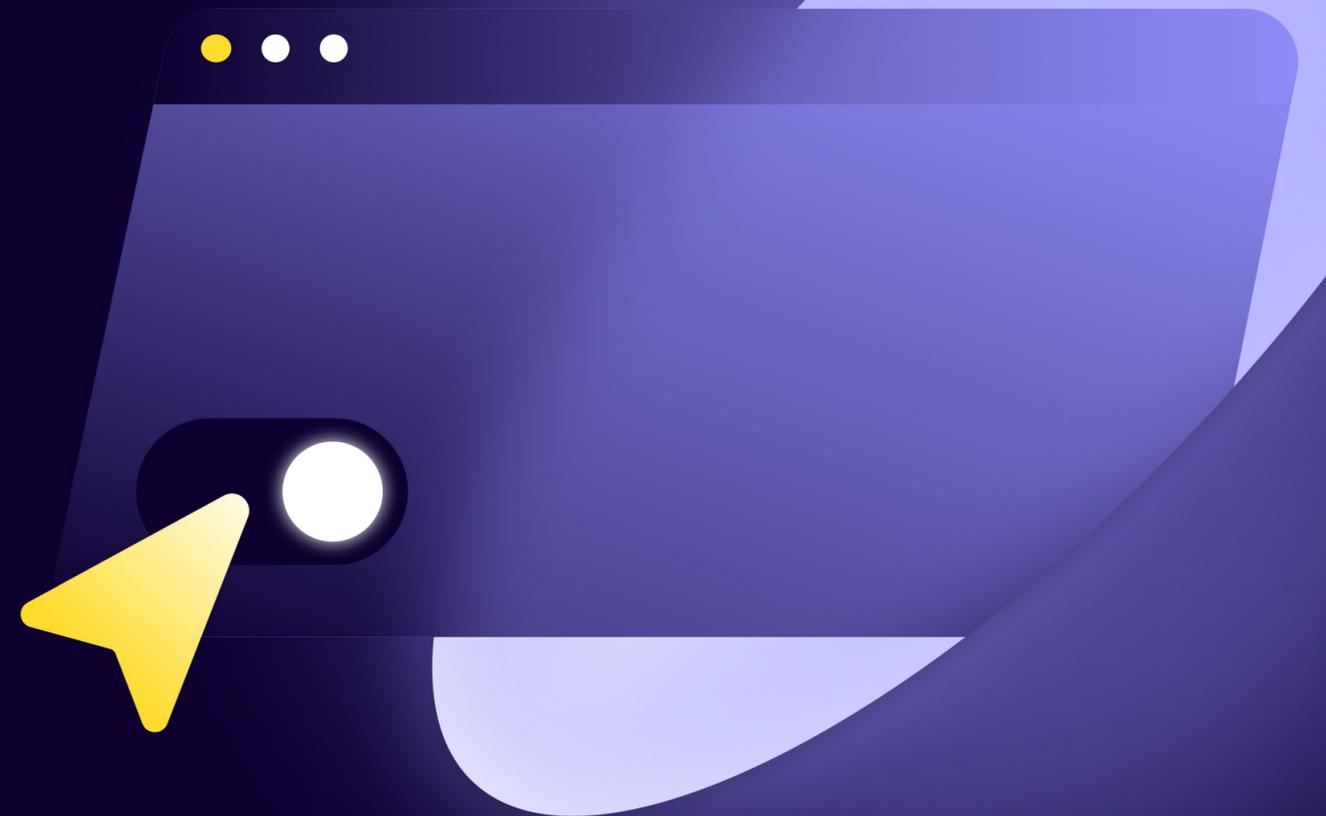


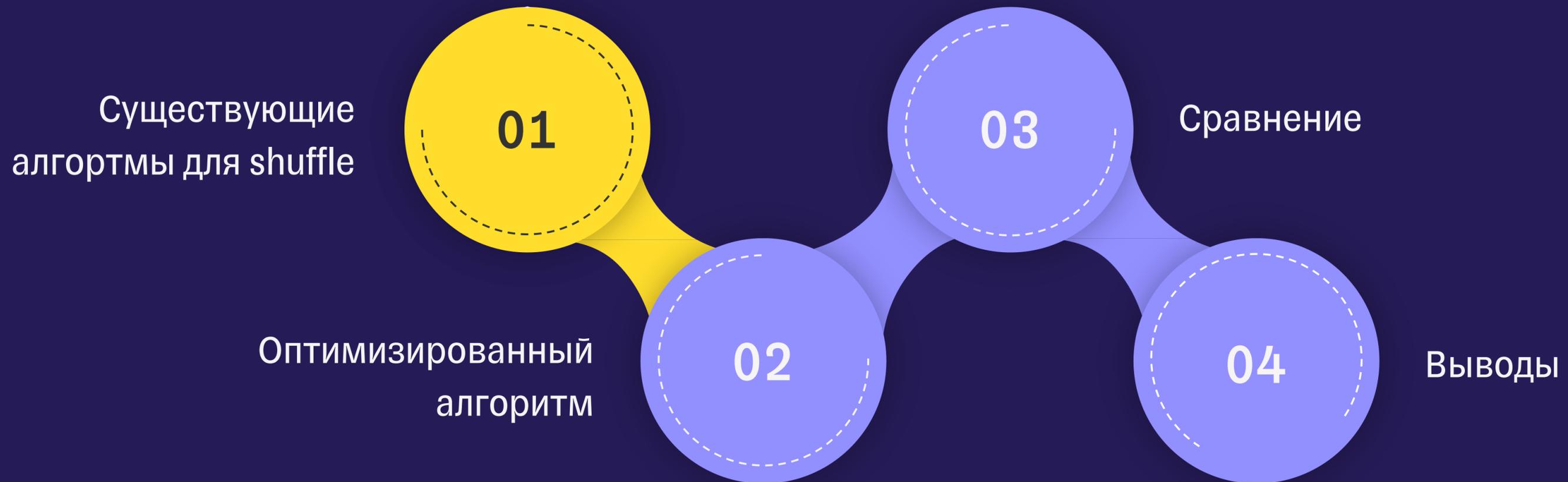


Оптимизация распределения партиций в последовательности задач распределенной обработки данных

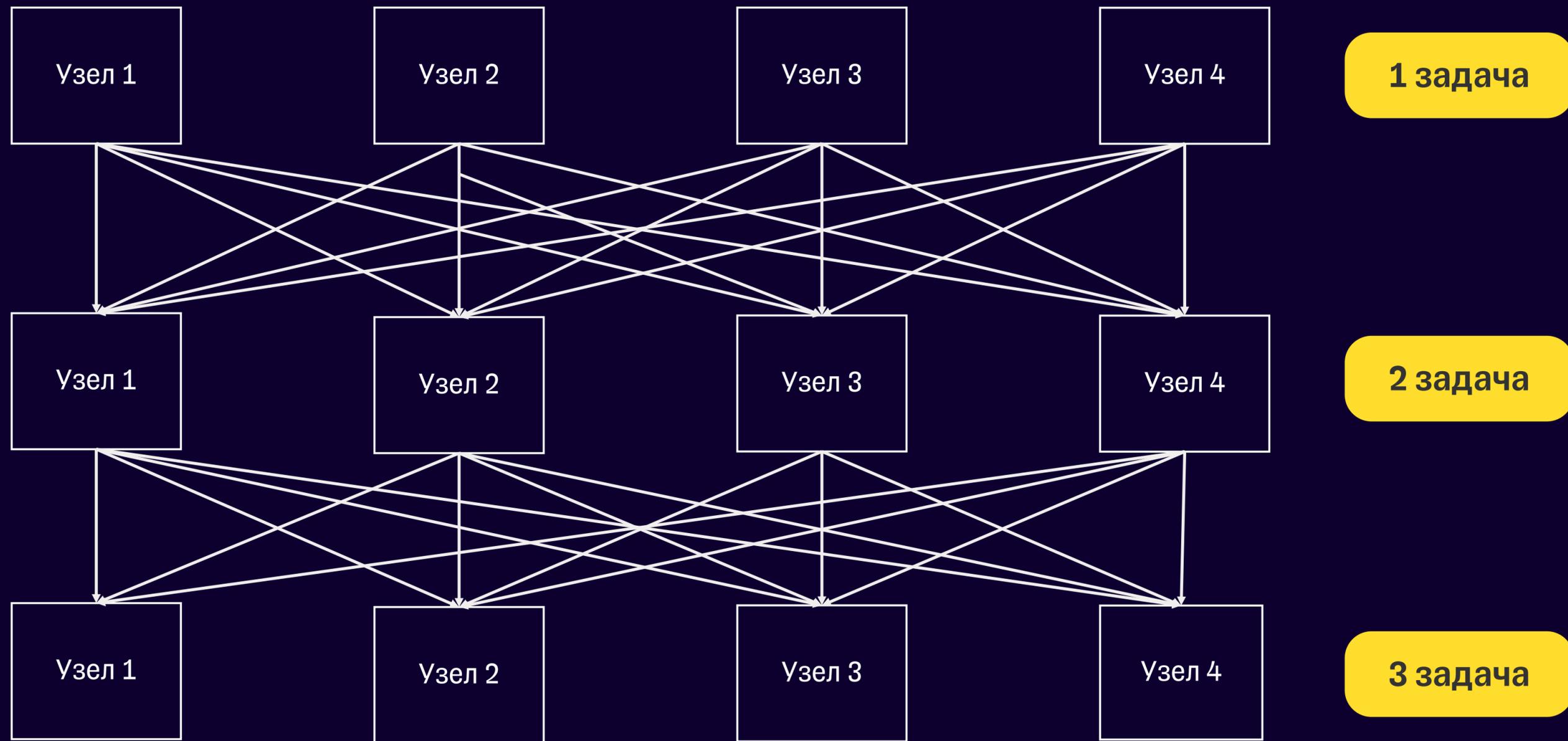
Булкина Милена
Data Engineer Т-Банк



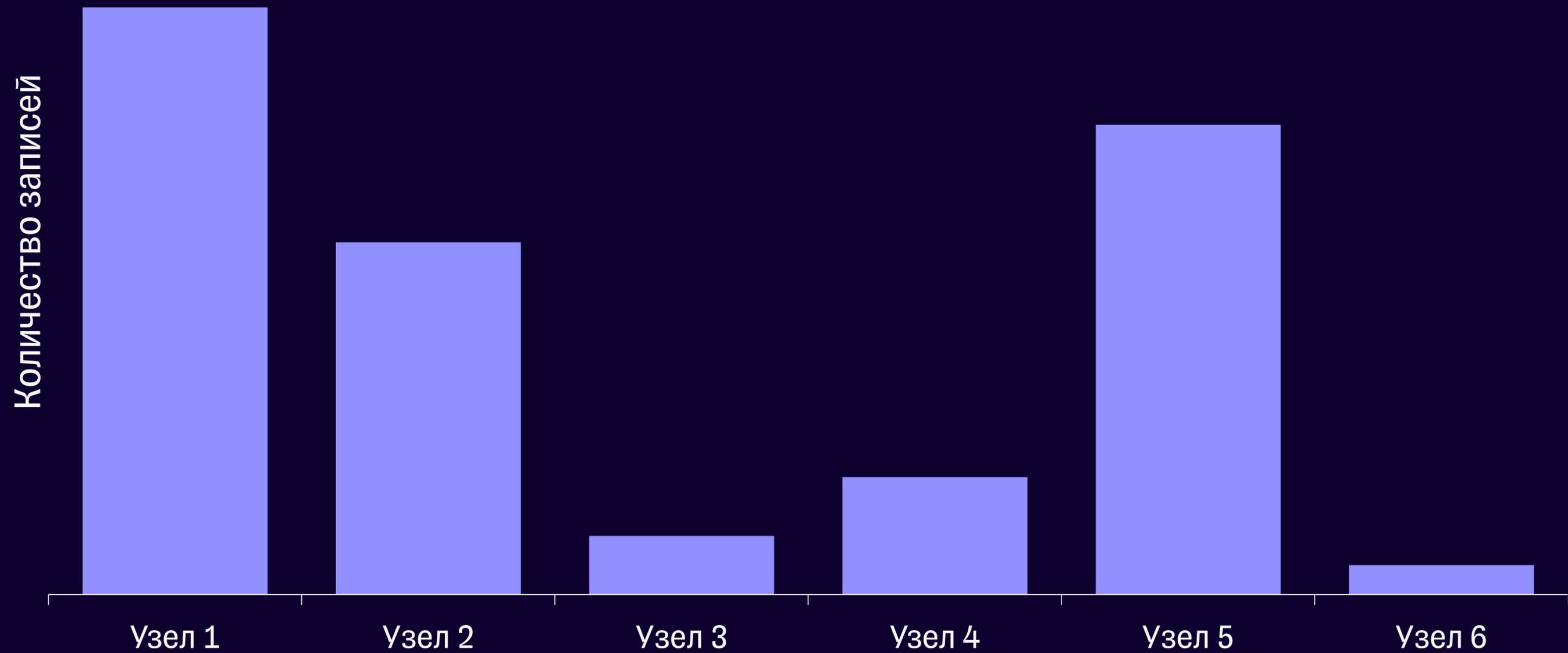
План доклада



Shuffle



Data skew



Мотивация

- ✎ Большое количество перераспределяемых данных приводит к существенному замедлению времени выполнения пайплайна
- ✎ Необходимо минимизировать объем shuffle-данных при минимальном data skew



Hash Shuffle

- ➔ Не возникает перекос на узлах, если равномерное распределение возможно
- ➔ Отсутствует оптимизация количества перемещаемых данных
- ➔ Задачи не учитываются как последовательность



Hash Shuffle

➔ Не возникает перекос на узлах, если равномерное распределение возможно

➔ Отсутствует оптимизация количества перемещаемых данных

➔ Задачи не учитываются как последовательность

Holistic Shuffle



➔ Часто возникает перекос после нескольких заданий

➔ Присутствует оптимизация shuffle-данных, но данные часто находятся на малом количестве узлов

➔ Задачи не учитываются как последовательность

1 – Holistic Shuffler for the Parallel Processing of SQL Window Functions

Fábio Coelho, José Pereira, Ricardo Vilaça, Rui Oliveira

DOI: 10.1007/978-3-319-39577-7_6

Основная идея Holistic Shuffle

Снижение объёма shuffle-данных
с помощью использования
статистики по ключам

A	NODE
1	1
1	2
1	3
1	3
2	1
2	1
2	3
3	4

Определение узла Holistic Shuffle

A	NODE
1	1
1	2
1	3
1	3

Партицирование по A



A	NODE
1	3
1	3
1	3
1	3

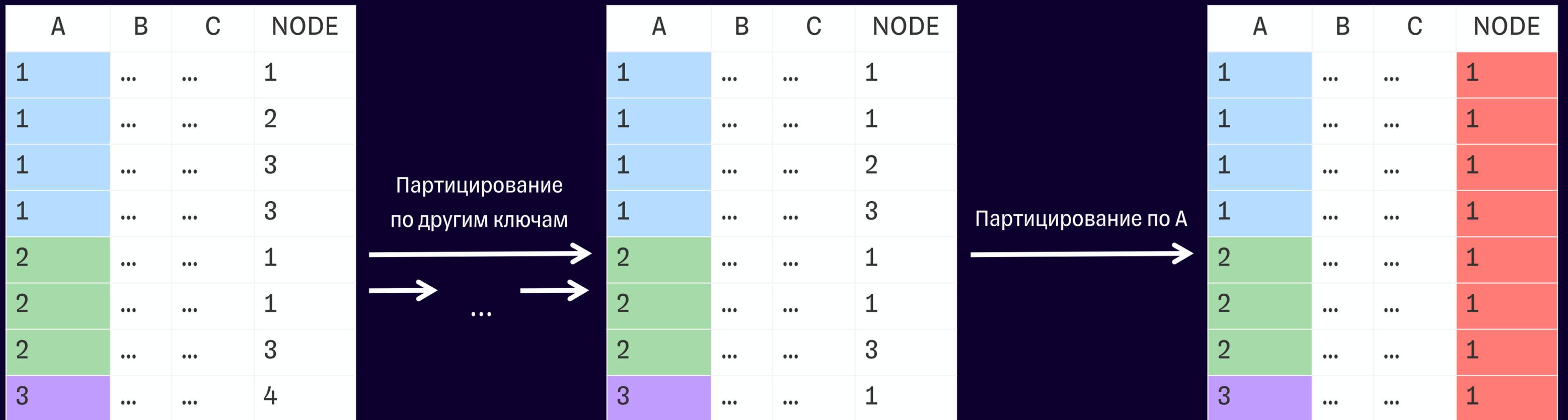
Определение узла Holistic Shuffle

A	NODE
1	1
1	2
1	3
1	3
2	1
2	1
2	3
3	4

Партицирование по A
→

A	NODE
1	3
1	3
1	3
1	3
2	1
2	1
2	1
3	4

Holistic Shuffle – data skew



Спустя несколько задач большое количество записей
могут оказаться на одном узле

Существующие
алгоритмы для shuffle

01

Оптимизированный
алгоритм

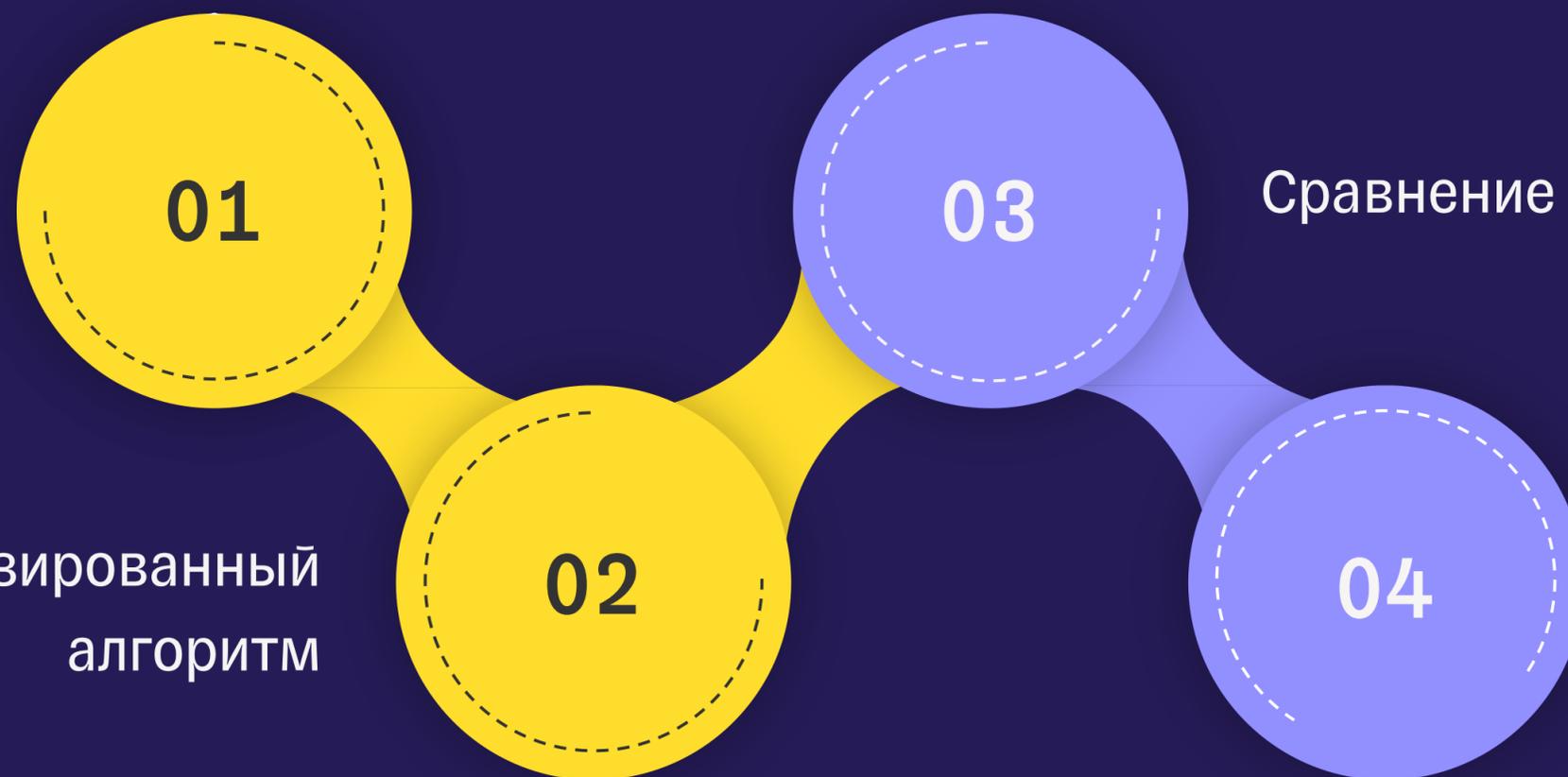
02

03

Сравнение

04

Выводы



Optimized Shuffle



В основе идея Holistic Shuffle — использование статистики для сокращения количества shuffle данных



Добавлен поиск общих ключей партицирования у задач, которые идут подряд (оптимизация ключей партицирования в пайплайне)



Ограничивается data skew

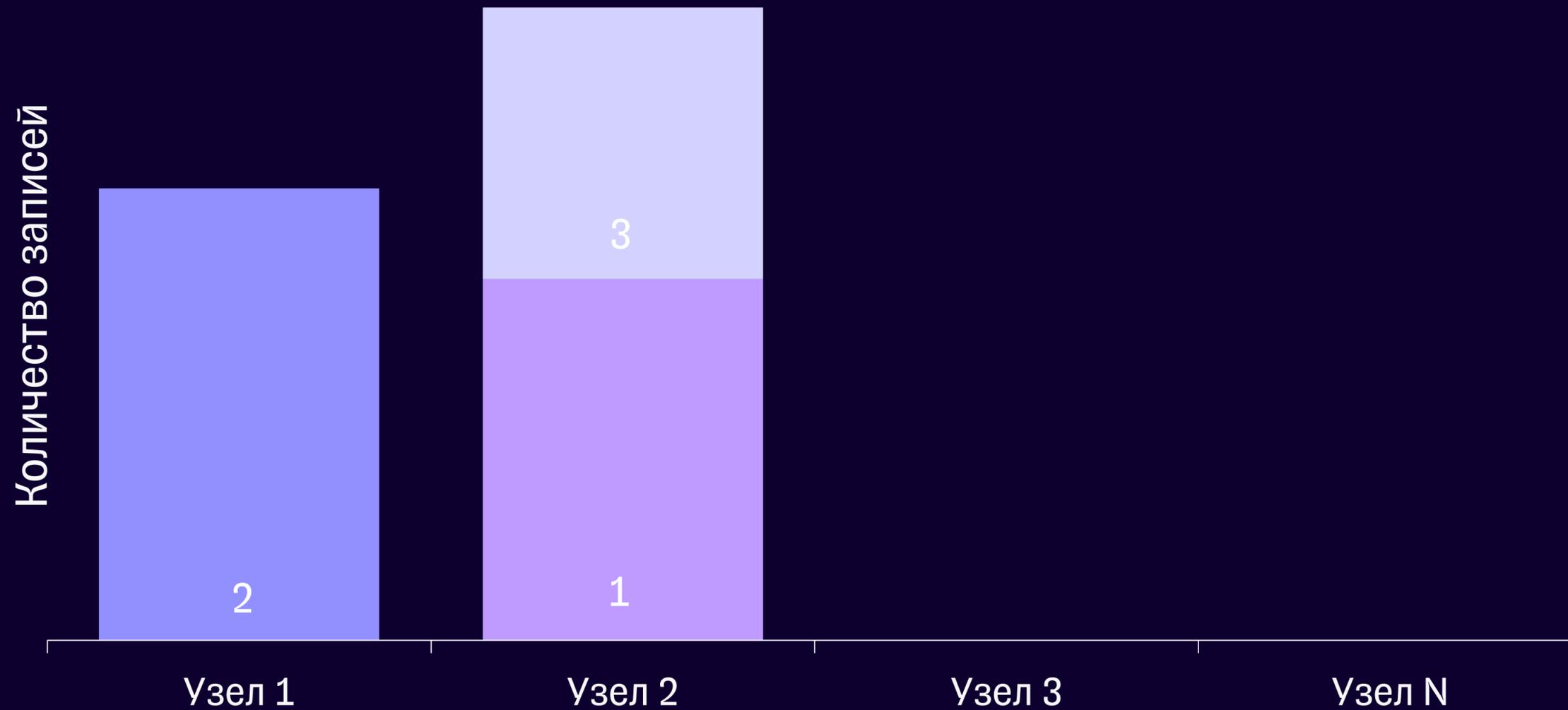
Data skew

- ✎ Входной аргумент – коэффициент, с помощью которого контролируется data skew
- ✎ Допускается перекося данных в заданных пределах

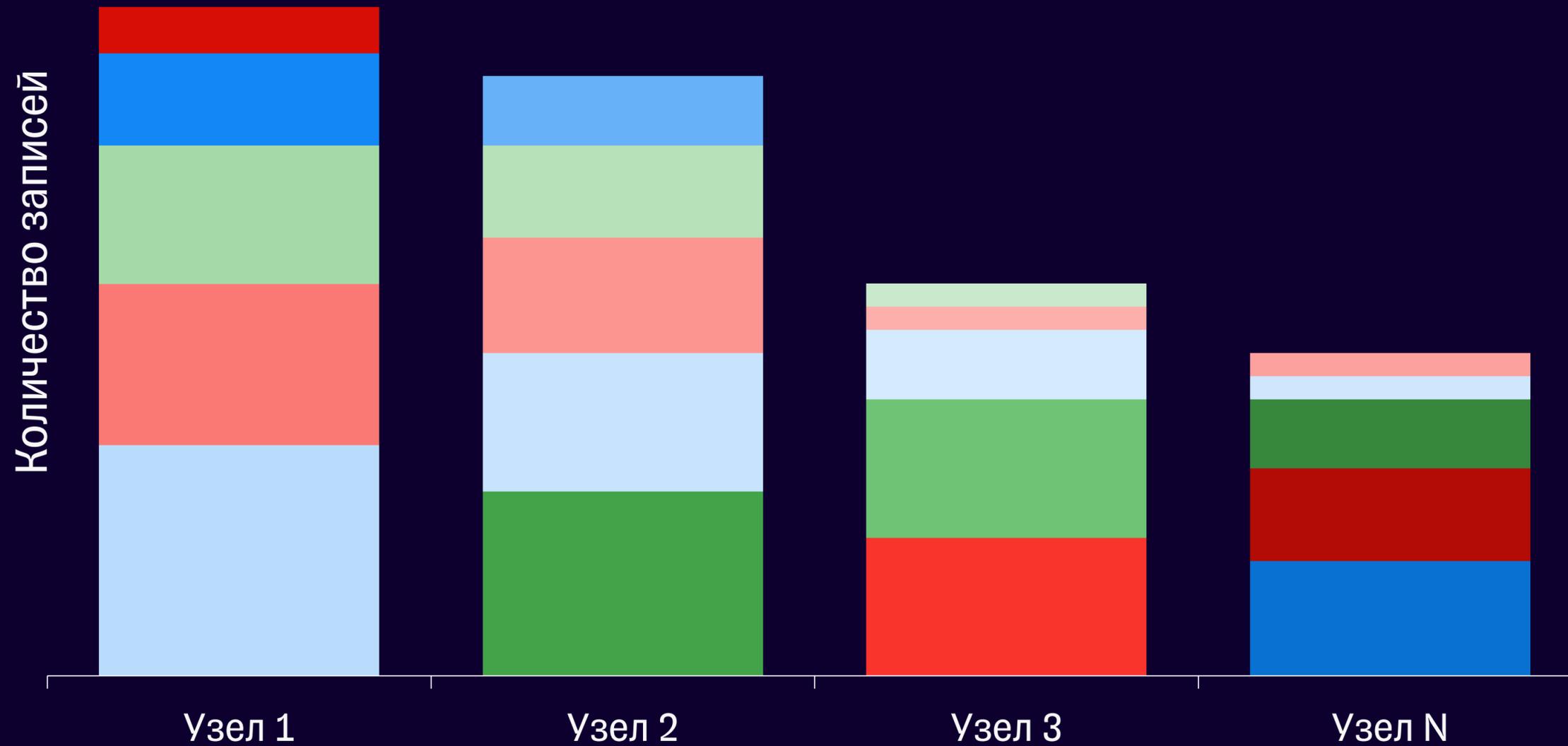
$$\text{rows_cnt} < k \times \frac{M}{N},$$

где M — общее количество записей, N — количество узлов

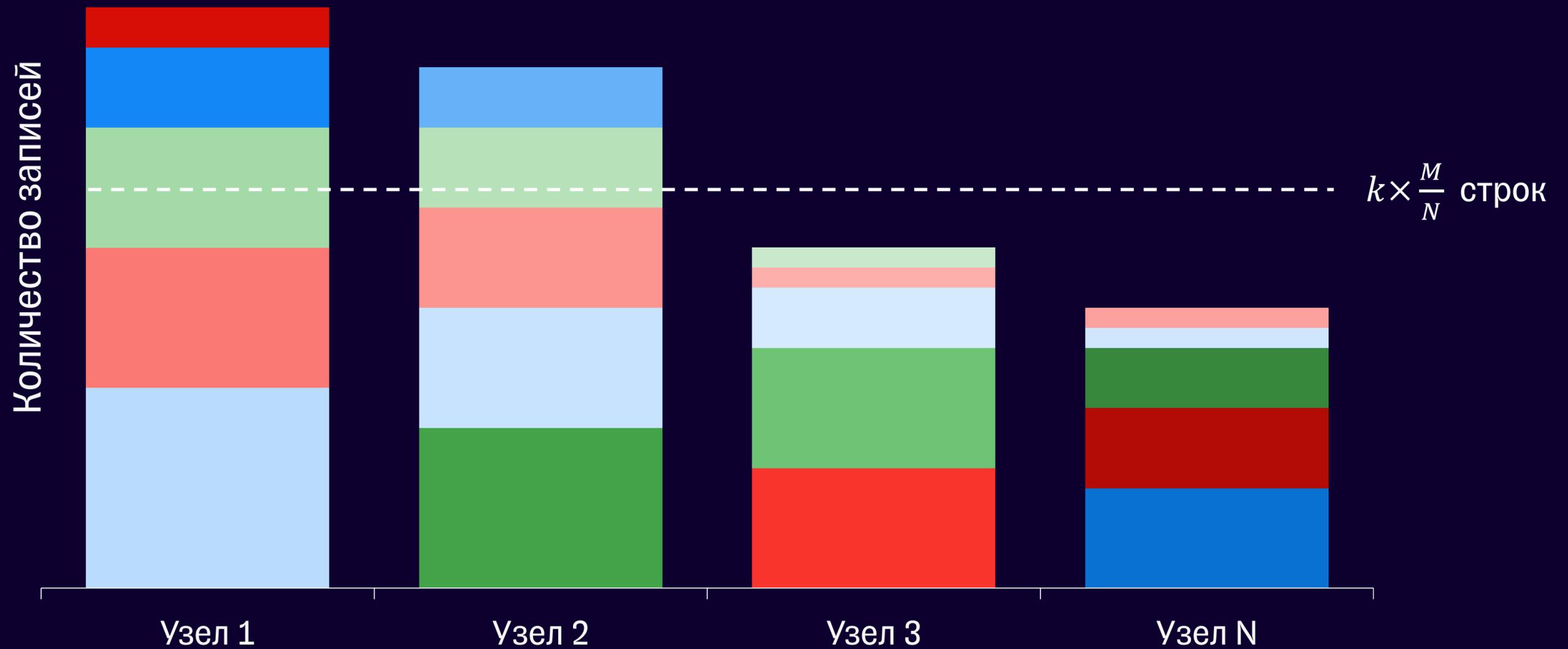
Выбор нового узла



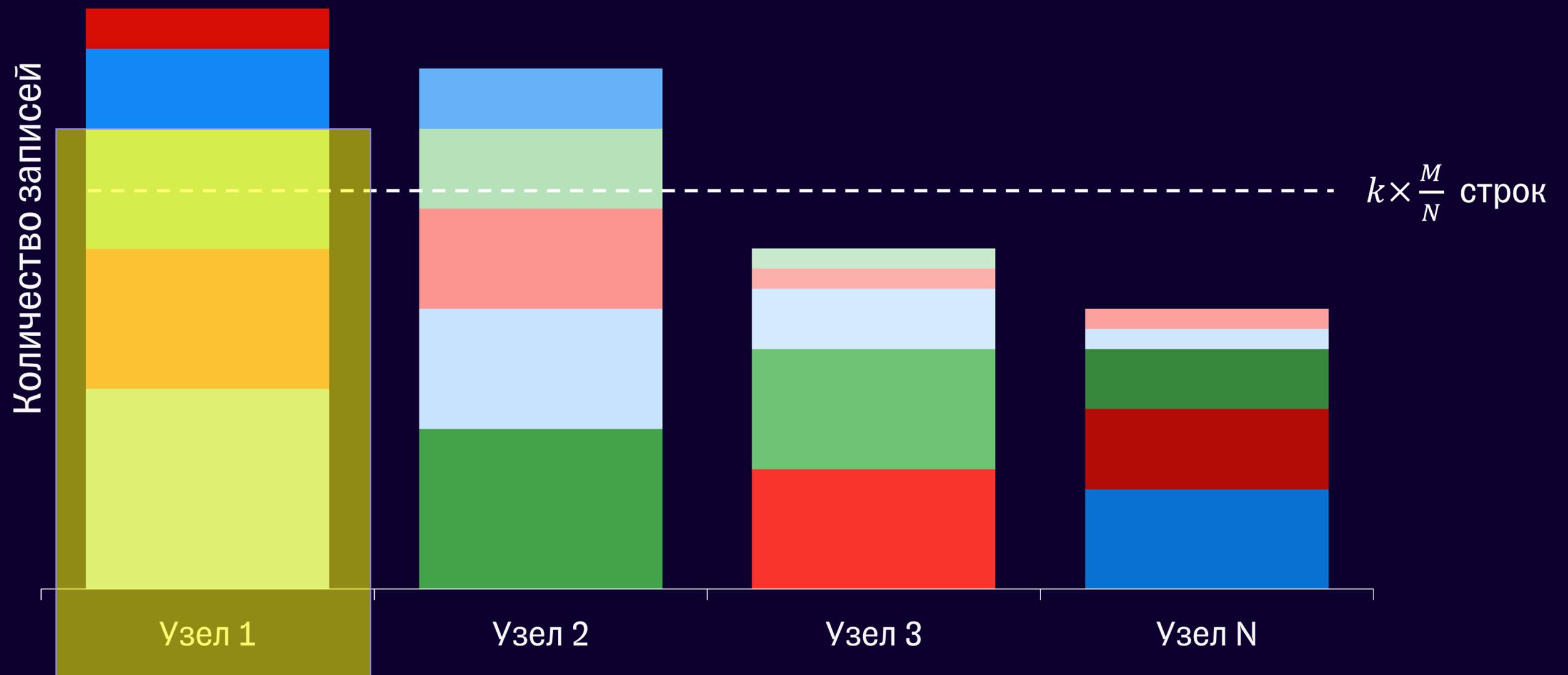
Выбор нового узла



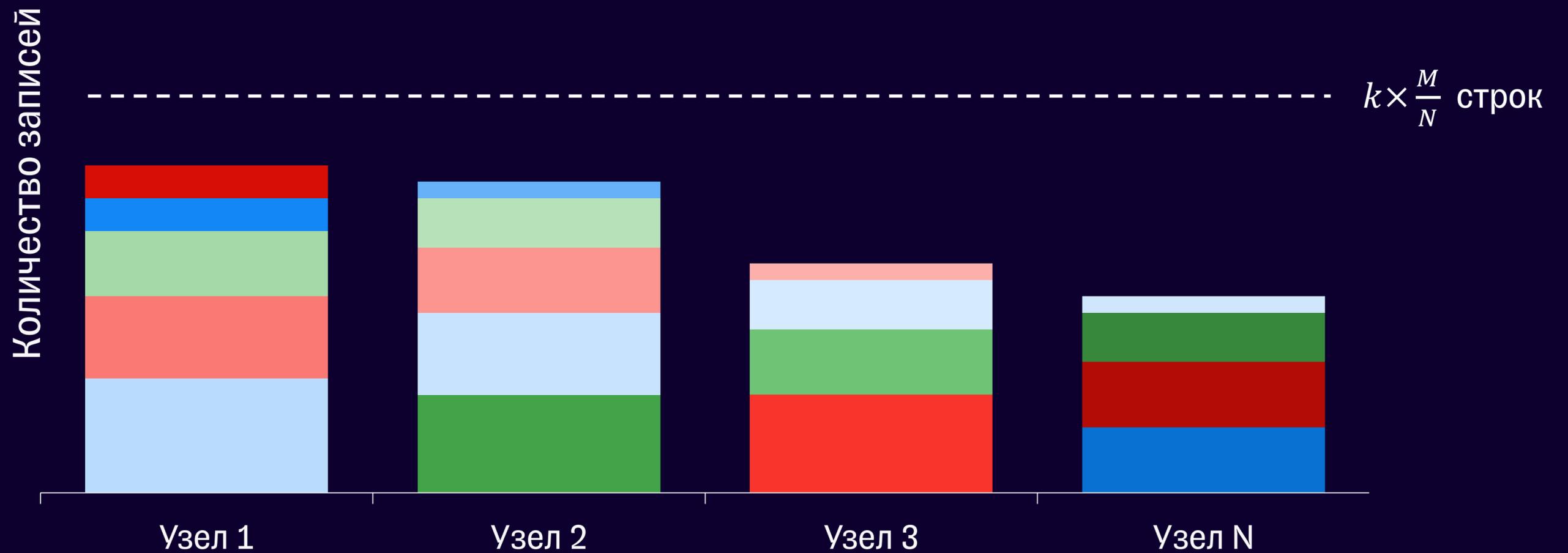
Недопустимый перекос



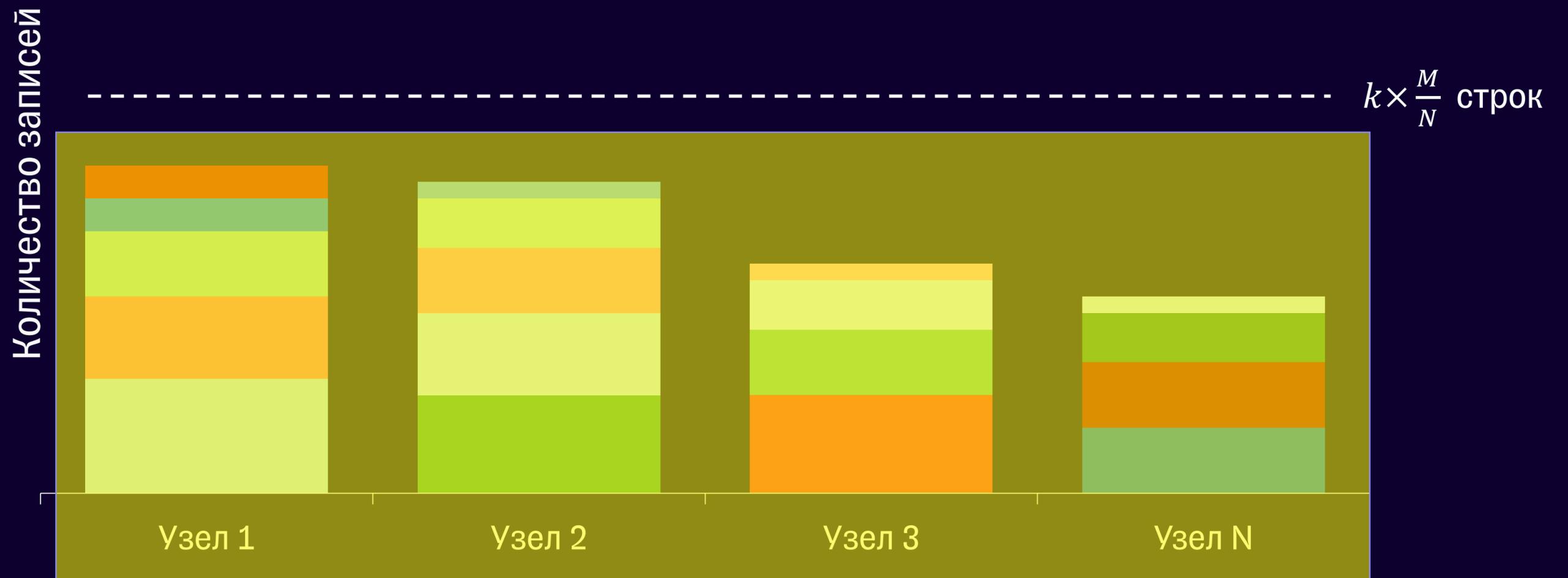
Новый узел



Новый узел



Новый узел



Сокращение количества ключей

- partition by a, b
- partition by b, c
- partition by b, d, e
- partition by f

При достаточной
кардинальности



- partition by b
- partition by f

Для задач, которые идут подряд и имеют общее подмножество ключей, можно сократить набор ключей при достаточной кардинальности подмножества ключей



Применимость алгоритма

Ограничения

- Объём ключей, используемых для репартиционирования, занимает незначительную часть всего датасета
- Только для колоночного формата данных

Пример данных, для которых может быть применен алгоритм

- Данные сейсморазведки
- Посты в блогах
- Широкие таблицы

Применимость алгоритма

Не используются в качестве
ключей репартиционирования



Существующие
алгоритмы для shuffle

01

03

Сравнение

Оптимизированный
алгоритм

02

04

Выводы

Тестовые данные

- Тест проводился на данных сейсморазведки
- Особенностью данных является большое количество атрибутов, малый процент из которых используется в качестве ключей партицирования
- 30 млн записей
- Сравнение проводилось на трех пайплайнах



Метрики

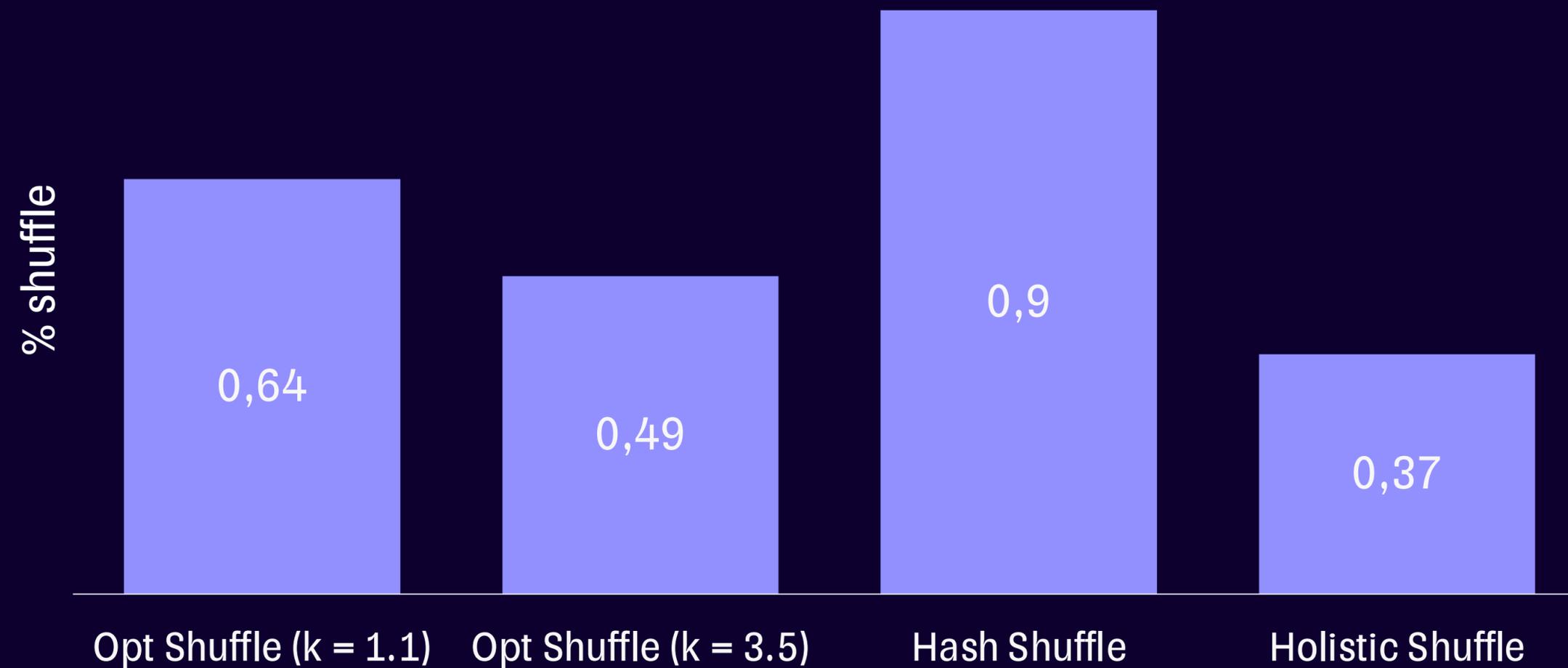
Метрики используемые для сравнения

- Доля передаваемых данных
- Коэффициент вариации

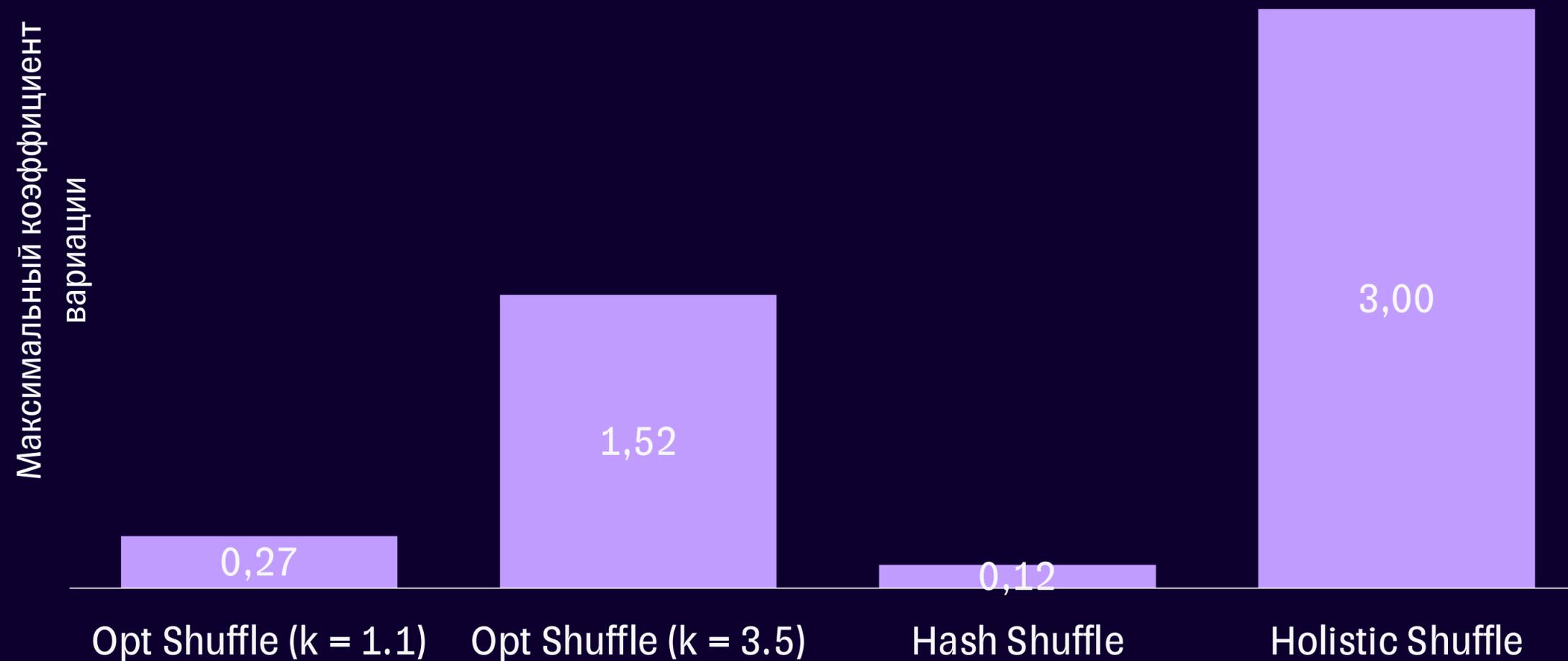
$$c_v = \frac{\sigma}{\mu}, \text{ где } \sigma \text{ — стандартное отклонение, } \mu \text{ — среднее}$$



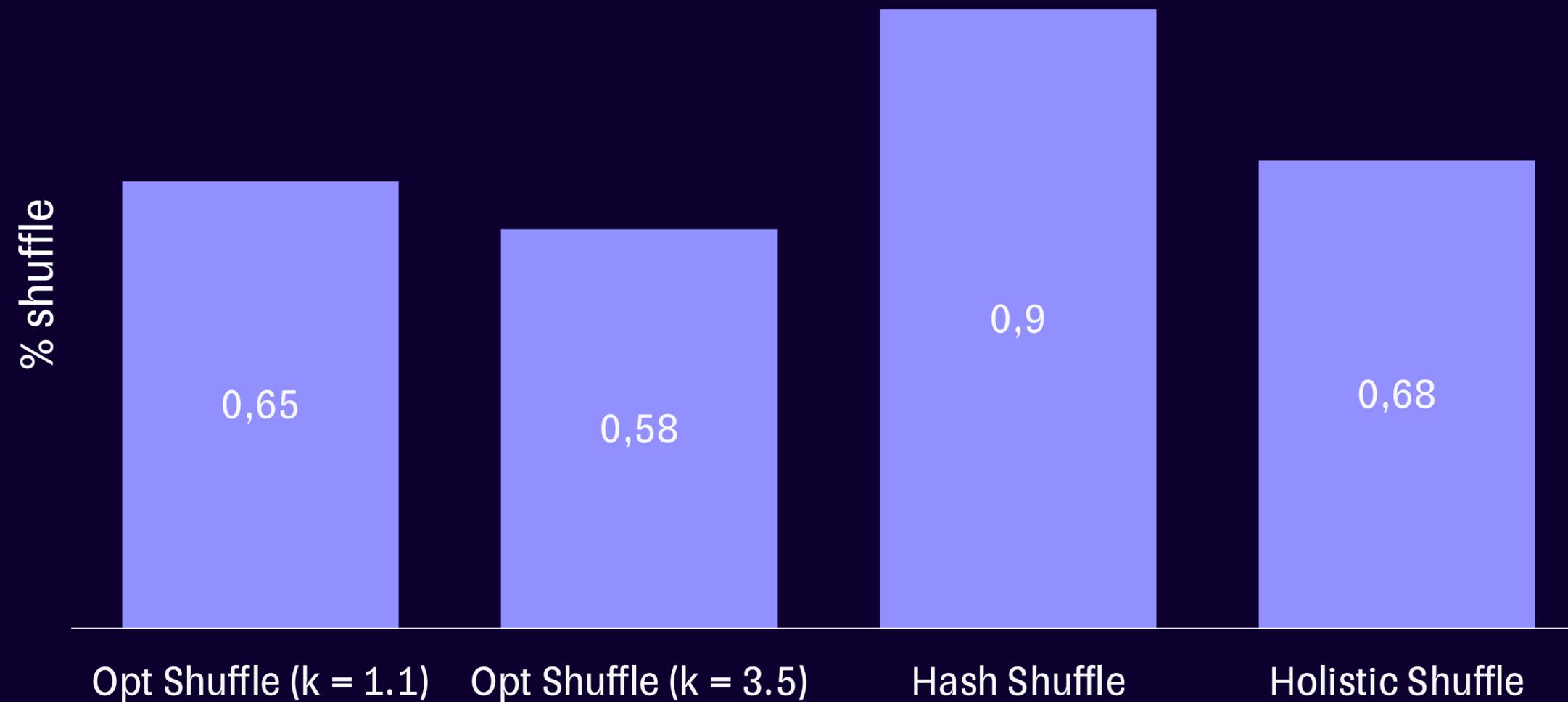
Сравнение – доля shuffle-данных



Сравнение – коэффициент вариации



Сравнение – доля shuffle-данных



Существующие
алгоритмы для shuffle

01

03

Сравнение

Оптимизированный
алгоритм

02

04

Выводы

Выводы

Для тестовых данных и $k = 1.1$

- Процент shuffle-данных — 64%
(Holistic — 45%, Hash — 90%)
- Коэффициент вариации не превышает 0.28 (Holistic — 3.0, Hash — 0.12)

Внедрение

Идея алгоритма была внедрена у индустриального партнера университета ИТМО, который занимается анализом данных сейсморазведки

Future Work



Поиск оптимального коэффициента k с помощью моделирования конфигурации системы (задача оптимизации, оптимизируем время выполнения пайплайна)



Масштабируемость



Определение максимального значения доли атрибутов, используемых в качестве ключей для репартиционирования при котором выгодно использовать алгоритм

Спасибо!

