

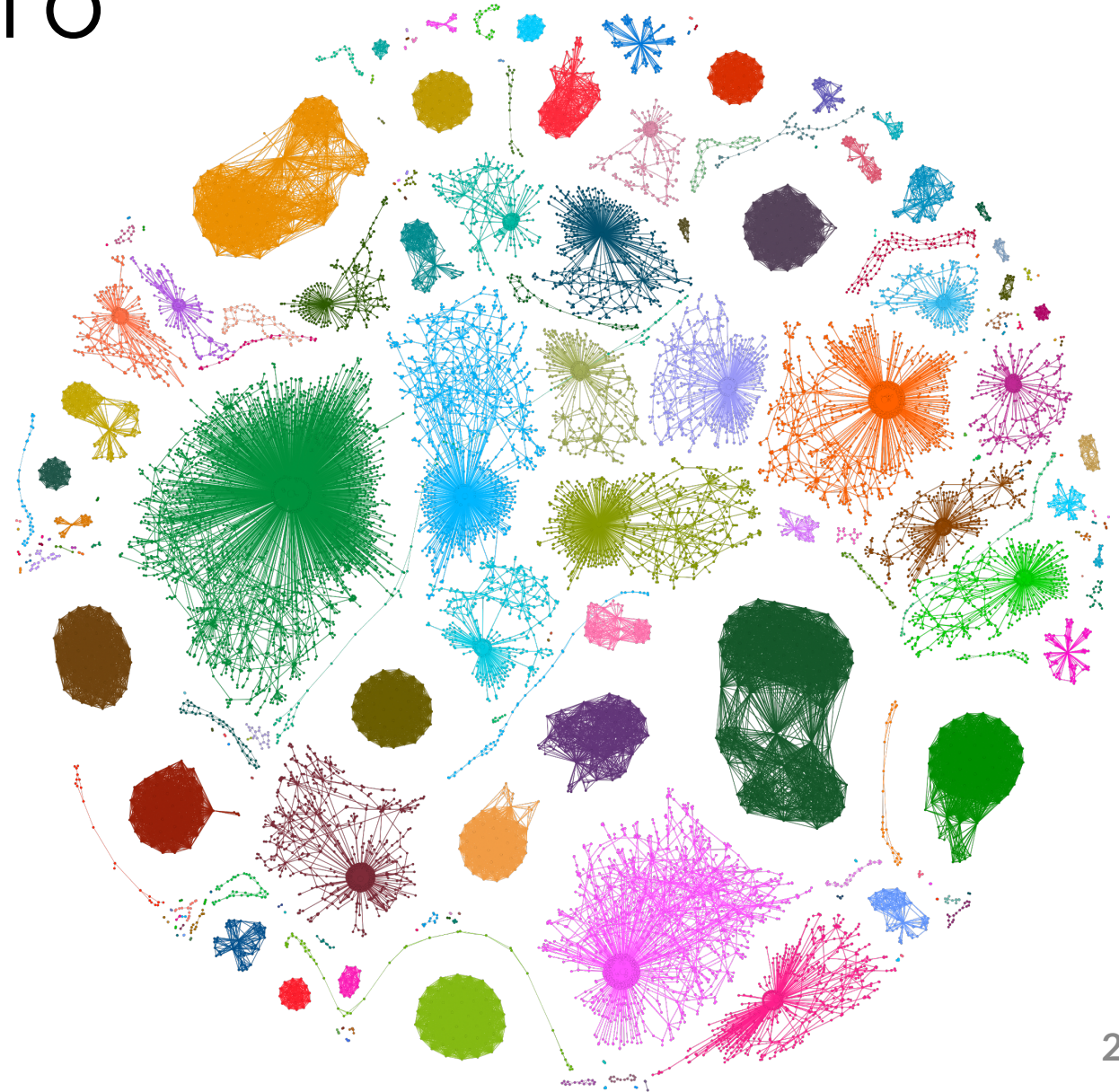
# Эмбеддинги графов без учителя

Антон Цицулин

# Графов в мире много

## Разные домены:

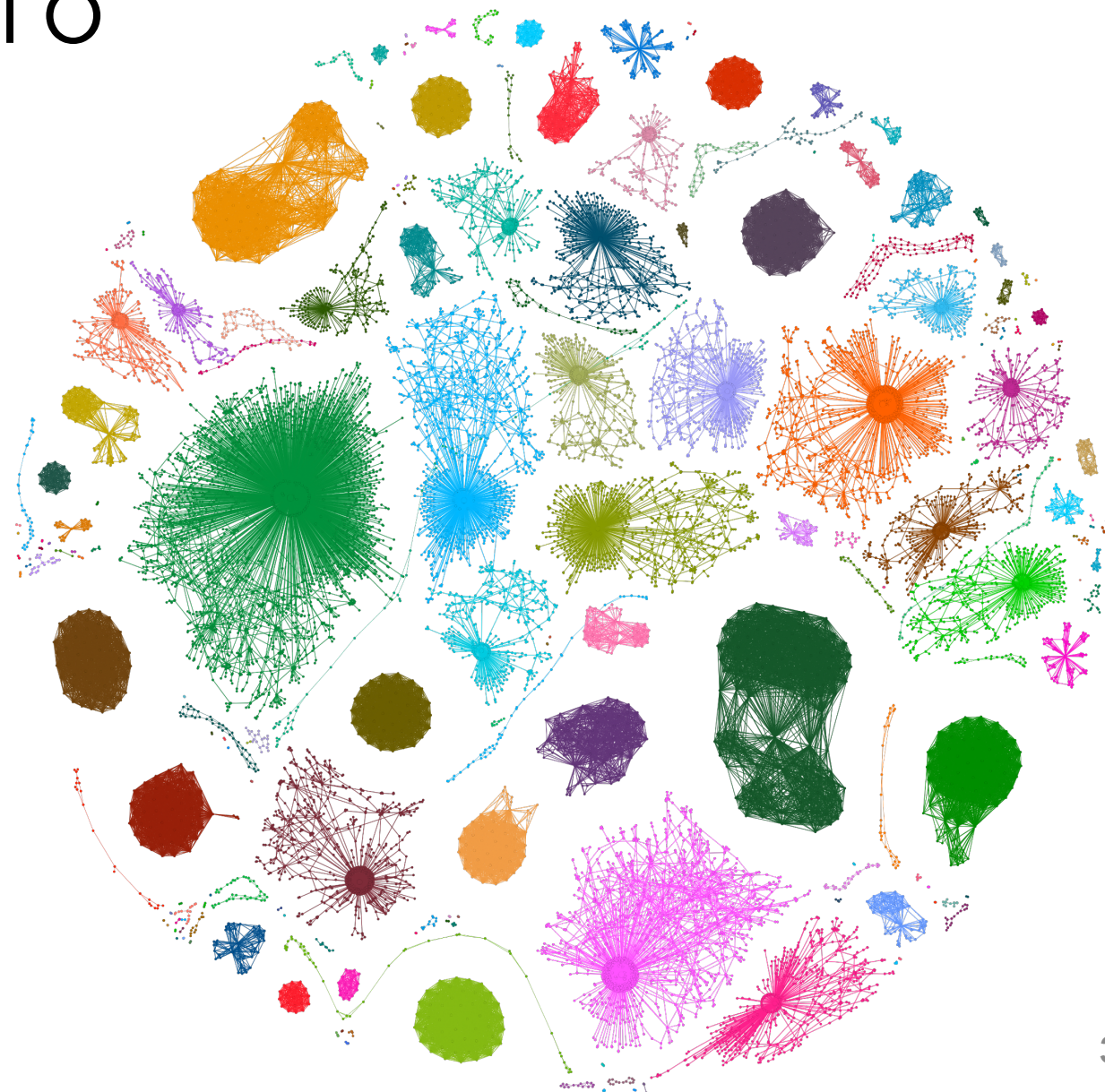
- Соцсети
- Биологические
- Транспортные сети
- ...



# Графов в мире много

## Разные типы графов:

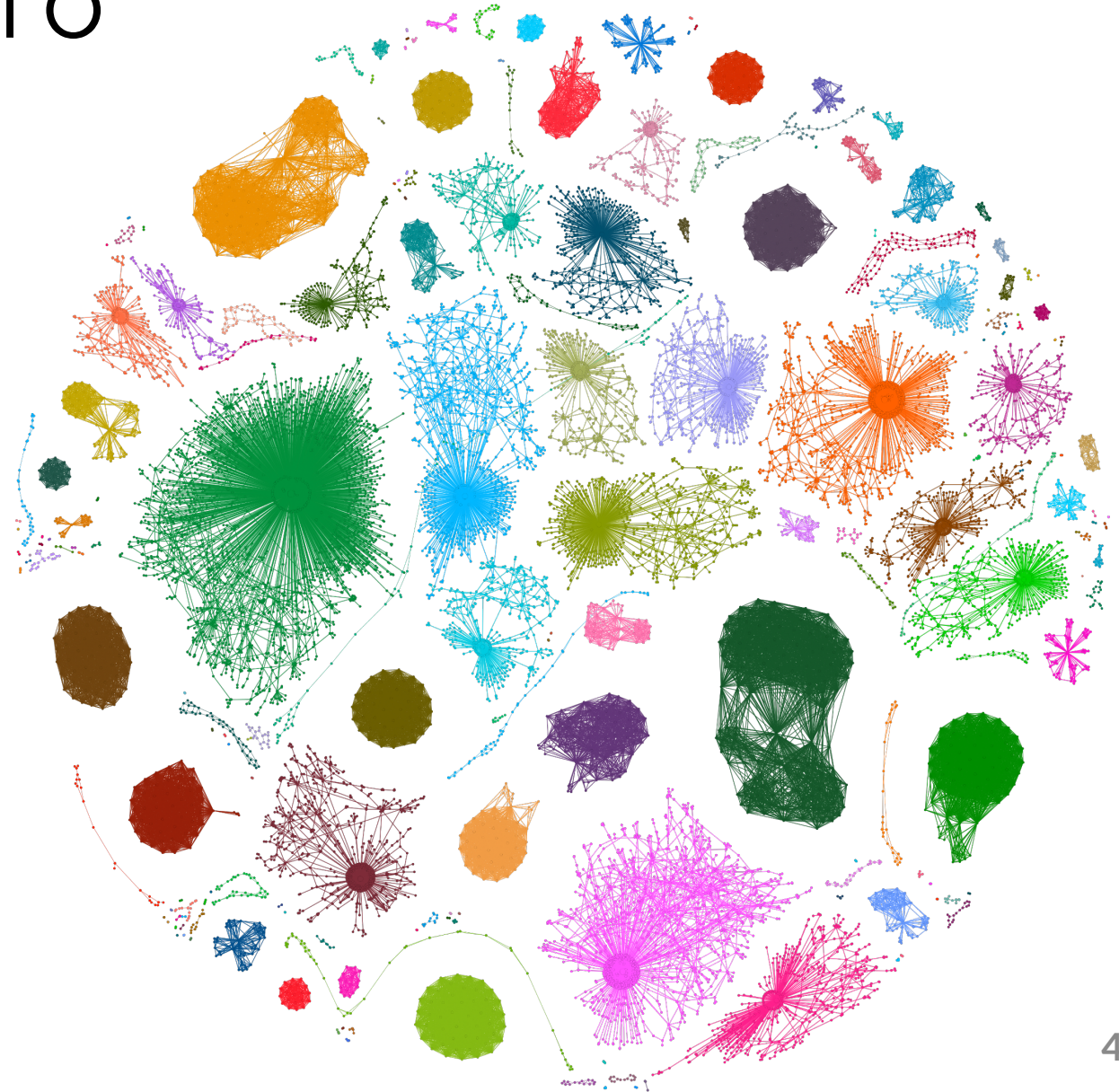
- (Не)направленные
- (Не)взвешенные
- С атрибутами
- Гетерогенные
- ...



# Графов в мире много

## Разные уровни анализа:

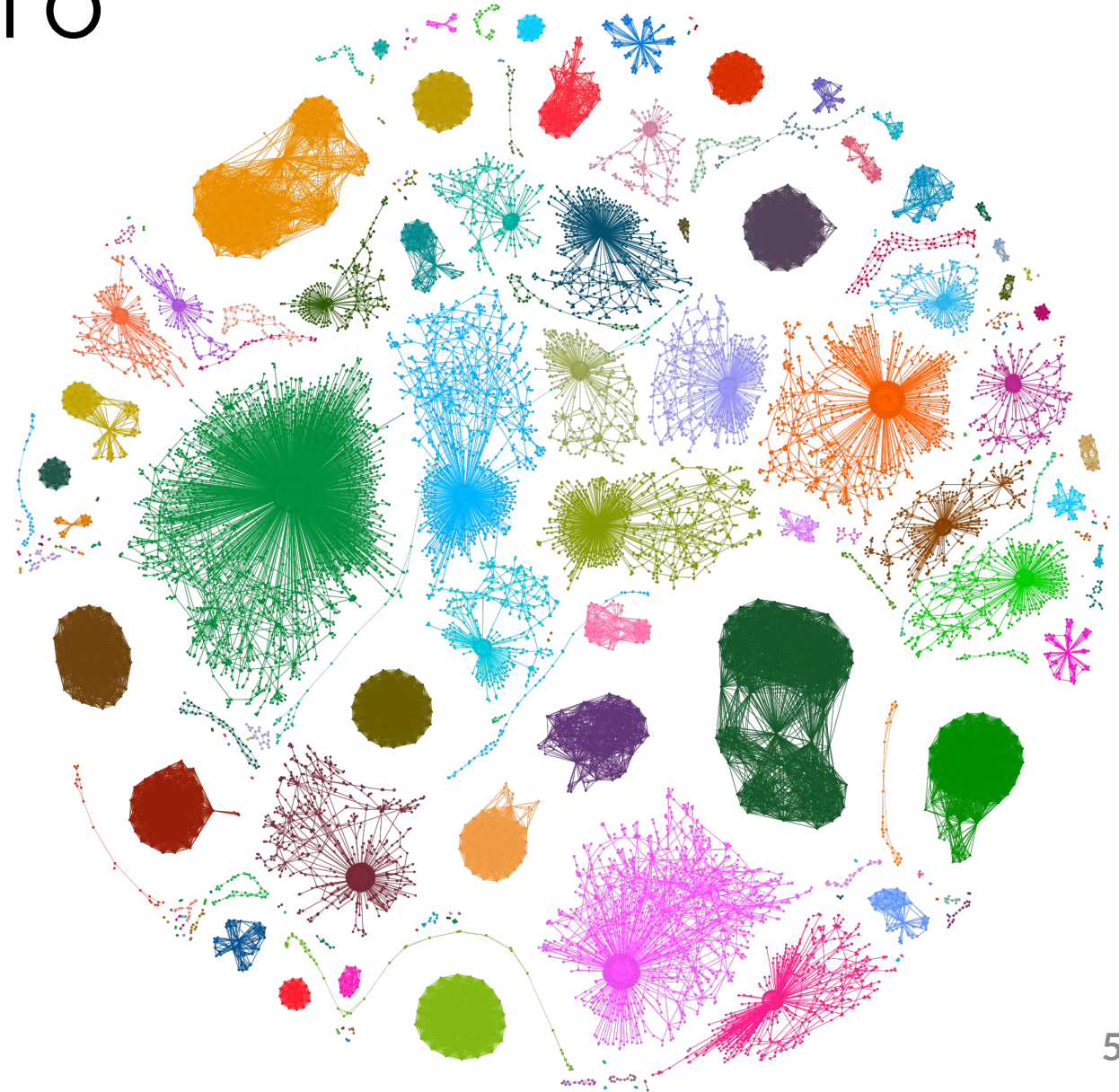
- Вершины
- Рёбра
- Подграфы
- Графы целиком
- ...



# Графов в мире много

## Разные задачи:

- Классификация
- Кластеризация
- Поиск аномалий
- ...



# Графов в мире много

## Домены

Соцсети  
Биологические  
Транспортные  
сети

## Типы

(Не)направленные  
(Не)взвешенные  
С атрибутами  
Гетерогенные

## Уровни

Вершины  
Рёбра  
Подграфы  
Графы целиком

## Задачи

Классификация  
Кластеризация  
Поиск  
аномалий

# Графов в мире много

## Домены

Соцсети  
Биологические  
Транспортные  
сети

## Типы

(Не)направленные  
(Не)взвешенные  
С атрибутами  
Гетерогенные

## Уровни

Вершины  
Рёбра  
Подграфы  
Графы целиком

## Задачи

Классификация  
Кластеризация  
Поиск  
аномалий



Эмбединги

# Графов в мире много

## Домены

Соцсети  
Биологические  
Транспортные  
сети

## Типы

(Не)направленные  
(Не)взвешенные  
С атрибутами  
Гетерогенные

## Уровни

Вершины  
Рёбра  
Подграфы  
Графы целиком

## Задачи

Классификация  
Кластеризация  
Поиск  
аномалий



Эмбединги



# Графов в мире много

## Домены

Соцсети  
Биологические  
Транспортные  
сети

## Типы

(Не)направленные  
(Не)взвешенные  
С атрибутами  
Гетерогенные

## Уровни

Вершины  
Рёбра  
Подграфы  
Графы целиком

## Задачи

Классификация  
Кластеризация  
Поиск  
аномалий



Эмбединги

# Графов в мире много

## Домены

Соцсети  
Биологические  
Транспортные  
сети

## Типы

(Не)направленные  
(Не)взвешенные  
С атрибутами  
Гетерогенные

## Уровни

Вершины  
Рёбра  
Подграфы  
Графы целиком

## Задачи

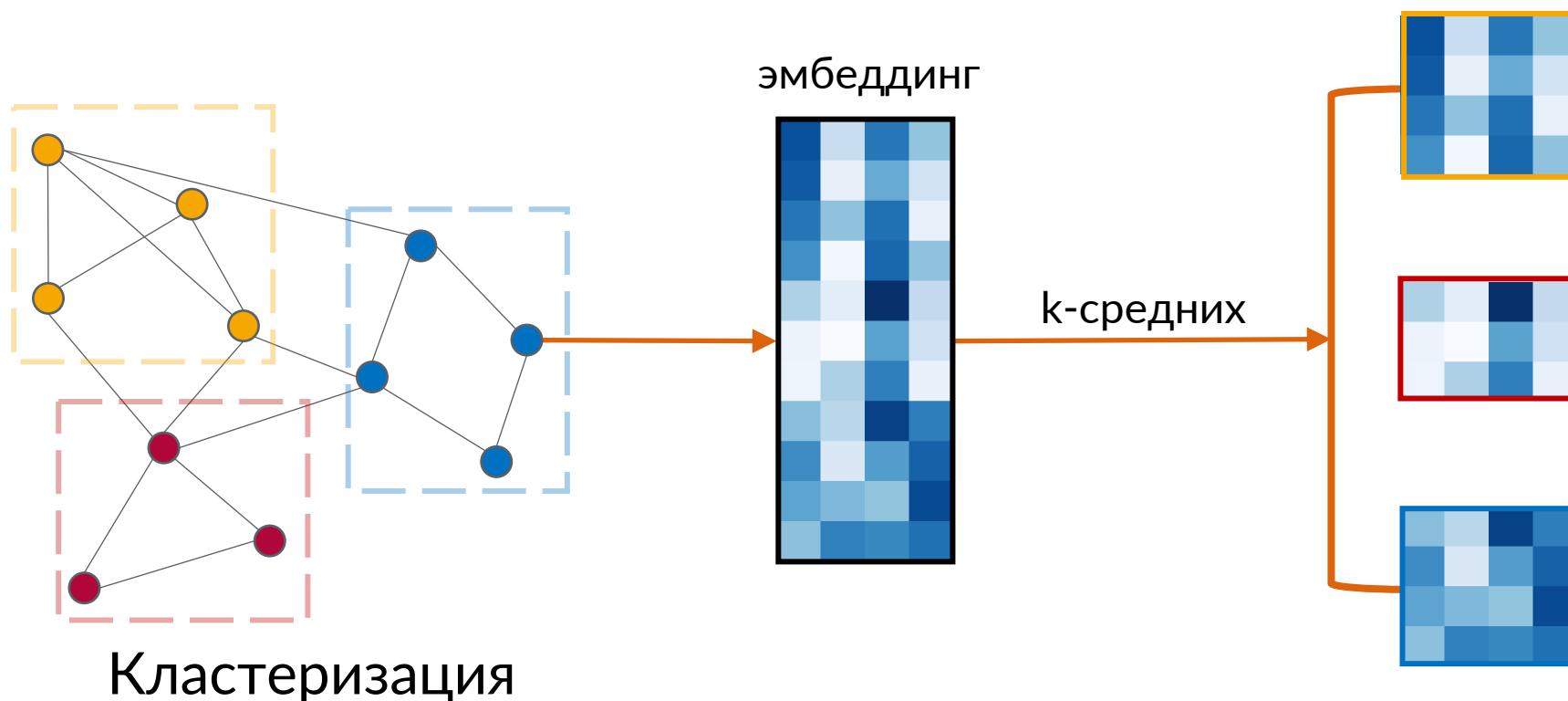
Классификация  
Кластеризация  
Поиск  
аномалий



Эмбединги

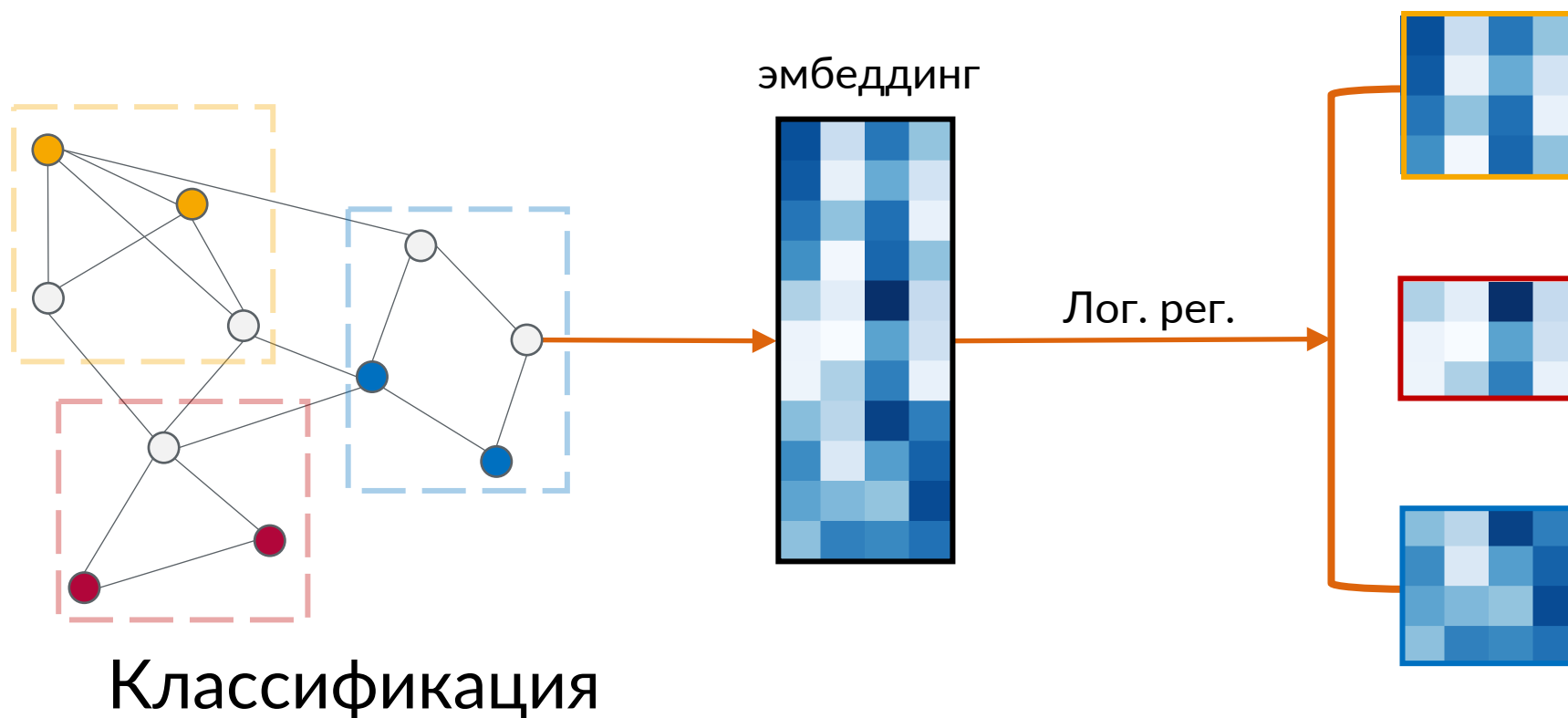
# Почему эмбединги?

У нас есть хорошие алгоритмы для анализа векторных данных ...



# Почему эмбединги?

У нас есть хорошие алгоритмы для анализа векторных данных ...

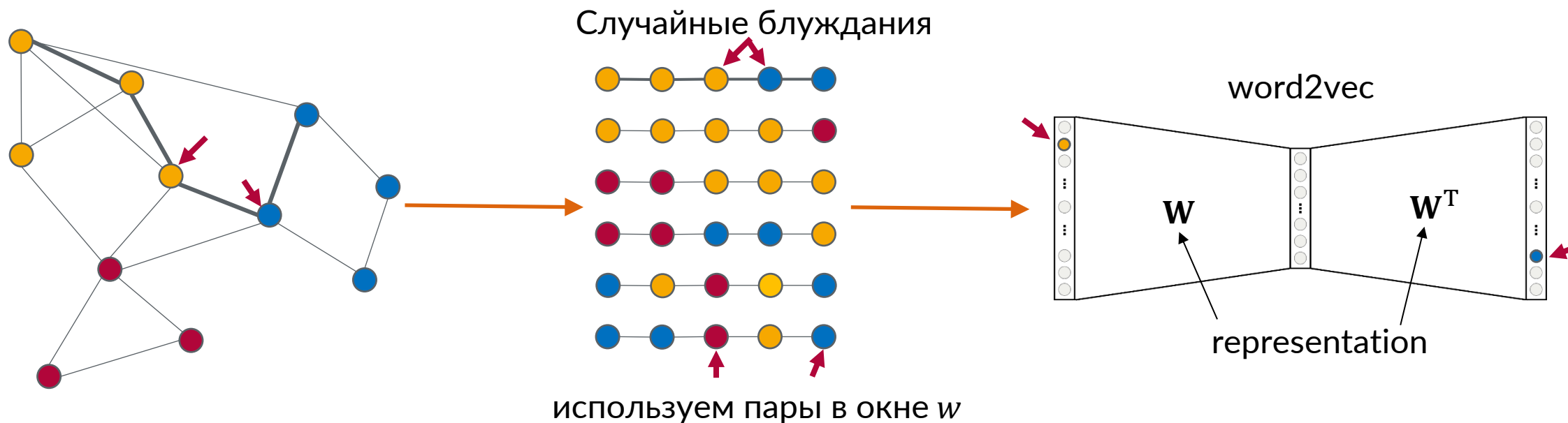


# Часть 1: графы без атрибутов

# DeepWalk

“Вершины в случ. блужданиях  $\approx$  словам в предложениях  $\rightarrow$  бахнем word2vec”

Начинаем  $\gamma$  блужданий длины  $t$  из каждой вершины



# DeepWalk: асимптотика и практика

На практике,  $\gamma = 80$ ,  $t = 80$ ,  $w = 10$ , получаем  $80 * 80 * n$  “текста”

**NB:** никогда не меняем  $w$

Если снижаете  $w$ , увеличивайте  $\gamma$  и  $t$

Как правильно тюнить параметры мы не знаем :(

# DeepWalk: асимптотика и практика

[Код на питоне](#) делает все блуждания и вызывает word2vec

Я написал версию на [C++](#) которая так не делает

Эпоха делается за  $O(d * \log n)$  :(

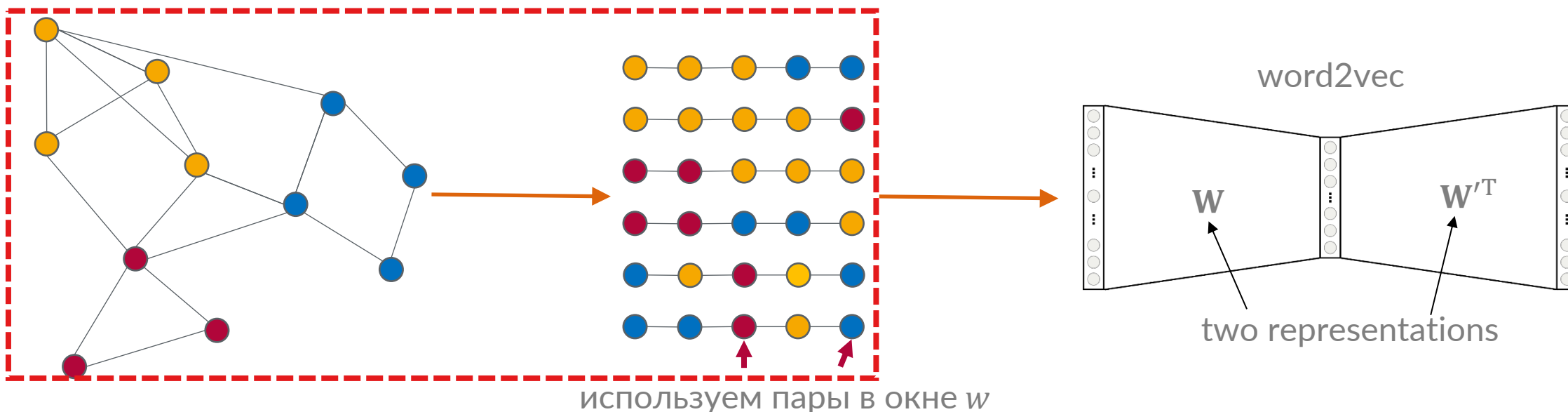
На практике: ~3М вершин



# Node2vec

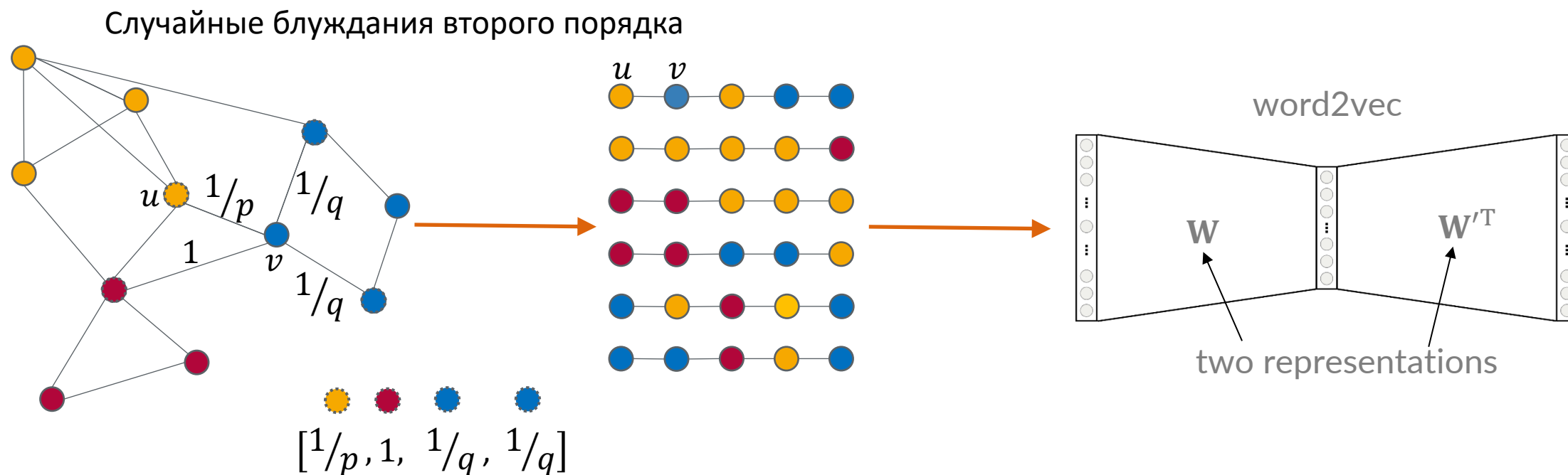
“Давайте-ка добавим еще два параметра ( $p, q$ ) в DeepWalk”

Начинаем  $\gamma$  блужданий длины  $t$  из каждой вершины



# Node2vec

“Давайте-ка добавим еще два параметра ( $p, q$ ) в DeepWalk”



# Node2vec: миф N°1

*Миф:* параметры  $(p, q)$  отвечают за BFS и DFS

*Реальность:* параметры  $(p, q)$  отвечают за треугольники  $\approx$  кластеры

Низкий  $q \rightarrow$  блуждатель ходит между кластерами

Высокий  $q \rightarrow$  блуждатель ходит внутри кластеров

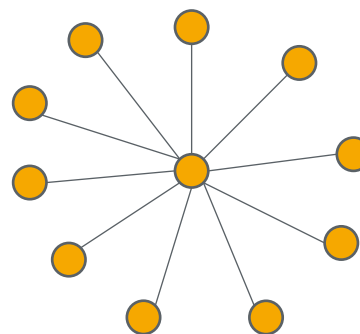
# Node2vec: миф N°2

*Миф:* node2vec быстро работает

*Реальность:* случайные блуждания второго порядка -  $O(n^2)$

Худший случай для графов-звёзд

(или любых графов с вершинами с высокой степенью)



# Node2vec: асимптотика и практика

**NB:** сравнения в статье обманчивы ( $\gamma = 10$  для всех)

В статье,  $\gamma = 10$ ,  $t = 80$ ,  $w = 10$ , перебор  $(p, q)$

$\gamma = 10$  даёт плохие результаты, используйте  $\gamma = 80$

Подбор  $(p, q)$  не даёт прироста на большинстве графов

# Node2vec: асимптотика и практика

[Код на питоне](#) делает все блуждания и вызывает word2vec

Я написал версию на [C++](#) которая так не делает

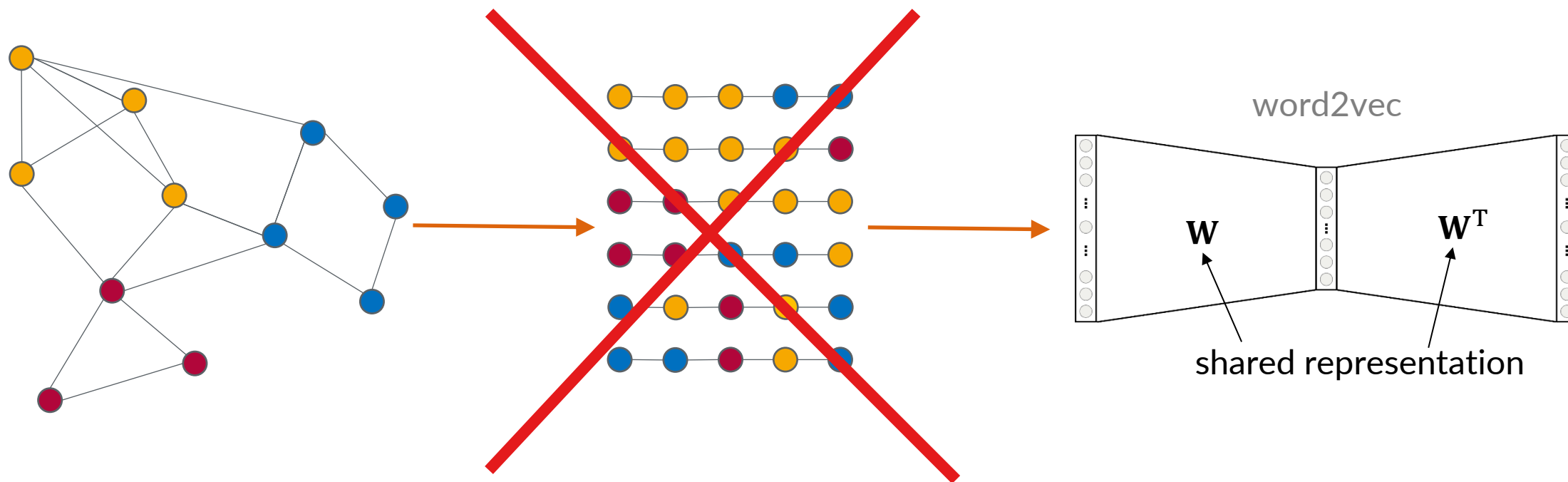
Препроцессинг до  $O(n^2)$  :(

Эпоха делается за  $O(d)$  :)

На практике: ~500т вершин если повезло с графом, если нет, ~50т

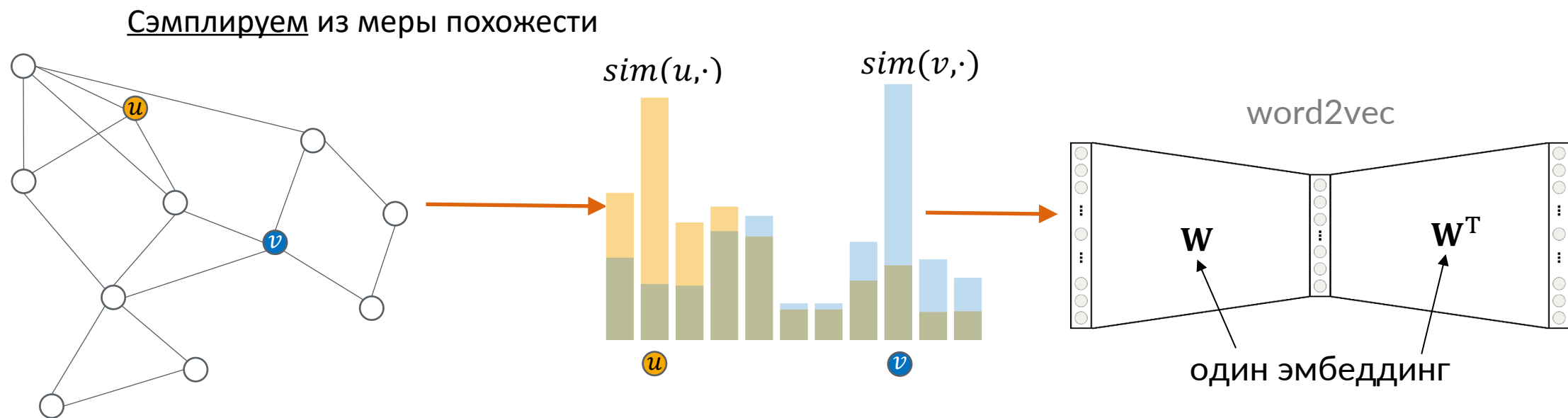
# VERSE

“Случайные блуждания – мера похожести вершин”



# VERSE

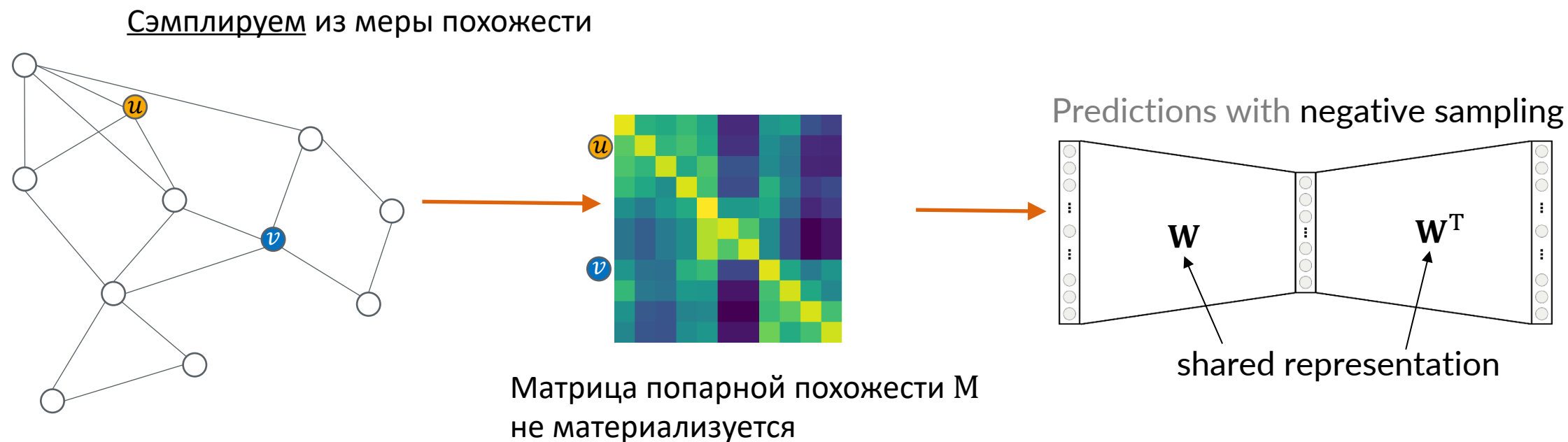
“Случайные блуждания – мера похожести вершин”





# VERSE

“Random walks define a similarity distribution”



# VERSE: интерпретация DeepWalk

Блуждания в DeepWalk  $\sim$  Personalized PageRank

Параметр PPR  $\alpha = \frac{w-2}{w+1}$  для  $w$  в DeepWalk

Мы можем мерить качество эмбеддингов :)

1 параметр вместо 3 или 5

# VERSE: асимптотика и практика

Простой и быстрый алгоритм, хорош на предсказании рёбер :)

Эпоха делается за  $O(d)$  :)

Код на [C++](#) работает хорошо

На практике: ~10М вершин

# Эмбединги для рёбер

Operator	Result
Average	$(\mathbf{a} + \mathbf{b})/2$
Concat	$[\mathbf{a}_1, \dots, \mathbf{a}_d, \mathbf{b}_1, \dots, \mathbf{b}_d]$
Hadamard	$[\mathbf{a}_1 * \mathbf{b}_1, \dots, \mathbf{a}_d * \mathbf{b}_d]$
Weighted L1	$[ \mathbf{a}_1 - \mathbf{b}_1 , \dots,  \mathbf{a}_d - \mathbf{b}_d ]$
Weighted L2	$[(\mathbf{a}_1 - \mathbf{b}_1)^2, \dots, (\mathbf{a}_d - \mathbf{b}_d)^2]$

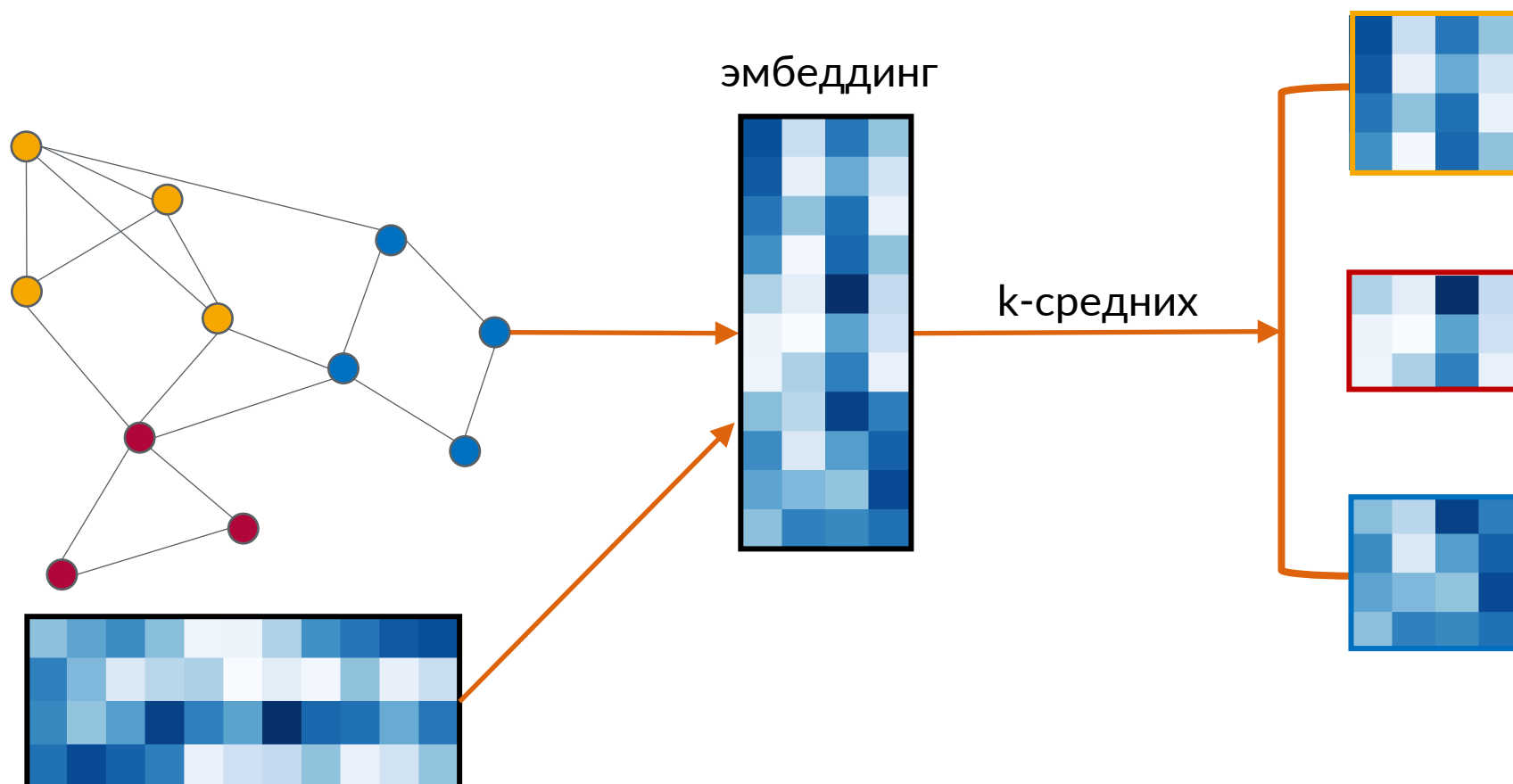
Table 2.3: Vector operators used for link-prediction task for each  $u, v \in V$  and corresponding embeddings  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^d$ .

**NB:** оператор выбираем в зависимости от алгоритма

# Часть 1: графы с атрибутами

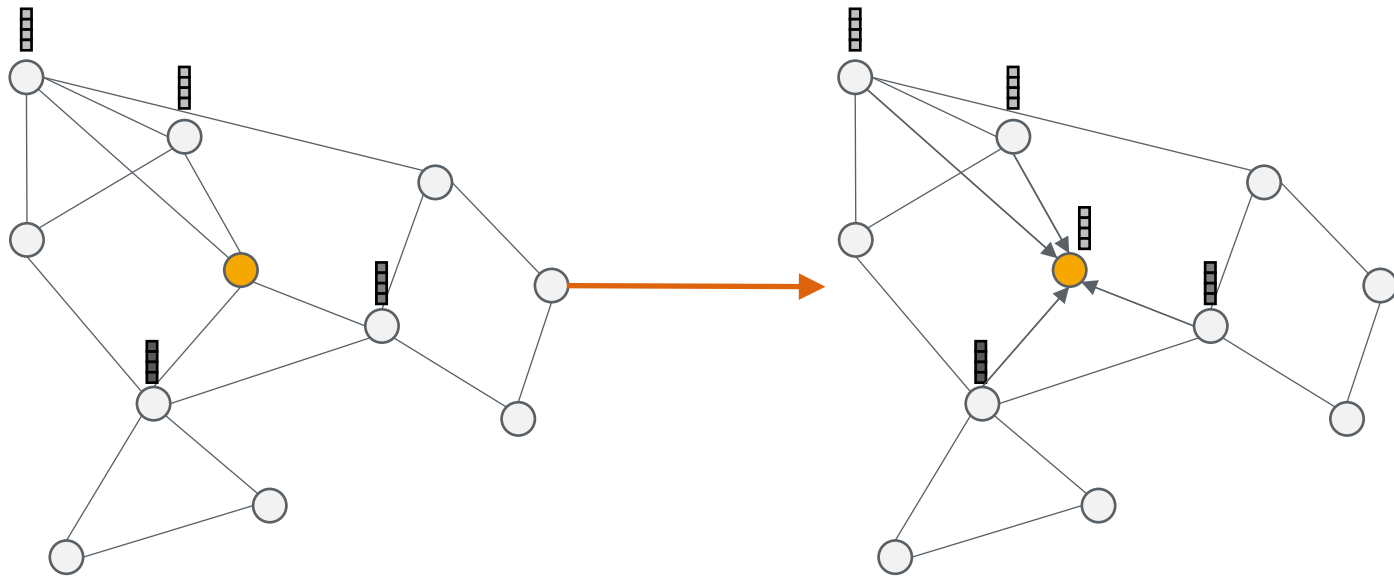
# Графы с атрибутами

На каждой вершине есть некоторый вектор с фичами



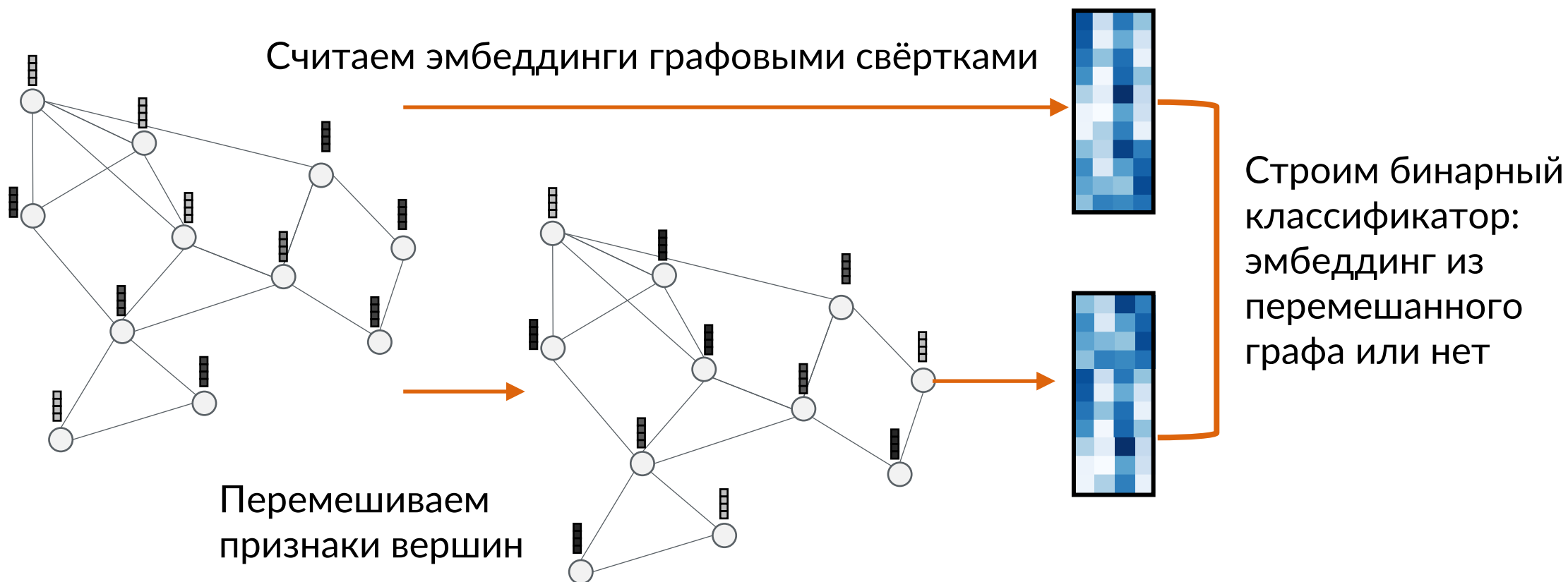
# Графовые свёртки

Скрытое представление вершины = среднее представление соседей



# Deep Graph Infomax

“Учимся отличать граф от него же с перемешанными признаками”





# DGI: асимптотика и практика

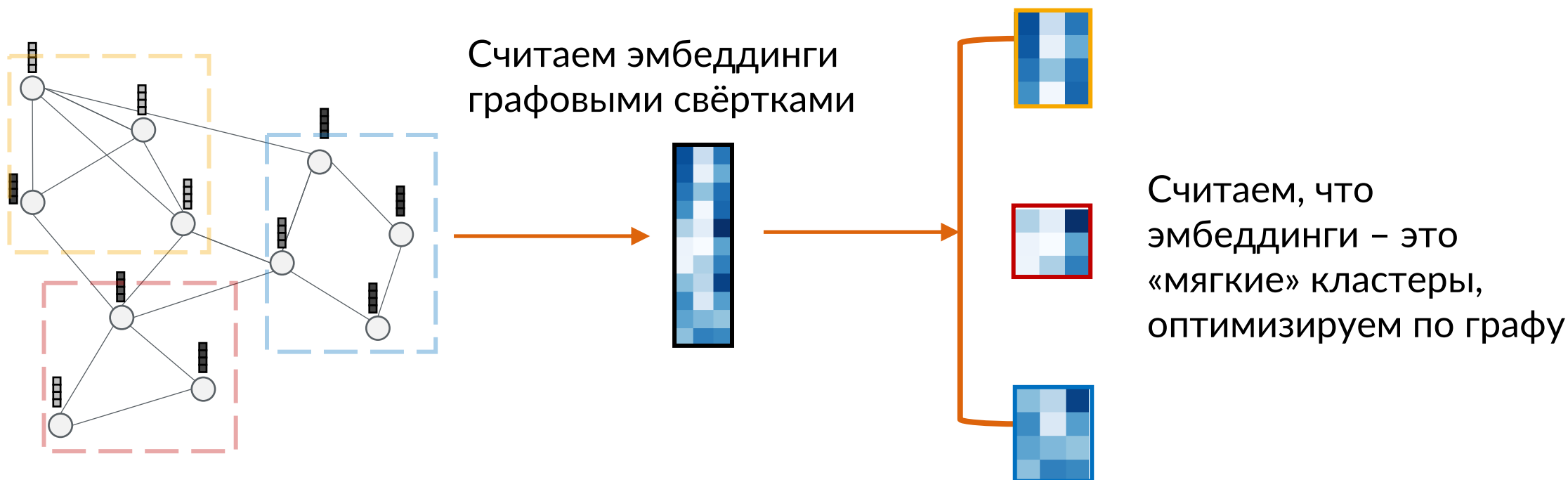
Можно делать для больших графов :)

«Размазывает» информацию по графу

Нетривиальное обучение – нужен early stopping :(

# Deep Modularity Networks

“Учимся кластеризовывать граф”



# DMoN: асимптотика и практика

Можно делать для средних графов (1М) :/

Группирует вершины графа по графу

Простое обучение

# Остались вопросы?

Twitter

Сайт

Пишите

[twitter.com/tsitsulin\\_](https://twitter.com/tsitsulin_)

[tsitsul.in/talks/datastart](http://tsitsul.in/talks/datastart)

[anton@tsitsul.in](mailto:anton@tsitsul.in)

← презентация