



Индексация в поисковой платформе Ozon

Денис Габайдулин, Search Runtime



Денис Габайдулин

- Staff software engineer в Ozon
- Ранее работал в ОК (big data), Yandex (YQL)
- Более 10 лет опыта в Infrastructure Development

- Архитектура и реализация индексации в движке O2*
- Архитектура записи документов в Apache Lucene

Поисковая платформа Ozon

Поисковая платформа Ozon


Текущий статус

Движок на базе
Apache Lucene

Десятки разных
индексов

Размеры индексов:
от 10К **до 160М**
документов

Размеры индексов:
от 10Mb **до 500Gb**

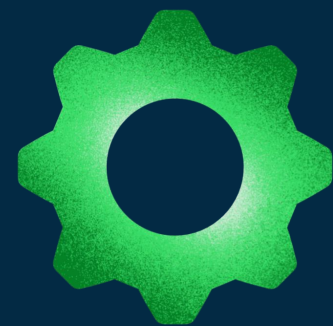
Максимальное 
количество полей
в документе: 1000+

Десятки тысяч RPS
поисковых запросов
в движок (без кэша)

Текущий статус



Успешно
прошла **два сезона**
распродаж



Рост основного
индекса **в 10+ раз**



RPS вырос
в 10+ раз



Но есть же
ElasticSearch!

Почему не Elasticsearch?

- Полный контроль над процессом ранжирования
- Полный контроль над процессами построения и репликации индекса
- Доступ к низкоуровневым оптимизациям на стадии поиска и фильтрации



Почему не Elasticsearch?

Полный контроль над процессом ранжирования

- Многоуровневое ранжирование
- Хотим менять базовую формулу ранжирования L1 (BM25)
- Хотим использовать ML в L1, L2 (learning to rank)

Почему не Elasticsearch?

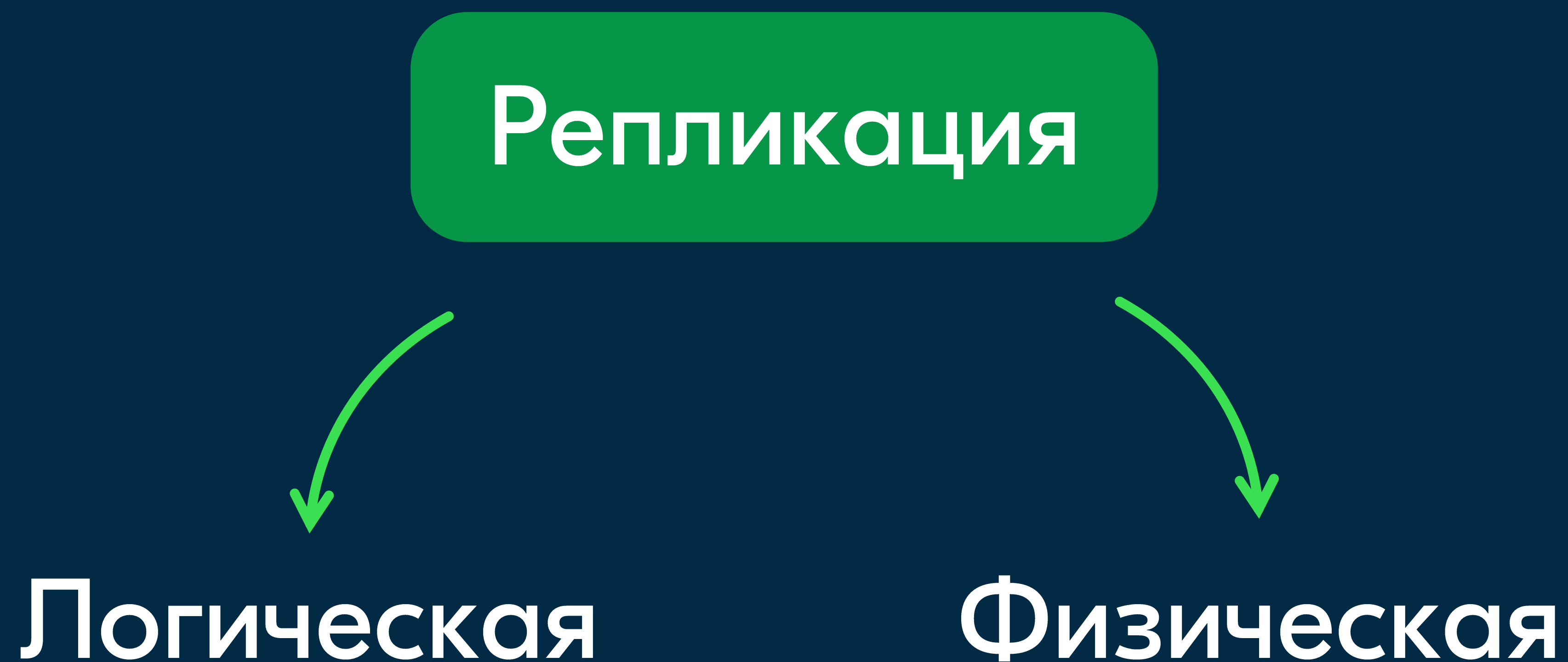
Полный контроль над процессом ранжирования

- Многоуровневое ранжирование
- Хотим менять базовую формулу ранжирования L1 (BM25)
- Хотим использовать ML в L1, L2 (learning to rank)

Почему не Elasticsearch?

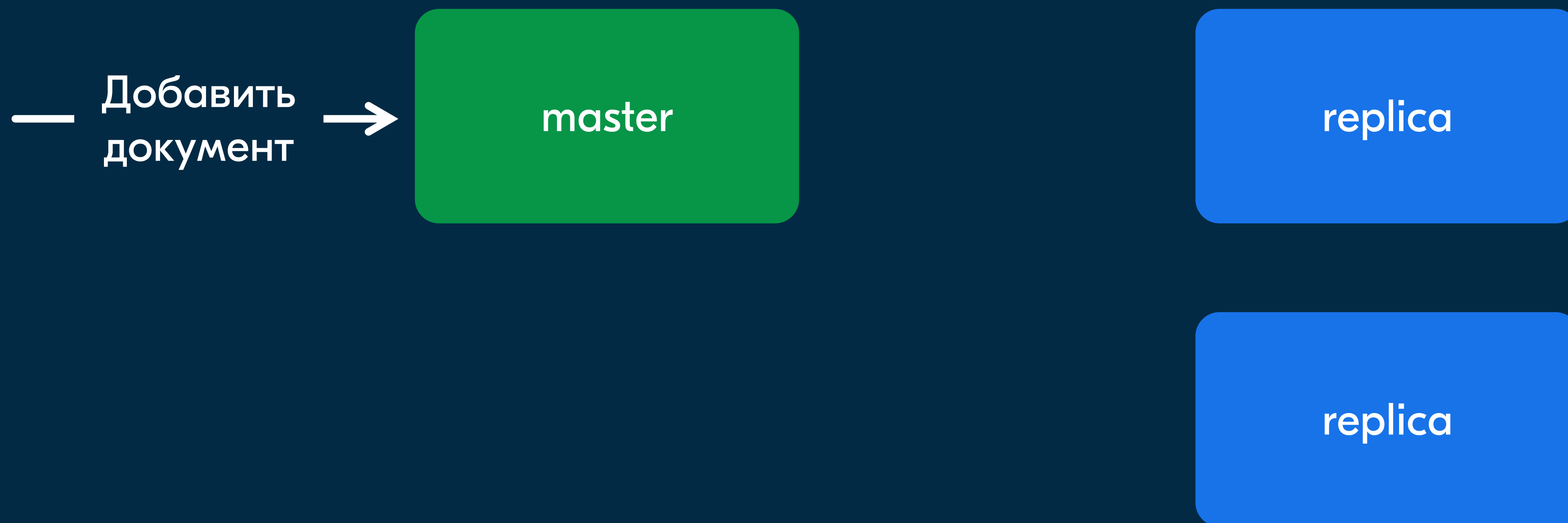
Полный контроль над процессом ранжирования

- Многоуровневое ранжирование
- Хотим менять базовую формулу ранжирования L1 (BM25)
- Хотим использовать ML в L1, L2 (learning to rank)



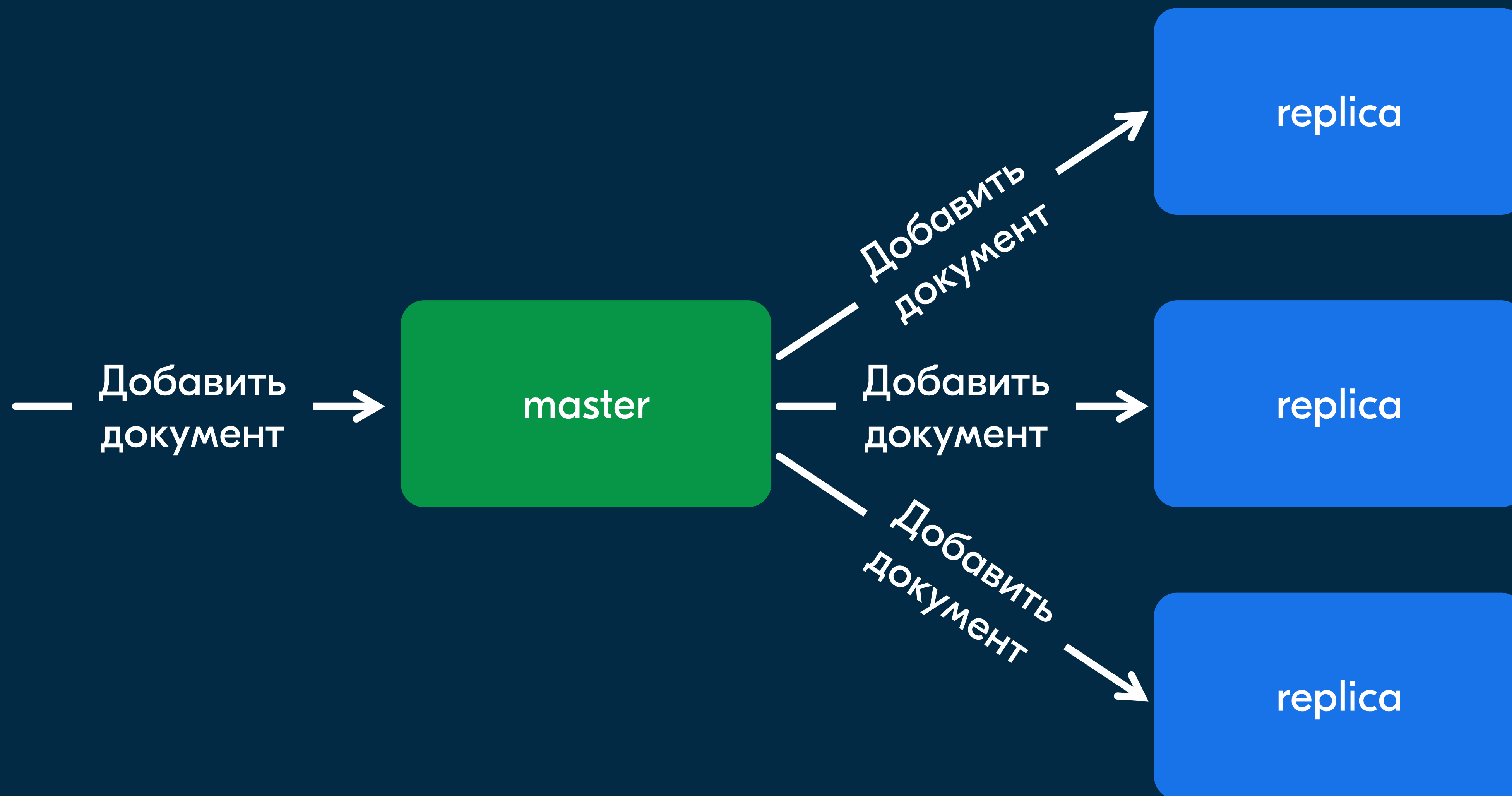
Почему не Elasticsearch?

Логическая репликация



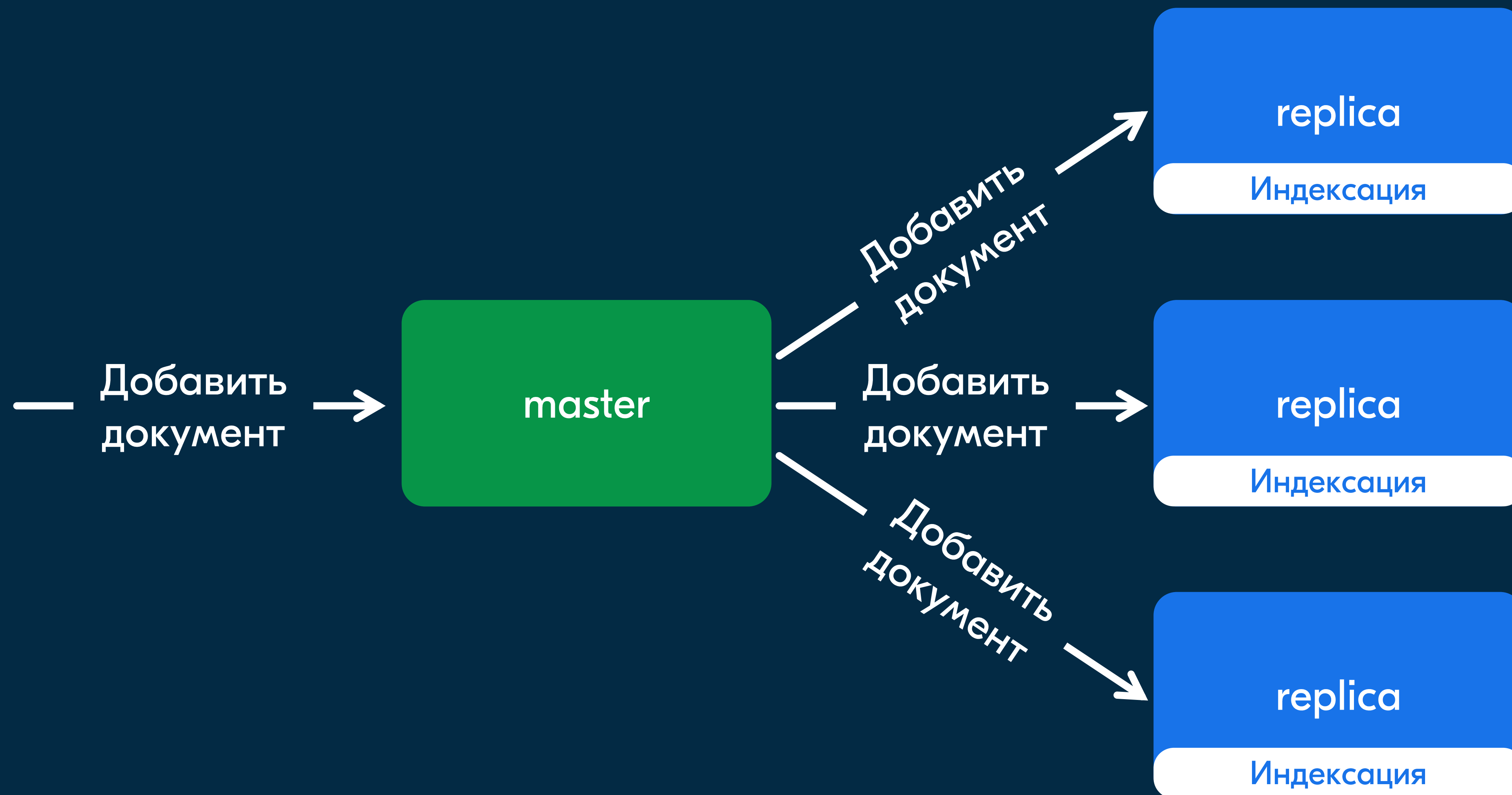
Почему не Elasticsearch?

Логическая репликация



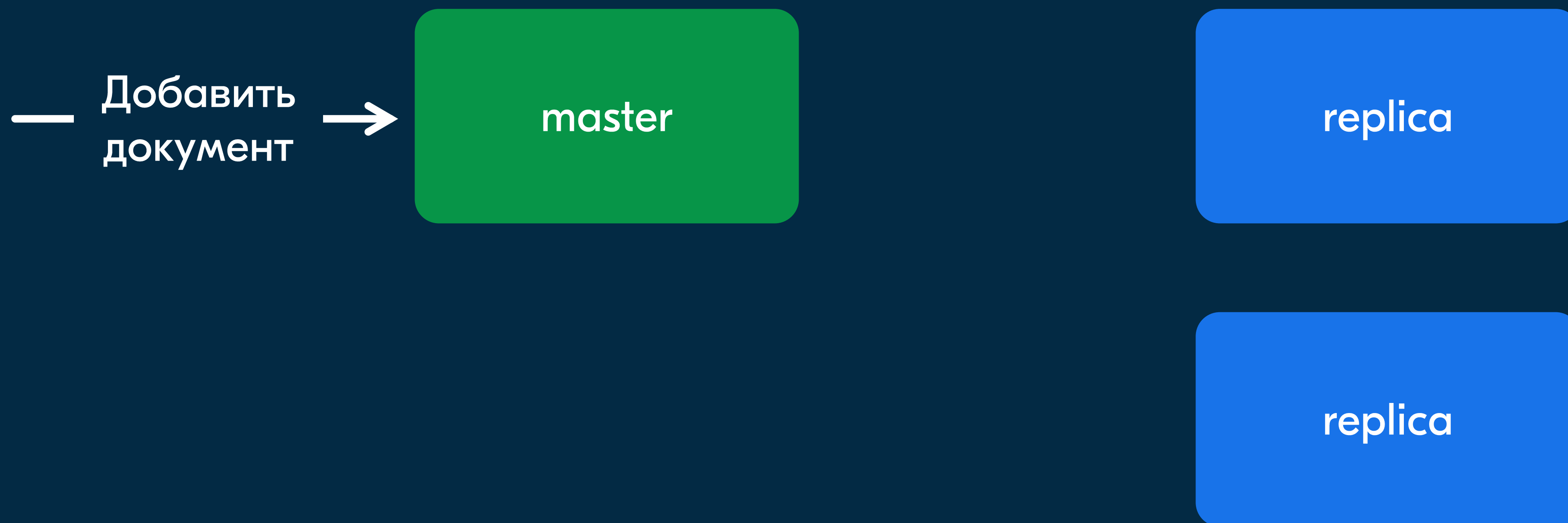
Почему не Elasticsearch?

Логическая репликация



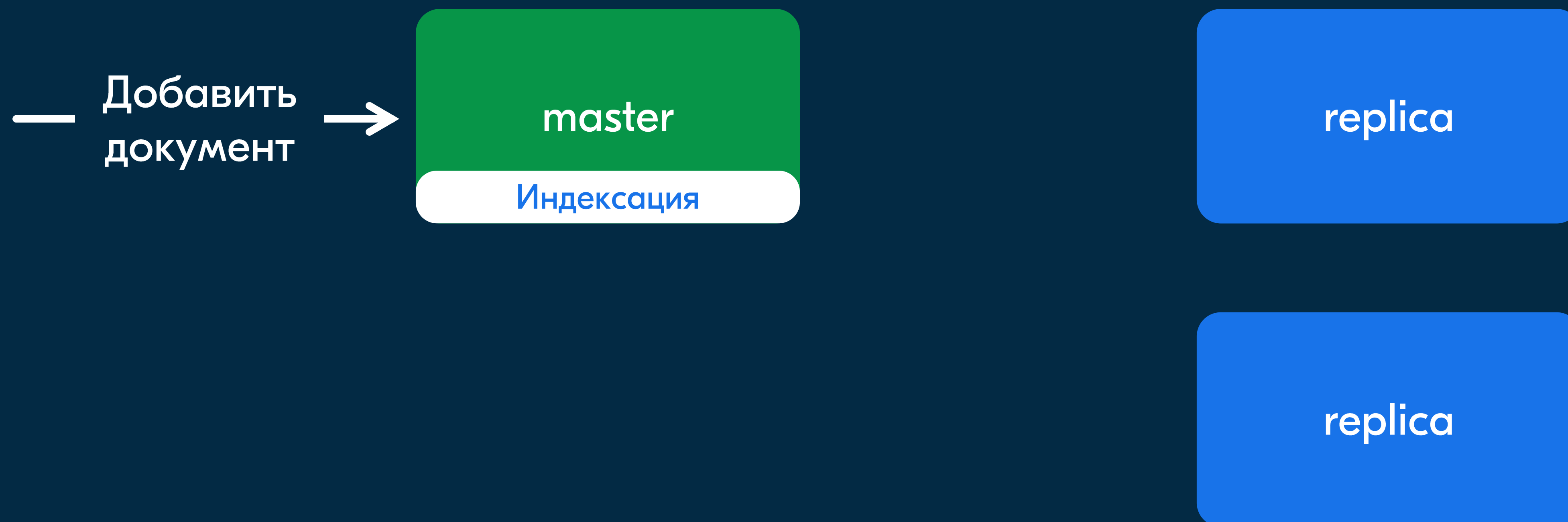
Почему не Elasticsearch?

Физическая репликация



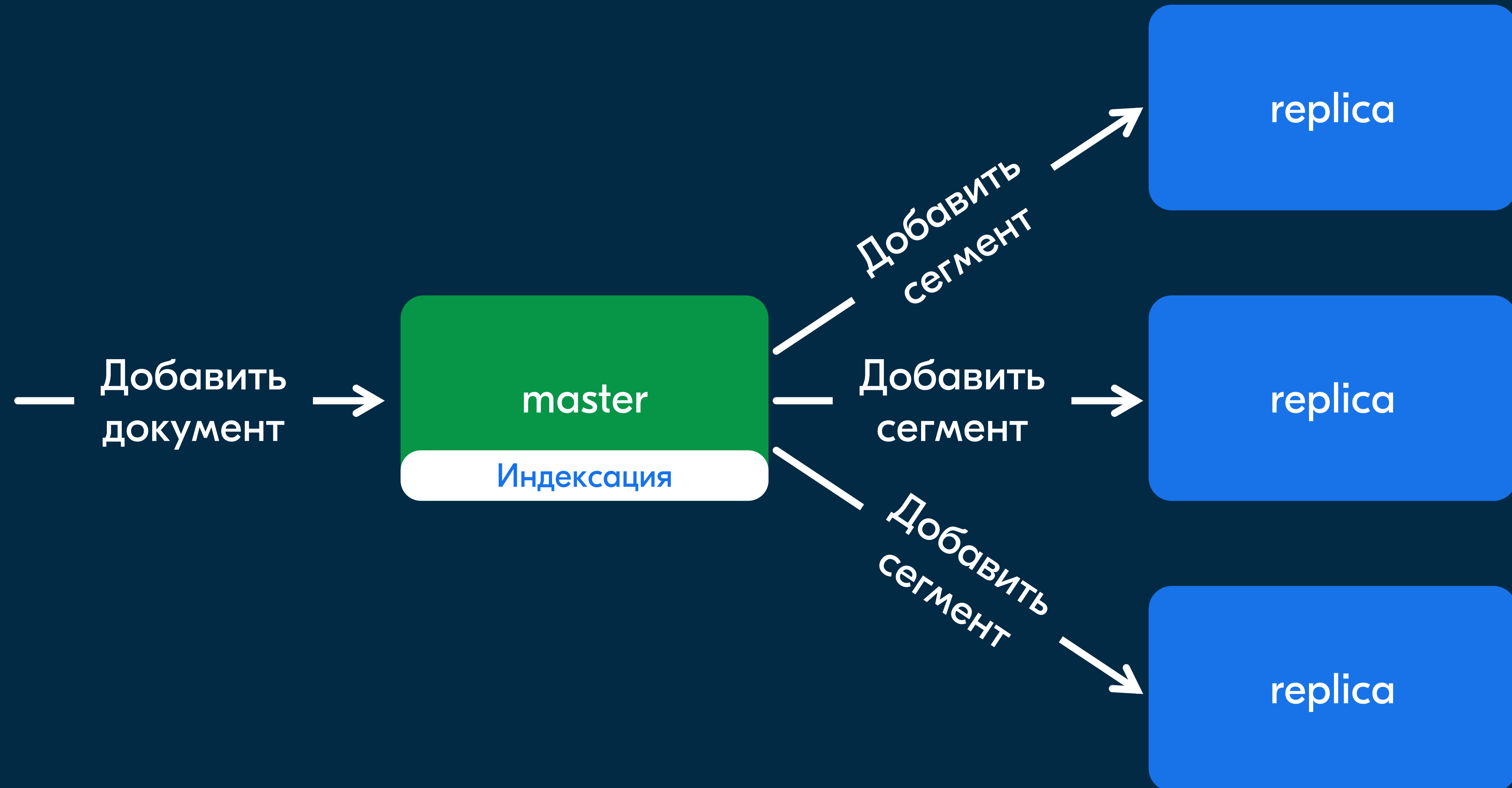
Почему не Elasticsearch?

Физическая репликация



Почему не Elasticsearch?

Физическая репликация



Почему не Elasticsearch?

Физическая репликация



Добавление готового сегмента
дешевле на порядки, чем индексация

Почему не Elasticsearch?

Физическая репликация

- Масштабируется на сотни поисковых реплик

- У нас **600+** реплик на данный момент

Почему не Elasticsearch?

Полный контроль над процессом построения индекса

- **Response time во многом зависит от геометрии индекса**
- Геометрия индекса — это размер и количество сегментов
- Управление геометрией индекса делается через `merge policy`
- В `elastic` нет возможности написать свою `merge policy`

Почему не Elasticsearch?

Полный контроль над процессом построения индекса

- Response time во многом зависит от геометрии индекса
- Геометрия индекса — это размер и количество сегментов
- Управление геометрией индекса делается через merge policy
- В elastic нет возможности написать свою merge policy

Почему не Elasticsearch?

Полный контроль над процессом построения индекса

- **Response time** во многом зависит от геометрии индекса
- Геометрия индекса — это размер и количество сегментов
- Управление геометрией индекса делается через **merge policy**
- В **elastic** нет возможности написать свою **merge policy**

Почему не Elasticsearch?

Полный контроль над процессом построения индекса

- **Response time** во многом зависит от геометрии индекса
- Геометрия индекса — это размер и количество сегментов
- Управление геометрией индекса делается через **merge policy**
- В **elastic** нет возможности написать свою **merge policy**

Почему не Elasticsearch?

Размеры сегментов (плохо )

Сегмент 6

Сегмент 5

Сегмент 3

Сегмент 4

Сегмент 2

Сегмент 1

Индекс с 6 файлами (сегментами)

Почему не Elasticsearch?

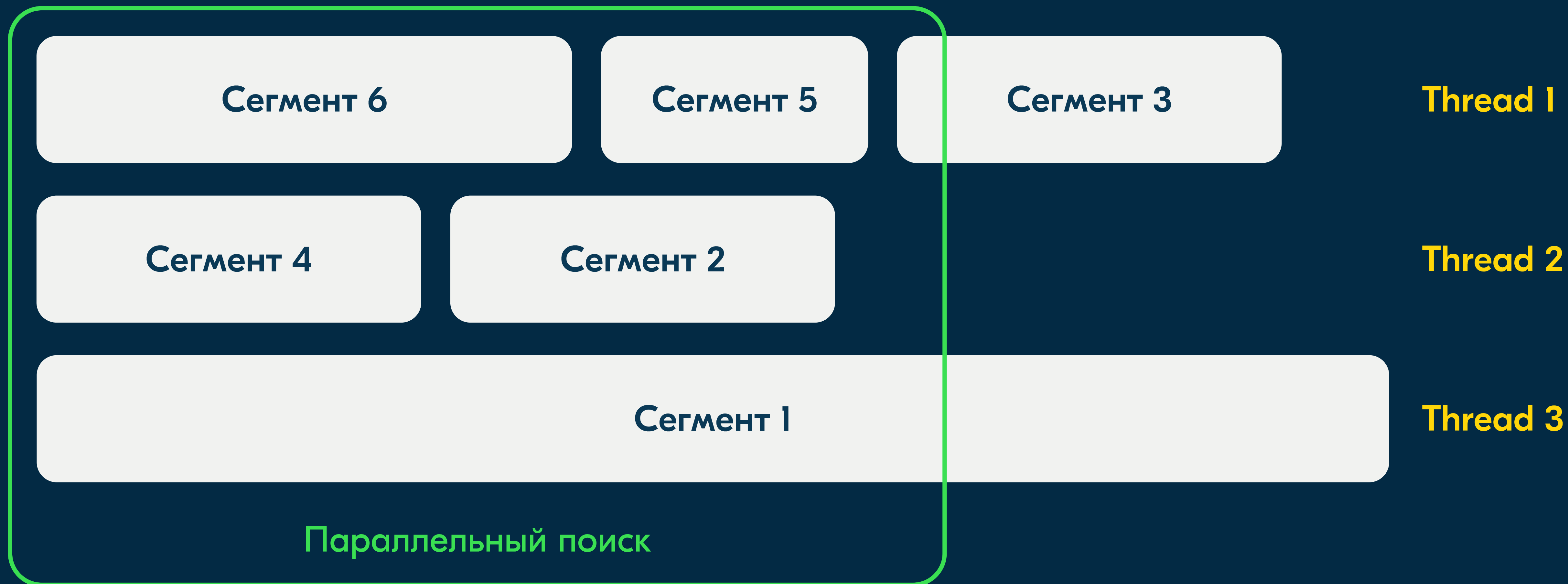
Размеры сегментов (плохо )



Поиск документов в 3-х потоках

Почему не Elasticsearch?

Размеры сегментов (плохо 🐔)



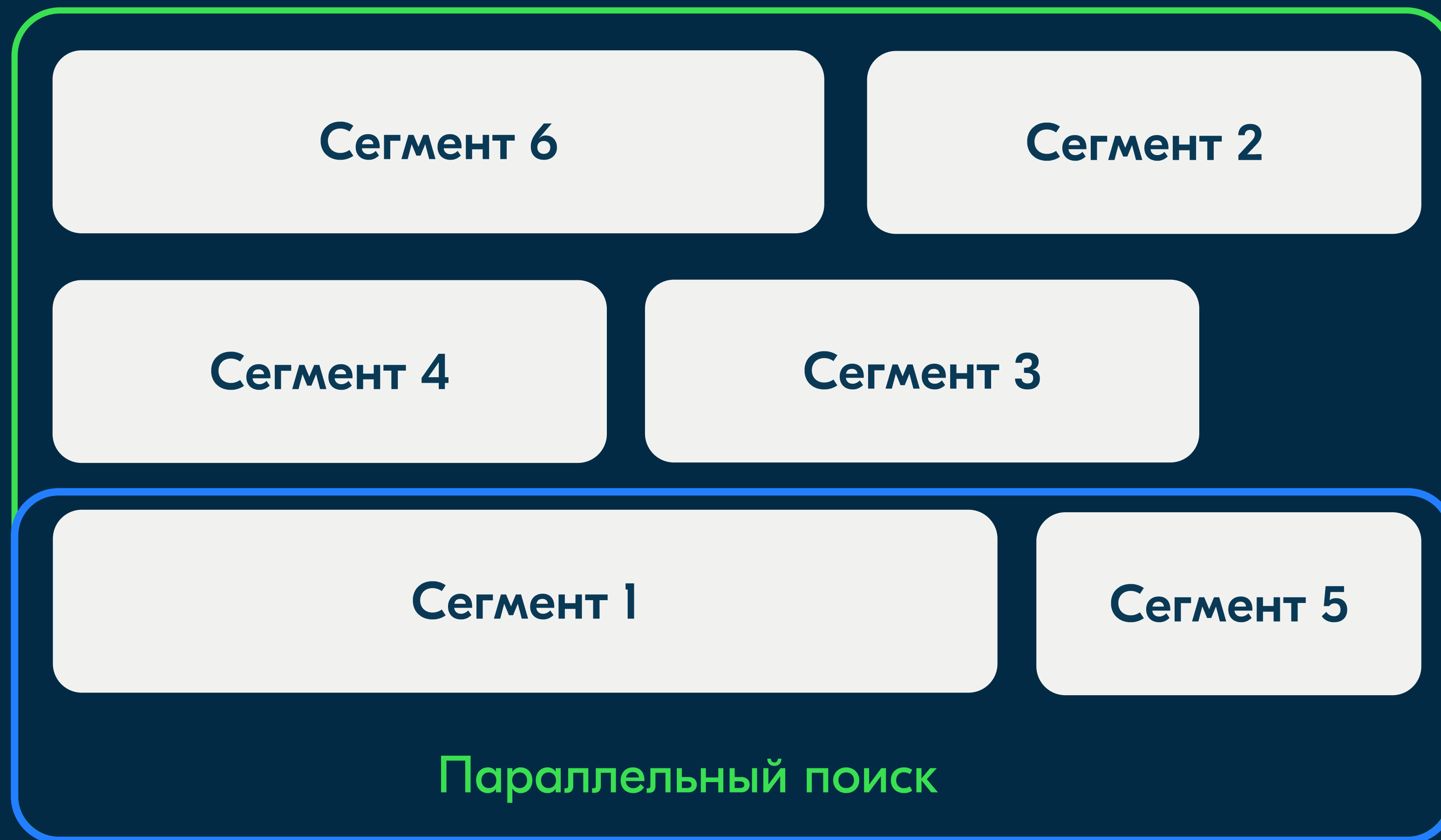
Почему не Elasticsearch?

Размеры сегментов (плохо )



Почему не Elasticsearch?

Размеры сегментов (хорошо 🐧)



Thread 1

Thread 2

Thread 3

Почему не Elasticsearch?

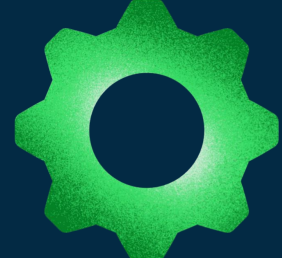
Эксплуатация

- Elasticsearch работал какое-то время (2017-2020 годы) [1]
- Время ответа в p95 было неприемлемо (1 секунда)
- O2 отвечает в p99: **250-300 ms**, p90: **90 ms**

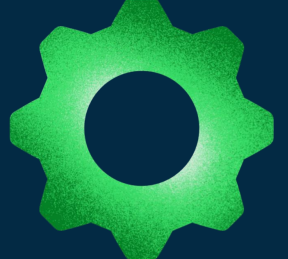
Поисковая платформа Ozon

Таймлайн разработки движка O2



 **Индексация** — это процесс построения специальных структур данных для поиска (inverted index, doc values, kd trees)

- На входе документы, на выходе поисковый индекс

 **O2 Master** — это сервис, который отвечает за управление индексами и непосредственно индексацию

Архитектура O2 Master

Наша цель

- Уменьшить время построения индекса за счет утилизации всех доступных ядер и серверов
- Перестраивать индекс полностью не менее чем раз в сутки



- Вертикальное масштабирование
- Горизонтальное масштабирование
- Отказоустойчивость (отказ ноды, отказ ЦОД)
- Throughput важнее latency

Ключевая метрика

Количество
проиндексированных
документов в секунду

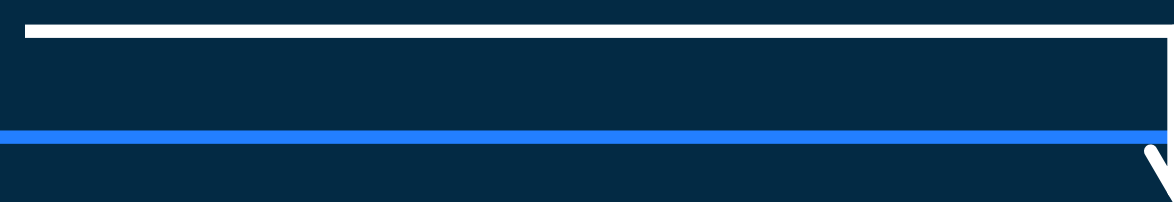


Или коротко:
**скорость
индексации**

Индексация в Apache Lucene

Архитектура индексации в Apache Lucene

Добавить документ

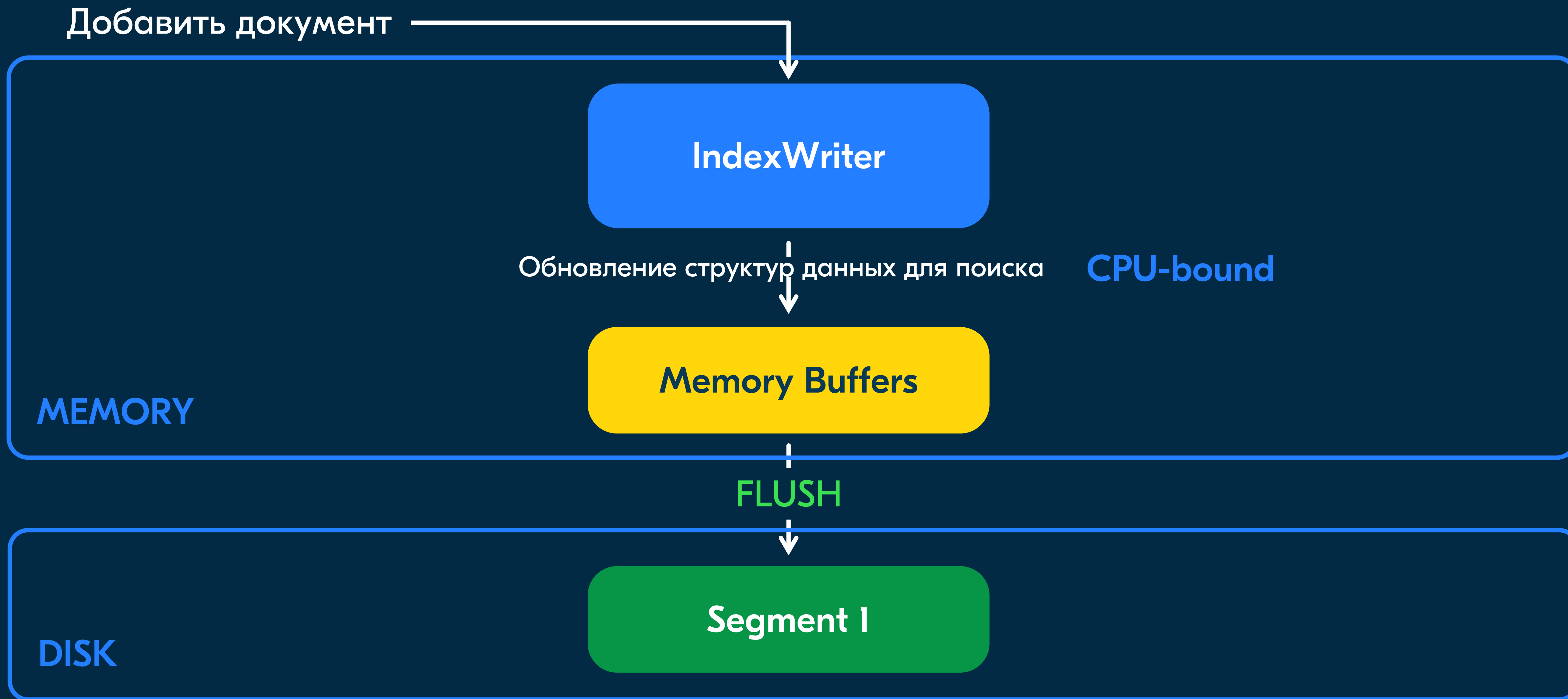


MEMORY

Архитектура индексации в Apache Lucene



Архитектура индексации в Apache Lucene

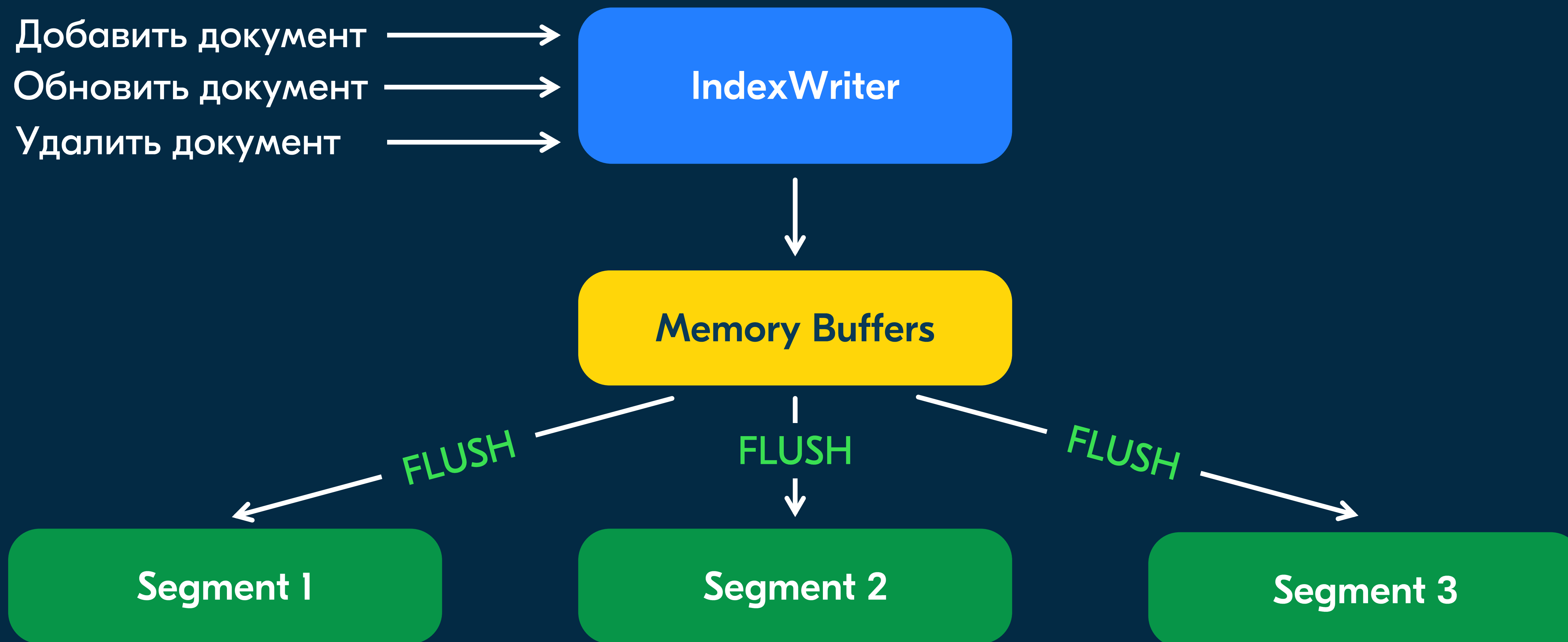


- CPU-bound-задача

Ресурсы

- CPU-bound:
 - Продвинутый анализ текста
 - Сложные структуры данных (FST, kd-trees, knn-vectors)
 - Много техник компрессии данных (doc values): delta encoding, gcd, vints, deduplication
 - Массивно параллельная обработка в памяти

Архитектура индексации в Apache Lucene



Index, version: 1

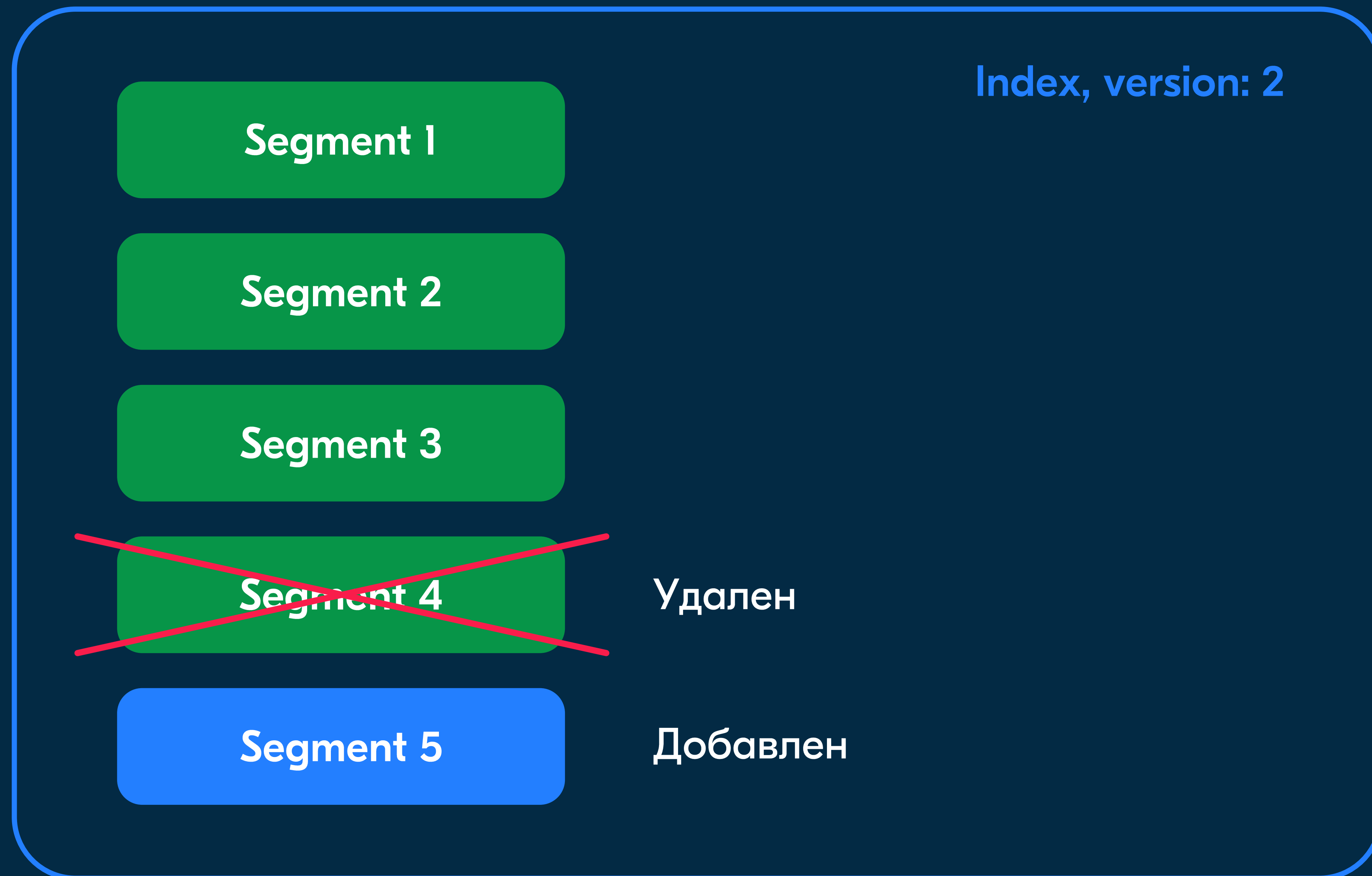
Segment 1

Segment 2

Segment 3

Segment 4

Сегменты
иммутабельны





А если сегментов
будет **слишком много?**

Архитектура индексации в Apache Lucene

Слишком много сегментов

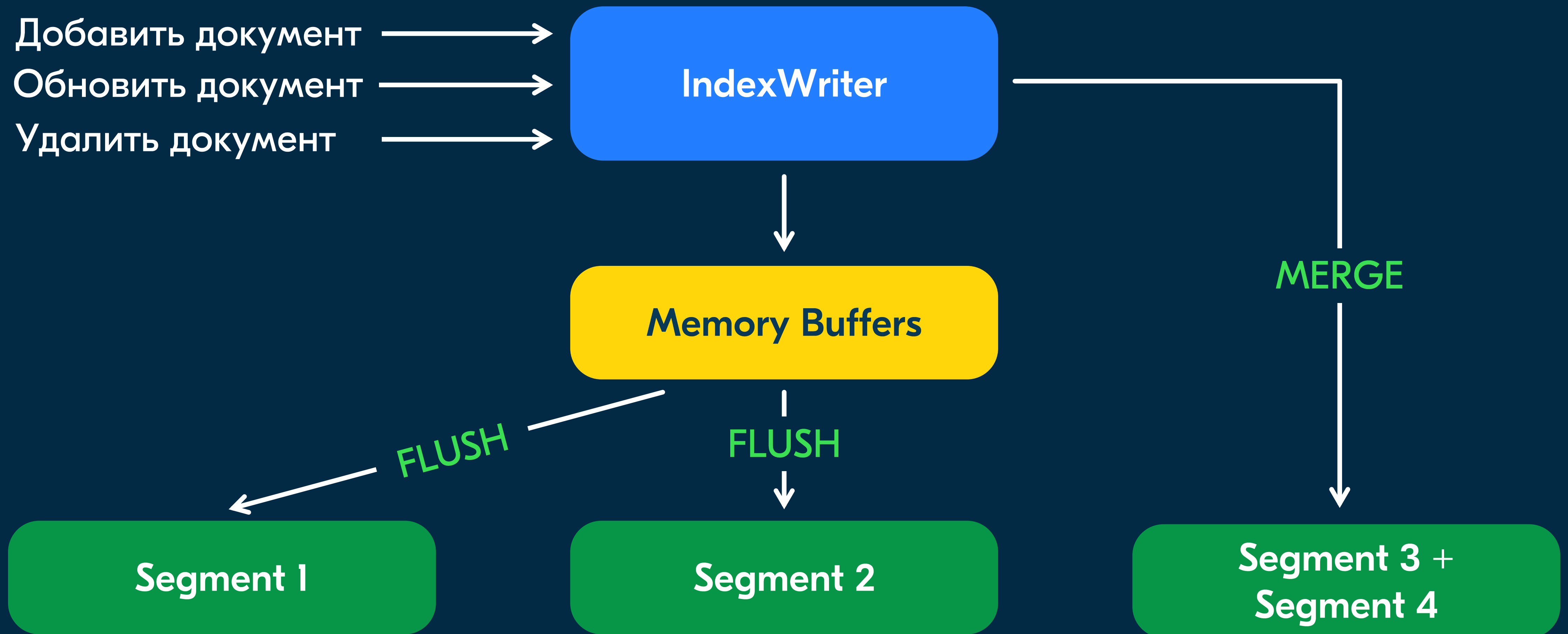
- Проблемы с местом на диске
- Не можем искать за приемлемое время

Архитектура индексации в Apache Lucene

Слишком много сегментов

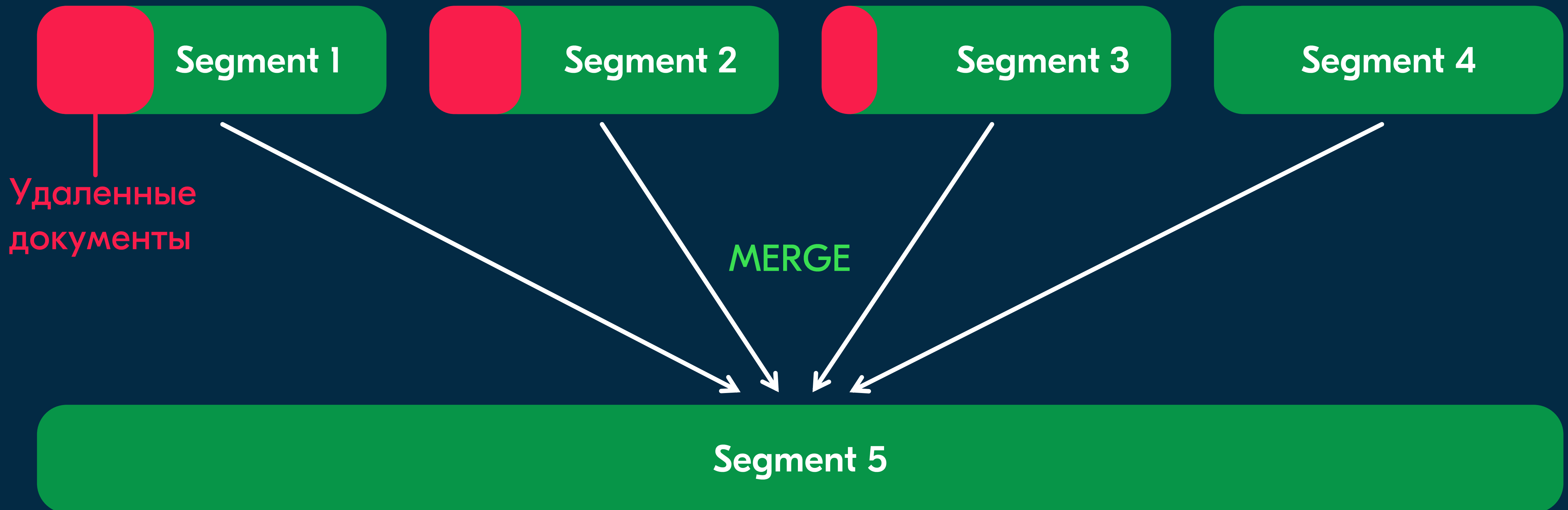
- Нужны операции `merge` и `compaction` как в `LSM-tree`

Архитектура индексации в Apache Lucene

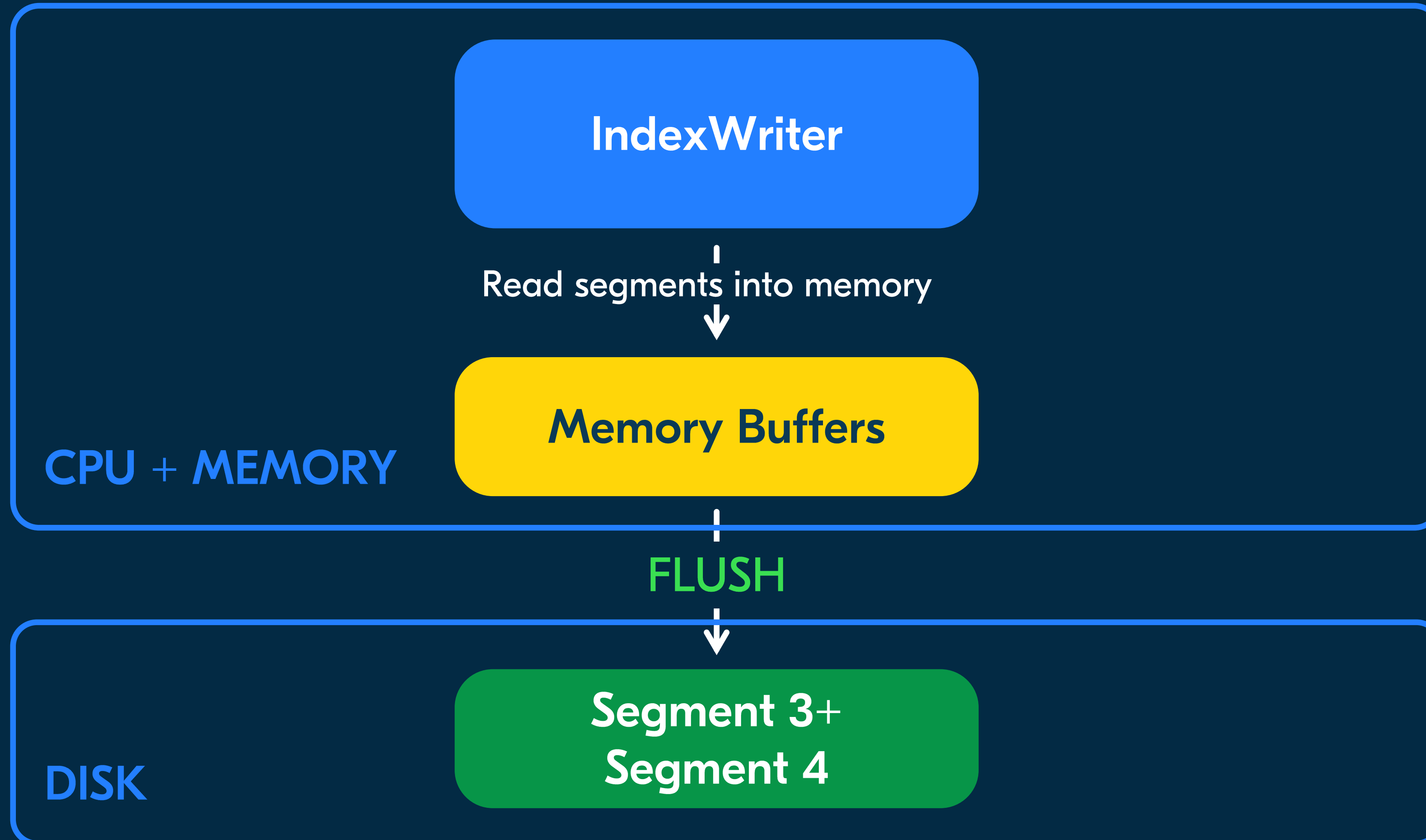


Архитектура индексации в Apache Lucene

Операция merge



Merge использует CPU, диск и память



Ресурсы

- CPU-bound-задача (INDEX, MERGE)
- Вторичные ресурсы: память и диск (MERGE, FLUSH)

Вертикальное масштабирование индексации

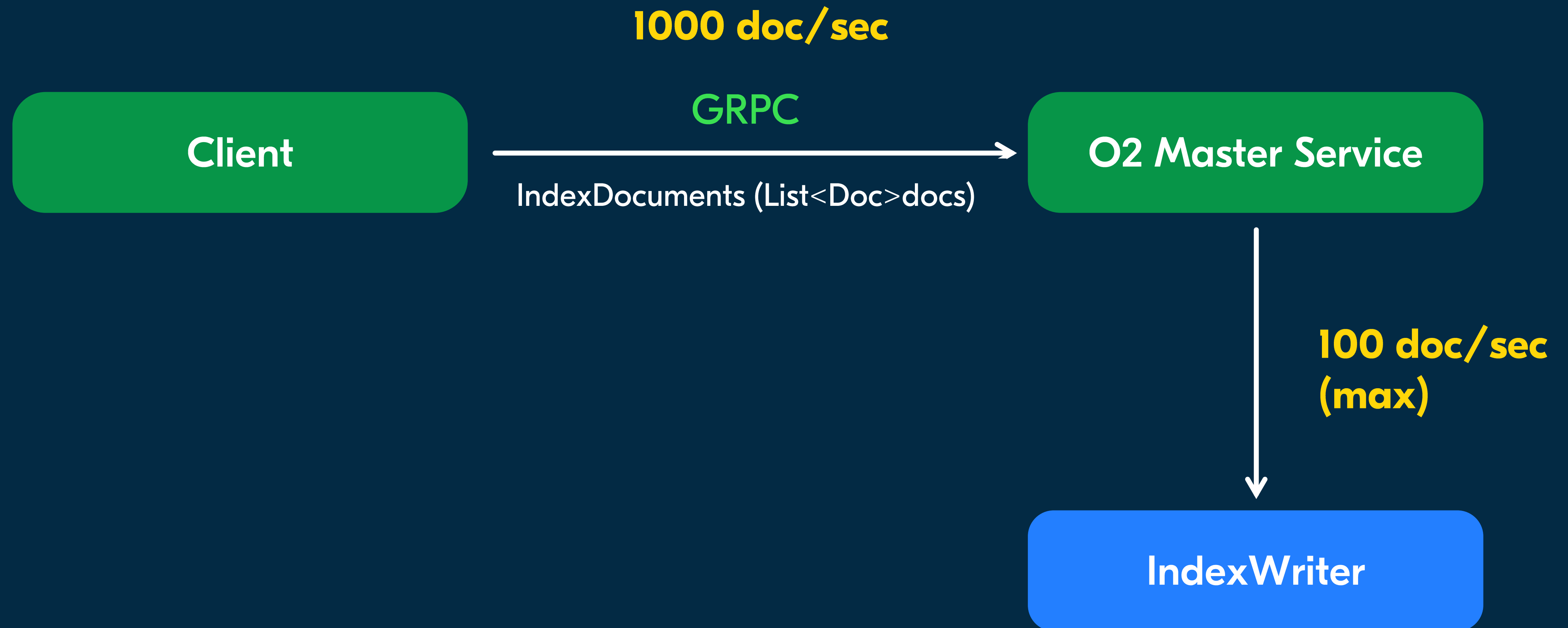
Архитектура кластера O2 Master (2021)

Первая версия архитектуры



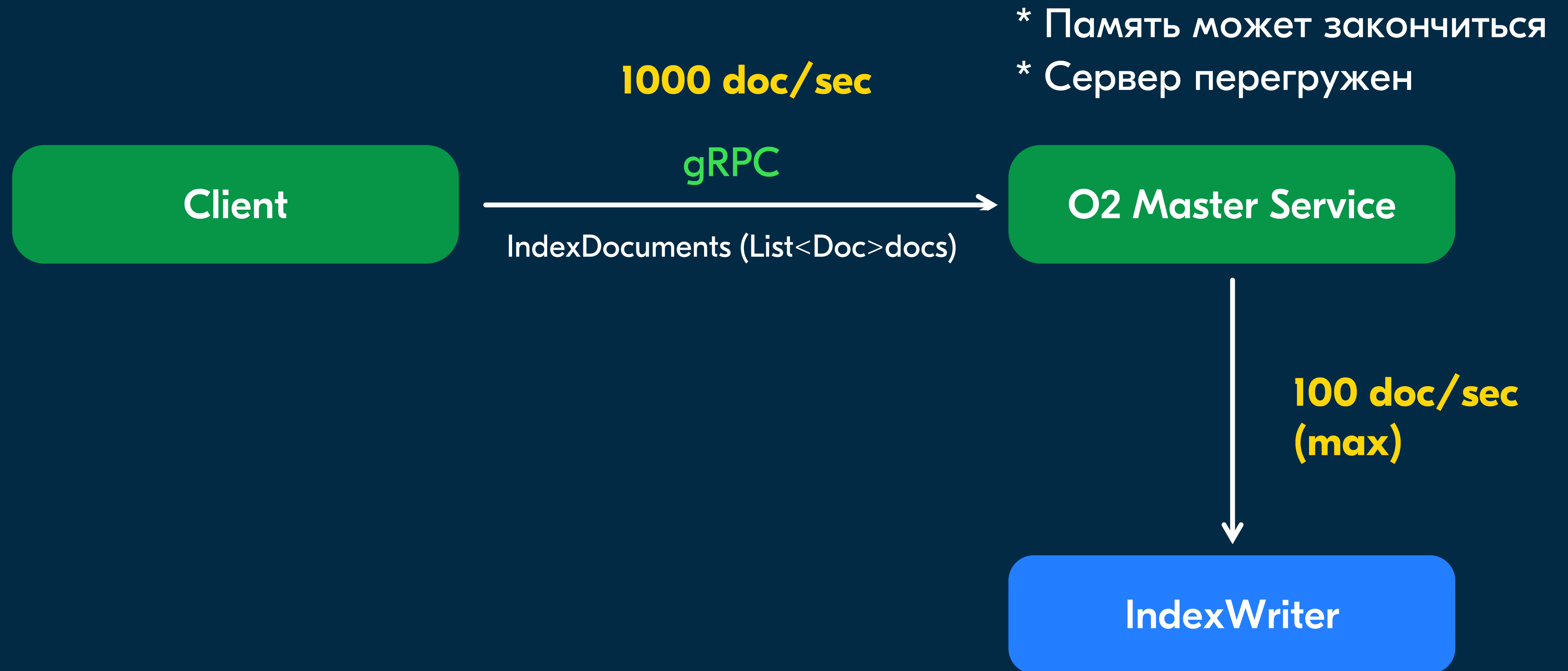
Архитектура кластера O2 Master (2021)

Первая версия архитектуры



Архитектура кластера O2 Master (2021)

Первая версия архитектуры



Архитектура кластера O2 Master (2021)

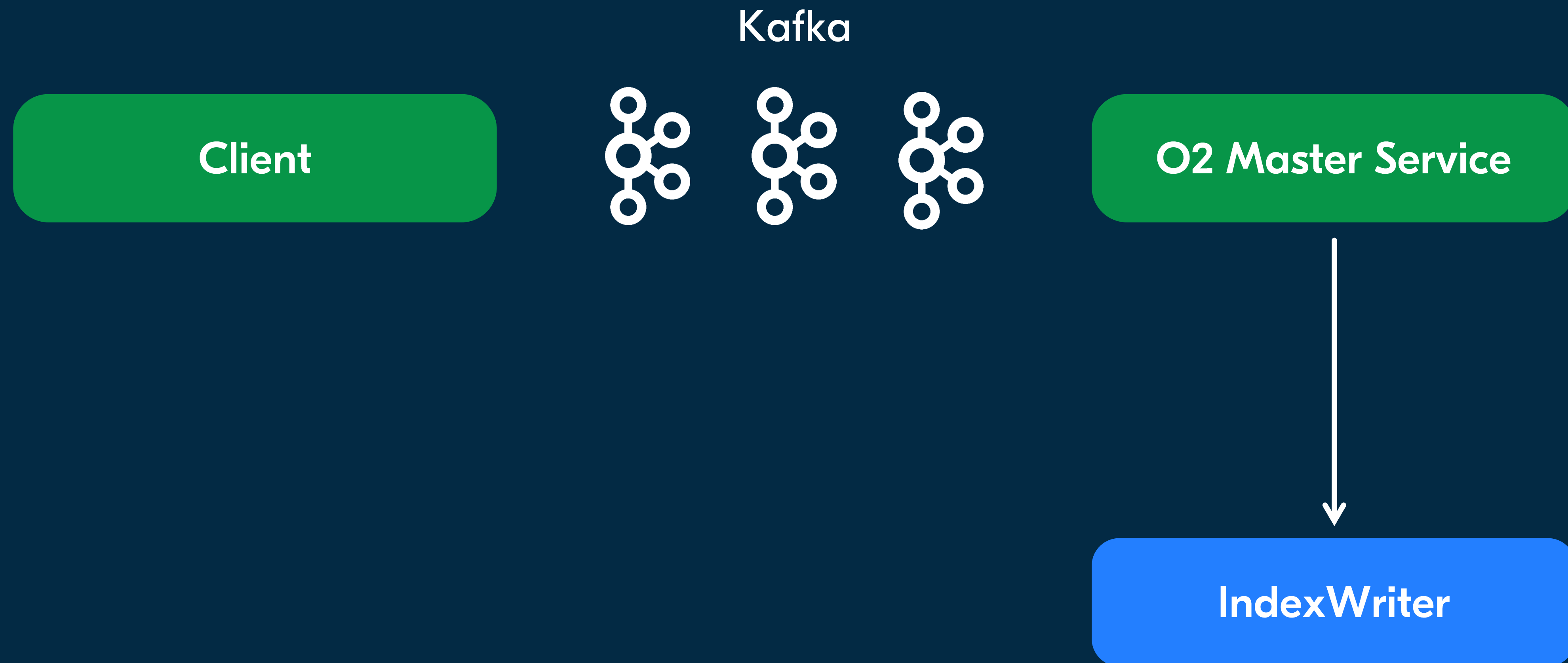
Первая версия архитектуры



Скорость индексации
не превышала **1000 doc/sec**

Архитектура кластера O2 Master (2021)

Вторая версия архитектуры



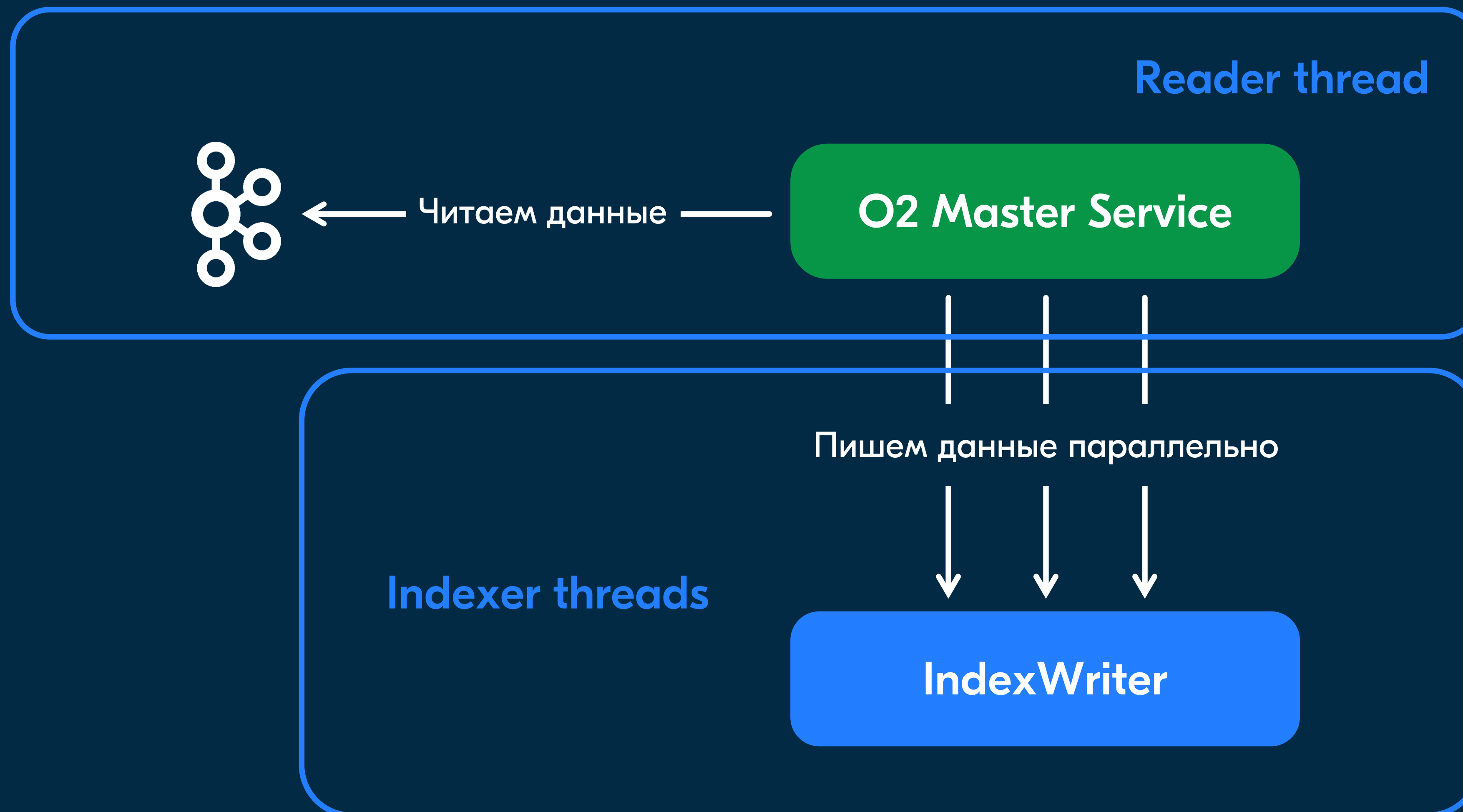
Архитектура кластера O2 Master (2021)

Вторая версия архитектуры

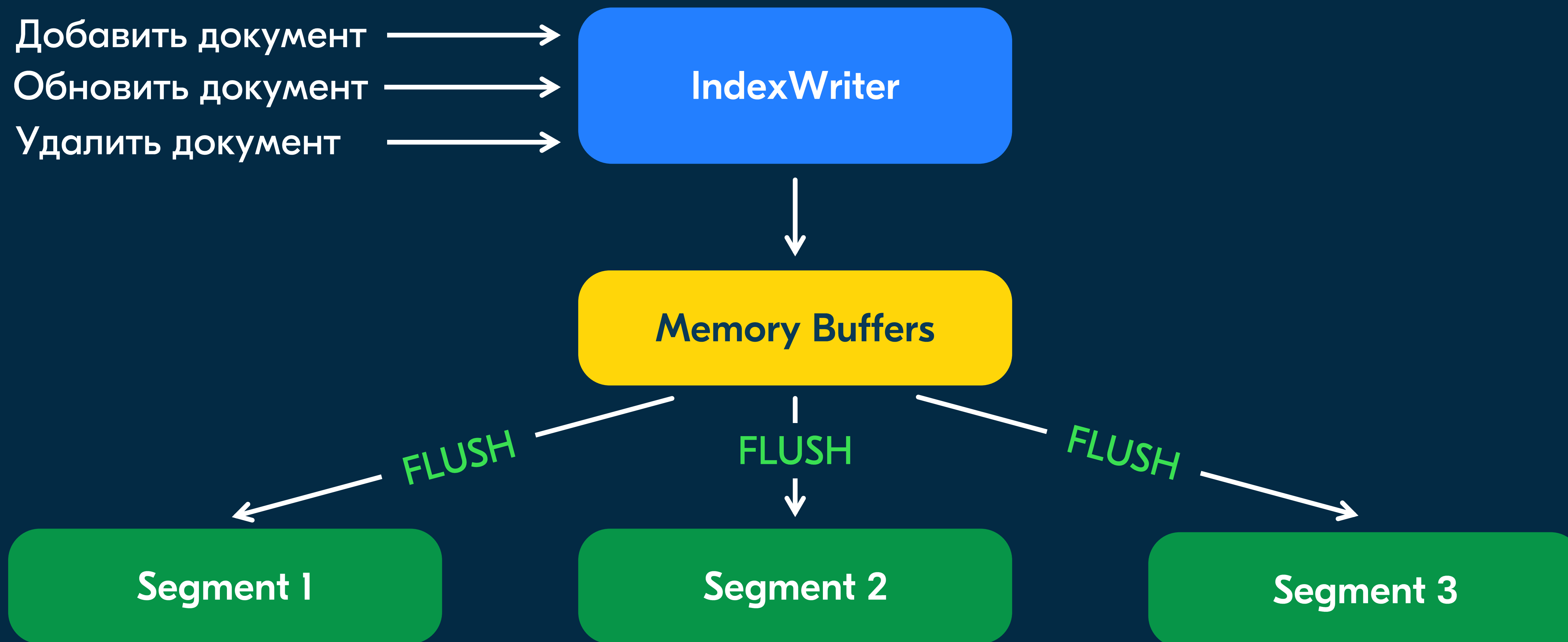
- Во второй версии заменили API-вставки на очередь (Kafka)
- Стало лучше:
 - Клиент и сервер больше не связаны напрямую
 - Скорость индексации выросла до 1500 doc/sec, и мы видели дальнейший потенциал

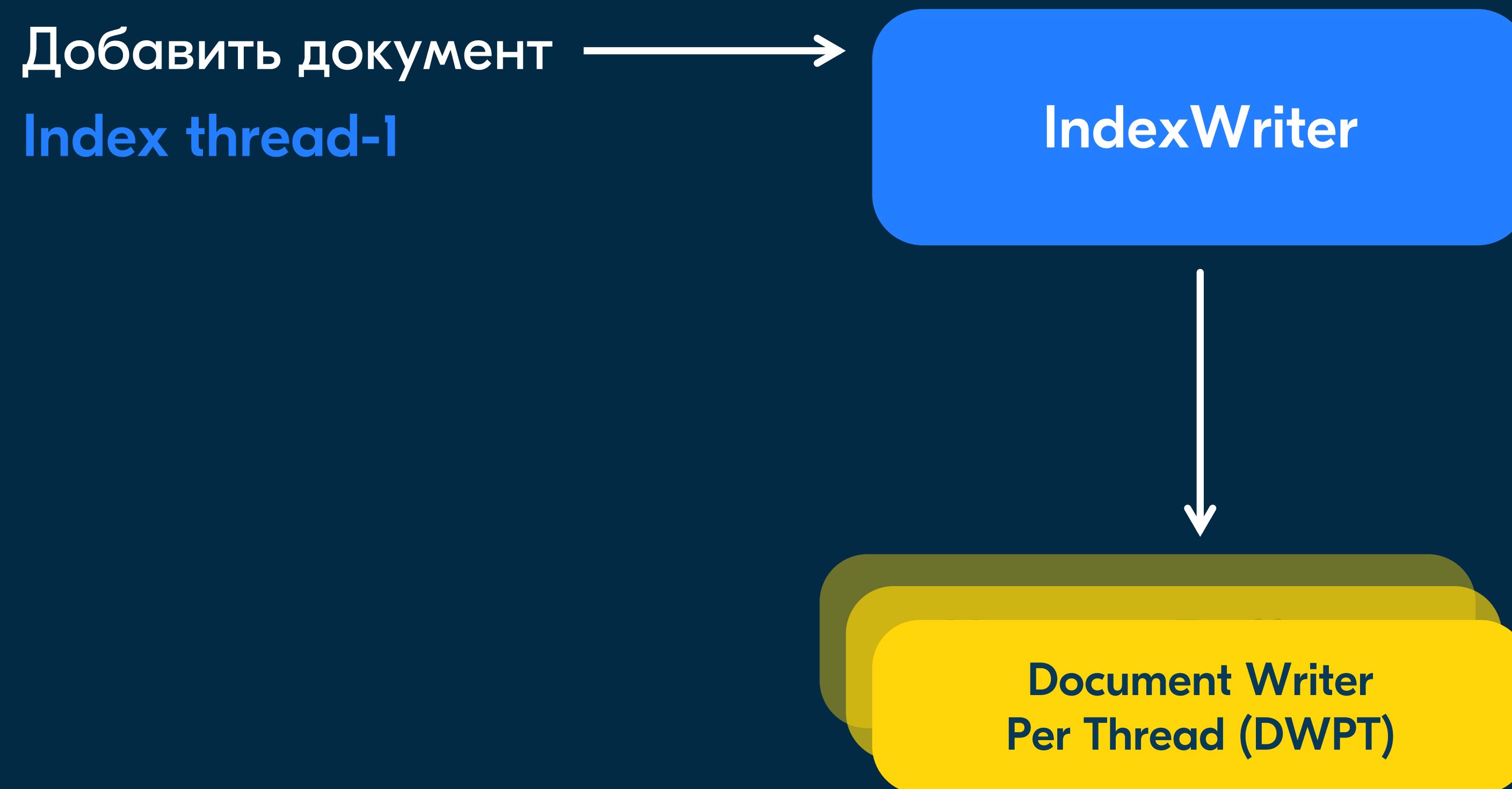
Архитектура кластера O2 Master (2021)

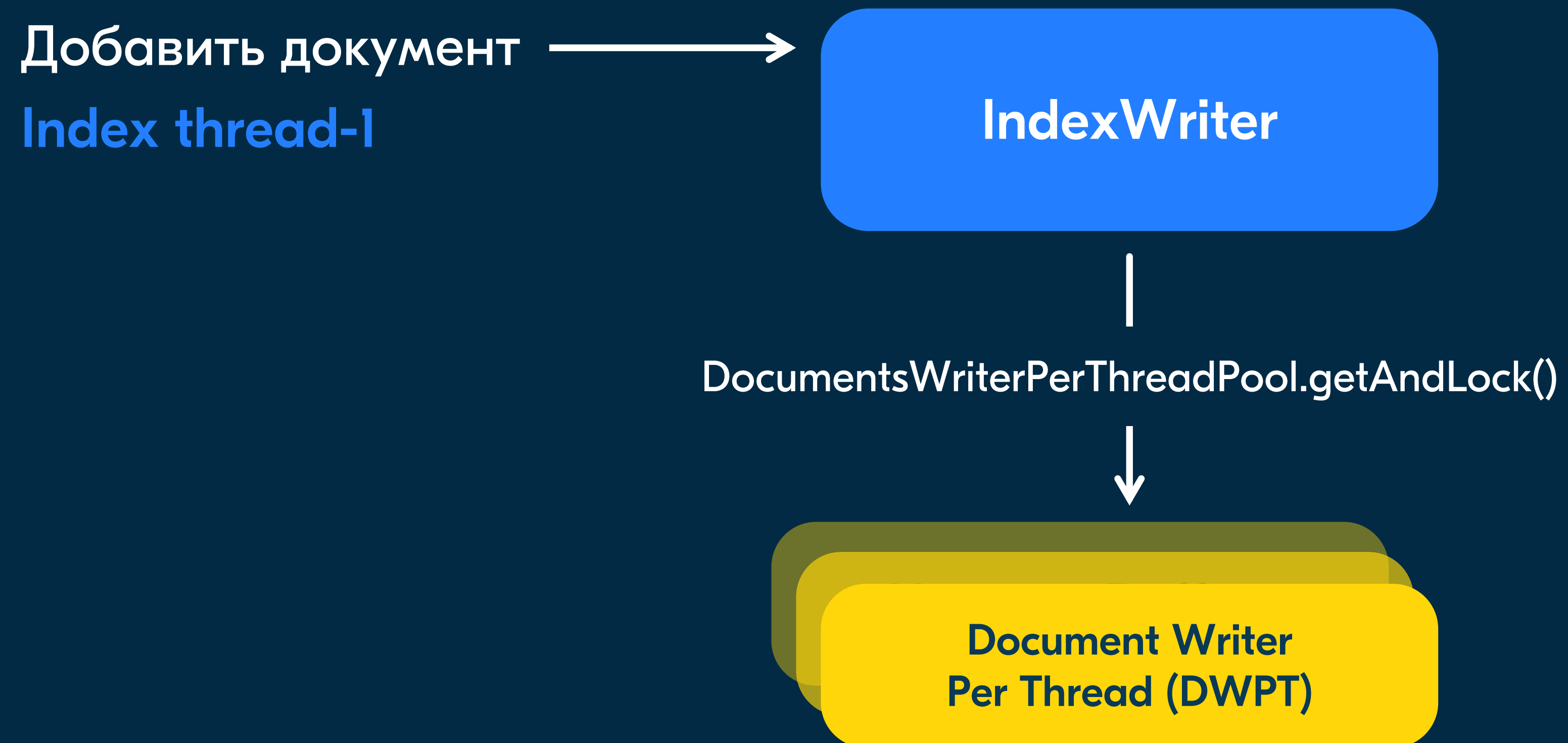
Вертикальное масштабирование индексации

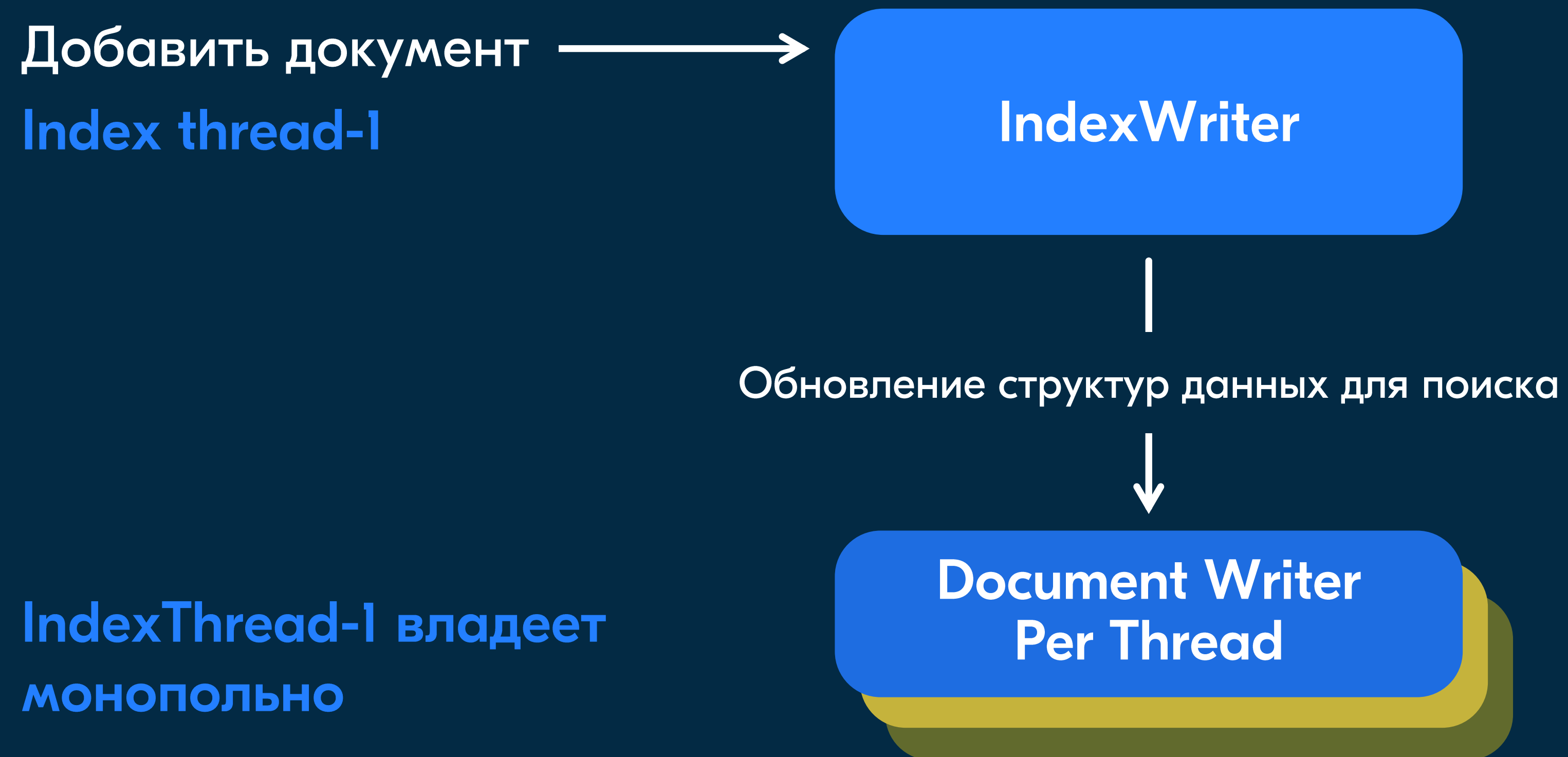


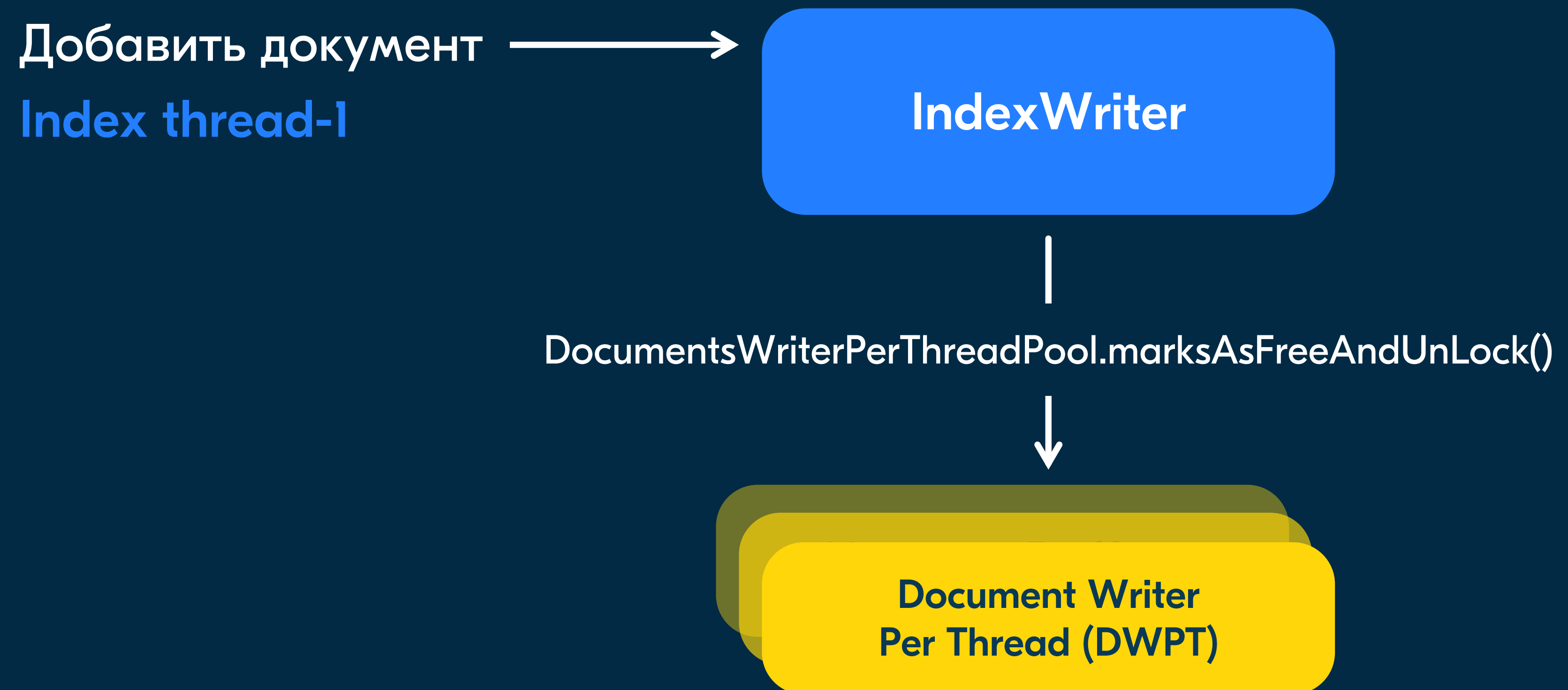
Архитектура индексации в Apache Lucene











Архитектура кластера O2 Master (2021)

Вертикальное масштабирование индексации

- **Итоги:**
 - Создаем пул потоков, равный количеству ядер
 - Получаем линейный рост целевой метрики?



Не все так просто

- Добавили ядер/потоков, а **throughput не растёт** 😞

Давайте посмотрим в профайлер...

```
[ 6] java.util.concurrent.locks.ReentrantLock.lock
[ 7] org.apache.lucene.index.DocumentsWriterFlushControl.markForFullFlush
[ 8] org.apache.lucene.index.DocumentsWriter.flushAllThreads
[ 9] org.apache.lucene.index.IndexWriter.getReader
[10] org.apache.lucene.index.StandardDirectoryReader.doOpenFromWriter
[11] org.apache.lucene.index.StandardDirectoryReader.doOpenIfChanged
[12] org.apache.lucene.index.StandardDirectoryReader.doOpenIfChanged
[13] org.apache.lucene.index.DirectoryReader.openIfChanged
[14] org.apache.lucene.search.SearcherManager.refreshIfNeeded
[15] org.apache.lucene.search.SearcherManager.refreshIfNeeded
[16] org.apache.lucene.search.ReferenceManager.doMaybeRefresh
[17] org.apache.lucene.search.ReferenceManager.maybeRefreshBlocking
[18] ru.ozon.search.o2.index.controller.O2IndexWriter.refreshAndAcquireSearcher
[19] ru.ozon.search.o2.index.controller.O2IndexWriter.search
[20] ru.ozon.search.o2.index.controller.O2IndexWriter.documentsCount
```

Понятнее не стало 😞

```
[ 6] java.util.concurrent.locks.ReentrantLock.lock
[ 7] org.apache.lucene.index.DocumentsWriterFlushControl.markForFullFlush
[ 8] org.apache.lucene.index.DocumentsWriter.flushAllThreads
[ 9] org.apache.lucene.index.IndexWriter.getReader
[10] org.apache.lucene.index.StandardDirectoryReader.doOpenFromWriter
[11] org.apache.lucene.index.StandardDirectoryReader.doOpenIfChanged
[12] org.apache.lucene.index.StandardDirectoryReader.doOpenIfChanged
[13] org.apache.lucene.index.DirectoryReader.openIfChanged
[14] org.apache.lucene.search.SearcherManager.refreshIfNeeded
[15] org.apache.lucene.search.SearcherManager.refreshIfNeeded
[16] org.apache.lucene.search.ReferenceManager.doMaybeRefresh
[17] org.apache.lucene.search.ReferenceManager.maybeRefreshBlocking
[18] ru.ozon.search.o2.index.controller.O2IndexWriter.refreshAndAcquireSearcher
[19] ru.ozon.search.o2.index.controller.O2IndexWriter.search
[20] ru.ozon.search.o2.index.controller.O2IndexWriter.documentsCount
```

Что странно?

- [6] `java.util.concurrent.locks.ReentrantLock.lock`
- [7] `org.apache.lucene.index.DocumentsWriterFlushControl.markForFullFlush`
- [8] `org.apache.lucene.index.DocumentsWriter.flushAllThreads`
- [9] `org.apache.lucene.index.IndexWriter.getReader`
- [10] `org.apache.lucene.index.StandardDirectoryReader.doOpenFromWriter`
- [11] `org.apache.lucene.index.StandardDirectoryReader.doOpenIfChanged`
- [12] `org.apache.lucene.index.StandardDirectoryReader.doOpenIfChanged`
- [13] `org.apache.lucene.index.DirectoryReader.openIfChanged`
- [14] `org.apache.lucene.search.SearcherManager.refreshIfNeeded`
- [15] `org.apache.lucene.search.SearcherManager.refreshIfNeeded`
- [16] `org.apache.lucene.search.ReferenceManager.doMaybeRefresh`
- [17] `org.apache.lucene.search.ReferenceManager.maybeRefreshBlocking`
- [18] `ru.ozon.search.o2.index.controller.O2IndexWriter.refreshAndAcquireSearcher`
- [19] `ru.ozon.search.o2.index.controller.O2IndexWriter.search`
- [20] `ru.ozon.search.o2.index.controller.O2IndexWriter.documentsCount`

Мы тут говорим про запись документов.
А в профайлере блокировка торчит из поиска?

[6] [java.util.concurrent.locks.ReentrantLock.lock](#)

[7] [org.apache.lucene.index.DocumentsWriterFlushControl.markForFullFlush](#)

[8] [org.apache.lucene.index.DocumentsWriter.flushAllThreads](#)

[9] [org.apache.lucene.index.IndexWriter.getReader](#)

[10] [org.apache.lucene.index.StandardDirectoryReader.doOpenFromWriter](#)

[11] [org.apache.lucene.index.StandardDirectoryReader.doOpenIfChanged](#)

[12] [org.apache.lucene.index.StandardDirectoryReader.doOpenIfChanged](#)

[13] [org.apache.lucene.index.DirectoryReader.openIfChanged](#)

[14] [org.apache.lucene.search.SearcherManager.refreshIfNeeded](#)

[15] [org.apache.lucene.search.SearcherManager.refreshIfNeeded](#)

[16] [org.apache.lucene.search.ReferenceManager.doMaybeRefresh](#)

[17] [org.apache.lucene.search.ReferenceManager.maybeRefreshBlocking](#)

[18] [ru.ozon.search.o2.index.controller.O2IndexWriter.refreshAndAcquireSearcher](#)

[19] [ru.ozon.search.o2.index.controller.O2IndexWriter.search](#)

73 [20] [ru.ozon.search.o2.index.controller.O2IndexWriter.documentsCount](#)

Поиск мешает индексации

- Проблема в методе **IndexWriter.getWriter**
- Метод может вызывать операцию flush all, которая сбрасывает все данные из памяти на диск
- Захватывает ту же блокировку, что и индексация (flush control)

Реализация O2 Master

Поиск мешает индексации

- Убрали все поисковые вызовы
- Скорость индексации достигла
~2500-3000 doc/sec



Выбор GC

- Какой GC выбрать?
- Может, ZGC?

Выбор GC

ZGC (JDK 14)

- `-XX:+UnlockExperimentalVMOptions -XX:+UseZGC`



Выбор GC

ZGC (JDK 14)

- И однажды...



Выбор GC

ZGC (JDK 14)

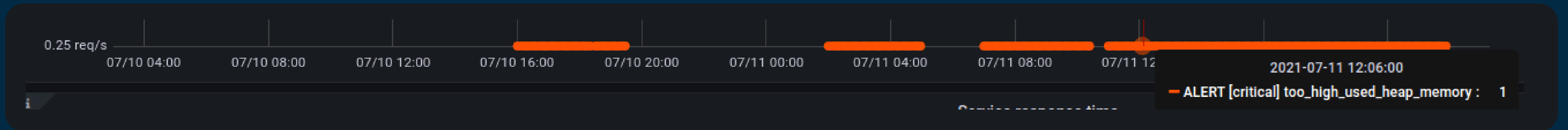
- Сервису стало совсем плохо 😞
- RT вырос **до 5 секунд** (max query timeout) в p90



Реализация O2 Master

ZGC (JDK 14)

- Сработал алерт на размер занятого хипа



Реализация O2 Master

ZGC (JDK 14)

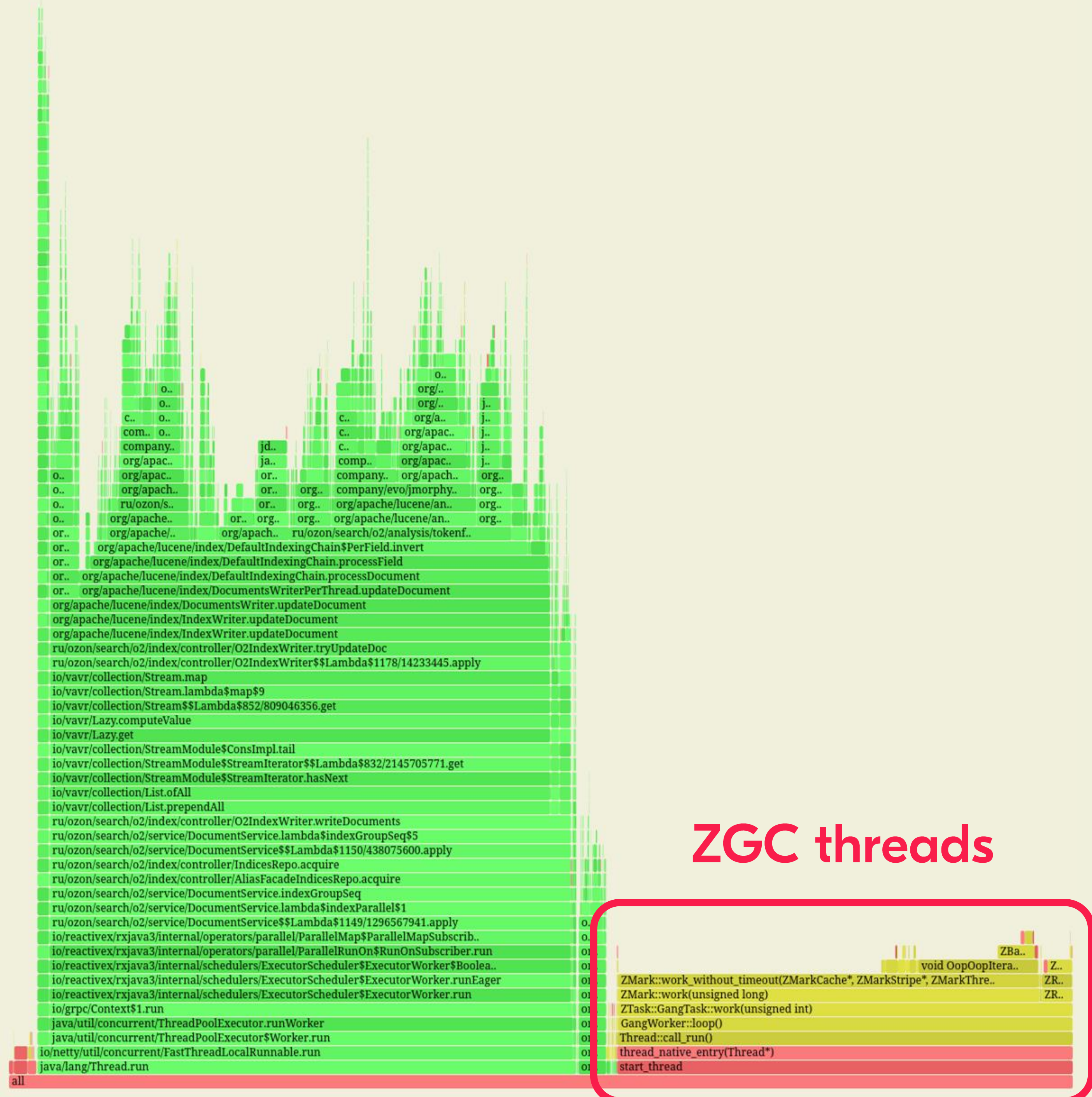
- Давайте посмотрим **в профайлер...**



Реализация O2 Master

ZGC (JDK 14)

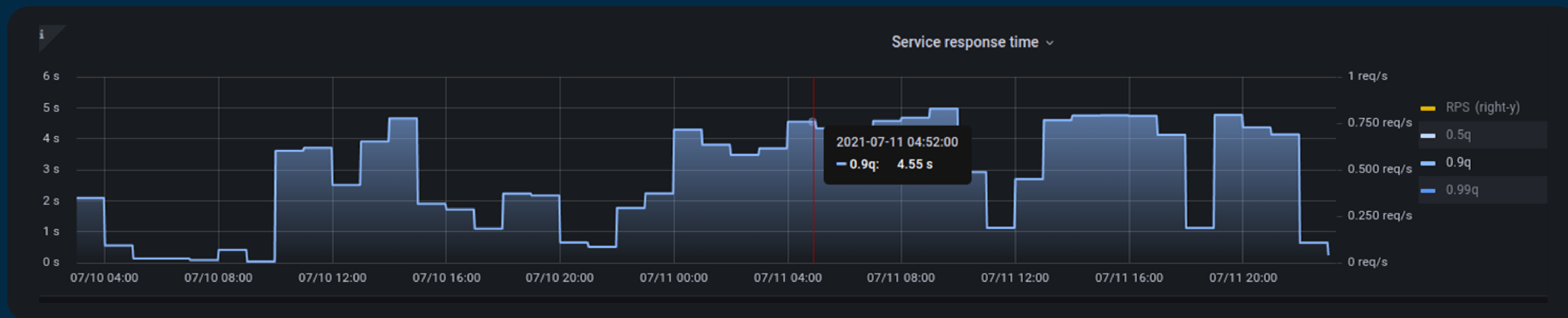
- Половину ядер занял ZGC (а мы CPU bound)



Выбор GC

ZGC (JDK 14)

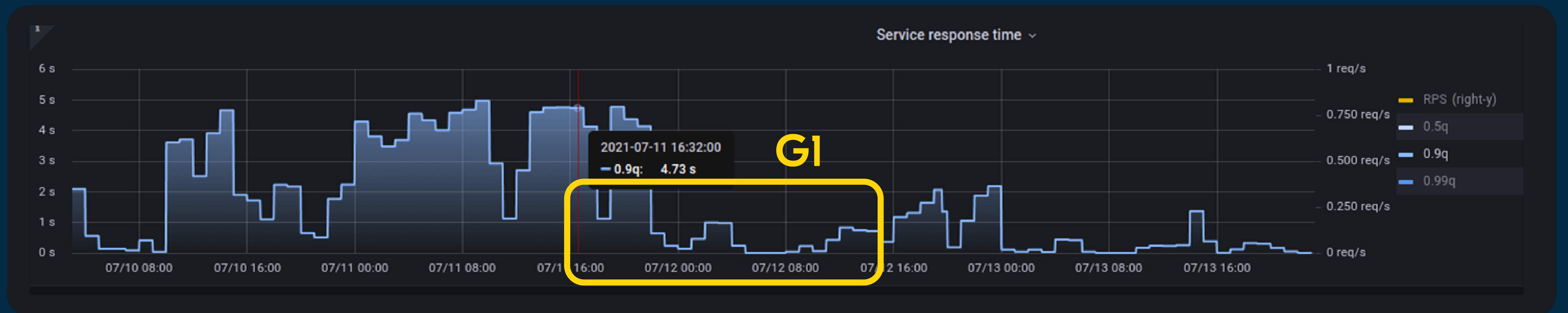
- В нашем случае ZGC занимал до 15-30% CPU 😞
- В какой-то момент GC перестал успевать очищать хип



Выбор GC

GI

- Переключились на GI



А вы пробовали
обновить java?



Выбор GC

Ситуация на сегодня



Прошло два года,
а что сегодня?

Выбор GC

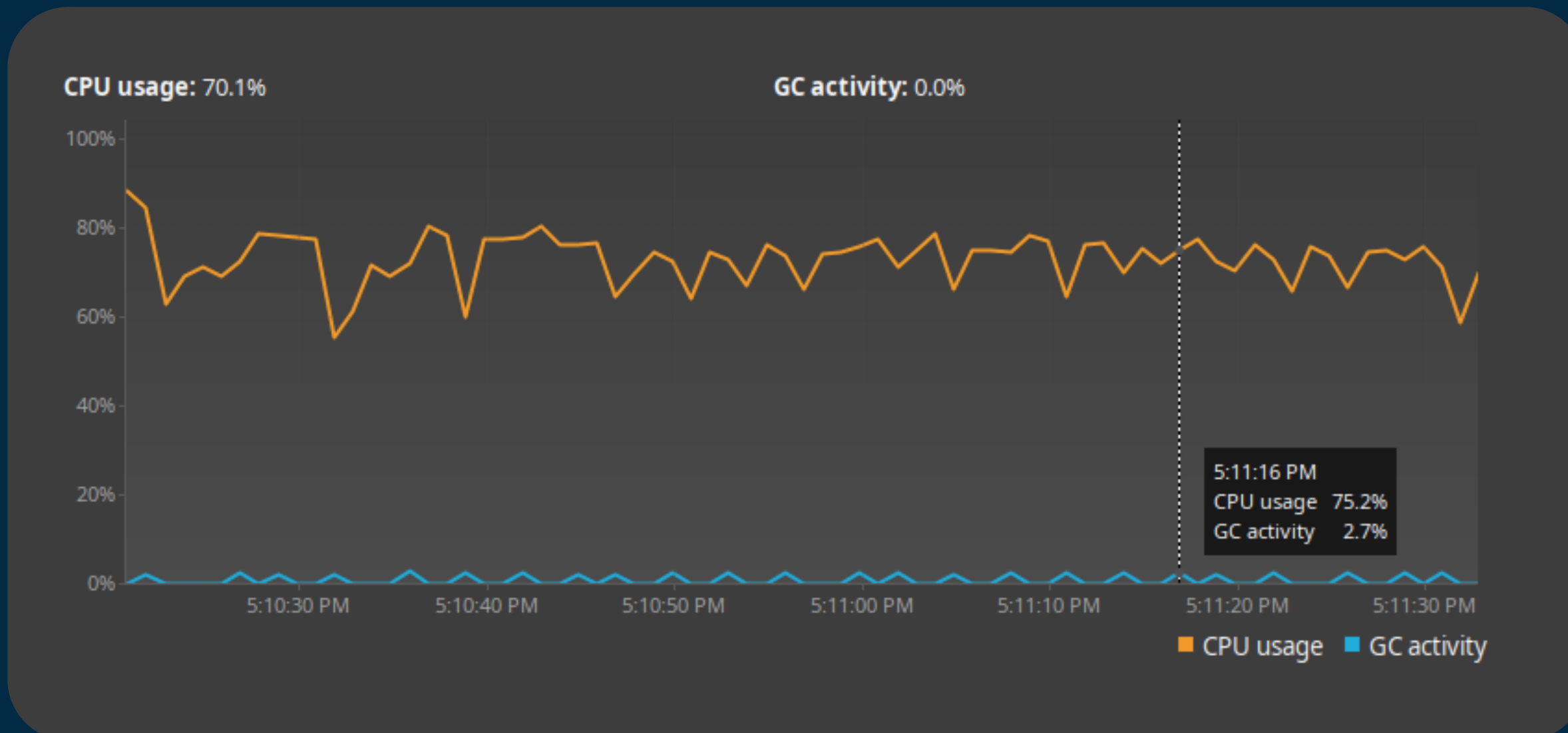
GI vs ZGC (JDK 17)

- Тест: приложение индексации документов
- Реальные данные
- 42 cores, 800 gb RAM
- Настройки GC дефолтные, -Xmx80g

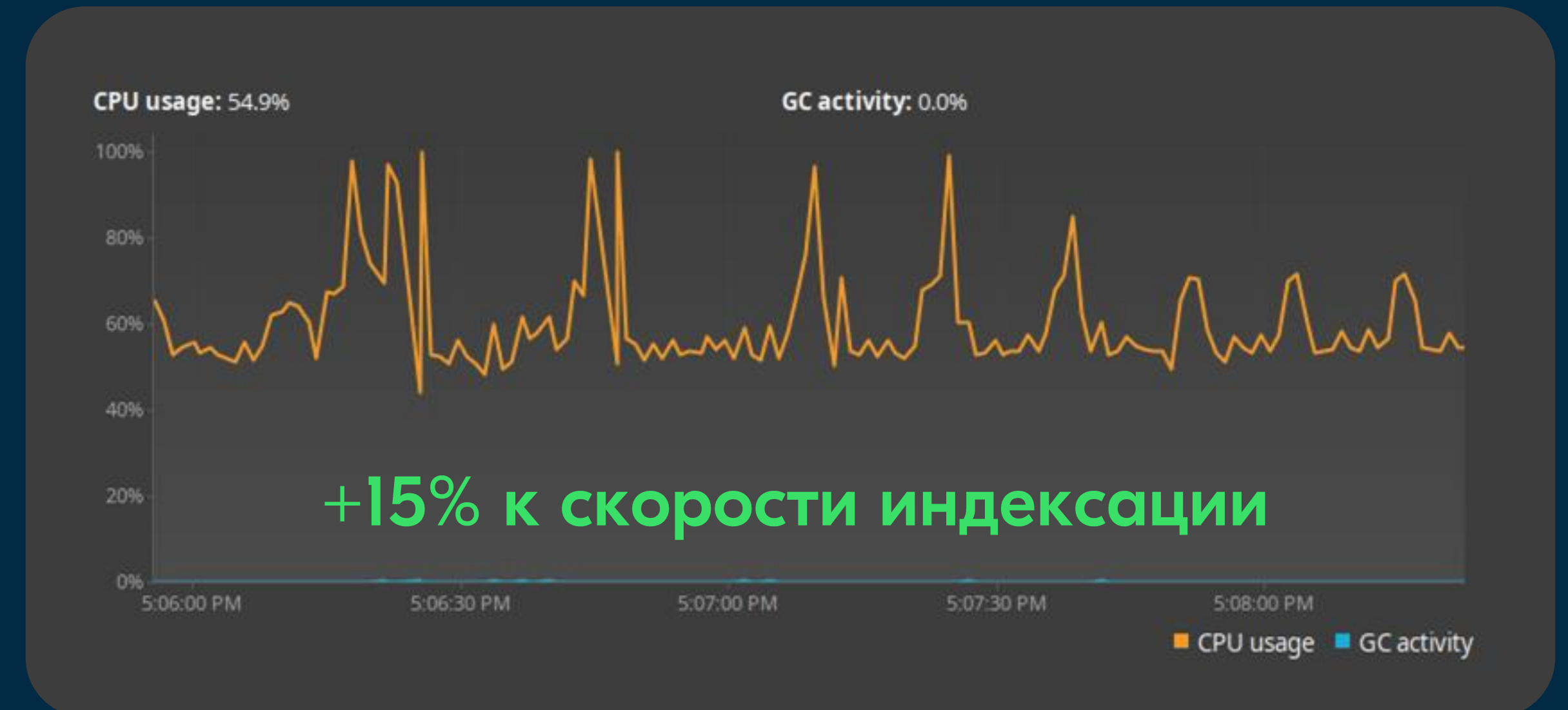
Выбор GC

GI vs ZGC (JDK 17)

- Нагрузка и данные одинаковы
- У ZGC размер пауз **меньше на порядок** за счет более частых сборок
- ZGC использует больше CPU



ZGC



GI

Выводы

- В каждом приложении необходимо выбирать GC через анализ метрик
- **ZGC** — хороший выбор для приложений, где важен latency
- **G1** все еще неплох для batch-нагрузки (CPU-bound)



Потребление памяти



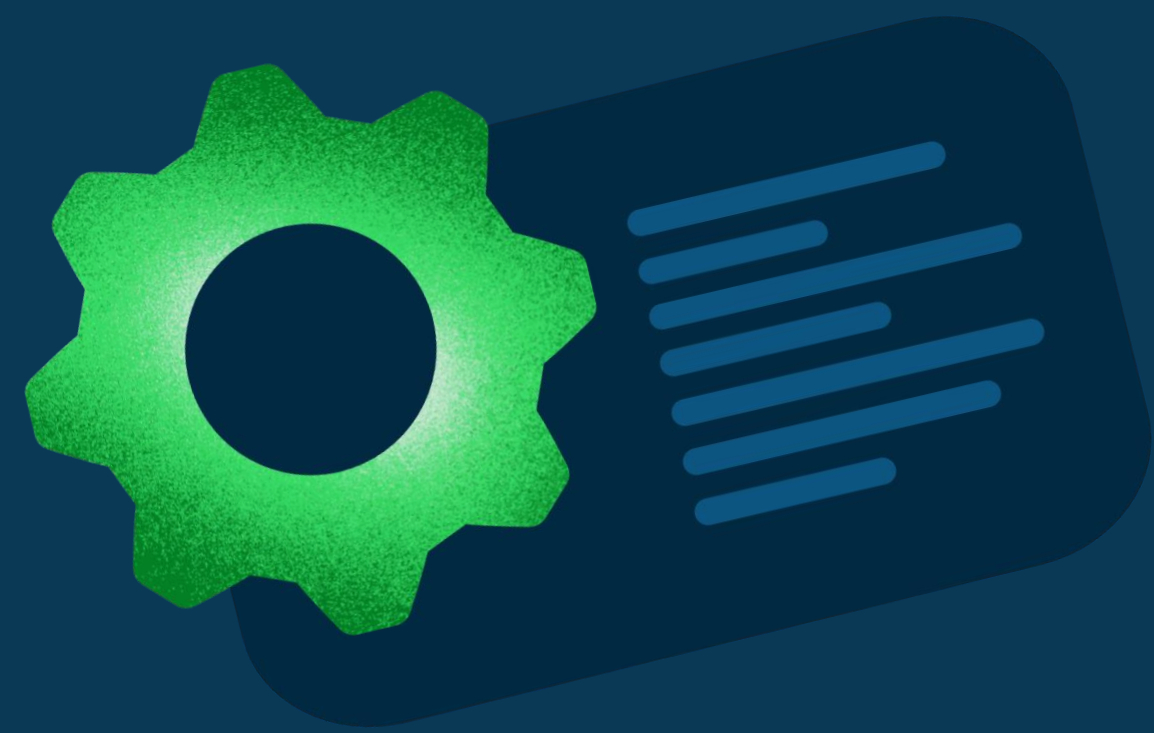
Сколько памяти нужно, чтобы
загрузить все потоки индексации?
(много)

Память, как ресурс

- На каждый документ надо иметь несколько мегабайт памяти
- Memory buffers для индексации
- Память под мержи
- Другие накладные расходы

Представление документа внутри

- bytes
- Protobuf
- `org.apache.lucene.document.Document` (объектное представление) 😞
- Конечные структуры данных внутри codec



В некоторых индексах у нас
500-1000 полей в документе
в среднем (!)

Как уменьшить потребление памяти java-приложения?

- Использовать примитивные коллекции (например, **fastutil**)
- Переиспользовать объекты
- Профилировать аллокации
- Использовать off-heap

Реализация O2 Master

Что используем мы?

JDK/guava

(больше в индексации)

fastutil

(больше в поиске)

Memory-mapped I/O

- Memory-mapped I/O (mmap) для неизменяемых данных

Memory-mapped I/O

- `org.apache.lucene.store.BufferedIndexInput`
- `one.nio.mem.MappedFile` + `Unsafe` [1]
- Можно работать с файлами > **2gb**

Project Panama

- Project Panama [1][2] (строим HNSW [3] на C++)

[1] <https://openjdk.org/projects/panama/>

[2] <https://denisgabaydulin.medium.com/first-look-into-panama-b1e016aa1a15>

[3] <https://arxiv.org/abs/1603.09320>

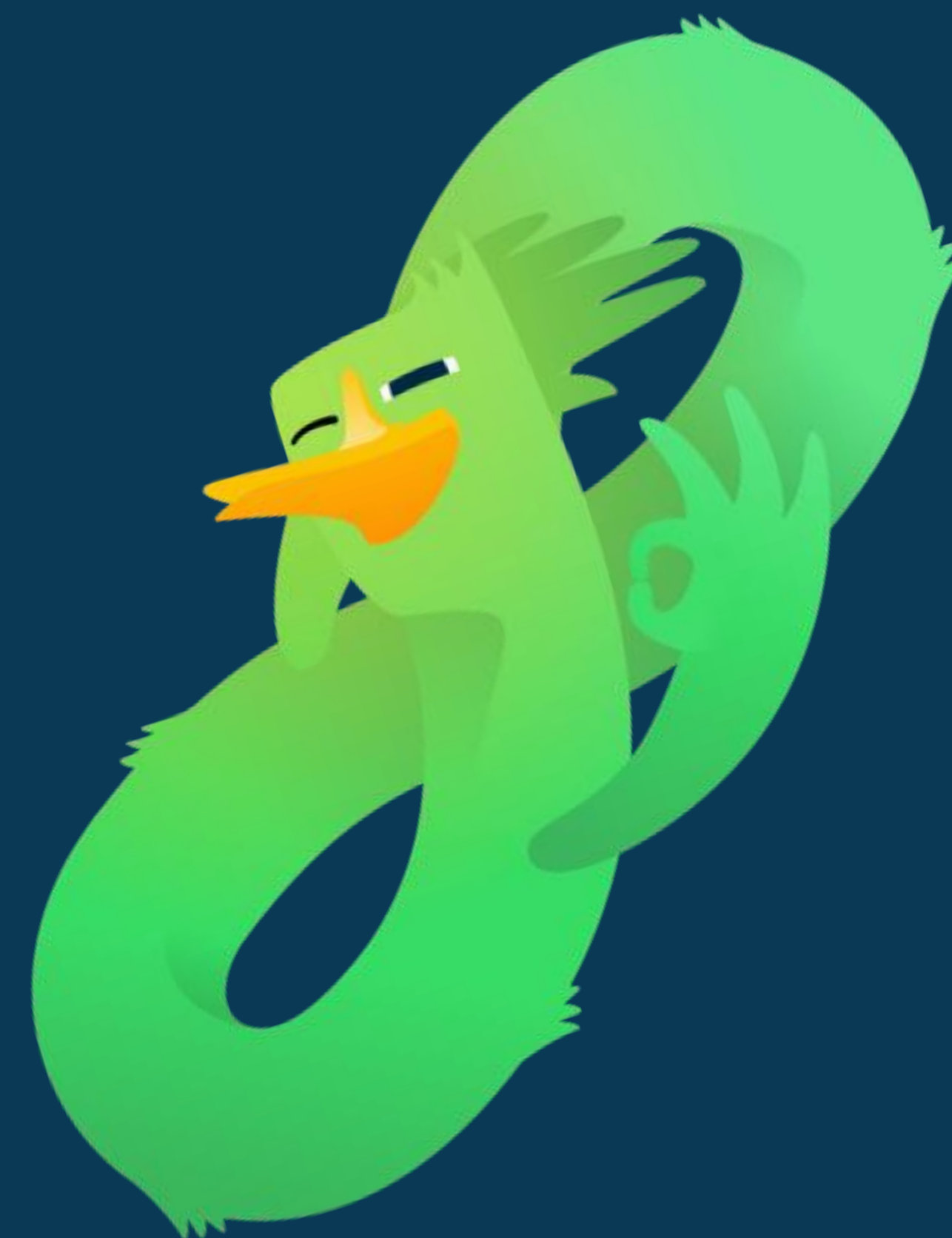
Project Panama

- Строим HNSW граф в hnswlib (C++)
- Копируем данные целиком и пишем в формате Lucene vector codec
- До **3x быстрее**, чем делает сам Lucene 9.2.0

Реализация O2 Master

Итоги

- Выбрали GC/heap size
- Уменьшили потребление памяти
- Смогли утилизировать **около 40 ядер**
- Скорость индексации достигла **~6000 doc/sec**



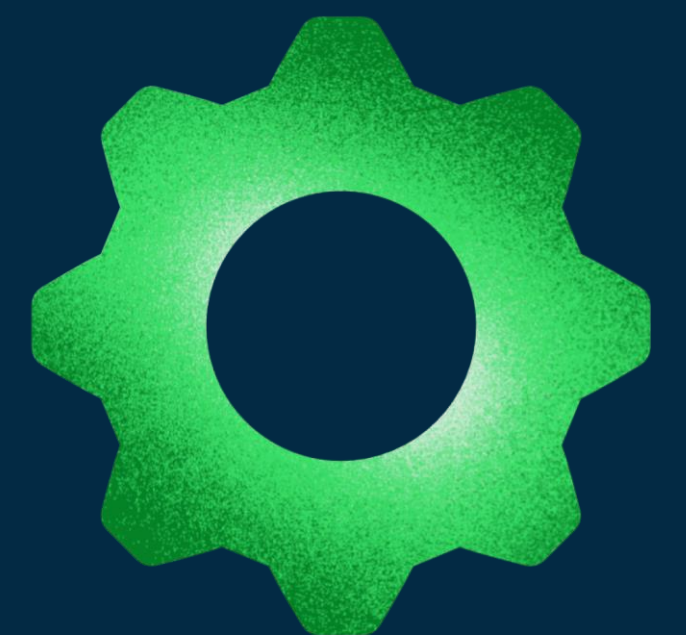
Индексация на максималках

Включили
индексацию
на максималках
~x2 ядер

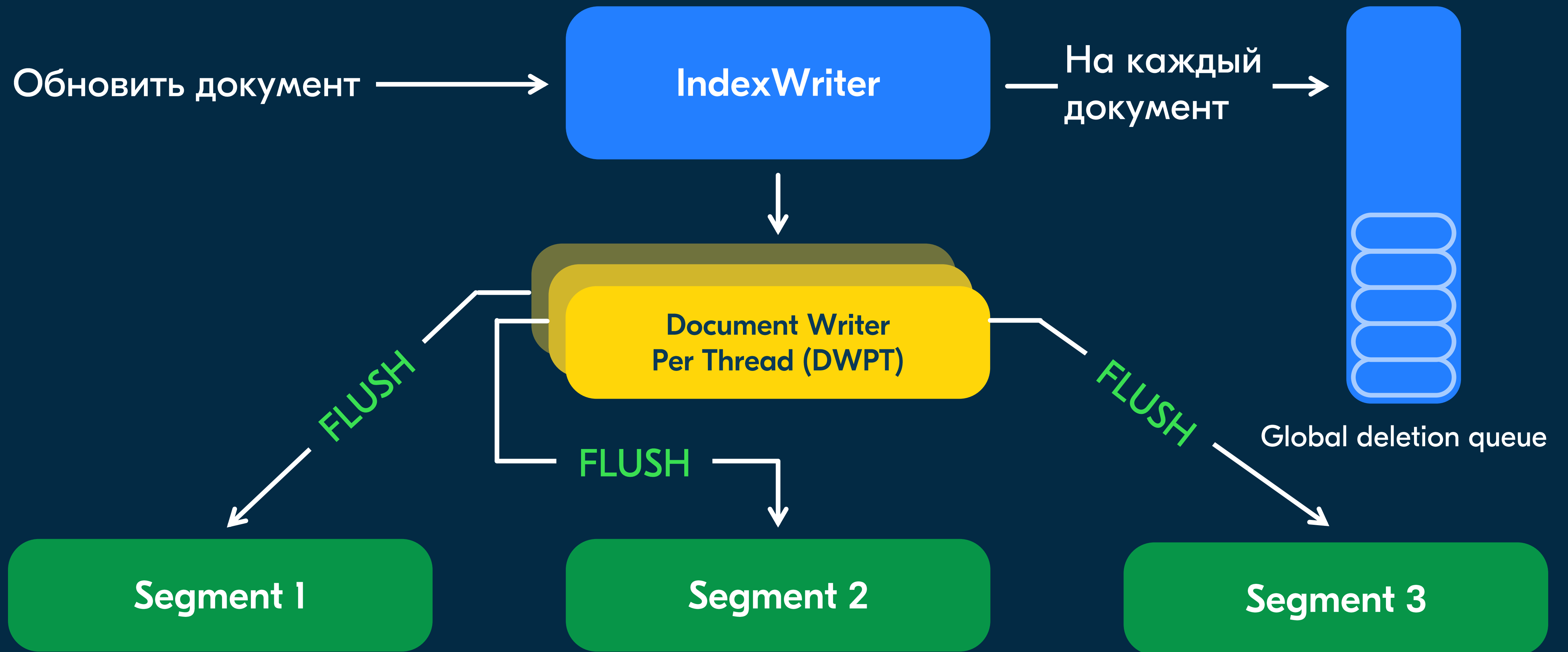
A throughput
вырос
НИЧТОЖНО 😞

Прирост меньше
15%

- Исторически у нас все вставки документов были, как апдейты
- Сигнатура метода:
updateDocument(Term id, Document doc)
- Update в Lucene: **insert + delete**,
только один документ

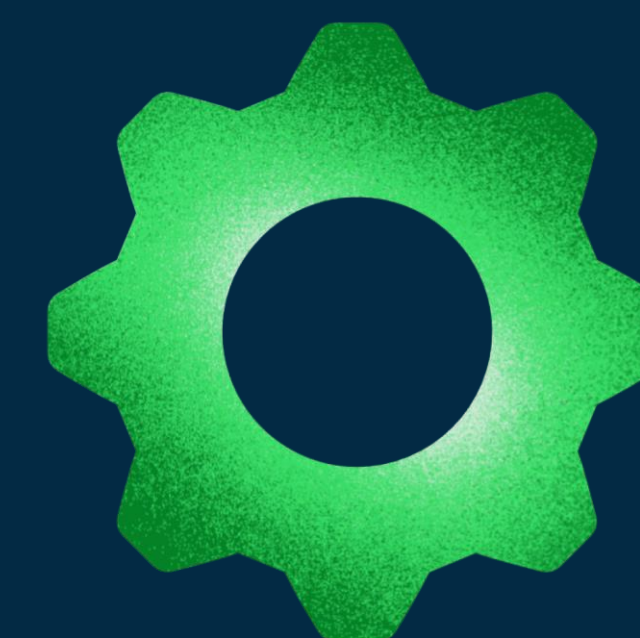


Архитектура индексации в Apache Lucene



Проблема апдейтов

- Сделали микробенчмарк
- **addDocument(Document doc)** работает на **20-30%** быстрее даже в один поток
- `profiler.sh -d 30 -e lock -o flamegraph=total`
- `profiler.sh -d 30 -e cpu -o flamegraph`



Проблема апдейтов #1

- Update имеет оверхед на каждую операцию
- Нет bulk-версии

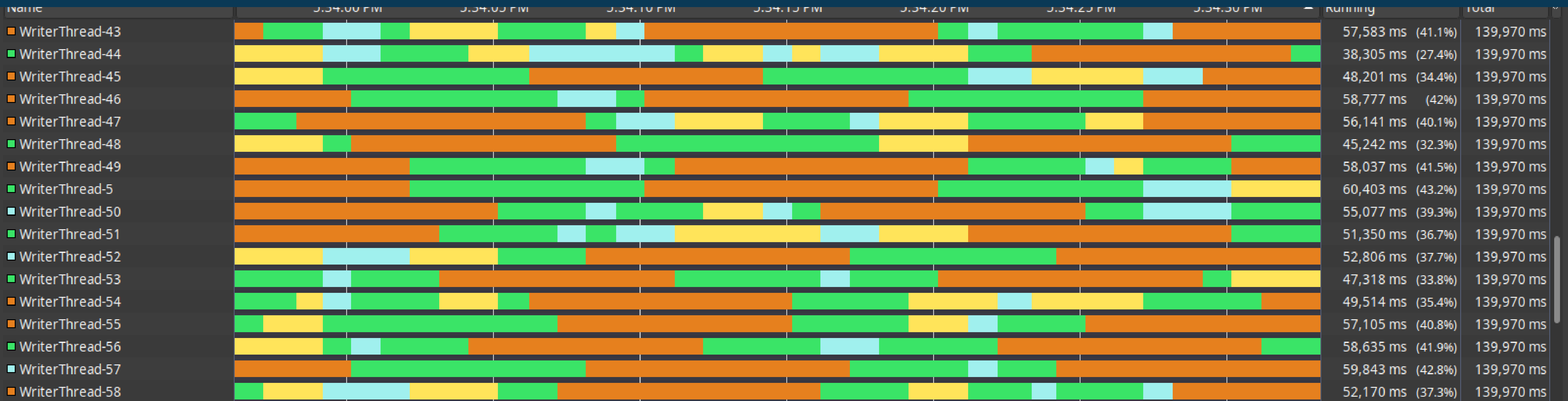
Проблема апдейтов #2

- При увеличении количества потоков flush становится узким местом в `updateDocument()`

Реализация O2 Master

Проблема апдейтов #2

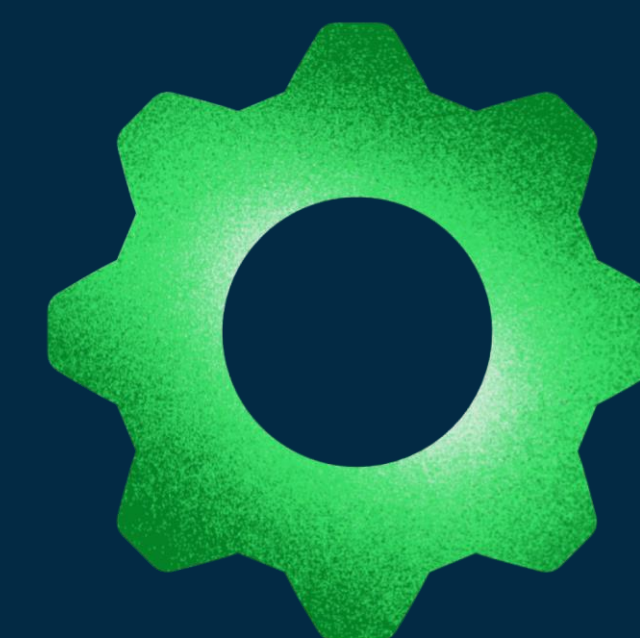
- Блокировки внутри потоков индексации на flush, выделены голубым



- Подробнее в отдельной статье [1]

Вставка документов вместо апдейтов

- Можно вставлять не по одному документу, а сразу пачкой, так еще быстрее
- Метод **`addDocuments(Iterable<Document> docs)`**
- Потребовалось сделать ручную дедупликацию документов



Реализация O2 Master

Дедупликация документов

Копии одного
документа



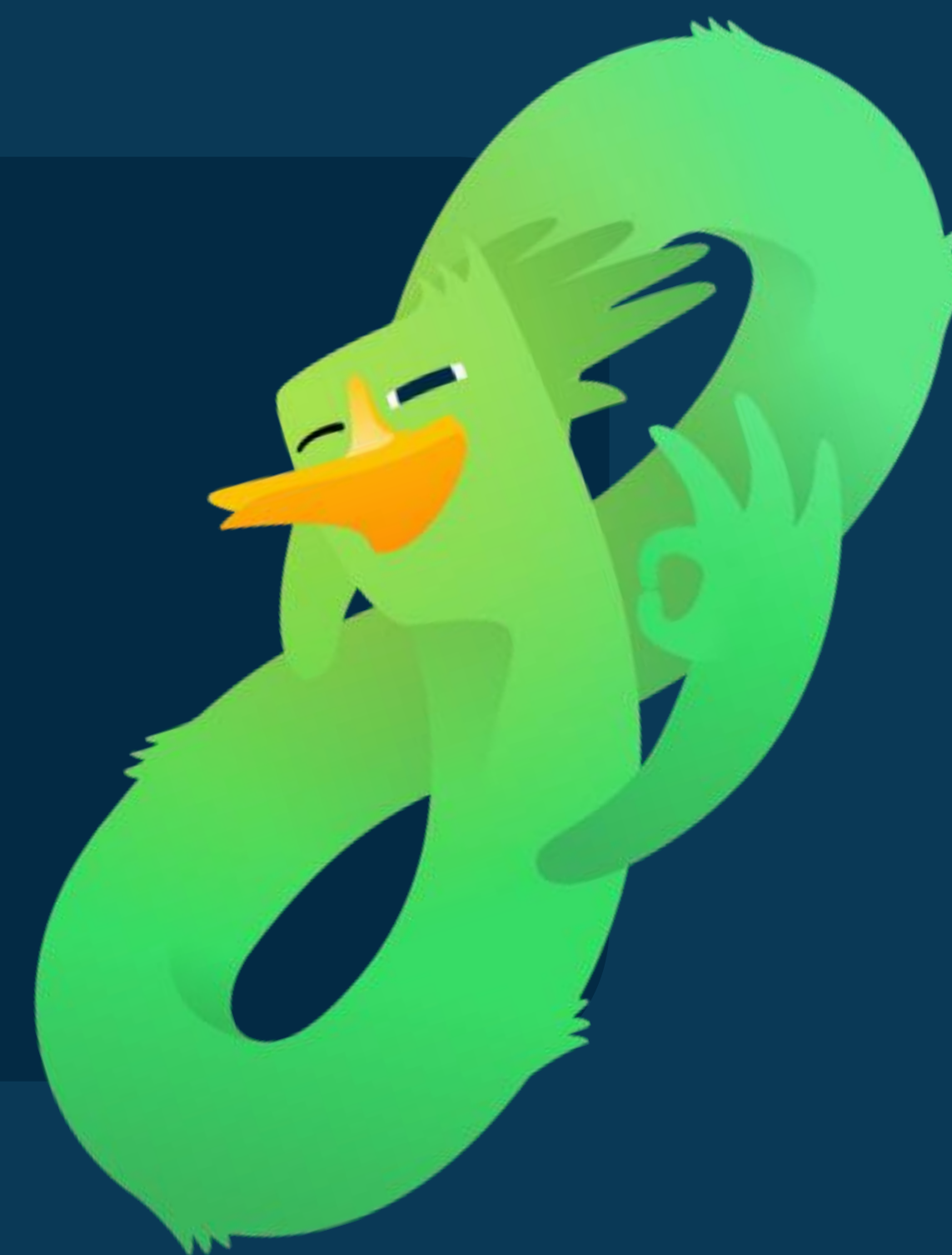
Deduplication

Удаленные
документы



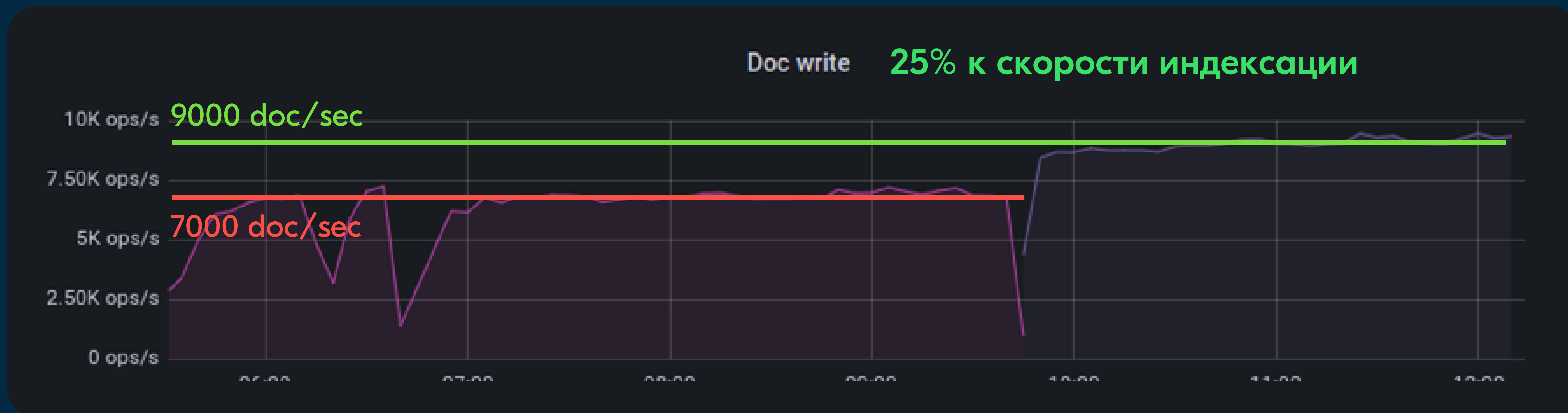
Итоги

Скорость индексации
достигла **~9000 doc/sec**



Реализация O2 Master

Результат замены апдейтов вставками



Архитектура кластера O2 Master

Результаты

- Утилизация ядер во время индексации
- Используется 70-80 ядер из 96



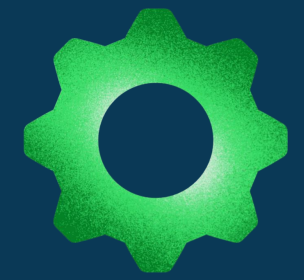
Реализация O2 Master

Вертикальное масштабирование индексации

- Спроектировали оптимальную архитектуру
- Разделили поиск и индексацию (**3000 doc/sec**)
- Снизили потребление памяти
- Подобрали GC/размер хипа (**6000 doc/sec**)
- Использовали bulk-операции, где это возможно (**9000 doc/sec**)

1. Вертикальное масштабирование дается **непросто**

2. Реализация так же важна, как архитектура



- **async-profiler**

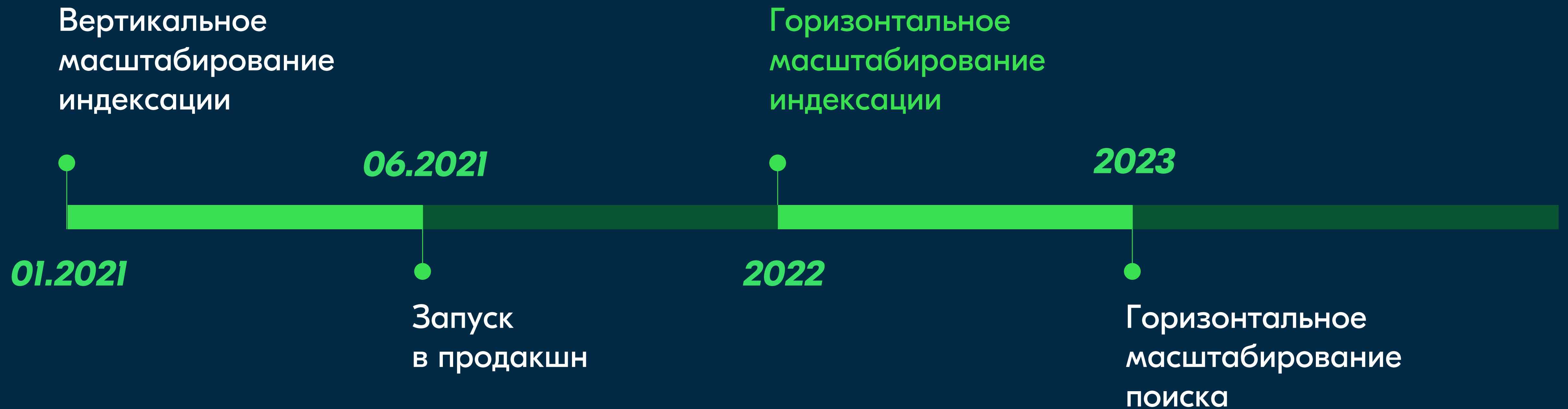
- **visual vm**

- **Java Microbenchmark Harness (JMH)**

- **Исходный код Lucene**

Поисковая платформа Ozon

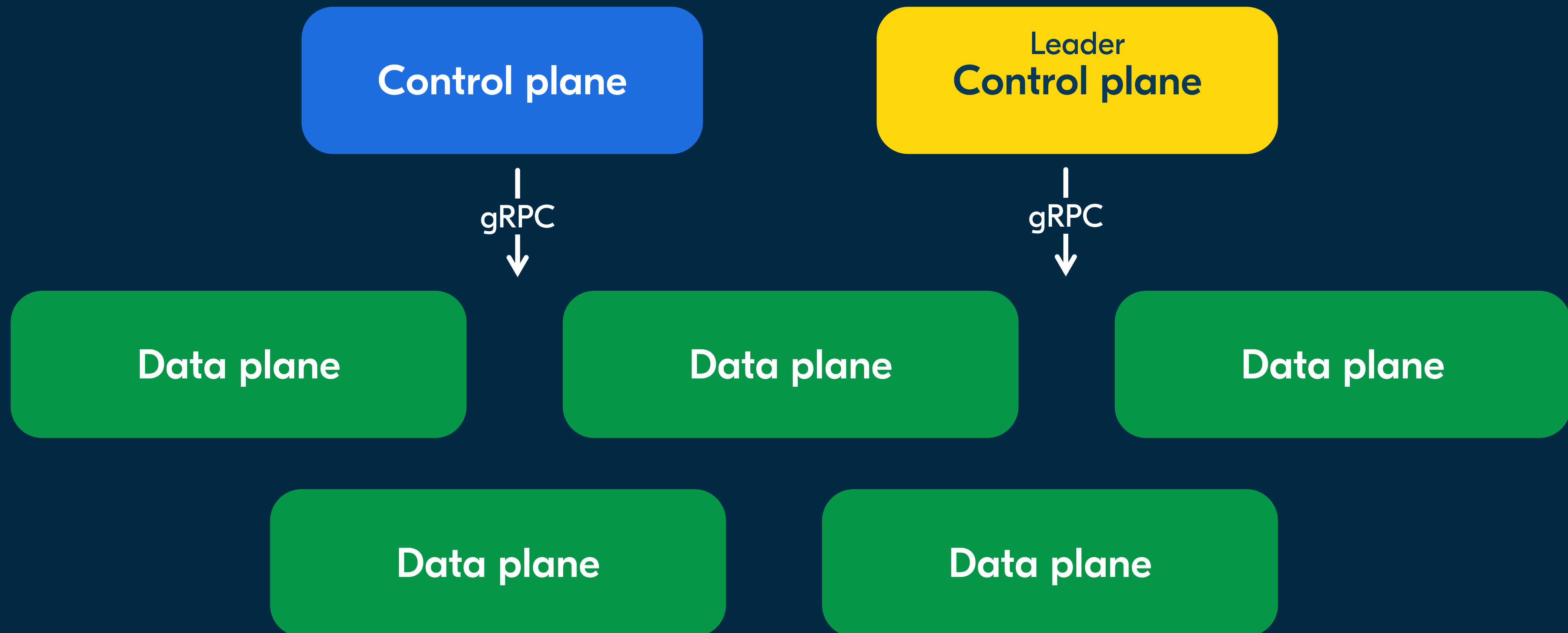
Таймлайн разработки движка O2



Горизонтальное масштабирование индексации (2022)

Архитектура кластера O2 Master (2022)

Горизонтальное масштабирование индексации



Горизонтальное масштабирование индексации

- Отдельный кластер в каждом ЦОД
- Два типа нод: control plane и data plane

Control plane — это API для клиента (достаточно двух-трех нод на кластер)

Data plane — это сервера индексации и публикации индекса (может быть несколько сотен нод в одном кластере)

Control plane

Горизонтальное масштабирование индексации

- Управляет метаданными индексов
- Предоставляет API для работы с индексами
- Выбирает лидера индекса
- Делает балансировку индексов по data plane-нодам
- Задает мэппинг индексов на реплики базового поиска
- Использует мало ресурсов (2-4 ядра, 8 GB RAM)

Горизонтальное масштабирование индексации

- Создает поисковые индексы
- Использует много ресурсов (80+ ядер, 700+ GB RAM)

Архитектура кластера O2 Master (2022)

Горизонтальное масштабирование индексации

Виртуальный индекс

Шард 1

Физический индекс Lucene

Шард 2

Физический индекс Lucene

Шард 3

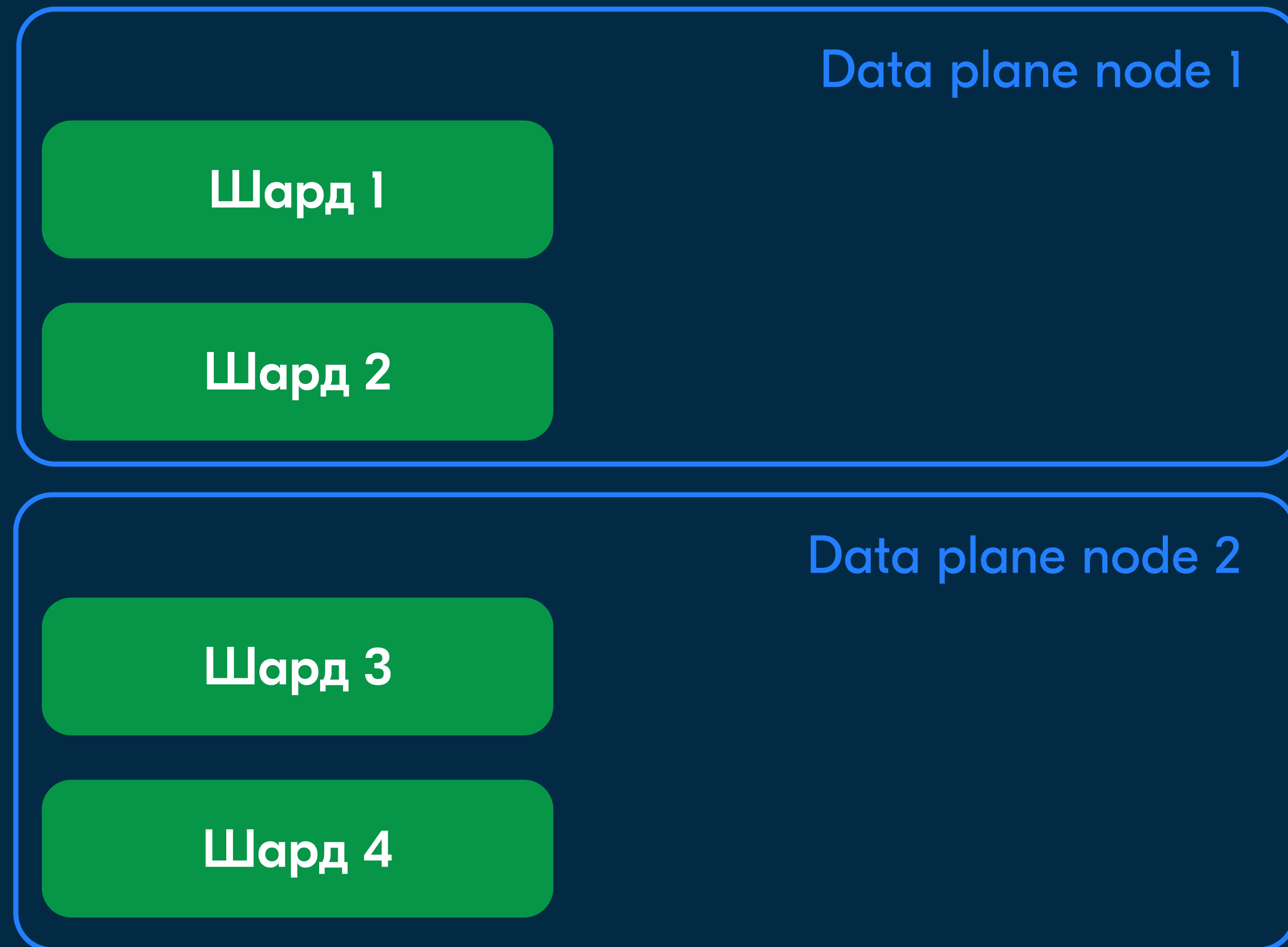
Физический индекс Lucene

Шард 4

Физический индекс Lucene

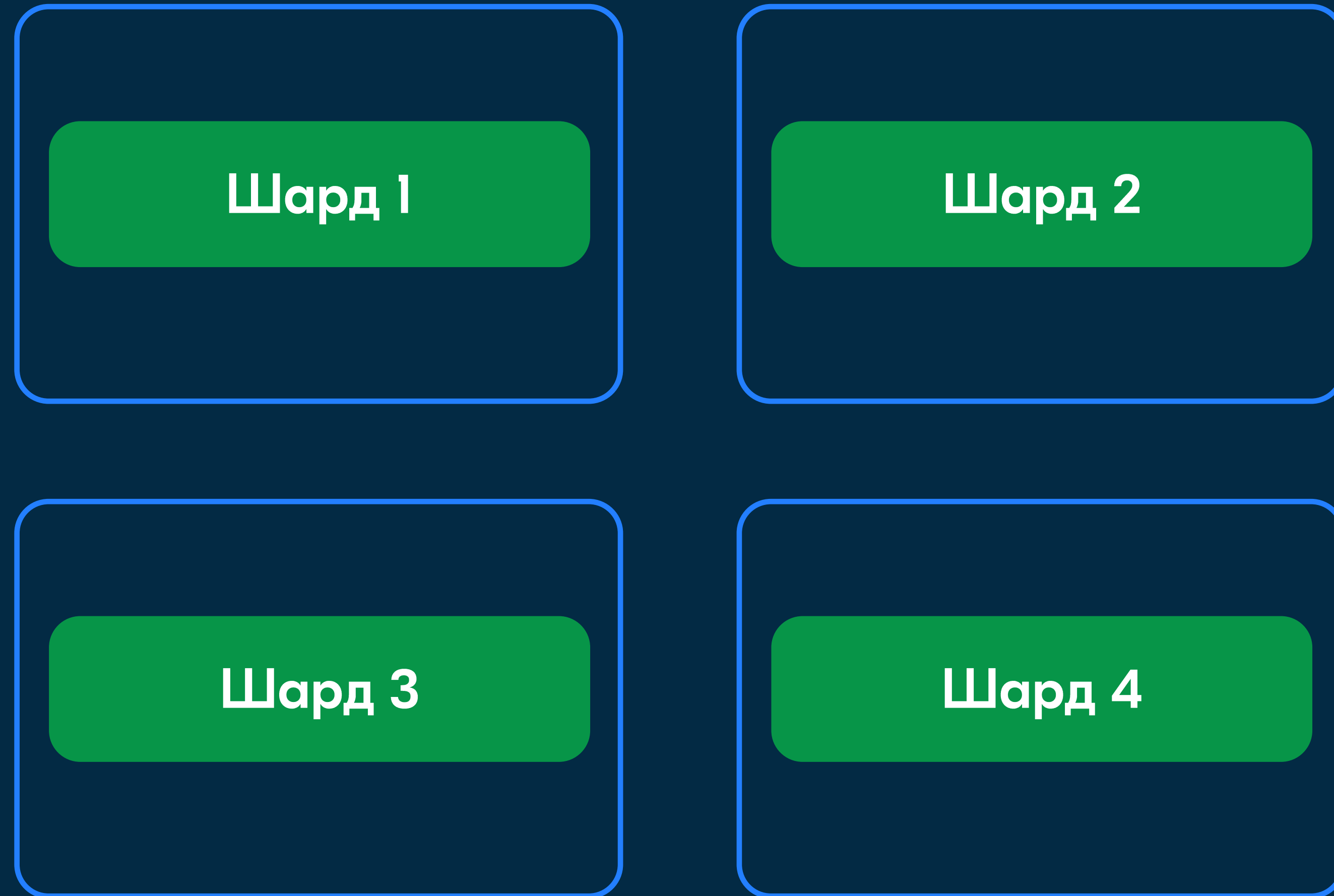
Архитектура кластера O2 Master (2022)

Горизонтальное масштабирование индексации



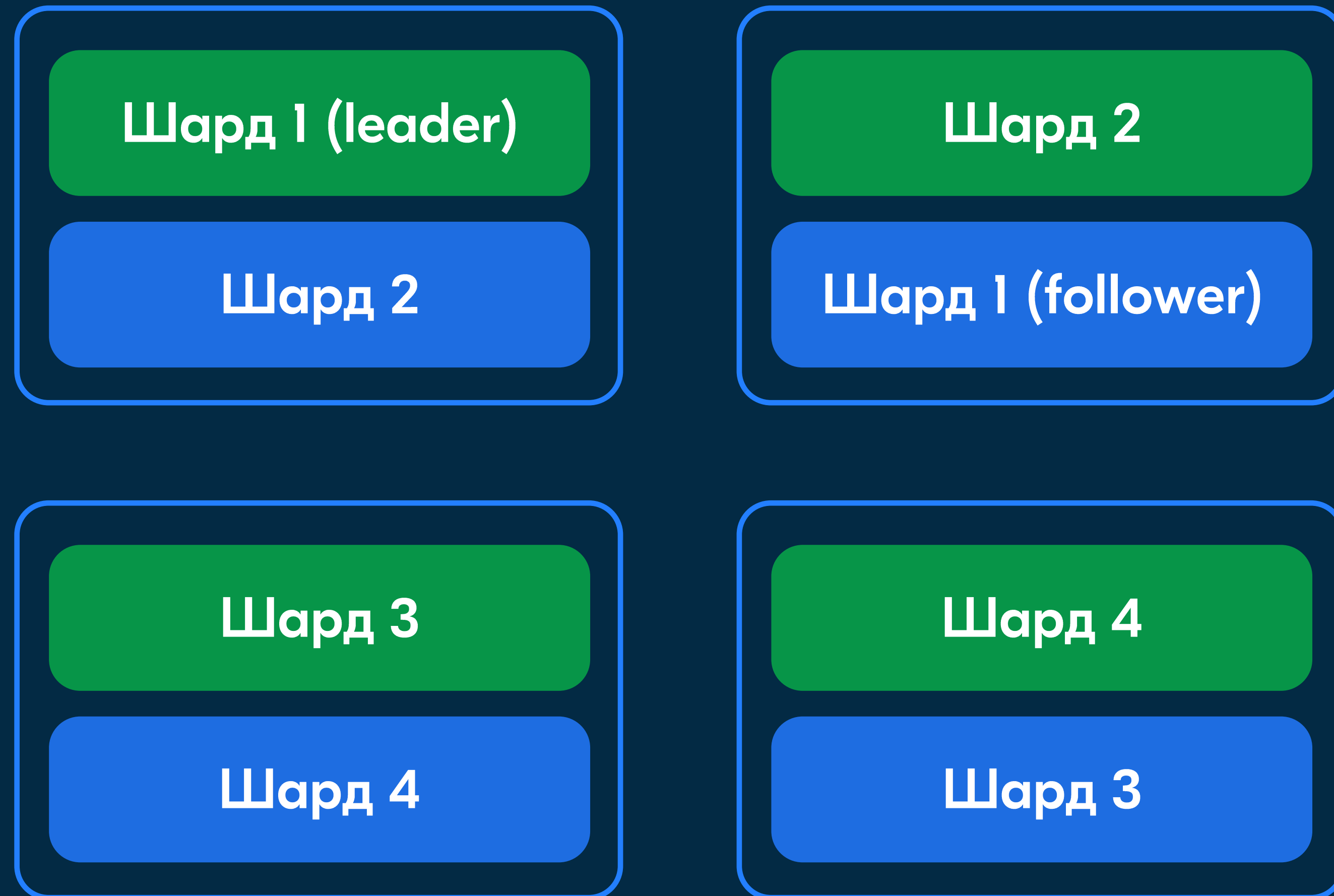
Архитектура кластера O2 Master (2022)

Динамическая балансировка



Архитектура кластера O2 Master (2022)

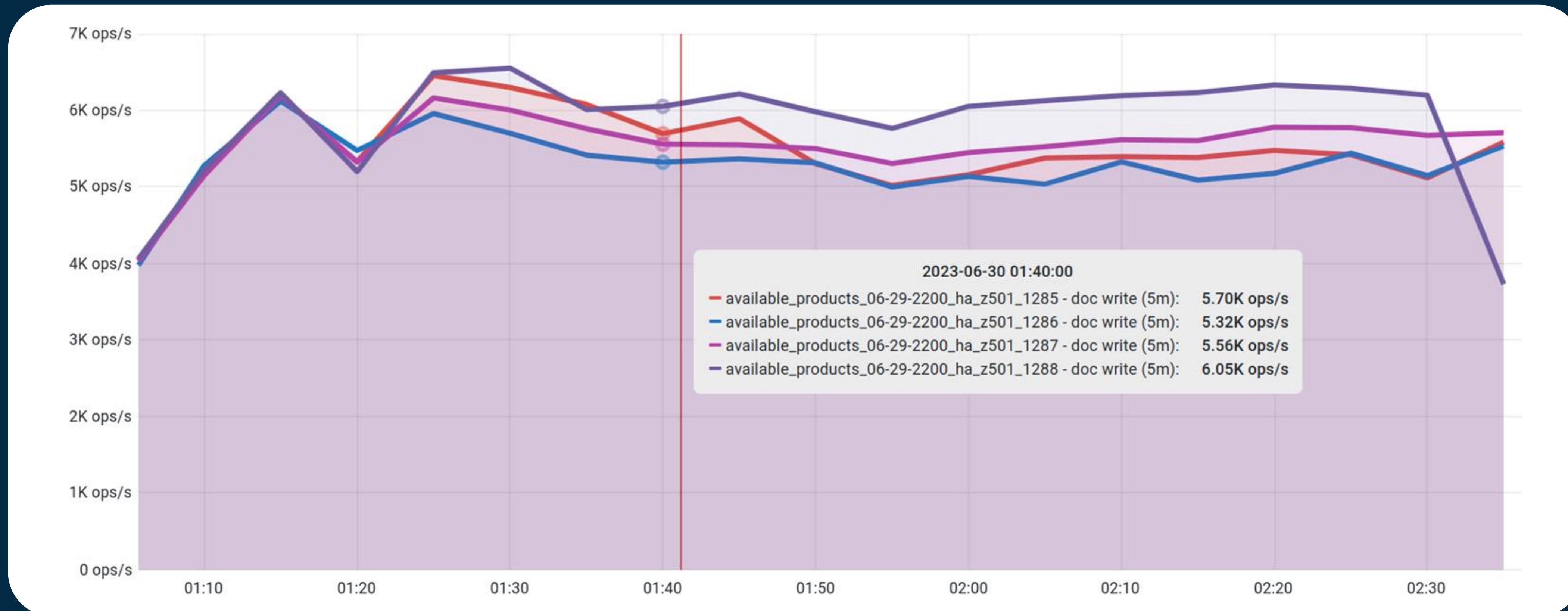
Hot standby



Архитектура кластера O2 Master (2022)

Результаты

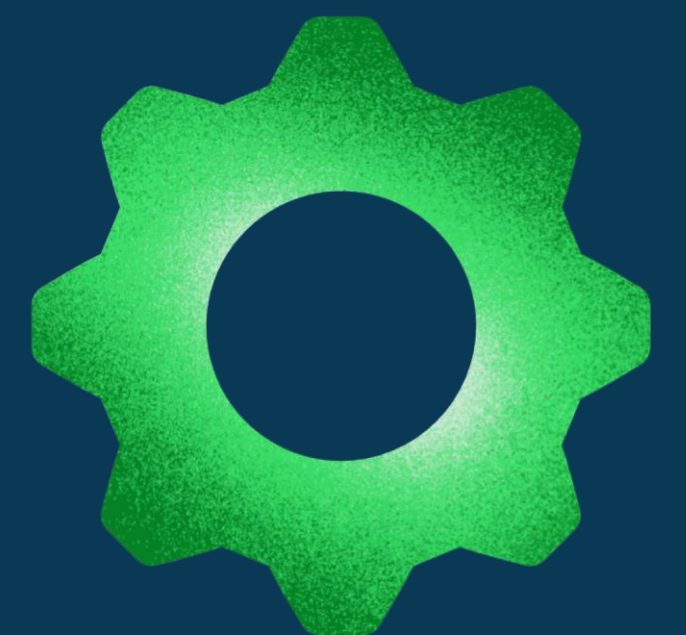
- Сборка шардированного индекса



Отказоустойчивость

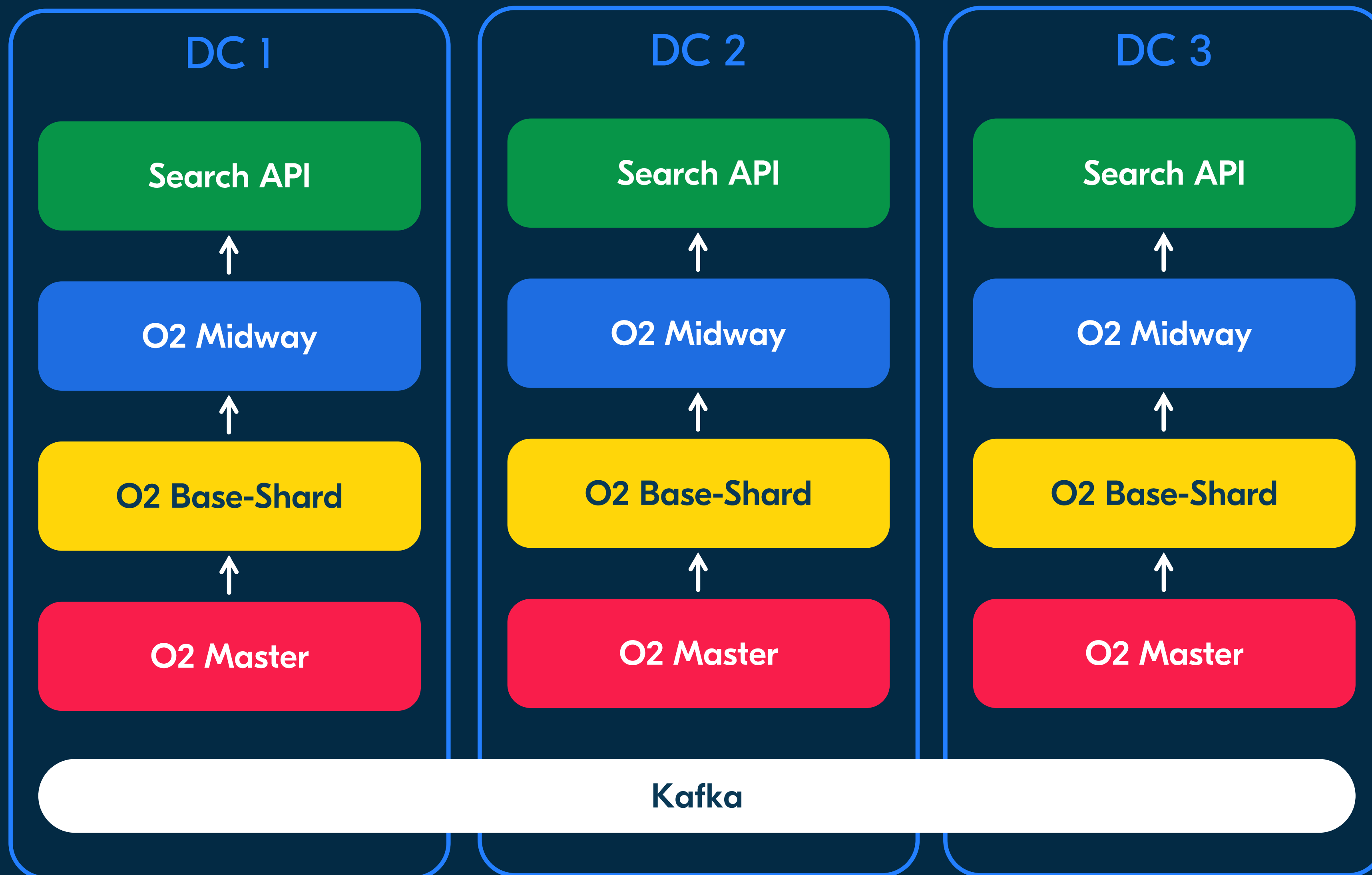
Отказоустойчивость

- Собираем несколько одинаковых индексов в каждом ЦОД
- Защищены от случайной порчи индекса (corrupted)



Поисковая платформа Ozon

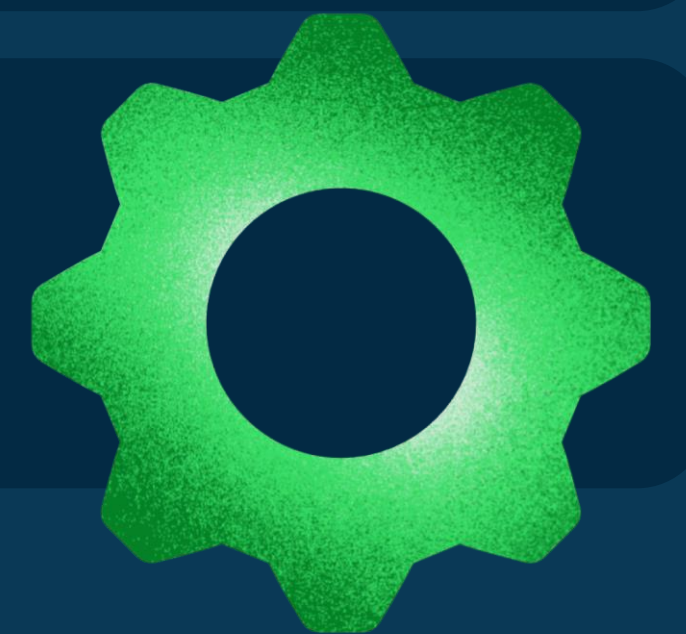
Отказоустойчивость



Архитектура кластера O2 Master

Отказоустойчивость

- Все типы нод поиска имеют более одного экземпляра
- Метаданные хранятся в PG (3 реплики)
- ETCD-кластер (5 реплик в трех ЦОД)
- Кафка в трех ЦОД



Архитектура кластера O2 Master

Результаты

1. Переживаем
отказ **любой**
ноды

2. Переживаем
отказ **стойки**

3. Переживаем
отказ **ЦОД**

В заключение

Цели достигнуты

- Индексация масштабируется вертикально и горизонтально
- Сервис отказоустойчив



В заключение

Как достичь производительности

- Знание стека и инструментов
- Изучение исходного кода открытых продуктов
- Реальное и виртуальное **общение** с коллегами в вашей области



Приглашение к дискуссии

Что не вошло в доклад

- Как работает наш поисковый кластер?
- Собственная реализация поисковых запросов
- Низкоуровневые оптимизации в Apache Lucene
- Кастомная политика мержа сегментов
- Распределенные алгоритмы коммита, ребалансировки



ozon{tech

Спасибо за внимание

Денис Габайдулин, Search Runtime

E-mail: gabaden@gmail.com

Telegram: [sherman81](https://t.me/sherman81)

