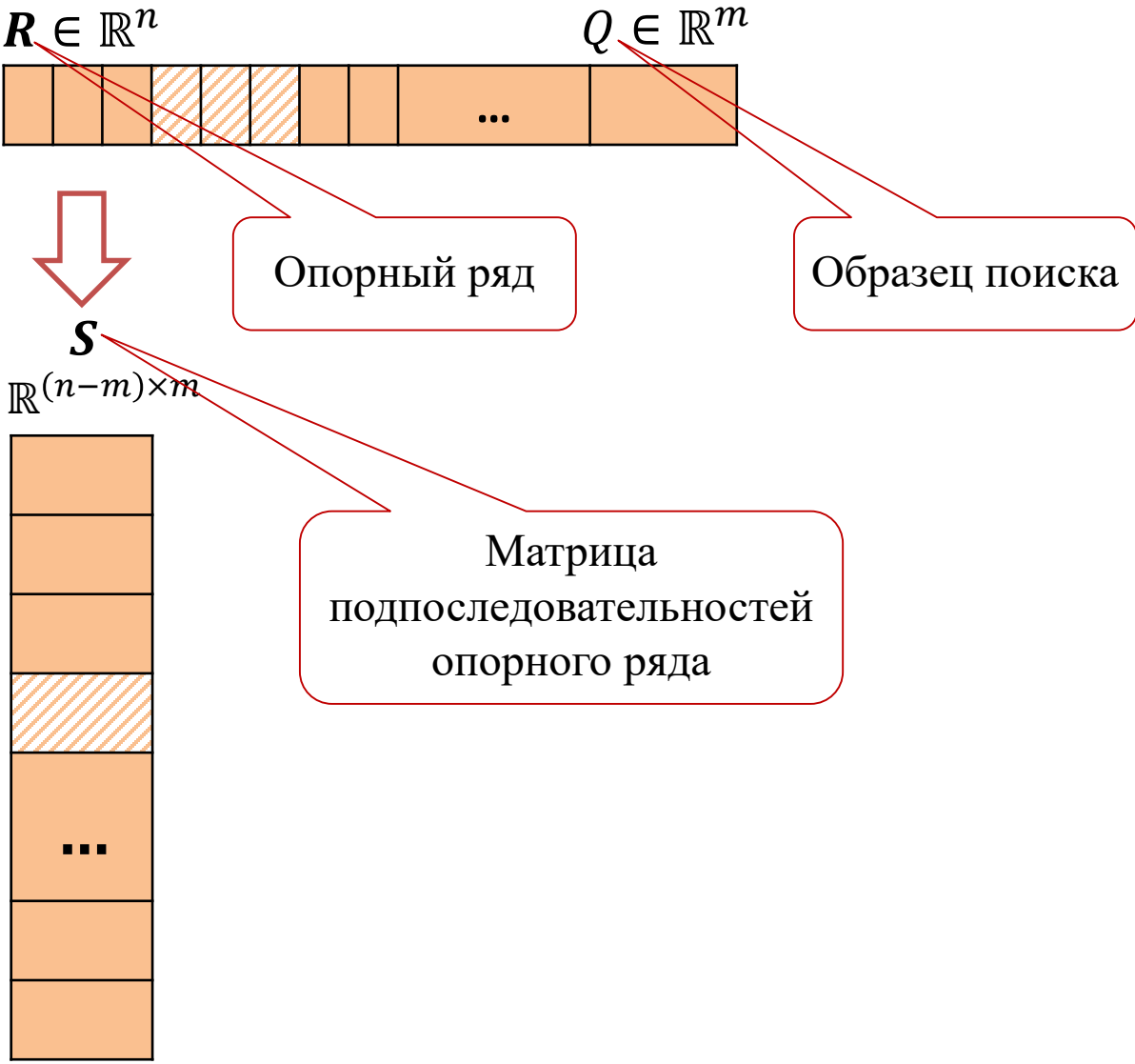
- Нижняя граница
 - функция $LB: \mathbb{R}^m \times \mathbb{R}^m \rightarrow \mathbb{R}_+$
с вычислительной сложностью менее $O(m^2)$
- Текущий минимум DTW
 - *bsf* (best-so-far)
- Отбрасывание соседей, далеких от образца, без вычисления DTW
 - если $LB(R[i:m], Q) > bsf$, то $DTW(R[i:m], Q) > bsf$
- Нормализация
 - Соседей и запрос нужно нормализовать

Отбрасывание: нижние границы



Нижняя граница	Сложность
$LB1(Q, C) = (q_1 - c_1)^2 + (q_m - c_m)^2$	$O(1)$
$LB2(Q, C) = \sum_{i=1}^m \begin{cases} (c_i - u_i)^2, & c_i > u_i \\ (c_i - \ell_i)^2, & c_i < \ell_i \\ 0, & \text{otherwise} \end{cases}$ $u_i = \max_{i-r \leq k \leq i+r} q_k$ $\ell_i = \min_{i-r \leq k \leq i+r} q_k$	$O(m)$
$LB3(Q, C) = LB2(C, Q)$	$O(m)$
$\mathbf{DTW}(Q, C)$	$O(mr)$

Параллельный поиск соседей по образцу

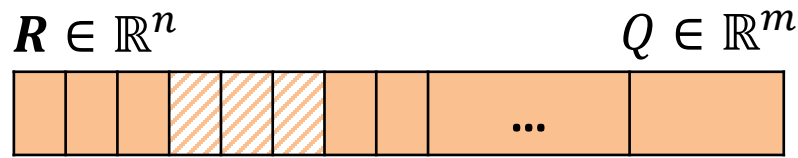


Опорный ряд

Образец поиска

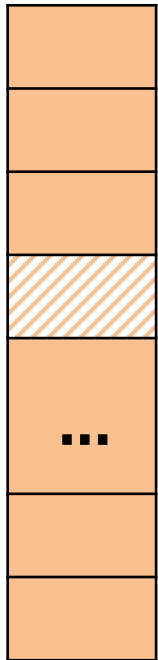
Матрица подпоследовательностей опорного ряда

Параллельный поиск соседей по образцу



S

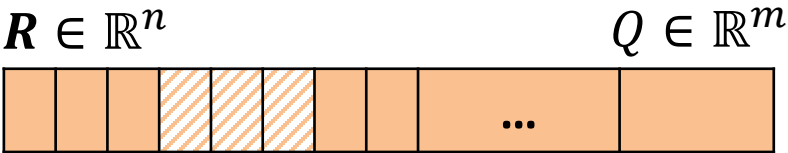
$\in \mathbb{R}^{(n-m) \times m}$



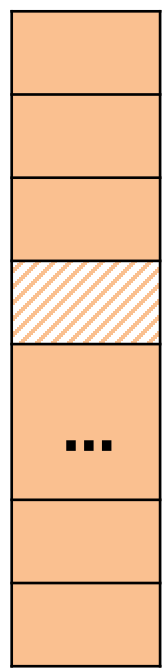
Нормализация

$$R[1:m] \Rightarrow \hat{R}[1:m], \hat{r}_i = \frac{r_i - \mu}{\sigma},$$
$$\mu = \frac{1}{m} \sum_{i=1}^m r_i, \sigma^2 = \frac{1}{m} \sum_{i=1}^m r_i^2 - \mu^2$$

Параллельный поиск соседей по образцу

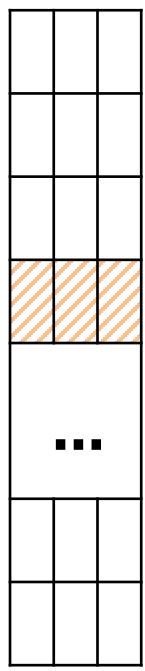


S
 $\in \mathbb{R}^{(n-m) \times m}$



LB
 $\in \mathbb{R}^{(n-m) \times 3}$

Матрица
нижних границ



Нормализация

Параллельный поиск соседей по образцу

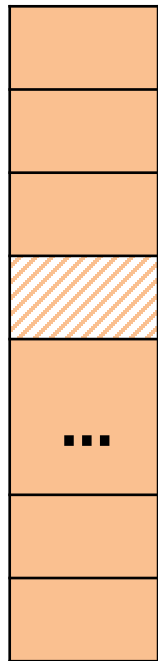
$R \in \mathbb{R}^n$

$Q \in \mathbb{R}^m$



S

$\in \mathbb{R}^{(n-m) \times m}$



LB

$\in \mathbb{R}^{(n-m) \times 3}$



$BitMap$

$\in \mathbb{B}^{n-m}$



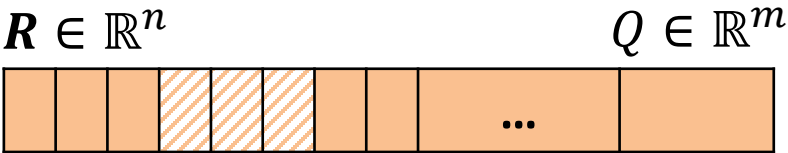
$bsf := DTW(Q, R[n-m:m])$

Битовая карта
подпоследовательностей

$$BitMap(i) := \bigwedge_{j=1}^3 LB(j) < bsf$$

Нормализация

Параллельный поиск соседей по образцу



$bsf := DTW(Q, R[n - m : m])$

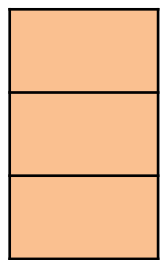
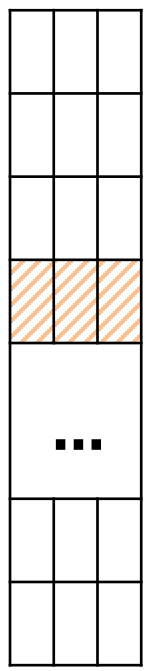
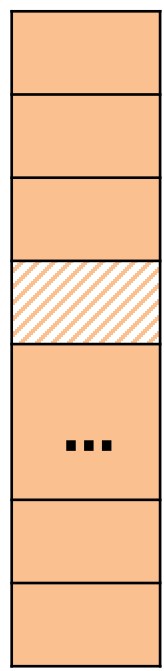


S
 $\in \mathbb{R}^{(n-m) \times m}$

LB
 $\in \mathbb{R}^{(n-m) \times 3}$

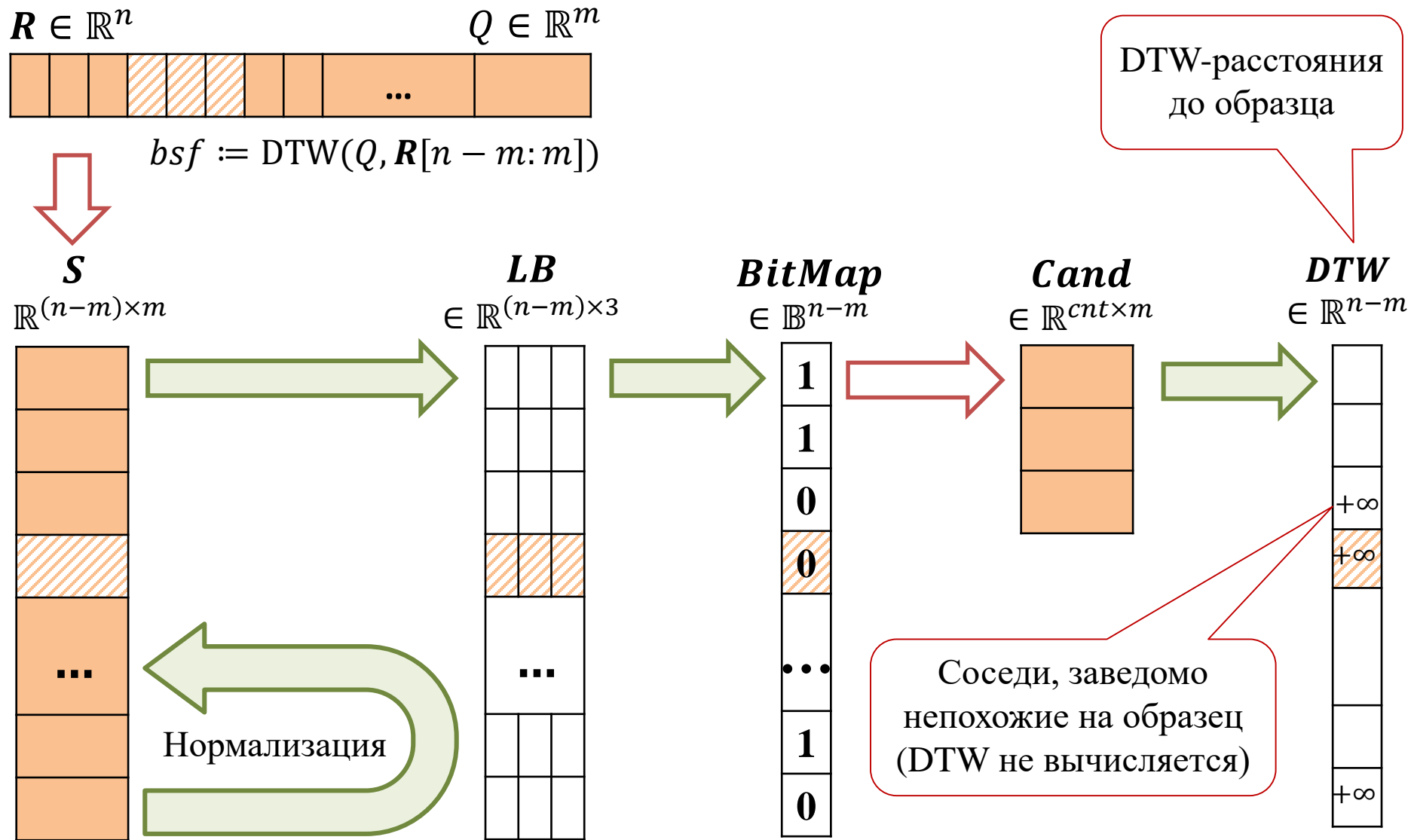
BitMap
 $\in \mathbb{B}^{n-m}$

Cand
 $\in \mathbb{R}^{cnt \times m}$



Матрица кандидатов в ближайшие соседи

Параллельный поиск соседей по образцу



Параллельный поиск соседей по образцу

$$R \in \mathbb{R}^n$$

$$Q \in \mathbb{R}^m$$

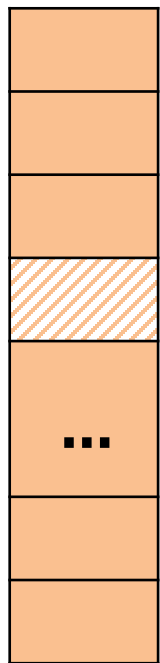


$$bsf := DTW(Q, R[n-m:m])$$

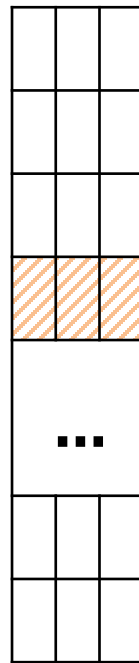
Уменьшение bsf для отбрасывания
большого числа соседей

$$bsf := \min(DTW)$$

$$S \in \mathbb{R}^{(n-m) \times m}$$



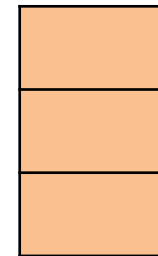
$$LB \in \mathbb{R}^{(n-m) \times 3}$$



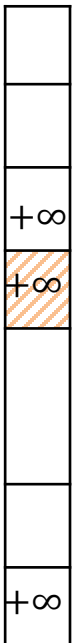
$$BitMap \in \mathbb{B}^{n-m}$$



$$Cand \in \mathbb{R}^{cnt \times m}$$

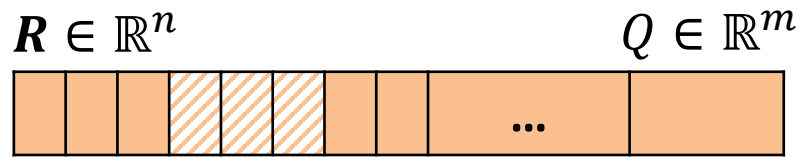


$$DTW \in \mathbb{R}^{n-m}$$



Нормализация

Параллельный поиск соседей по образцу



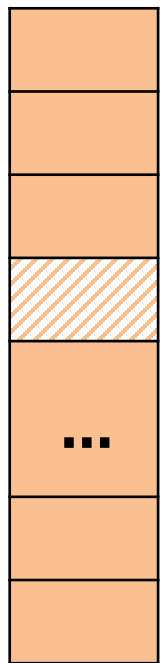
$$bsf := DTW(Q, R[n - m : m])$$

$\frac{1}{DTW(i) + \epsilon}$ используется
в скоринге интервалов



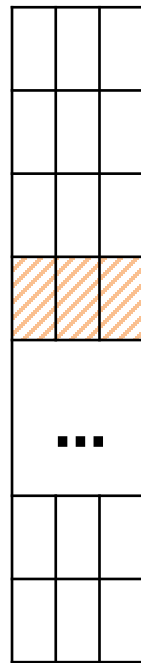
S

$$\in \mathbb{R}^{(n-m) \times m}$$



LB

$$\in \mathbb{R}^{(n-m) \times 3}$$



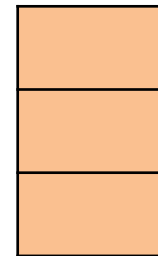
BitMap

$$\in \mathbb{B}^{n-m}$$



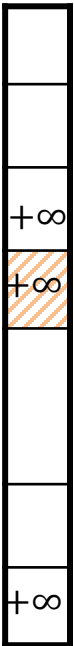
Cand

$$\in \mathbb{R}^{cnt \times m}$$



DTW

$$\in \mathbb{R}^{n-m}$$



Эксперименты: платформа и данные

- Аппаратная платформа:

Процессор	Intel Xeon Gold 6254
Ядра	18 ядер по 4 GHz
Память	192 Gb

- Точность восст-я:

$$RMSE = \sqrt{\frac{1}{h} \sum_{i=1}^h (x_i - \tilde{x}_i)^2}$$

$h = 1000$

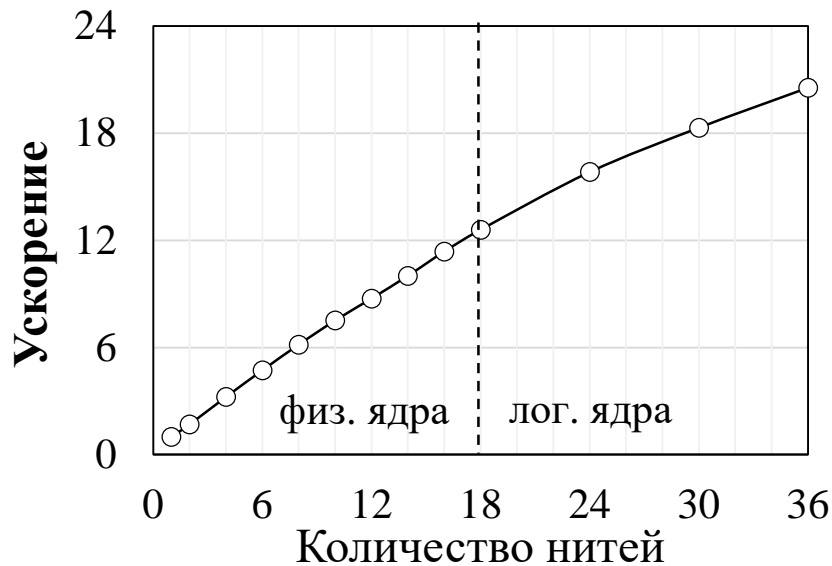
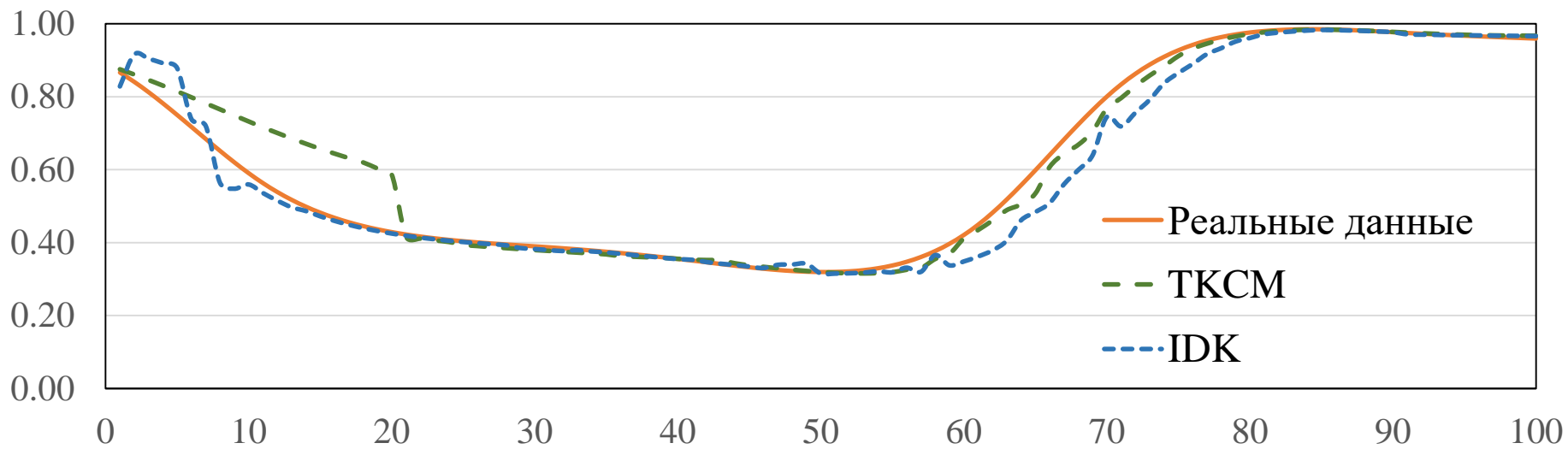
- Наборы данных:

Набор	Длина ряда, n	К-во опорных рядов, d	К-во ближ. соседей, k	Длина запроса, m	Полоса Сако–Чиба, r
Chlorine	4 209	3	3	100	25
Marel Carnot	50 000	3	3	100	25
Walk run	210 000	3	3	50	12

- Аналог: ТКСМ*

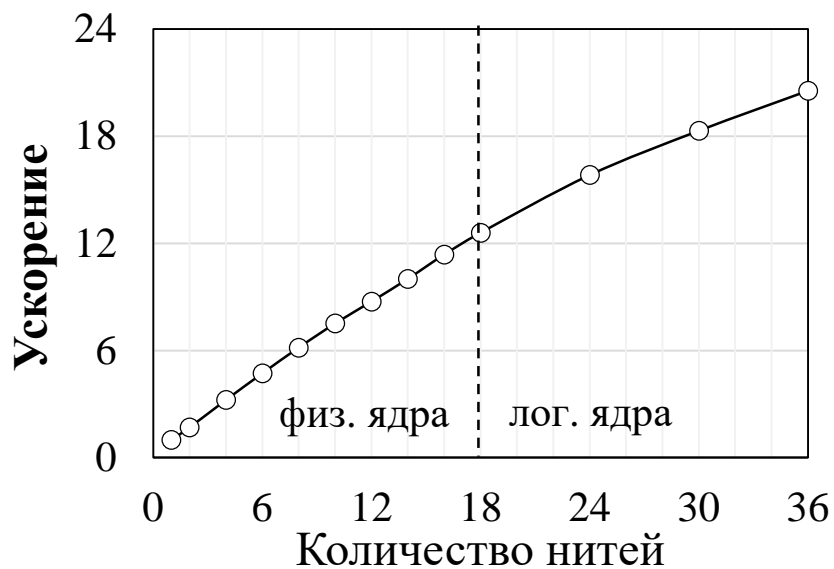
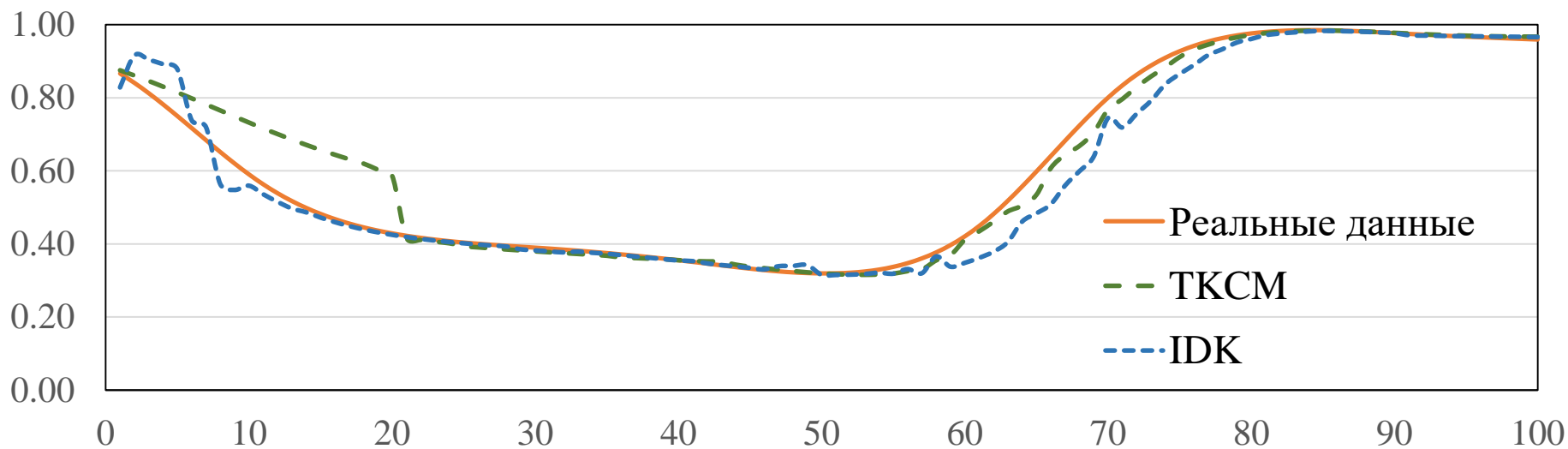
* Wellenzohn K. *et al.* Continuous Imputation of Missing Values in Streams of Pattern-Determining Time Series. Proc. of EDBT 2017. P. 330–341. DOI: 10.5441/002/edbt.2017.30

Эксперименты: набор Chlorine



Алгоритм	RMSE	Runtime, s
TKCM	0.062	0.013
IDK	0.049	0.006

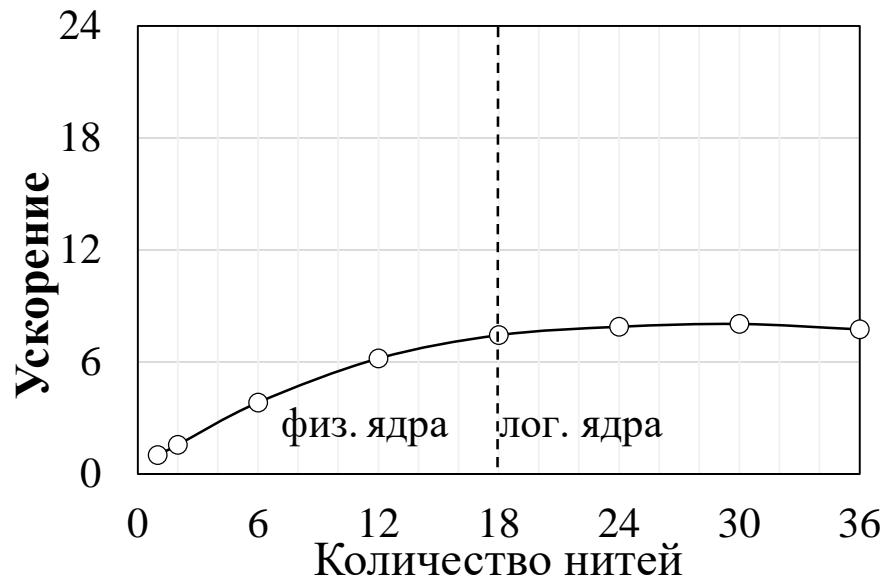
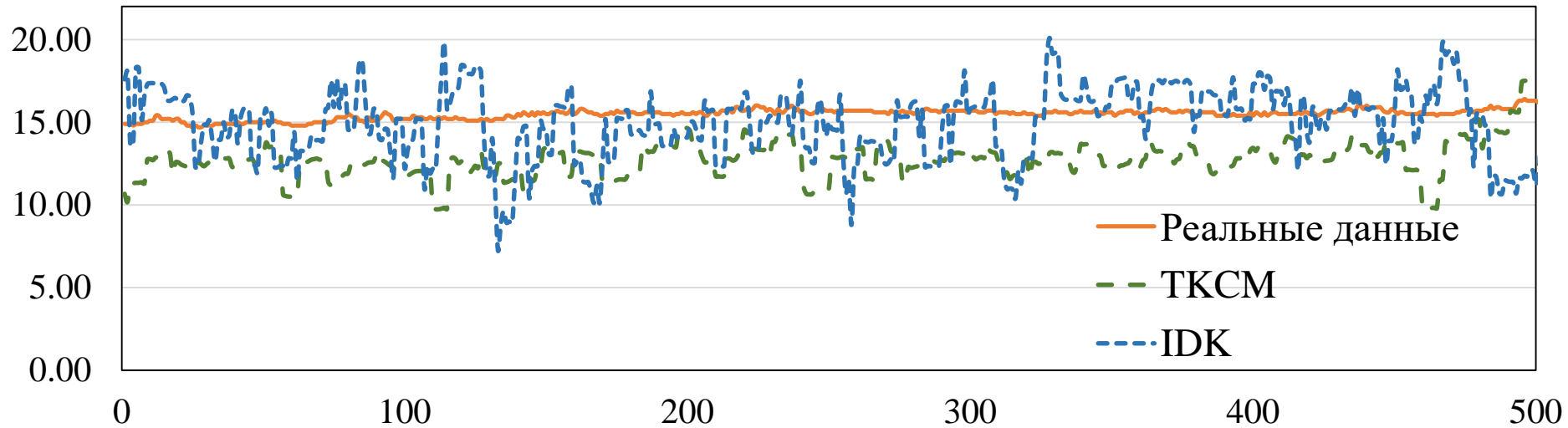
Эксперименты: набор Chlorine



Алгоритм	RMSE	Runtime, s
TKCM	0.062	0.013
IDK	0.049	0.006

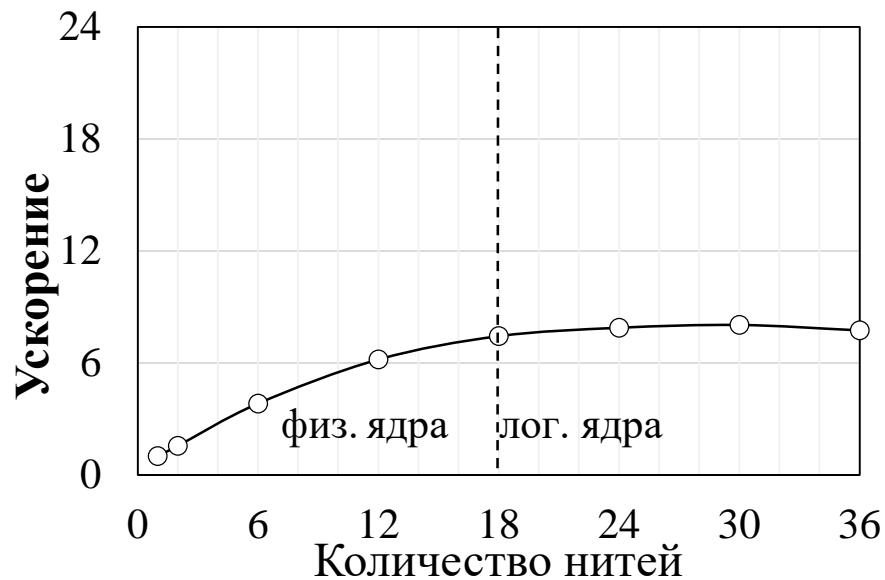
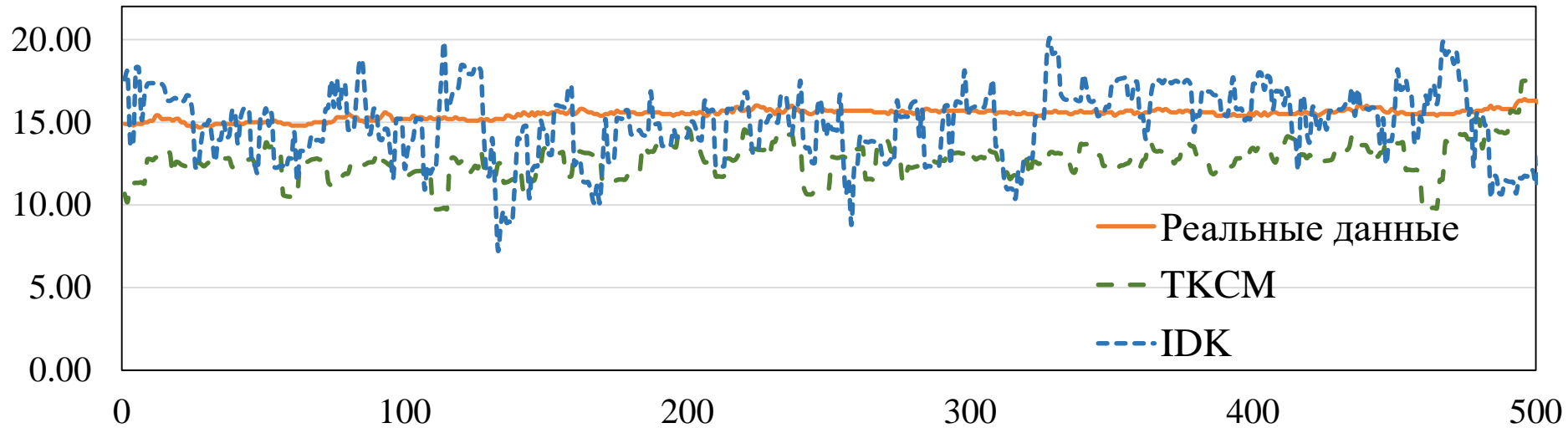
Отброшено 98% соседей
без вычисления DTW

Эксперименты: набор Marel Carnot



Алгоритм	RMSE	Runtime, s
TKCM	2.497	0.057
IDK	2.434	0.054

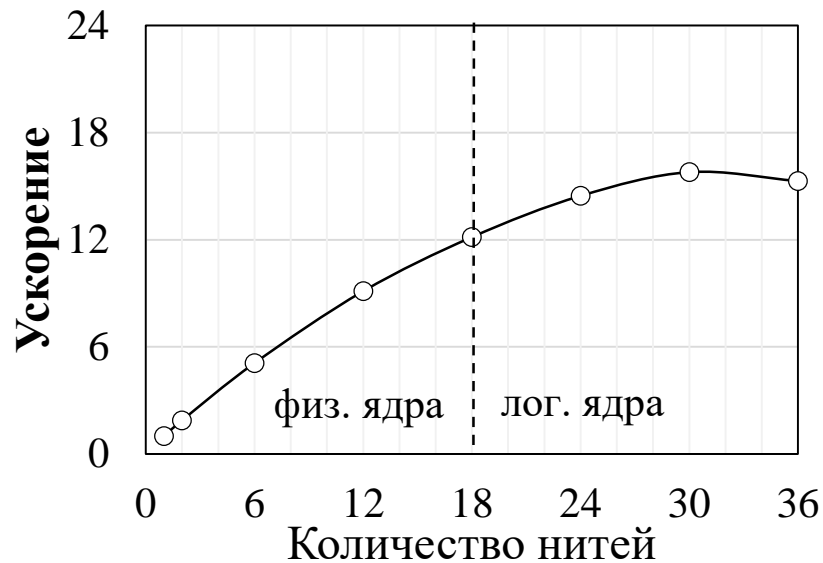
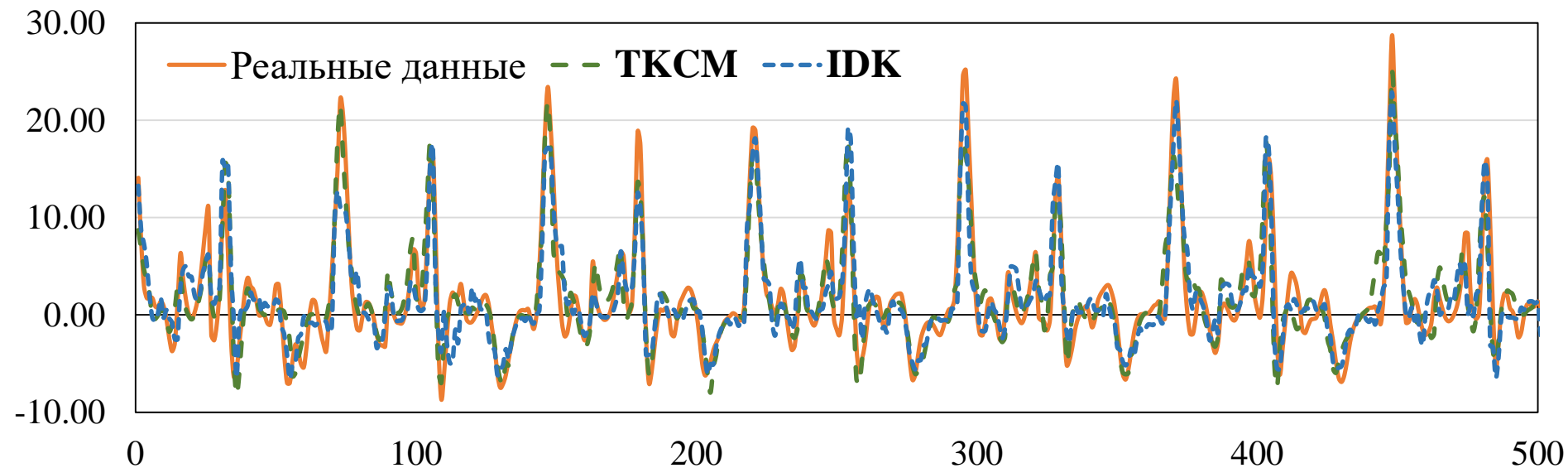
Эксперименты: набор Marel Carnot



Алгоритм	RMSE	Runtime, s
TKCM	2.497	0.057
IDK	2.434	0.054

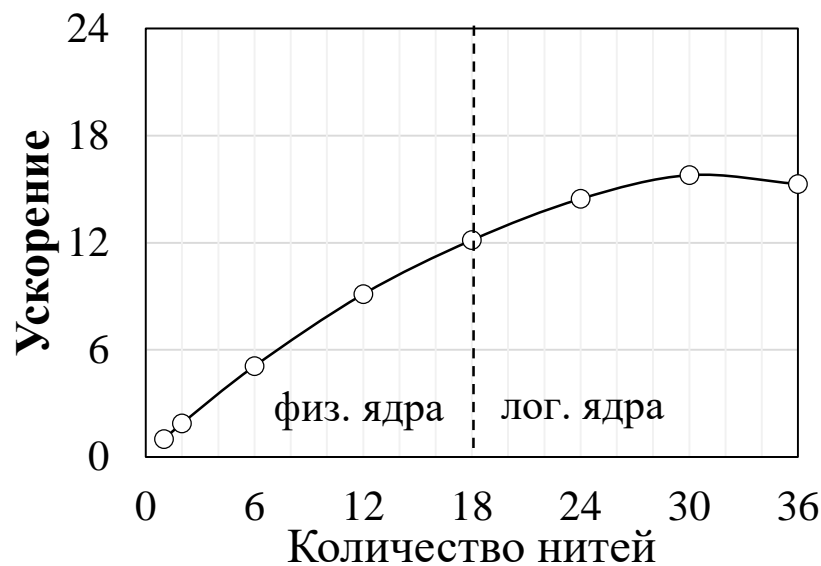
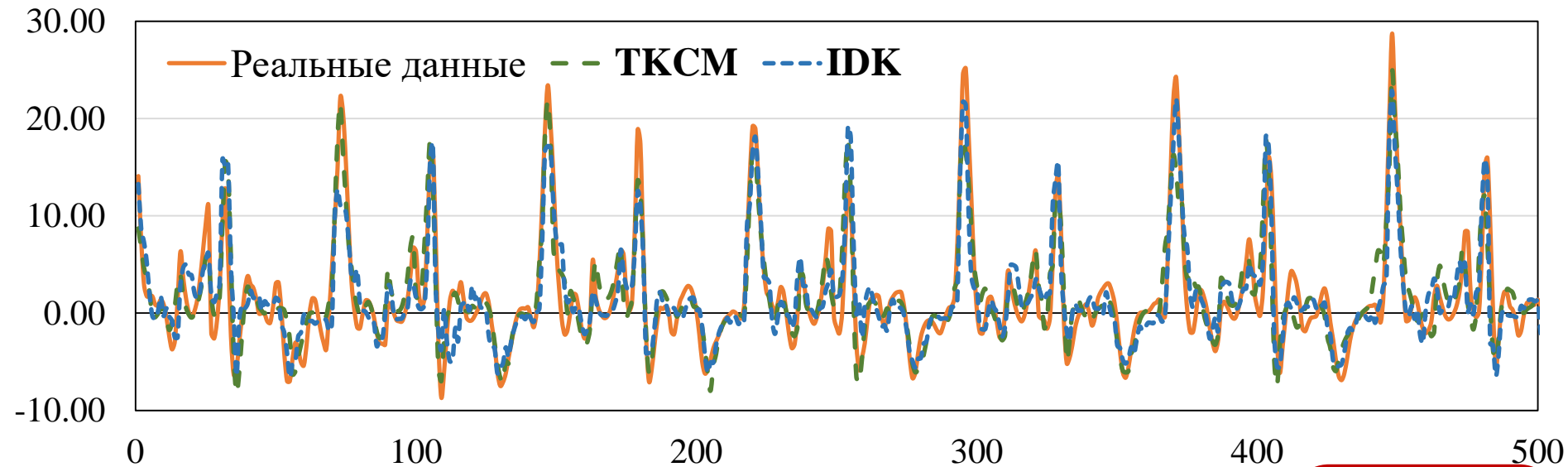
Отброшено 75% соседей
без вычисления DTW

Эксперименты: набор Walk run



Алгоритм	RMSE	Runtime, s
TKCM	2.785	0.242
IDK	2.937	0.080

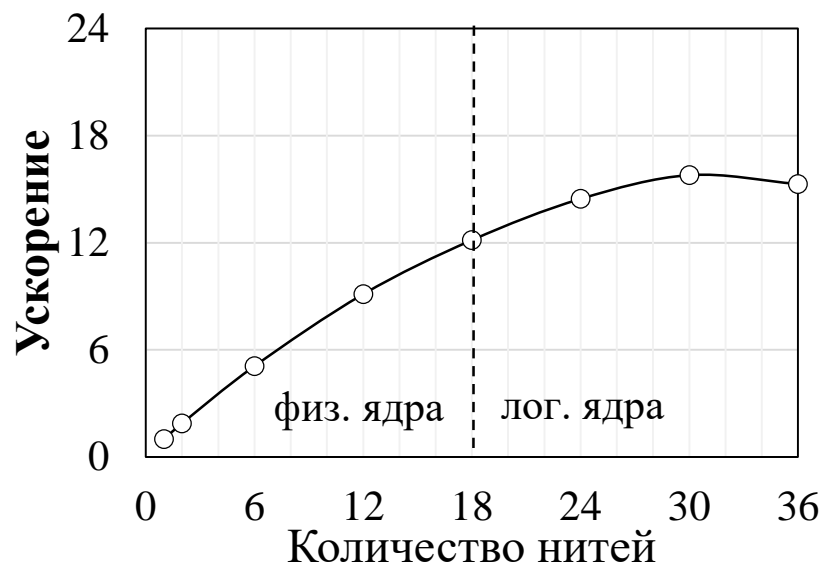
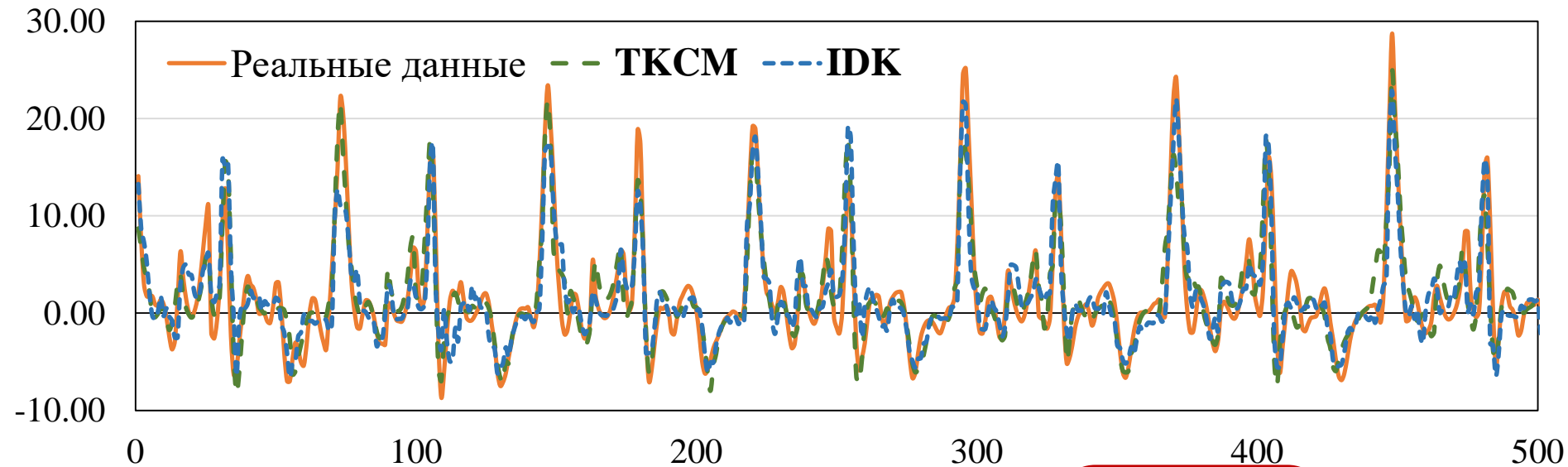
Эксперименты: набор Walk run



Алгоритм	RMSE	Runtime, s
TKCM	2.785	0.242
IDK	2.937	0.080

Отброшено 97% соседей
без вычисления DTW

Эксперименты: набор WalkRun



Алгоритм	RMSE	Runtime, s
TKCM	2.785	0.242
IDK	2.937	0.080

DTW лучше определяет схожесть по форме
Нормализация – ошибаемся по амплитуде

Заключение

- IDK – новый параллельный алгоритм восстановления пропущенных значений потокового временного ряда в режиме реального времени на многоядерном процессоре
- Эксперименты показывают хорошие масштабируемость и точность восстановления IDK
- Будущие исследования: версия IDK для GPU

Спасибо за внимание! Вопросы?

Андрей Николаевич Полуянов andrey.poluyanov@gmail.com

Михаил Леонидович Цымблер mzym@susu.ru