

# Реактивный раздатчик

[ok.ru/music](http://ok.ru/music)

Вадим Цесько



одноклассники

Joker<?>

# Одноклассники

- **71М пользователей (MAU)**
- **7000 машин в 4 ДЦ**

---

<sup>1</sup>Без учёта CDN



# Одноклассники

- **71М пользователей (MAU)**
- **7000 машин в 4 ДЦ**
- **2 Тб/с<sup>1</sup>**

---

<sup>1</sup>Без учёта CDN



# Одноклассники

- **71М пользователей** (MAU)
- **7000 машин** в 4 ДЦ
- **2 Тб/с<sup>1</sup>**
- До **100К запросов/с** на ноду (4 ядра)
- **99% < 100 мс**

---

<sup>1</sup>Без учёта CDN



# Одноклассники

- **71М пользователей** (MAU)
- **7000 машин** в 4 ДЦ
- **2 Тб/с<sup>1</sup>**
- До **100К запросов/с** на ноду (4 ядра)
- **99% < 100 мс**
- **Java**

---

<sup>1</sup>Без учёта CDN



# ok.ru/music

The screenshot displays the OK.ru music player interface. At the top, there is a navigation bar with icons for Messages, Discussions, Notifications, Friends, Guests, Events, Music, and Video, along with a search bar. The main player area shows a track titled "Múm - The Colorful Stabwound" with a progress bar at 0:18 / 3:40. Below the player, a search bar is active with the text "Search for music". The search results are displayed in a list format, showing the track name, duration, and a play button. The track "Múm - The Colorful Stabwound" is highlighted. To the left of the search results, there is a sidebar with a "My music" section containing a list of albums and their durations. Below the search results, there is a "Similar to your tracks" section with a list of recommended tracks. The interface is dark-themed.

ОДНОКЛАССНИКИ

Messages Discussions Notifications Friends Guests Events Music Video Search

0:18 / 3:40

Now playing uploaded by Helina Post as status Added

Upload your tracks from VK Popular music Collections Upload music Purchased My radio My music 323 Create a collection

M83 100 Grease 13 Lamb 77 Groove Armada 100 Bonobo 92 The Album Leaf 35 Ambient (ACS) 200 Roxy Music 85 Ghost In The Shell 5

Search for music Search

My music

Múm

- Twin Peaks - Twin Peaks Theme 5:04
- M83 - Midnight City 4:03
- Múm - The Colorful Stabwound 3:40
- Godspeed You! Black Emperor - providence 29:02
- Bryan Ferry - Knockin' On Heaven's Door 6:14
- Coldplay - Viva La Vida 4:01
- Kenji Kawai - Making Of Cyborg 4:28
- R.E.M. - Losing My Religion 4:29
- A Silver Mt. Zion - 1,000,000 Died To Make This Sound 14:42
- A Silver Mt. Zion - Track12 0:11
- Bonobo - Trouble Man - Change Is What We Need 6:42
- Bonobo - Bonobo - Recurring 4:52
- Bonobo - Bonobo - Change Down The Sugar Rhyme 4:10
- Bonobo - Jazz Juice - Marra Bossa 4:19
- Bonobo - Muro Maldoro - I Know A Little Cuban 2:47

Upload music Edit

Similar to your tracks

- Bryan Ferry - Slave To Love 4:26
- M83 - We Own The Sky 5:02
- Skillet - Hero 3:06

Show more tracks

# Раздатчик ok.ru/music

Name	× Headers Preview Response Timing
<input type="checkbox"/> stream?id=v0_10039836541_1_1&...	
<input type="checkbox"/> stream?id=v0_10039836541_1_1&...	
<input checked="" type="checkbox"/> stream?id=v0_10039836541_1_1&...	<b>General</b> <b>Request URL:</b> <a href="https://musicd.mycdn.me/v0/stream?id=v0_10039836541_1_1&amp;cid=v0_10039836541_1_1&amp;c4feb7b00abc485bed1f350fbd4c1080288bc11e55e8dd80db7a390addce662d&amp;ts=1538391797113&amp;md5=f34d94fe3501688219866">https://musicd.mycdn.me/v0/stream?id=v0_10039836541_1_1&amp;cid=v0_10039836541_1_1&amp;c4feb7b00abc485bed1f350fbd4c1080288bc11e55e8dd80db7a390addce662d&amp;ts=1538391797113&amp;md5=f34d94fe3501688219866</a> <b>Request Method:</b> <u>GET</u> <b>Status Code:</b> <span style="color: green;">●</span> <u>206 Partial Content</u> <b>Remote Address:</b> 217.20.155.15:443 <b>Referrer Policy:</b> origin
<input type="checkbox"/> play30;jsessionid=JsNuCw-IRC2Nv...	<b>Response Headers</b> <a href="#">view source</a> <b>Accept-Ranges:</b> bytes <b>Cache-Control:</b> max-age=86400 <b>Content-Length:</b> 8570377 <b>Content-Range:</b> <u>bytes 262144-8832520/8832521</u> <b>Content-Type:</b> audio/mpeg
	<b>Request Headers</b> <a href="#">view source</a> <b>Accept:</b> */* <b>Accept-Encoding:</b> identity;q=1, *;q=0 <b>Accept-Language:</b> en-US,en;q=0.9,ru;q=0.8 <b>Cache-Control:</b> no-cache

# Раздатчик ok.ru/music

- Раздаёт музыку с 2010

---

<sup>2</sup><http://profyclub.ru/docs/174>

<sup>3</sup><http://2017.jokerconf.com/2017/talks/4fasygftomyekgaqaaiuo/>





# Раздатчик ok.ru/music

- Раздаёт музыку с 2010
- $\leq$  **100 Гб/с**
- $\leq$  **500К соединений**

---

<sup>2</sup><http://profyclub.ru/docs/174>

<sup>3</sup><http://2017.jokerconf.com/2017/talks/4fasygftomyekgaqaaiuo/>

# Раздатчик ok.ru/music

- Раздаёт музыку с 2010
- $\leq$  **100 Гб/с**
- $\leq$  **500К соединений**
- Кеширующий фронтенд перед  
one-blob-storage<sup>2</sup> + one-cold-storage<sup>3</sup>

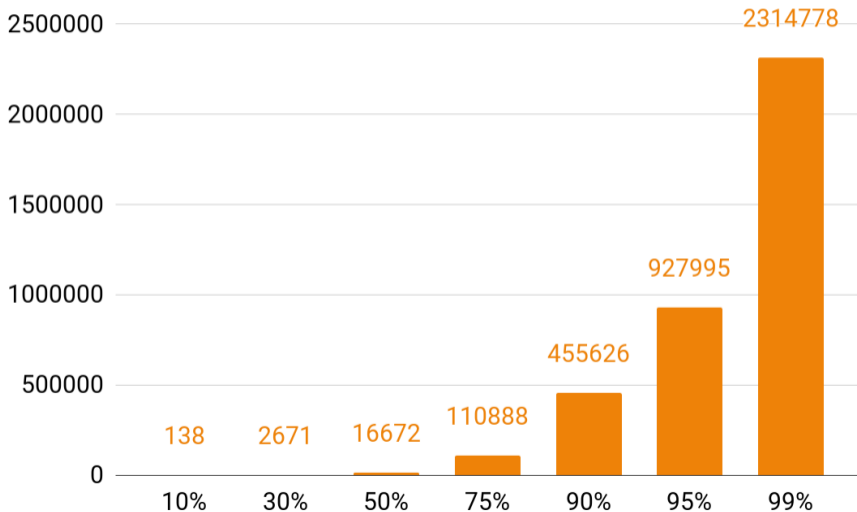
---

<sup>2</sup><http://profyclub.ru/docs/174>

<sup>3</sup><http://2017.jokerconf.com/2017/talks/4fasygftomyekgaqaaiuo/>



# Длинный хвост



# Must have

- До 100К «медленных» клиентов на ноду
- Масштабируемость
- Отказоустойчивость



# Масштабировать

Что:



# Масштабировать

Что:

- **Пропускную способность** кластера



# Масштабировать

Что:

- **Пропускную способность** кластера
- **Суммарную ёмкость** кеша кластера



# Масштабировать

Что:

- **Пропускную способность** кластера
- **Суммарную ёмкость** кеша кластера

Как:





# Масштабировать

Что:

- **Пропускную способность** кластера
- **Суммарную ёмкость** кеша кластера

Как:

- **Горизонтально:** ДЦ и машины



# Масштабировать

Что:

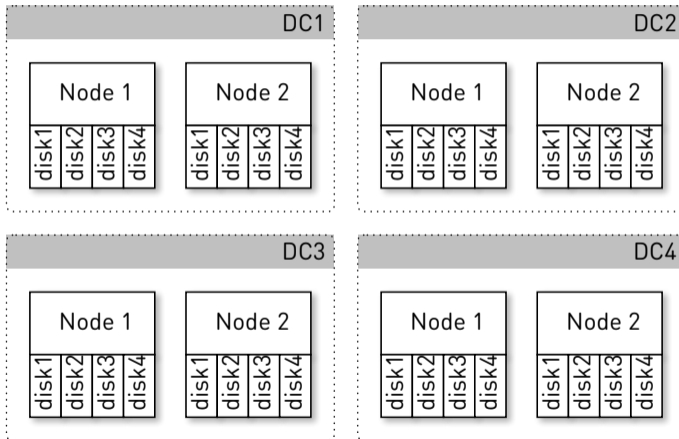
- **Пропускную способность** кластера
- **Суммарную ёмкость** кеша кластера

Как:

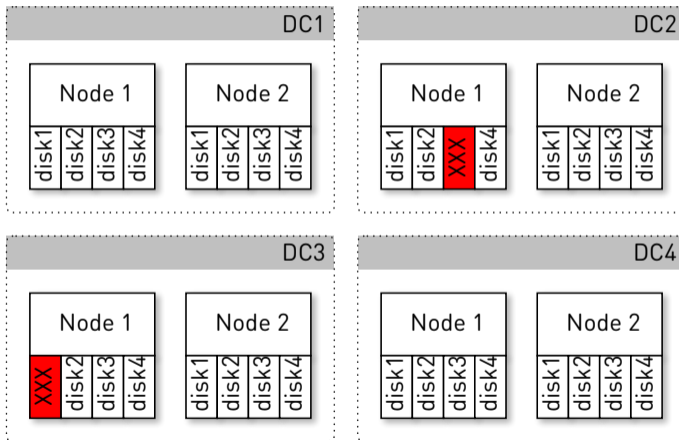
- **Горизонтально:** ДЦ и машины
- **Вертикально:** диски и RAM



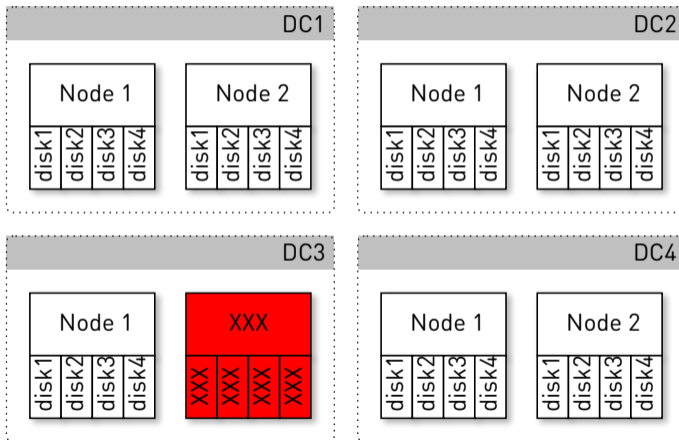
# Эффект масштаба



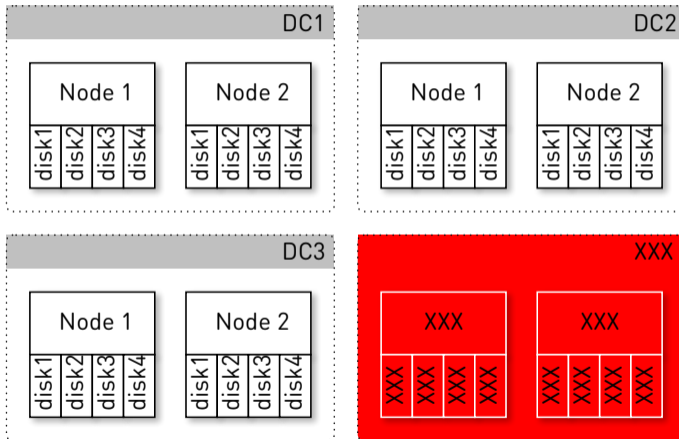
# Отказ диска



# Отказ машины



# Отказ ДЦ



# Всё (с)ломается<sup>4</sup>

- Прорабатываем сценарии отказов
  - ДЦ
  - Машины
  - Диски

---

<sup>4</sup><https://techtrain.ru/talks/1pqtojnosesumeo00ayccm/>



# Всё (с)ломается<sup>4</sup>

- Прорабатываем сценарии отказов
  - ДЦ
  - Машины
  - Диски
- Резервируем

---

<sup>4</sup><https://techtrain.ru/talks/1pqtojnosesumeo00ayccm/>





# Всё (с)ломается<sup>4</sup>

- Прорабатываем сценарии отказов
  - ДЦ
  - Машины
  - Диски
- Резервируем
  - Реплицируем данные

---

<sup>4</sup><https://techtrain.ru/talks/1pqtojnosesumeo00ayccm/>



# Всё (с)ломается<sup>4</sup>

- Прорабатываем сценарии отказов
  - ДЦ
  - Машины
  - Диски
- Резервируем
  - Реплицируем данные
  - Запас по пропускной способности

---

<sup>4</sup><https://techtrain.ru/talks/1pqtojnosesumeo00ayccm/>



# Балансировка запросов

- По **ДЦ**
  - Манёвры
  - Аварии

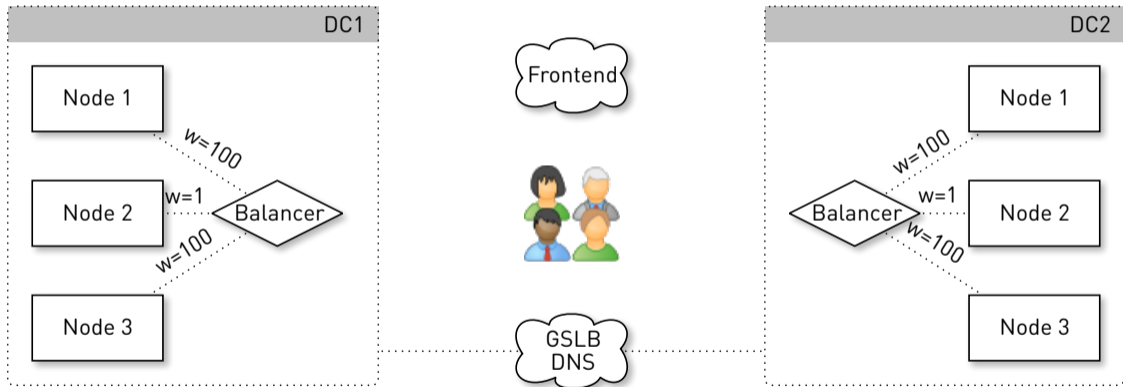


# Балансировка запросов

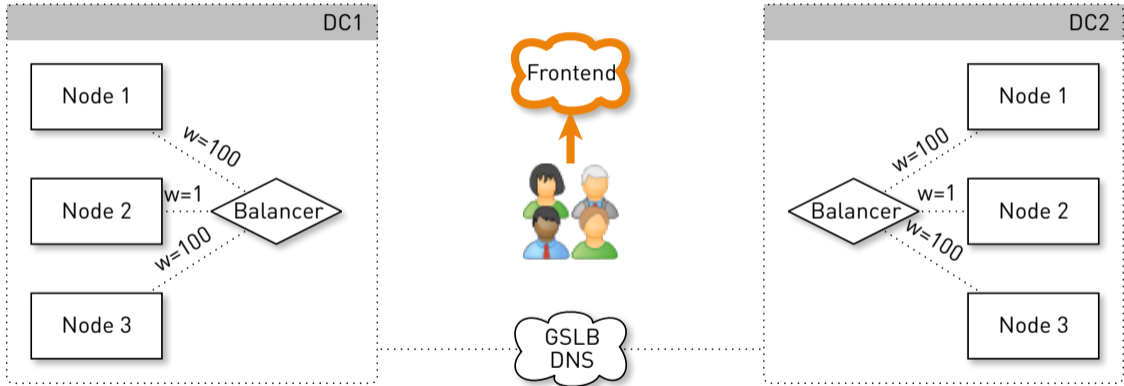
- По **ДЦ**
  - Манёвры
  - Аварии
- По **нодам** внутри ДЦ
  - С **весами**
  - Плавный ввод в ротацию
  - Нагрузочное тестирование



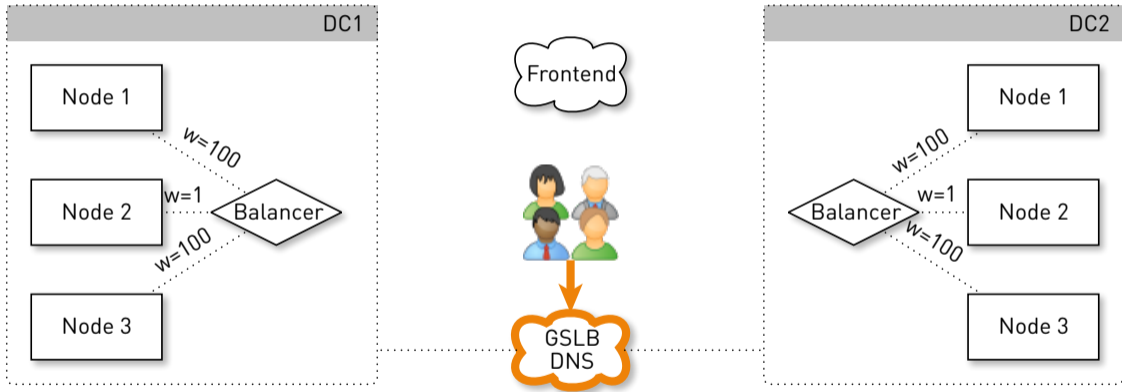
# Путь запроса



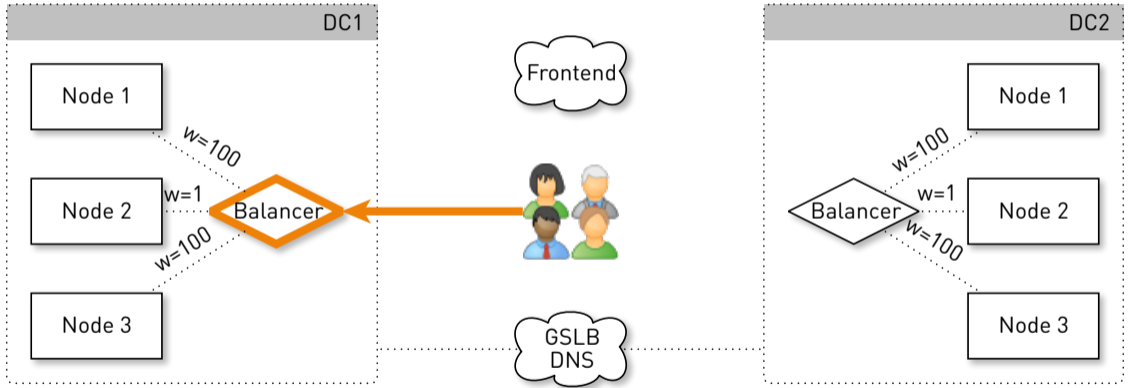
<https://ok.ru/music>



musicd.mycdn.me/v0/stream?id=...

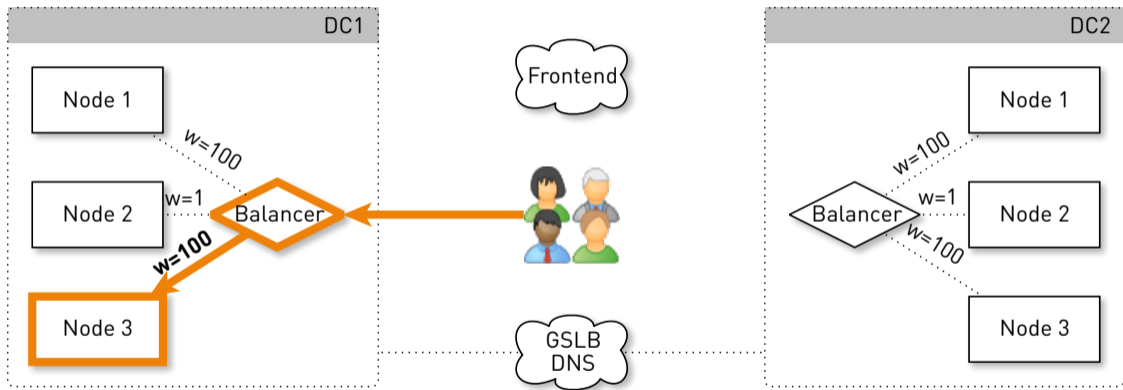


# 5.61.23.6/v0/stream?id=...





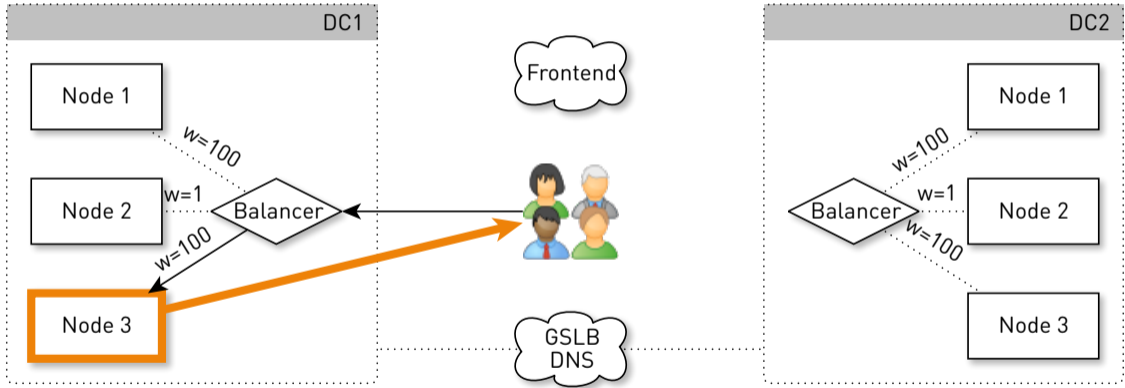
# Соединение с нодой<sup>5</sup>



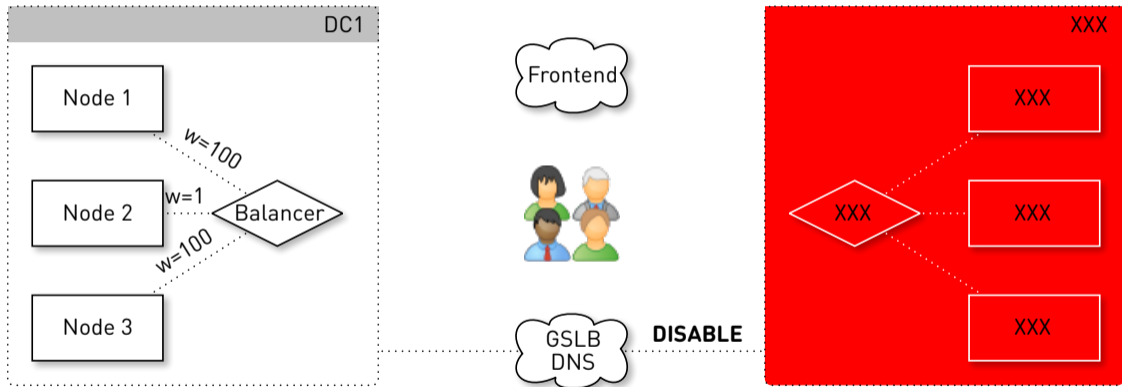
<sup>5</sup><http://www.highload.ru/2017/abstracts/2858.html>



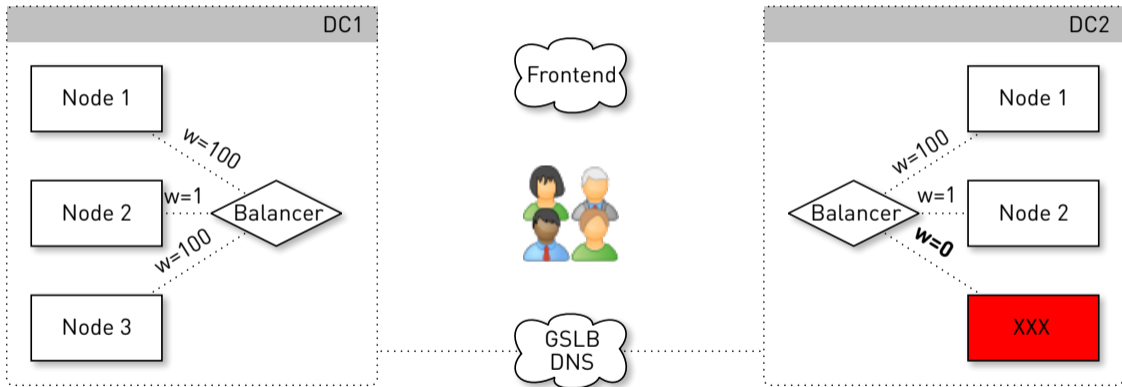
# Direct Server Return



# Сбой ДЦ



# Сбой ноды



# Балансировка треков по нодам

- Независимые ДЦ



# Балансировка треков по нодам

- Независимые ДЦ
- Равномерно по машинам



# Балансировка треков по нодам

- Независимые ДЦ
- Равномерно по машинам
- Репликация треков на случай отказов



# Балансировка треков по нодам

- Независимые ДЦ
- Равномерно по машинам
- Репликация треков на случай отказов
- Перераспределение нагрузки при отказах





# Балансировка треков по нодам

- Независимые ДЦ
- Равномерно по машинам
- Репликация треков на случай отказов
- Перераспределение нагрузки при отказах
- **Consistent hashing**<sup>6</sup>

---

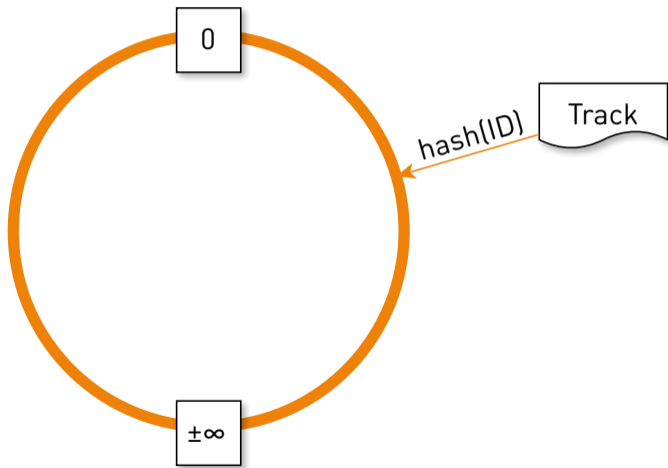
<sup>6</sup><https://medium.com/@dgryski/consistent-hashing-algorithmic-tradeoffs-ef6b8e2fcae8>



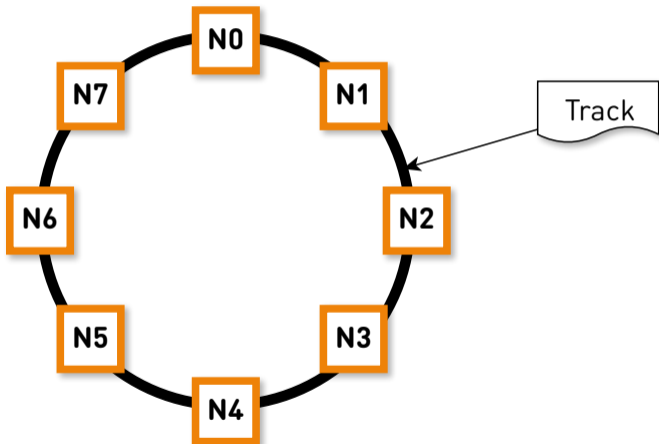
# Consistent hashing



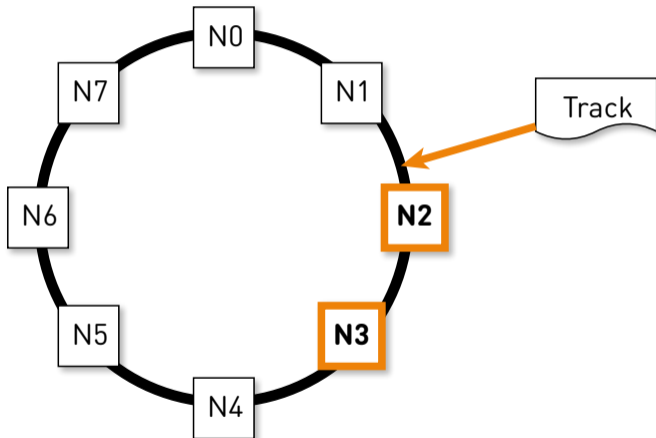
# Строим кольцо хешей



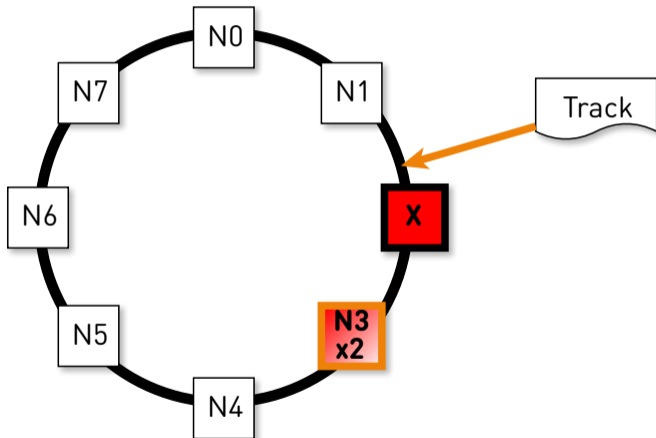
# Делим интервалы между нодами



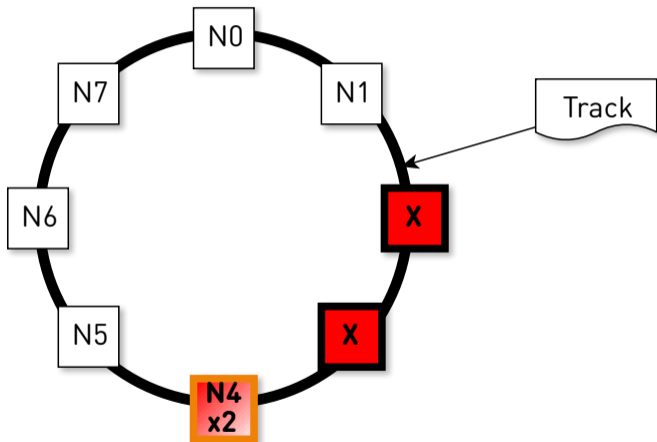
# Определяем реплики



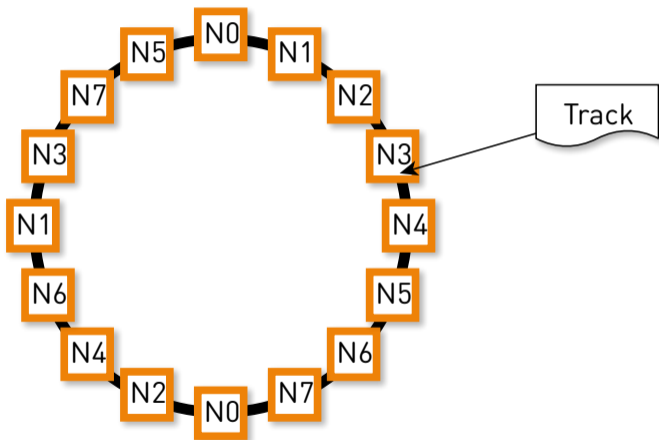
# Отказ ноды



# Каскадный отказ

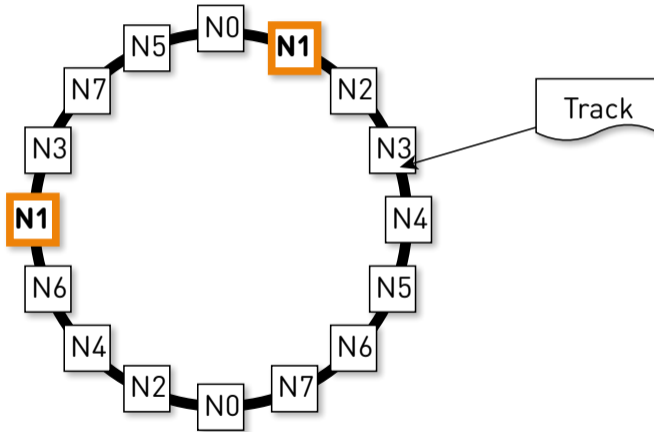


На практике  $vnodes \ggg nodes$

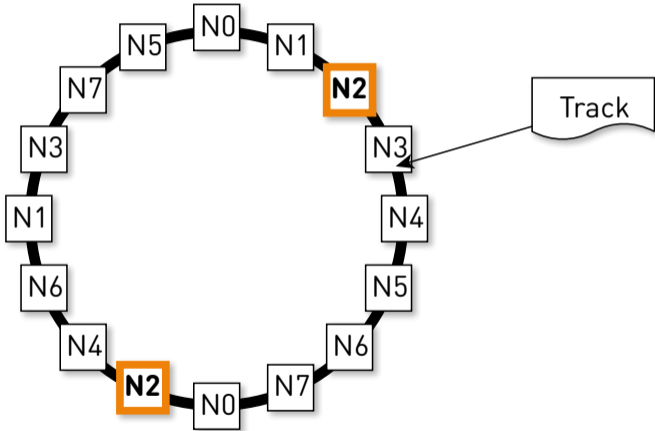




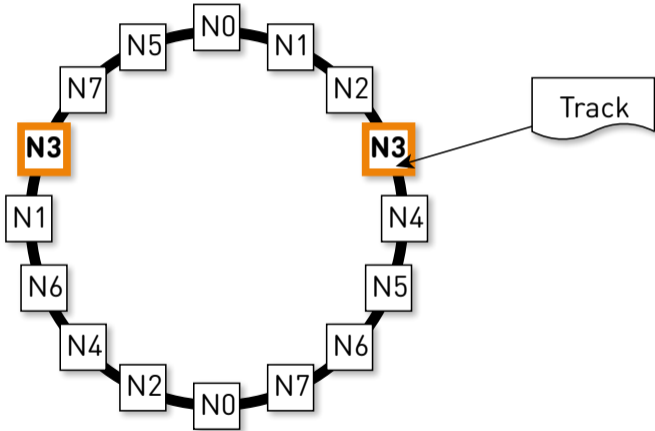
# N1



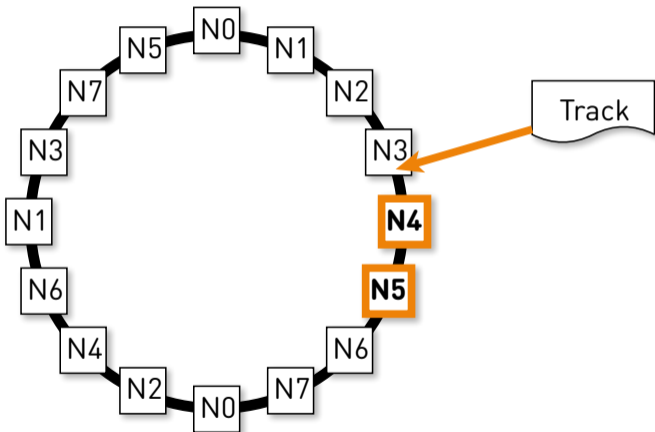
# N2



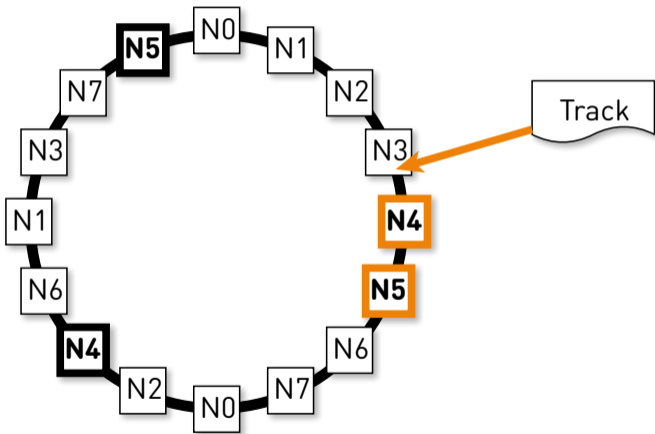
# N3



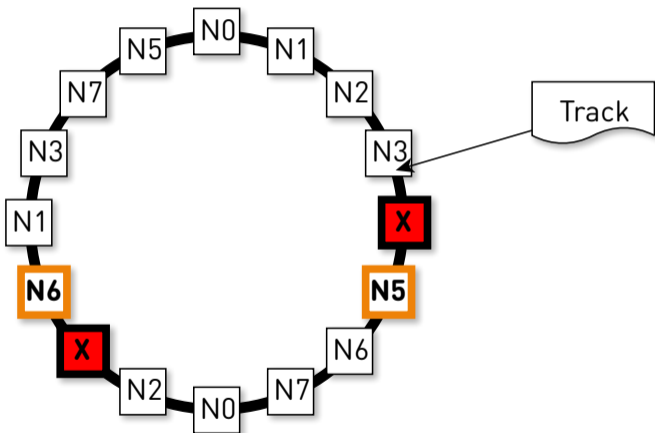
# Так же определяем реплики



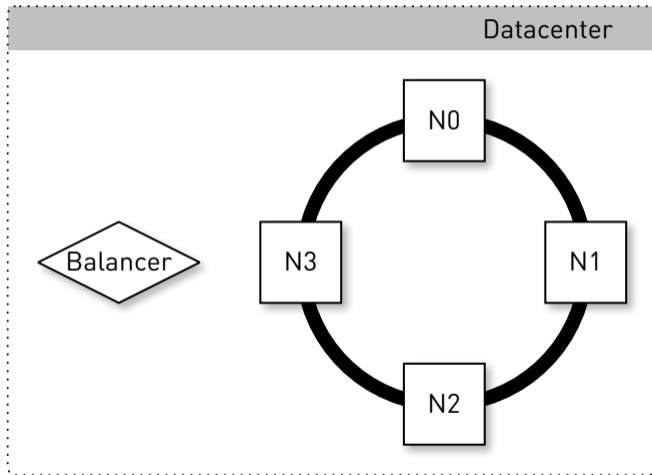
# N4 + N5



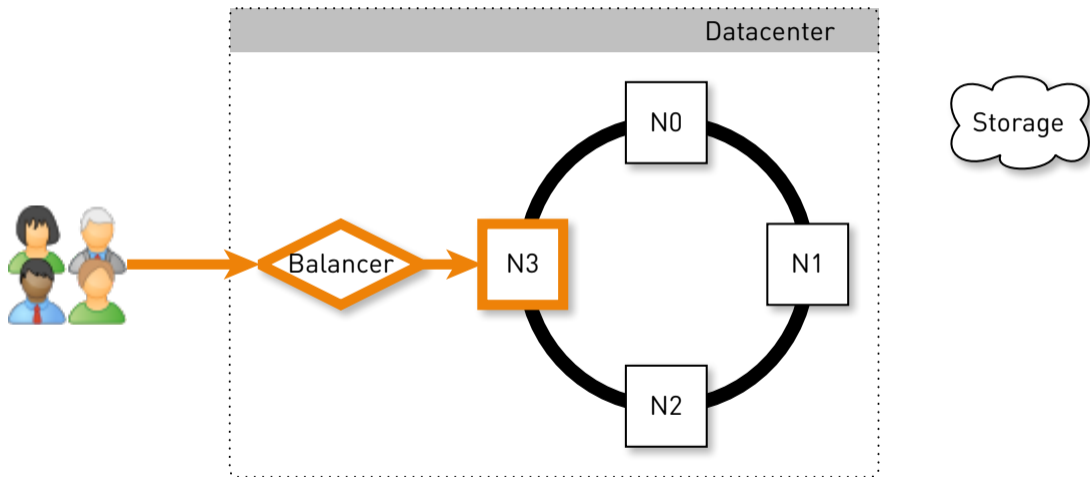
# Отказ ноды



# Отдача трека

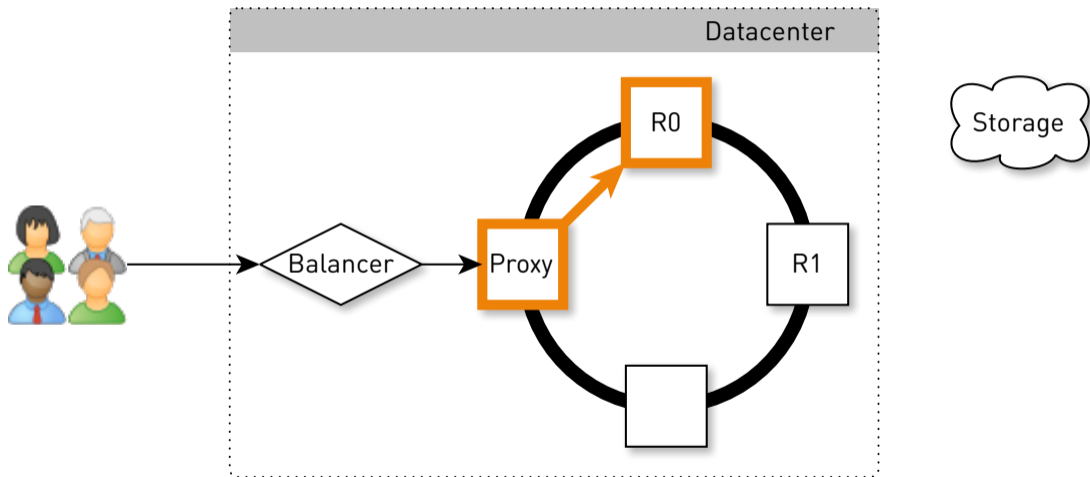


# Случайная нода

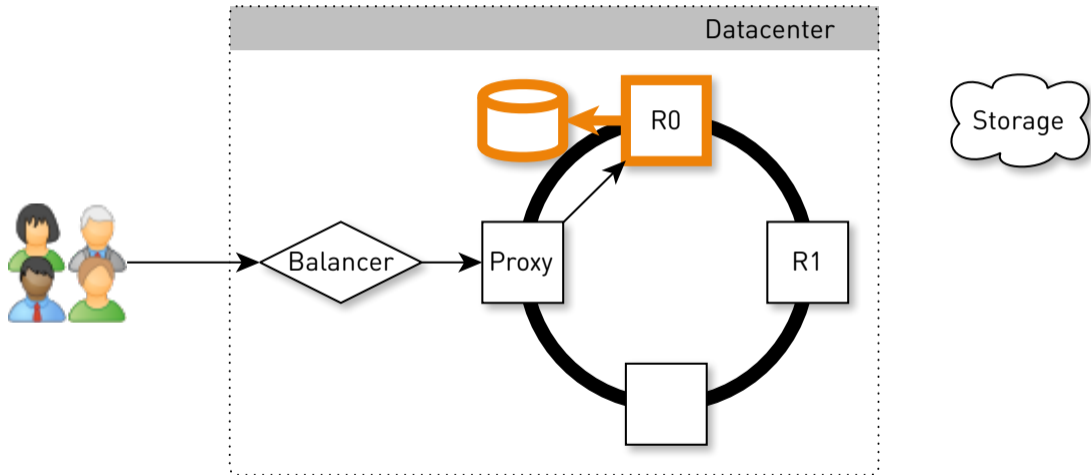




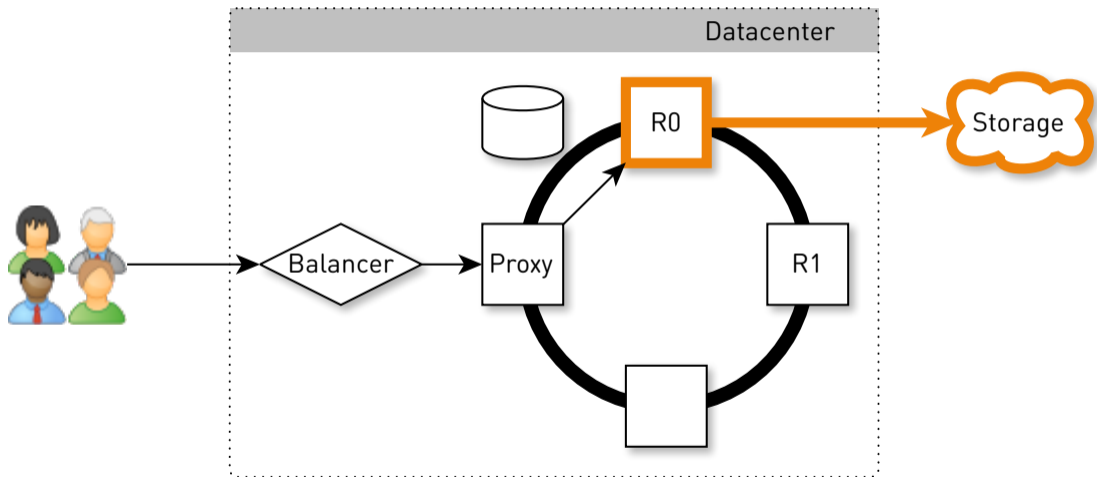
# Проксируем



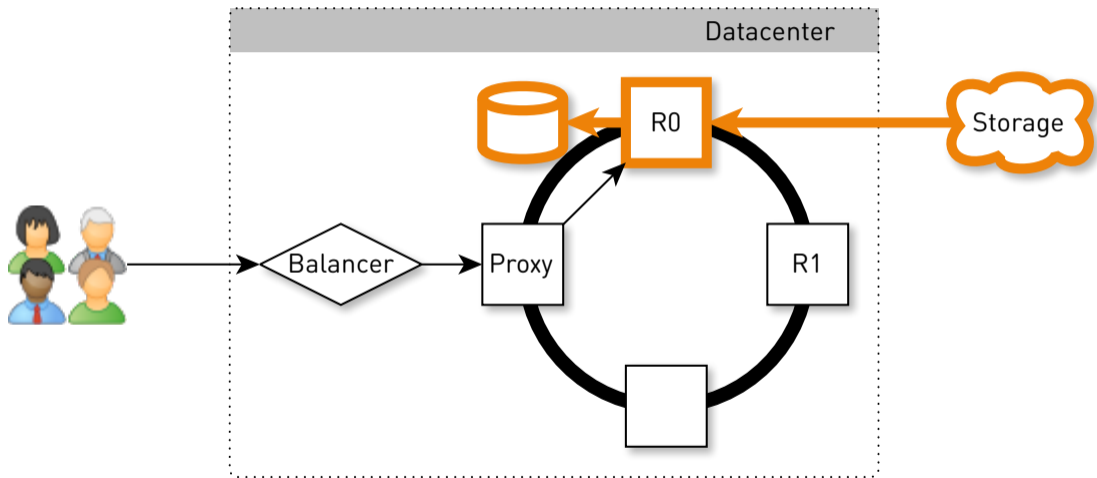
# Ищем локально



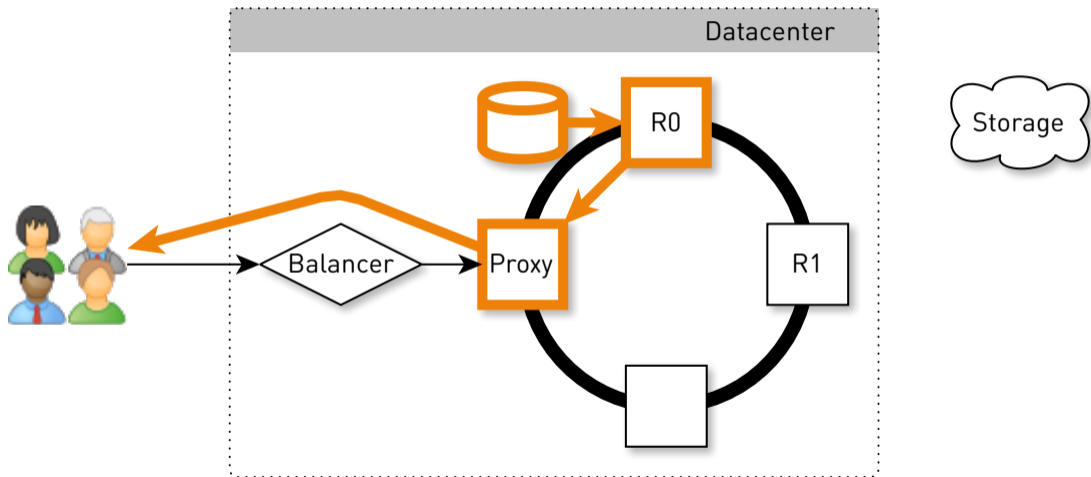
# Подтягиваем данные



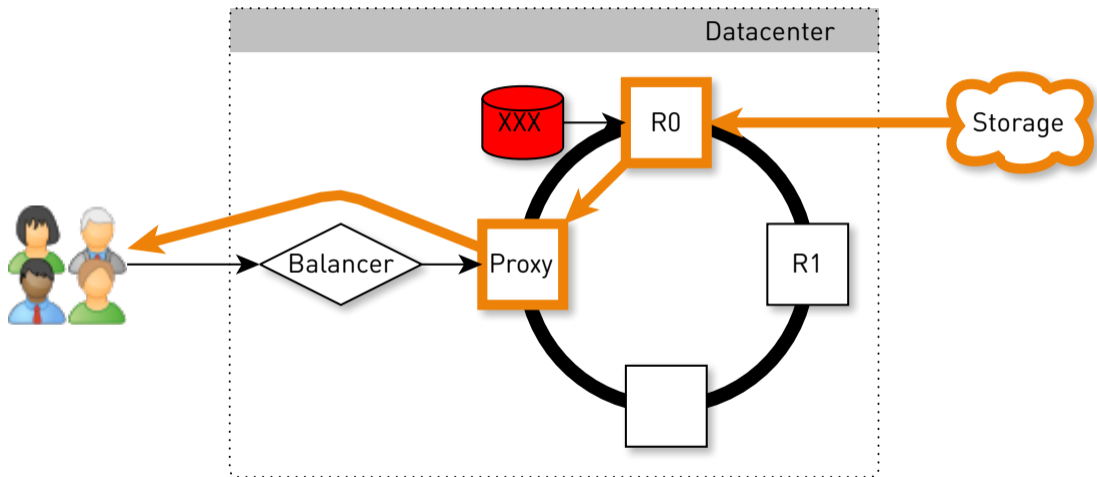
# Сохраняем



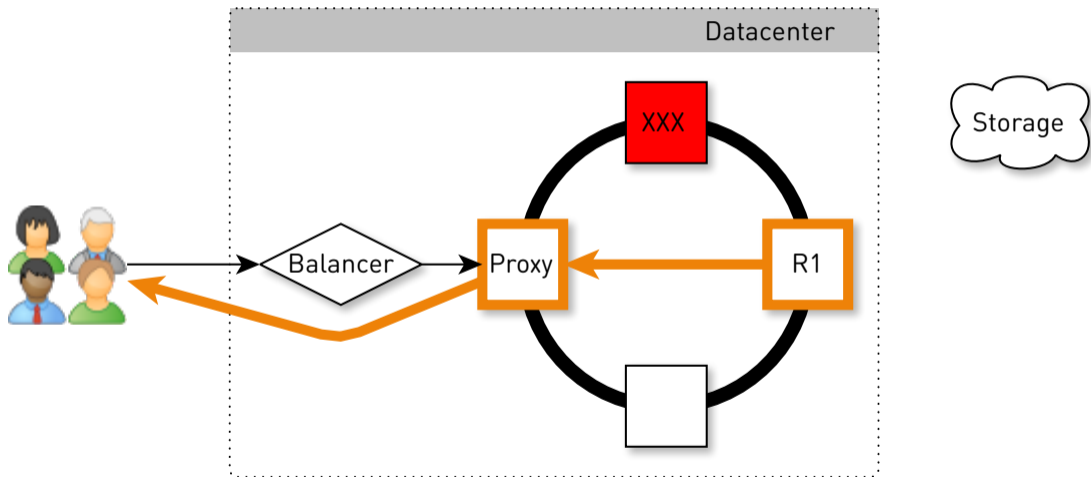
# Отдаём с диска



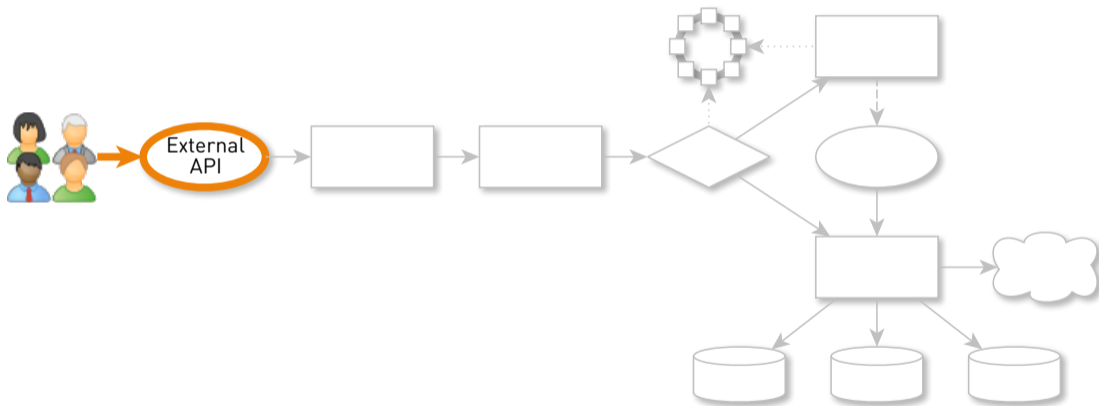
# Сбой диска



# Сбой реплики

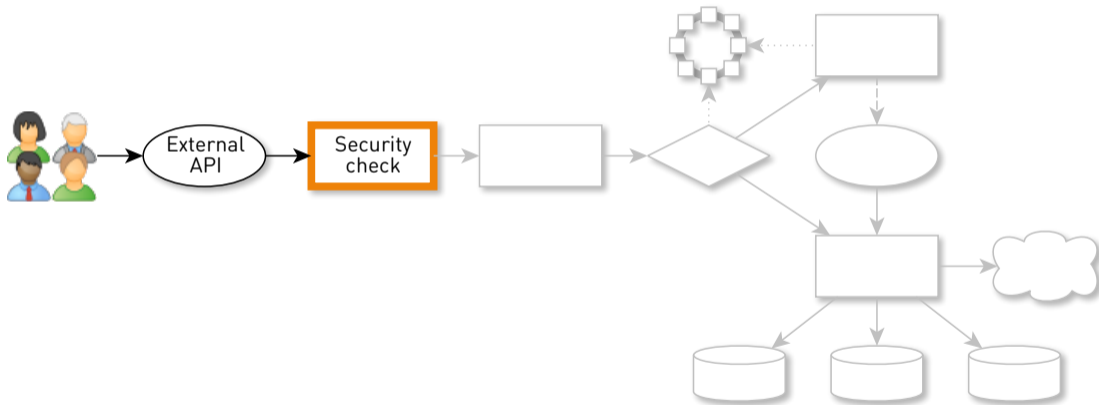


# Нода изнутри

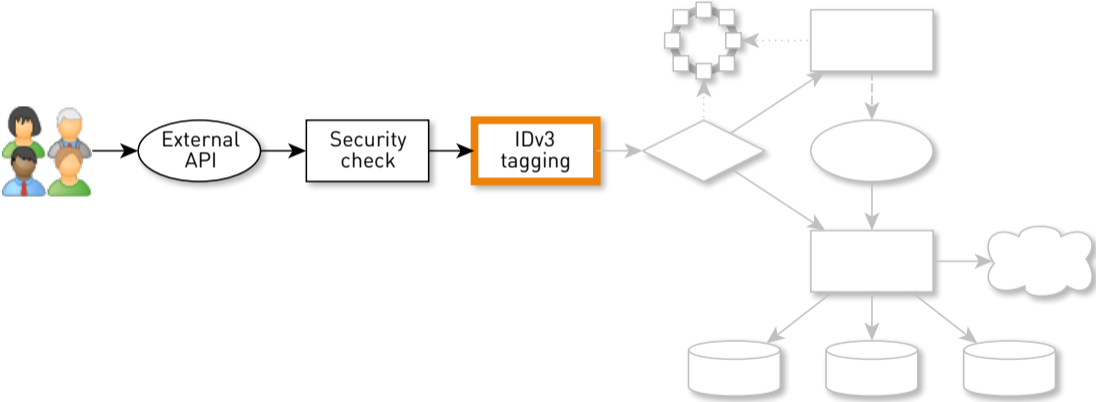




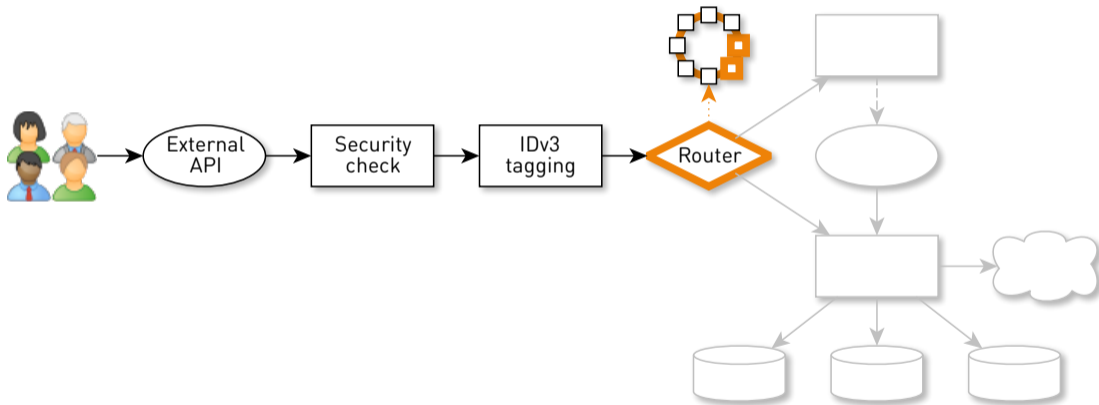
# Проверка подписи



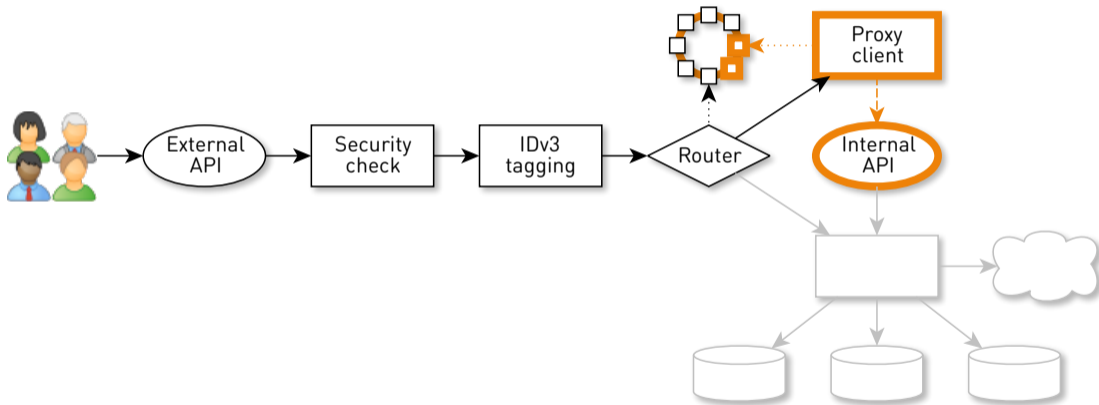
# Добавление IDv3



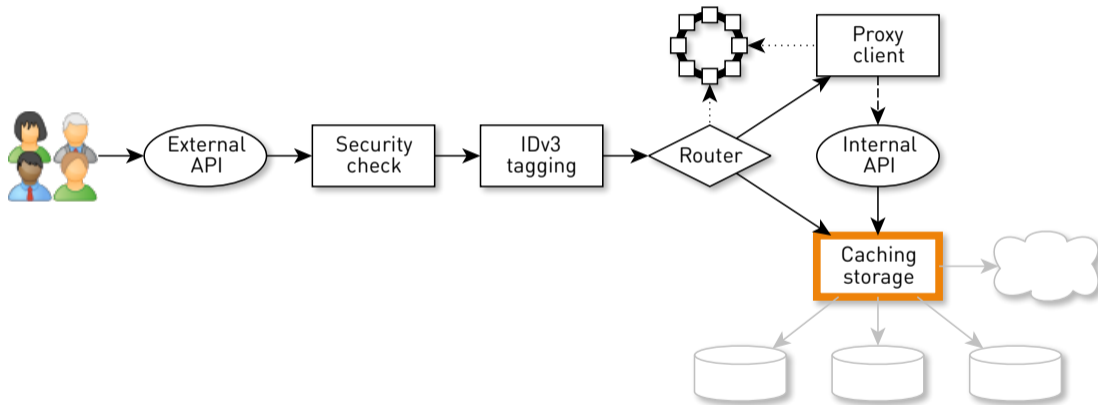
# Определение реплик



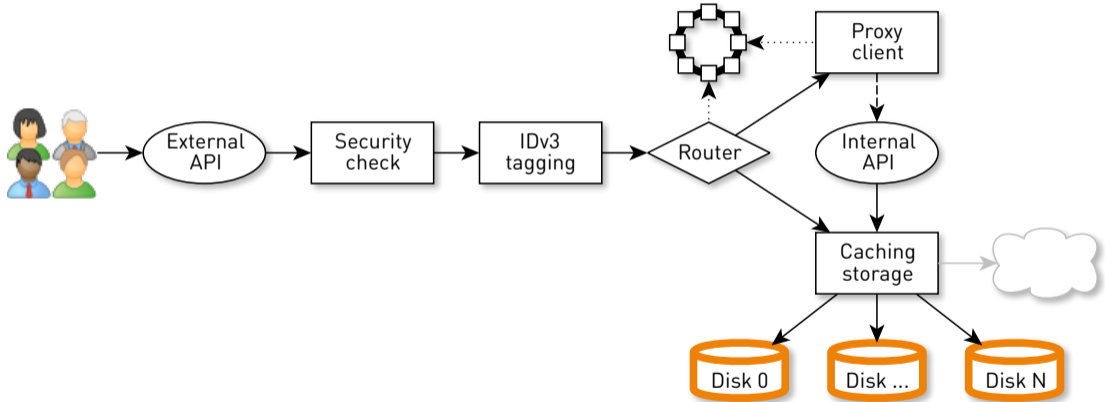
# Проксируем с соседа



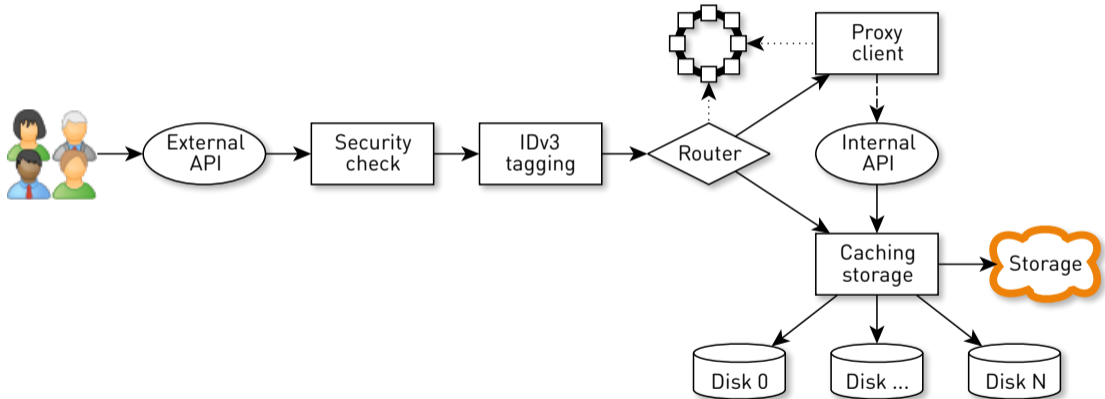
# Отдаём локально



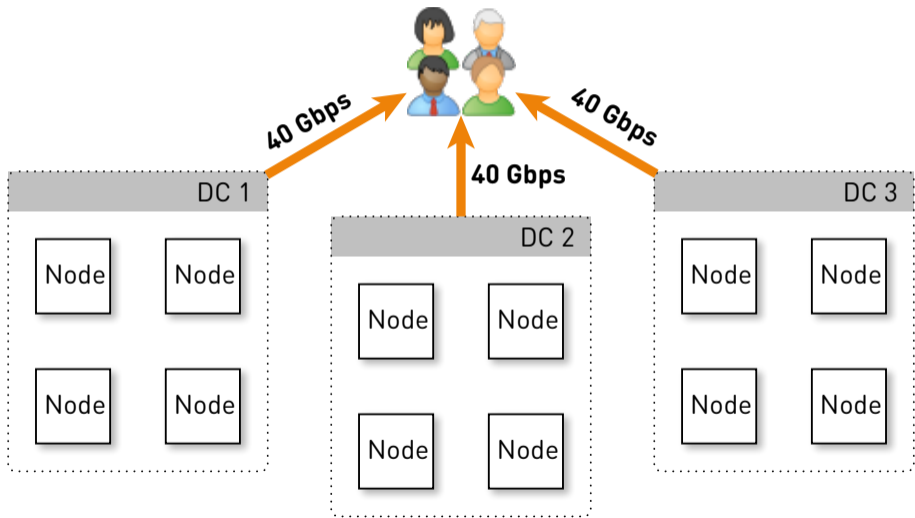
# Hit



# Miss

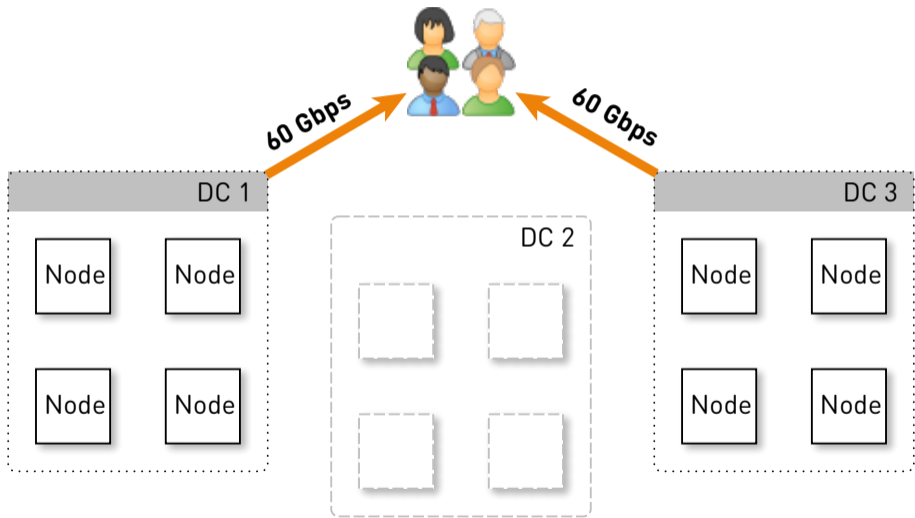


Пусть 120 Гб/с...

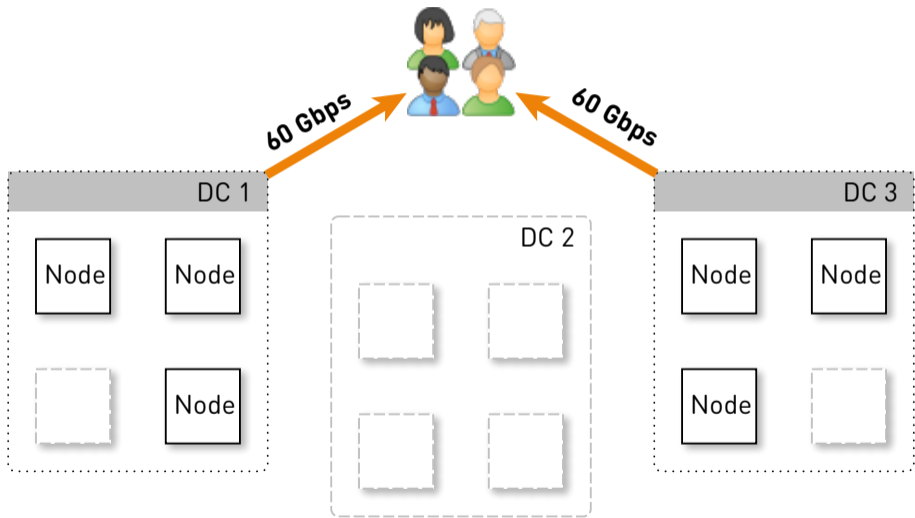




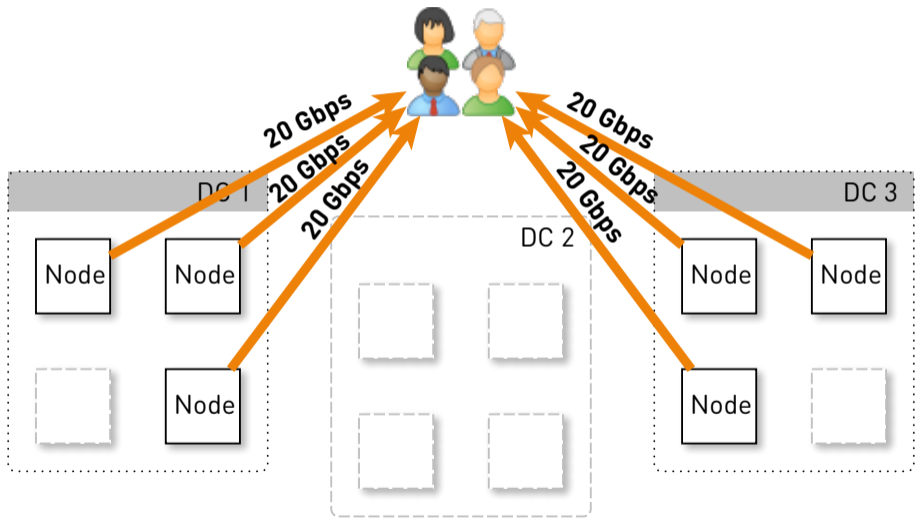
# -1 ДЦ



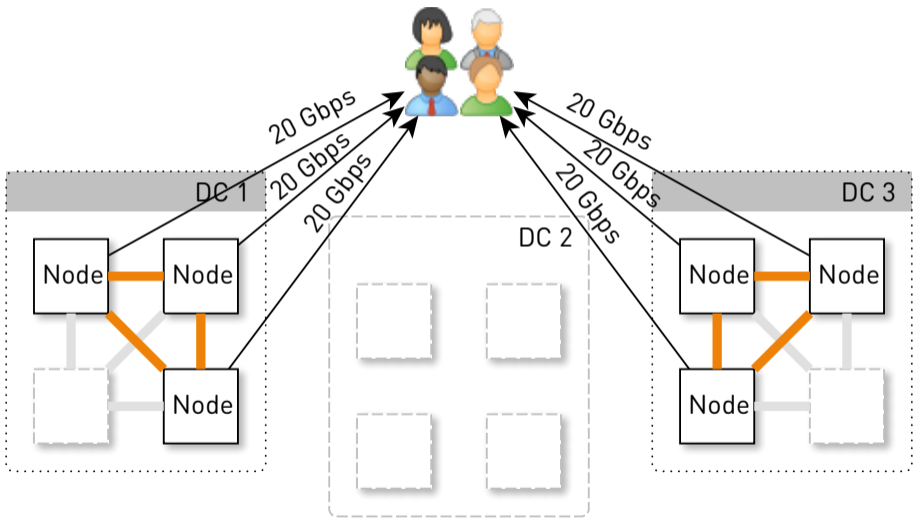
# Ещё и релиз



# С ноды



$RF = 2 \Rightarrow 50\%$  проксируется



# Итого одна нода

- Отдаёт пользователям **20 Гб/с**



# Итого одна нода

- Отдаёт пользователям **20 Гб/с**
- Из них **10 Гб/с проксирует**



# Итого одна нода

- Отдаёт пользователям **20 Гб/с**
- Из них **10 Гб/с проксирует**
- И отдаёт **соседям-прокси 10 Гб/с**



# Итого одна нода

- Отдаёт пользователям **20 Гб/с**
- Из них **10 Гб/с проксирует**
- И отдаёт **соседям-прокси 10 Гб/с**
- В итоге **наружу 30 Гб/с**, из них **20 Гб/с с дисков**





# Итого одна нода

- Отдаёт пользователям **20 Гб/с**
- Из них **10 Гб/с проксирует**
- И отдаёт **соседям-прокси 10 Гб/с**
- В итоге **наружу 30 Гб/с**, из них **20 Гб/с с дисков**
- **512 ГБ RAM** вместит  $\approx$  **50К** горячих треков



# Итого одна нода

- Отдаёт пользователям **20 Гб/с**
- Из них **10 Гб/с проксирует**
- И отдаёт **соседям-прокси 10 Гб/с**
- В итоге **наружу 30 Гб/с**, из них **20 Гб/с с дисков**
- **512 ГБ RAM** вместит  $\approx$  **50К** горячих треков
- $\approx$  **30-40%** хвоста провалится в диски



# Итого одна нода

- Отдаёт пользователям **20 Гб/с**
- Из них **10 Гб/с проксирует**
- И отдаёт **соседям-прокси 10 Гб/с**
- В итоге **наружу 30 Гб/с**, из них **20 Гб/с с дисков**
- **512 ГБ RAM** вместит  $\approx$  **50К** горячих треков
- $\approx$  **30-40%** хвоста провалится в диски
- $\approx$  **8 Гб/с** с дисков



# Как хранить данные?



# Как хранить данные?

- Трек = файл?
  - **Десятки ТБ**
  - **Миллионы файлов**



# Как хранить данные?

- Трек = файл?
  - **Десятки ТБ**
  - **Миллионы файлов**
  - Накладные расходы на **метаданные**



# Как хранить данные?

- Трек = файл?
  - **Десятки ТБ**
  - **Миллионы файлов**
  - Накладные расходы на **метаданные**
  - Тянуть треки **только сначала**



# Как хранить данные?

- Трек = файл?
  - **Десятки ТБ**
  - **Миллионы файлов**
  - Накладные расходы на **метаданные**
  - Тянуть треки **только сначала**
  - Хранить треки **только целиком**



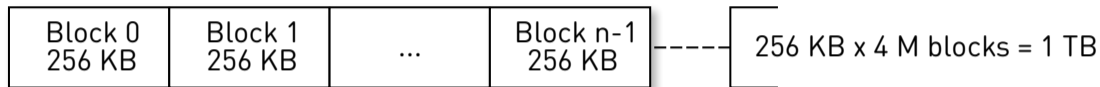


# Как хранить данные?

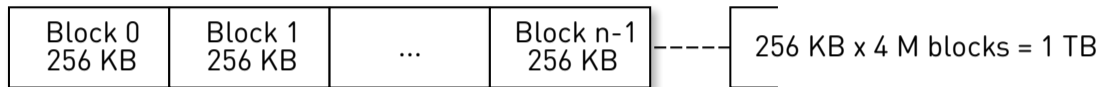
- Трек = файл?
  - **Десятки ТБ**
  - **Миллионы файлов**
  - Накладные расходы на **метаданные**
  - Тянуть треки **только сначала**
  - Хранить треки **только целиком**
- **Так мы не делаем**



# Трек = $N \times 256 \text{ KB}$ блоков (SSD)



# Трек = $N \times 256$ КБ блоков (SSD)



- Каждый диск — независимое хранилище



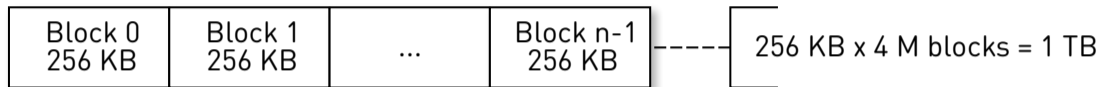
# Трек = $N \times 256$ КБ блоков (SSD)



- Каждый диск — независимое хранилище
- `hash(track, block_num)` чтобы **размазать блоки**



# Трек = $N \times 256$ КБ блоков (SSD)



- Каждый диск — независимое хранилище
- `hash(track, block_num)` чтобы **размазать блоки**
- **Page cache** ОС



# Как хранить блоки?

- Один файл на XFS размером с диск
  - `xfswalk -c resvsp 0 ...`



# Как хранить блоки?

- Один файл на XFS размером с диск
  - `xfs_io -c resvsp 0 ...`
- Сырое блочное устройство
  - `new RandomAccessFile("/dev/sdc", "rw")`



# Как хранить блоки?

- Один файл на XFS размером с диск
  - `xfstool -c resvsp 0 ...`
- Сырое блочное устройство
  - `new RandomAccessFile("/dev/sdc", "rw")`

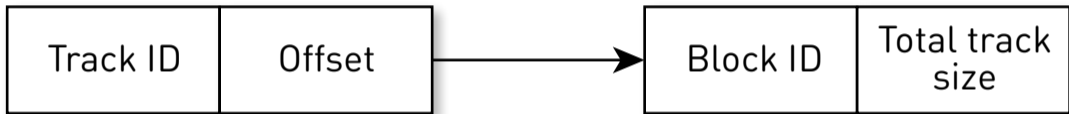
## **Блочное устройство vs файл на XFS:**

- **LA в 2 раза ниже**
- **Запись в 1.5 раза быстрее**
- **Время ответа в 2-3 раза ниже**

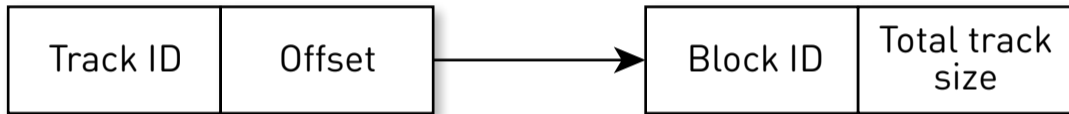




# Индекс



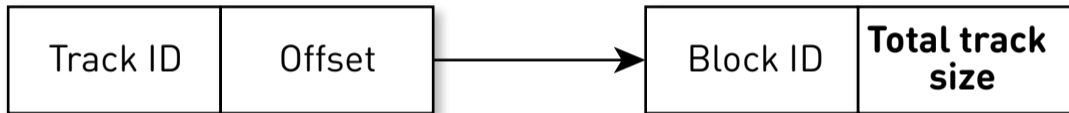
# Индекс



- **29 байт** на блок ( $\approx$  1 ГБ на 10 ТБ блоков)



# Total track size



- **29 байт** на блок ( $\approx$  1 ГБ на 10 ТБ блоков)



# Total track size

## ▼ Response Headers [view source](#)

**Accept-Ranges:** bytes

**Cache-Control:** max-age=86400

**Content-Length:** 8570377

**Content-Range:** bytes 262144-8832520/8832521

**Content-Type:** audio/mpeg

# Как хранить?

- In-memory

---

<sup>7</sup><https://github.com/ben-manes/caffeine/wiki/Efficiency>



# Как хранить?

- In-memory
- Компактно

---

<sup>7</sup><https://github.com/ben-manes/caffeine/wiki/Efficiency>



# Как хранить?

- In-memory
- Компактно
- Персистентно

---

<sup>7</sup><https://github.com/ben-manes/caffeine/wiki/Efficiency>



# Как хранить?

- In-memory
- Компактно
- Персистентно
- Вытеснение LRU<sup>7</sup>

---

<sup>7</sup><https://github.com/ben-manes/caffeine/wiki/Efficiency>





# Как хранить?

- In-memory
- Компактно
- Персистентно
- Вытеснение LRU<sup>7</sup>
- Thread-safe

---

<sup>7</sup><https://github.com/ben-manes/caffeine/wiki/Efficiency>



# Как хранить?

- In-memory
- Компактно
- Персистентно
- Вытеснение LRU<sup>7</sup>
- Thread-safe
  
- SharedMemoryFixedMap<sup>8</sup> на tmpfs

---

<sup>7</sup><https://github.com/ben-manes/caffeine/wiki/Efficiency>

<sup>8</sup><https://github.com/odnoklassniki/one-nio>



# Durability

- **Рестарт** машины — теряем tmpfs



# Durability

- **Рестарт** машины — теряем tmpfs
- **Падение** процесса — неконсистентность



# Durability

- **Рестарт** машины — теряем `tmpfs`
- **Падение** процесса — неконсистентность
- **Snapshot**

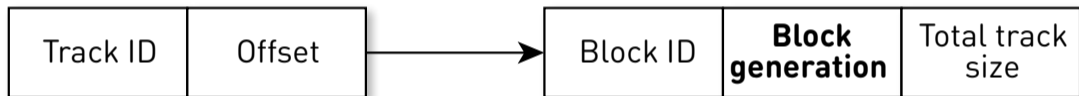


# Durability

- **Рестарт** машины — теряем tmpfs
- **Падение** процесса — неконсистентность
- **Snapshot**
- **WAL**

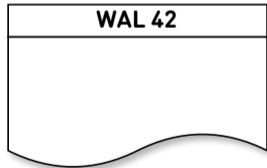
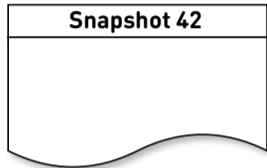


# Порядок операций



# Snapshot + WAL


Track ID	Offset	Block ID	Generation	Track size
1.mp3	0	3	17	10 MB
1.mp3	1	7	18	10 MB
2.mp3	0	9	21	15 MB





# Пишем в snapshot

Track ID	Offset	Block ID	Generation	Track size
<b>1.mp3</b>	<b>0</b>	<b>3</b>	<b>17</b>	<b>10 MB</b>
1.mp3	1	7	18	10 MB
2.mp3	0	9	21	15 MB




Snapshot 42
1.mp3,0 → 3,17,10MB*

WAL 42
--------



# Пишем в snapshot

Track ID	Offset	Block ID	Generation	Track size
1.mp3	0	3	17	10 MB
<b>1.mp3</b>	<b>1</b>	<b>7</b>	<b>18</b>	<b>10 MB</b>
2.mp3	0	9	21	15 MB



Snapshot 42
1.mp3,0 → 3,17,10MB 1.mp3,1 → 7,18,10MB*

WAL 42
--------



# Вытеснение в WAL

Track ID	Offset	Block ID	Generation	Track size
1.mp3	0	3	17	10 MB
		<b>7</b>	<b>18</b>	
2.mp3	0	9	21	15 MB

Snapshot 42
1.mp3,0 → 3,17,10MB 1.mp3,1 → 7,18,10MB

WAL 42
7,18 dirty*



# Вытеснение в WAL

Track ID	Offset	Block ID	Generation	Track size
1.mp3	0	3	17	10 MB
		7	18	
		<b>9</b>	<b>21</b>	

Snapshot 42
1.mp3,0 → 3,17,10MB 1.mp3,1 → 7,18,10MB

WAL 42
7,18 dirty 9,21 dirty*



# Переиспользование блока

Track ID	Offset	Block ID	Generation	Track size
1.mp3	0	3	17	10 MB
		7	18	
<b>3.mp3</b>	<b>0</b>	<b>9</b>	<b>22</b>	<b>8 MB</b>

Snapshot 42
1.mp3,0 → 3,17,10MB 1.mp3,1 → 7,18,10MB

WAL 42
7,18 dirty 9,21 dirty



# Заканчиваем snapshot

Track ID	Offset	Block ID	Generation	Track size
1.mp3	0	3	17	10 MB
		7	18	
<b>3.mp3</b>	<b>0</b>	<b>9</b>	<b>22</b>	<b>8 MB</b>

Snapshot 42
1.mp3,0 → 3,17,10MB
1.mp3,1 → 7,18,10MB
3.mp3,0 → 9,22,8MB*

WAL 42
7,18 dirty
9,21 dirty



# В итоге

Track ID	Offset	Block ID	Generation	Track size
1.mp3	0	3	17	10 MB
3.mp3	0	9	22	8 MB

Snapshot 42
1.mp3,0 → 3,17,10MB
1.mp3,1 → 7,18,10MB
3.mp3,0 → 9,22,8MB

WAL 42
7,18 dirty
9,21 dirty



# Восстанавливаем

Track ID	Offset	Block ID	Generation	Track size
1.mp3	0	3	17	10 MB
3.mp3	0	9	22	8 MB

Track ID	Offset	Block ID	Generation	Track size
----------	--------	----------	------------	------------

Snapshot 42
1.mp3,0 → 3,17,10MB
1.mp3,1 → 7,18,10MB
3.mp3,0 → 9,22,8MB

WAL 42
7,18 dirty
9,21 dirty





# Сканируем WAL

Track ID	Offset	Block ID	Generation	Track size
1.mp3	0	3	17	10 MB
3.mp3	0	9	22	8 MB

Track ID	Offset	Block ID	Generation	Track size
----------	--------	----------	------------	------------

Snapshot 42	
1.mp3,0	→ 3,17,10MB
1.mp3,1	→ 7,18,10MB
3.mp3,0	→ 9,22,8MB

Dirty	
Block ID	Generation
7	18
9	21



# Обычная запись

Track ID	Offset	Block ID	Generation	Track size
1.mp3	0	3	17	10 MB
3.mp3	0	9	22	8 MB

Track ID	Offset	Block ID	Generation	Track size
1.mp3	0	3	17	10 MB

Snapshot 42	
1.mp3,0	→ 3,17,10MB*
1.mp3,1	→ 7,18,10MB
3.mp3,0	→ 9,22,8MB

Dirty	
Block ID	Generation
7	18
9	21



# Вытесненная запись

Track ID	Offset	Block ID	Generation	Track size
1.mp3	0	3	17	10 MB
3.mp3	0	9	22	8 MB

Track ID	Offset	Block ID	Generation	Track size
1.mp3	0	3	17	10 MB

Snapshot 42	
1.mp3,0	→ 3,17,10MB
1.mp3,1	→ 7,18,10MB*
3.mp3,0	→ 9,22,8MB

Dirty	
Block ID	Generation
7	18
9	21



# Переиспользованный блок

Track ID	Offset	Block ID	Generation	Track size
1.mp3	0	3	17	10 MB
3.mp3	0	9	22	8 MB

Track ID	Offset	Block ID	Generation	Track size
1.mp3	0	3	17	10 MB
<b>3.mp3</b>	<b>0</b>	<b>9</b>	<b>22</b>	<b>8 MB</b>

Snapshot 42	
1.mp3,0	→ 3,17,10MB
1.mp3,1	→ 7,18,10MB
3.mp3,0	→ 9,22,8MB*

Dirty	
Block ID	Generation
7	18
<b>9</b>	<b>21</b>



# Успех

Track ID	Offset	Block ID	Generation	Track size
1.mp3	0	3	17	10 MB
3.mp3	0	9	22	8 MB

Track ID	Offset	Block ID	Generation	Track size
1.mp3	0	3	17	10 MB
3.mp3	0	9	22	8 MB

Snapshot 42	
1.mp3,0	→ 3,17,10MB
1.mp3,1	→ 7,18,10MB
3.mp3,0	→ 9,22,8MB

Dirty	
Block ID	Generation
7	18
9	21



# Уберите детей от экранов

## Disclaimer

Your mileage **will** vary.



# Page cache

- < **20%** попаданий<sup>9</sup>

---

<sup>9</sup>fs/cachestat from <https://github.com/brendangregg/perf-tools>

<sup>10</sup><https://github.com/odnoklassniki/one-nio>



# Page cache

- < **20%** попаданий<sup>9</sup>
- Read ahead?

---

<sup>9</sup>fs/cachestat from <https://github.com/brendangregg/perf-tools>

<sup>10</sup><https://github.com/odnoklassniki/one-nio>





# Page cache

- < **20%** попаданий<sup>9</sup>
- Read ahead?

`Mem.posix_fadvise(..., POSIX_FADV_RANDOM)`<sup>10</sup>:

---

<sup>9</sup>fs/cachestat from <https://github.com/brendangregg/perf-tools>

<sup>10</sup><https://github.com/odnoklassniki/one-nio>



# Page cache

- < **20%** попаданий<sup>9</sup>
- Read ahead?

`Mem.posix_fadvise(..., POSIX_FADV_RANDOM)`<sup>10</sup>:

- > **70%** попаданий
- **В 2 раза** меньше чтений с диска
- HTTP latency **на 20% ниже**

---

<sup>9</sup>fs/cachestat from <https://github.com/brendangregg/perf-tools>

<sup>10</sup><https://github.com/odnoklassniki/one-nio>



# -XX : +UseHugeTLBFS

- -Xmx24g



# -XX : +UseHugeTLBFS

- -Xmx24g
- GC Time/Safepoint Total Time **на 20-30% ниже**
- Более **равномерная** загрузка ядер



## -XX : +UseHugeTLBFS

- -Xmx24g
- GC Time/Safepoint Total Time **на 20-30% ниже**
- Более **равномерная** загрузка ядер
- Но **нет эффекта** на HTTP latency



# Инцидент

- 1 В саппорте жалобы



# Инцидент

- 1 В саппорте жалобы
- 2 Сузили круг поиска до одной машины



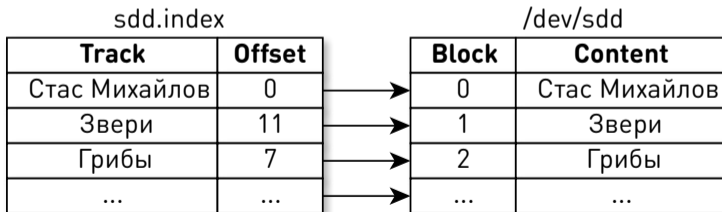
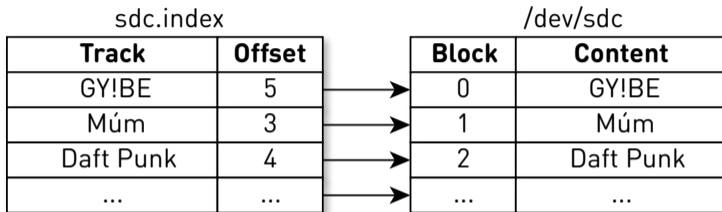
# Инцидент

- 1 В саппорте жалобы
- 2 Сузили круг поиска до одной машины
- 3 Машину недавно перезапускали...

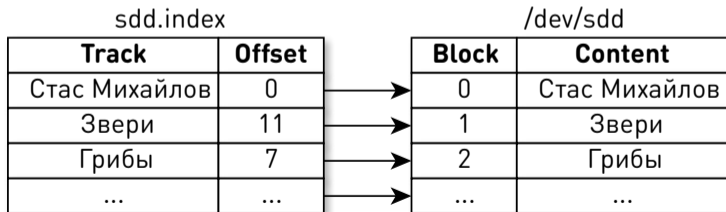
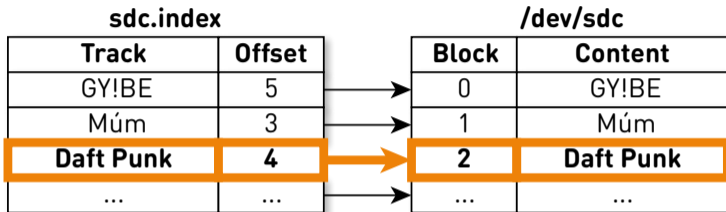




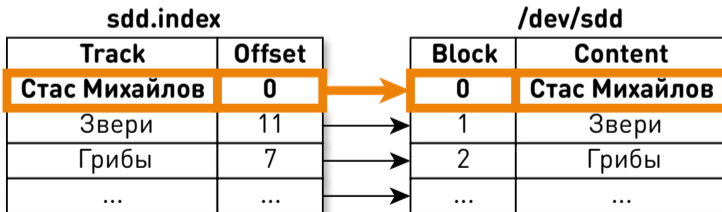
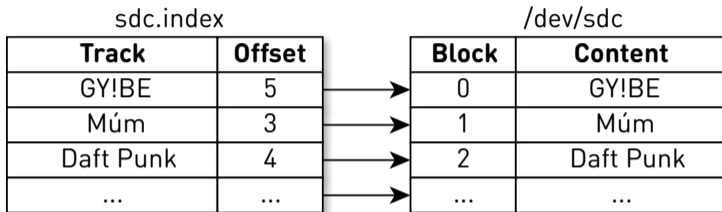
# До перезагрузки



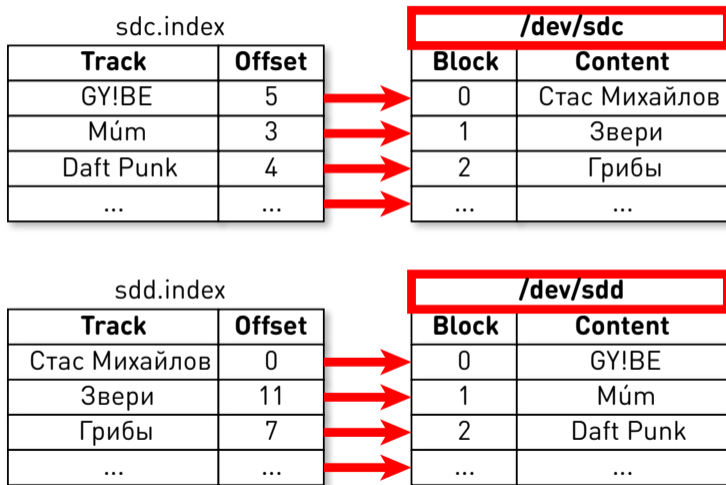
# До перезагрузки



# До перезагрузки

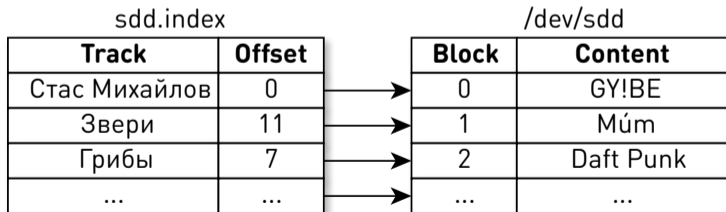
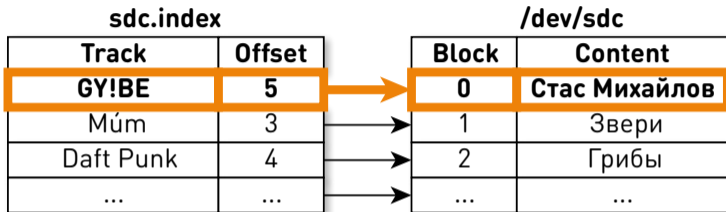


# После перезагрузки<sup>11</sup>

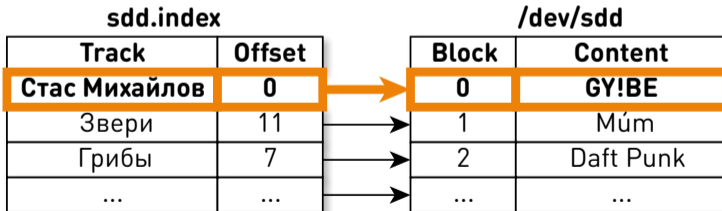
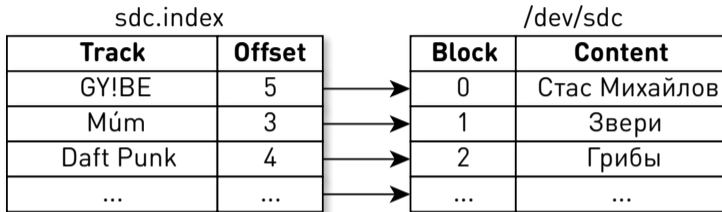


<sup>11</sup>[https://wiki.archlinux.org/index.php/Persistent\\_block\\_device\\_naming](https://wiki.archlinux.org/index.php/Persistent_block_device_naming)

# Party mode



# Party mode



# Fix

- **Персистентный World Wide Name**<sup>12</sup>

---

<sup>12</sup>[https://en.wikipedia.org/wiki/World\\_Wide\\_Name](https://en.wikipedia.org/wiki/World_Wide_Name)



# Fix

- **Персистентный World Wide Name**<sup>12</sup>
- `/dev/disk/by-id/wwn-*` в качестве **ID**
  - Индексов
  - Snapshot
  - WAL

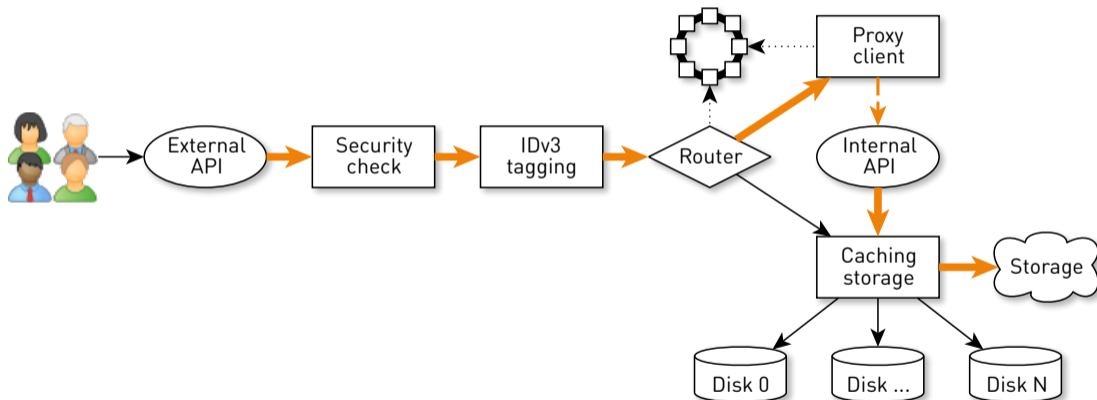
---

<sup>12</sup>[https://en.wikipedia.org/wiki/World\\_Wide\\_Name](https://en.wikipedia.org/wiki/World_Wide_Name)





# Распределённая отладка



# TraceID<sup>13</sup>-метка в запросах

```
1 GET /v0/stream?id=v0_10000051227 HTTP/1.1
2 Range: bytes=10878976-11141119
3 X-OK-Trace-ID: 570752372832:bec:v0_10000051227
```

---

<sup>13</sup>See [OpenTracing/Zipkin](#)



# При отладке TraceID в логах

```
1 2018-09-13 17:40:40,147 DEBUG
2 [cache-provider-4] CachingDataProvider:
3 -> [570752372832:bec:v0_10000051227]
4 Serving chunk from cache:
5 Chunk[10878976-11010047/16346382]
```



# Как отправлять данные?

```
1 ByteBuffer buffer = ByteBuffer.allocate(size);
2 int count = fileChannel.read(buffer, position);
3 if (count <= 0) {
4     // ...
5 }
6 buffer.flip();
7 socketChannel.write(buffer);
```



## SYNOPSIS

```
#include <sys/sendfile.h>
```

```
ssize_t sendfile(int out_fd, int in_fd, off_t *offset, size_t count);
```

## DESCRIPTION

**sendfile()** copies data between one file descriptor and another. Because this copying is done within the kernel, **sendfile()** is more efficient than the combination of **read(2)** and **write(2)**, which would require transferring data to and from user space.

# sendfile()

```
1 final RandomAccessFile blocks;  
2  
3 int write(Socket socket) {  
4     if (socket.getSslContext() == null) {  
5         socket.sendFile(blocks, offset, size);  
6     }
```

- **Zero-copy**, но только не SSL



# В случае SSL

```
1 ByteBuffer buffer = ByteBuffer.allocate(size);
2 int count = fileChannel.read(buffer, position);
3 if (count <= 0) {
4     // ...
5 }
6 buffer.flip();
7 socketChannel.write(buffer);
```



# IOUtil

```
1 ByteBuffer bb =
    Util.getTemporaryDirectBuffer(dst.remaining());
2 try {
3     int n = readIntoNativeBuffer(fd, bb, position, nd);
4     bb.flip();
5     if (n > 0)
6         dst.put(bb);
7     return n;
8 } finally {
9     Util.offerFirstTemporaryDirectBuffer(bb);
10 }
```





# one-nio: DirectMemory и друзья

```
1 final Allocator allocator =  
2     new MallocMT(size, concurrency);  
3  
4  
5  
6  
7  
8  
9  
10 . . .
```



# Если SSL...

```
1 final Allocator allocator =  
2     new MallocMT(size, concurrency);  
3  
4 int write(Socket socket) {  
5     if (socket.getSslContext() != null) {  
6  
7  
8  
9  
10    ...
```



# Выделяем нативный буфер

```
1 final Allocator allocator =  
2     new MallocMT(size, concurrency);  
3  
4 int write(Socket socket) {  
5     if (socket.getSslContext() != null) {  
6         long address = allocator.malloc(size);  
7  
8  
9  
10    ...
```



# Оборачиваем в ByteBuffer

```
1 final Allocator allocator =
2     new MallocMT(size, concurrency);
3
4 int write(Socket socket) {
5     if (socket.getSslContext() != null) {
6         long address = allocator.malloc(size);
7         ByteBuffer buf =
8             DirectMemory.wrap(address, size);
9
10    ...
```



# Читаем

```
1 final Allocator allocator =
2     new MallocMT(size, concurrency);
3
4 int write(Socket socket) {
5     if (socket.getSslContext() != null) {
6         long address = allocator.malloc(size);
7         ByteBuffer buf =
8             DirectMemory.wrap(address, size);
9         int available = channel.read(buf, offset);
10    ...
```

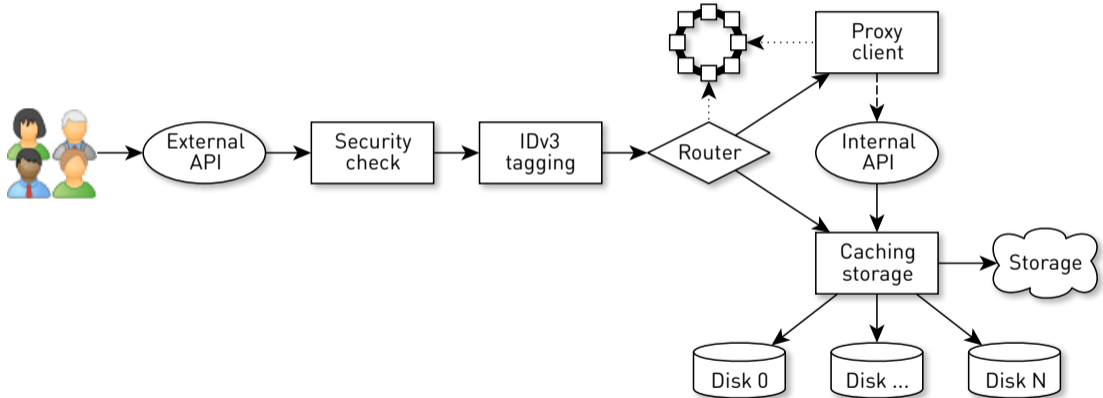


# Пишем

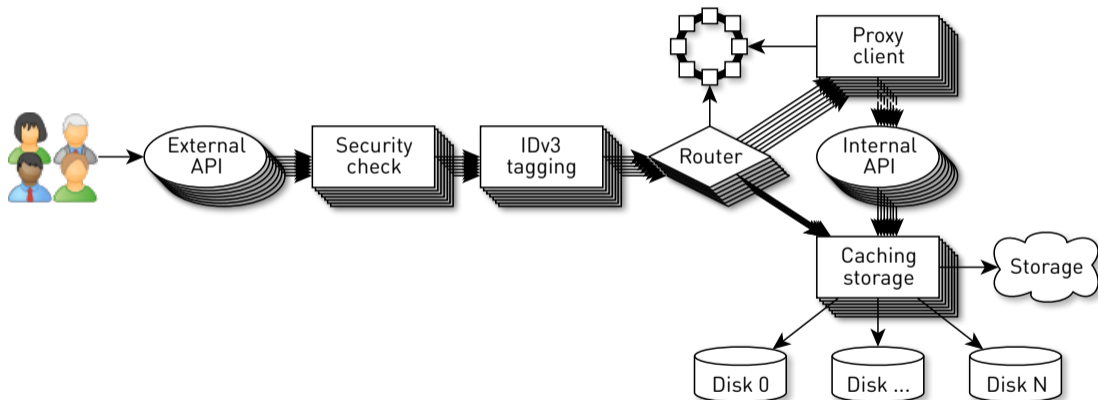
```
1 final Allocator allocator =
2     new MallocMT(size, concurrency);
3
4 int write(Socket socket) {
5     if (socket.getSslContext() != null) {
6         long address = allocator.malloc(size);
7         ByteBuffer buf =
8             DirectMemory.wrap(address, size);
9         int available = channel.read(buf, offset);
10        socket.writeRaw(address, available, flags);
```



# One more thing...



# Очень много клиентов





# Варианты<sup>14</sup>

Клиент = Thread:

- **Синхронно**

---

<sup>14</sup><https://thestrangeloop.com/2017/zuuls-journey-to-non-blocking.html>



# Варианты<sup>14</sup>

Клиент = Thread:

- **Синхронно**
- **100К** ПОТОКОВ?!

---

<sup>14</sup><https://thestrangeloop.com/2017/zuuls-journey-to-non-blocking.html>



# Варианты<sup>14</sup>

Клиент = Thread:

- **Синхронно**
- **100К** потоков?!
- **Нежизнеспособно**

---

<sup>14</sup><https://thestrangeloop.com/2017/zuuls-journey-to-non-blocking.html>



# Варианты<sup>14</sup>

Клиент = Thread:

- **Синхронно**
- **100К** потоков?!
- **Нежизнеспособно**

Клиент = конвейер:

- **Асинхронно**

---

<sup>14</sup><https://thestrangeloop.com/2017/zuuls-journey-to-non-blocking.html>



# Варианты<sup>14</sup>

Клиент = Thread:

- **Синхронно**
- **100К** потоков?!
- **Нежизнеспособно**

Клиент = конвейер:

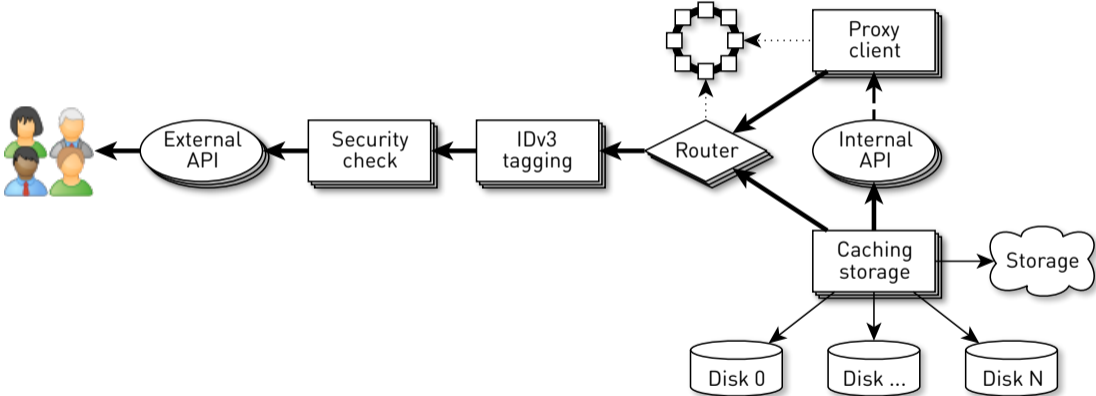
- **Асинхронно**
- Стадия — пул потоков
- **Push** данных клиенту

---

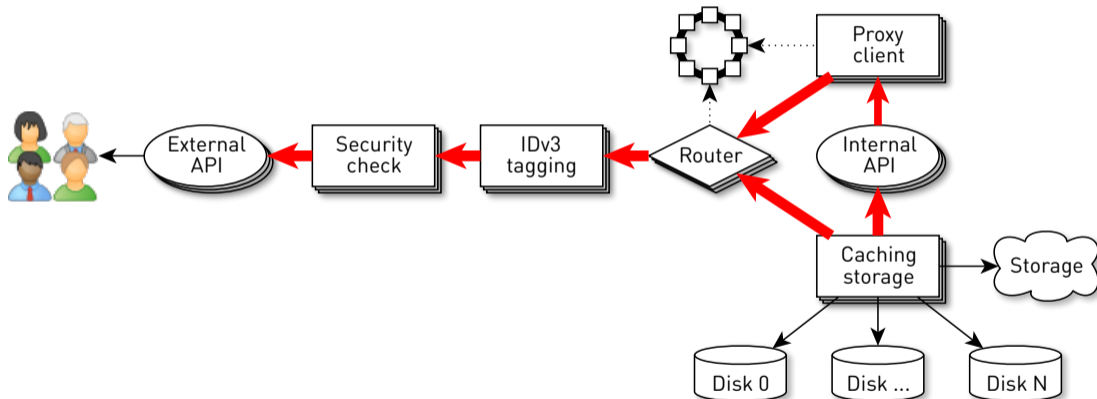
<sup>14</sup><https://thestrangeloop.com/2017/zuuls-journey-to-non-blocking.html>



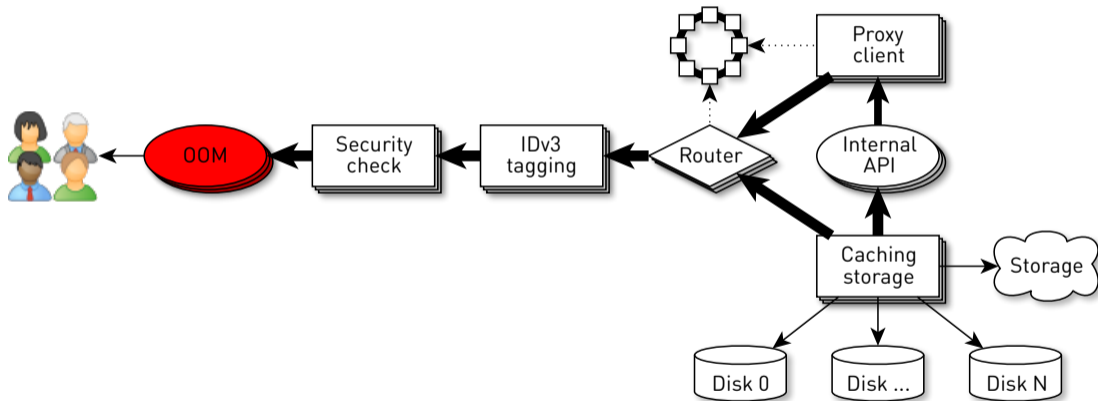
# Асинхронный push



# Бекенды гораздо быстрее

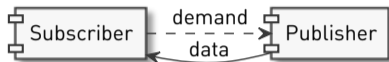


# Если неограниченные очереди





# Reactive streams<sup>15</sup> to the rescue

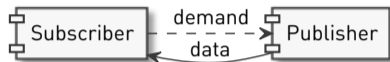


- **Subscriber** управляет скоростью
- Память ограничена **demand**

---

<sup>15</sup><http://www.reactive-streams.org>

# Reactive streams<sup>15</sup> to the rescue

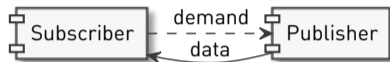


- **Subscriber** управляет скоростью
- Память ограничена **demand**
- Subscriber быстрее — **push**

---

<sup>15</sup><http://www.reactive-streams.org>

# Reactive streams<sup>15</sup> to the rescue

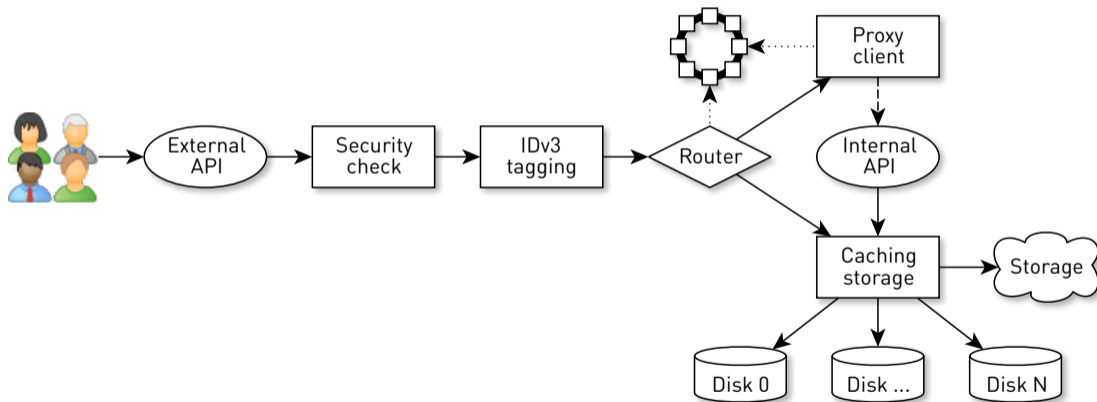


- **Subscriber** управляет скоростью
- Память ограничена **demand**
- Subscriber быстрее — **push**
- Publisher быстрее — **pull**

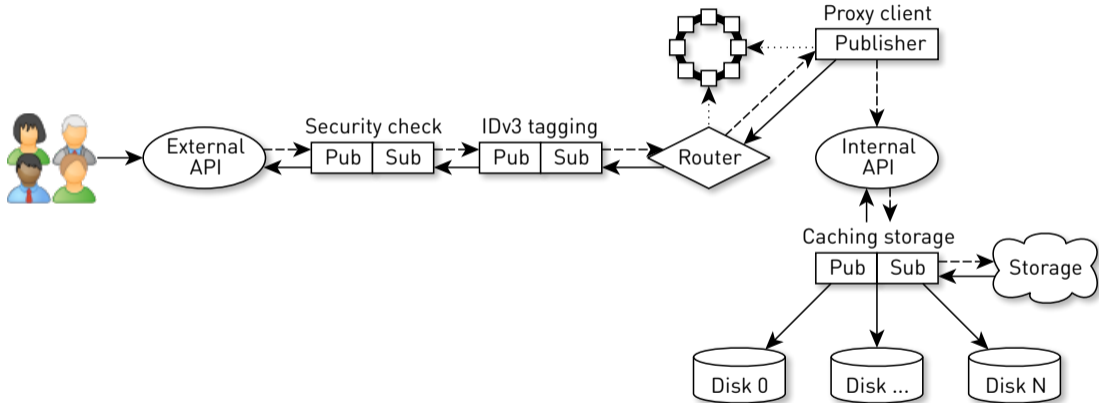
---

<sup>15</sup><http://www.reactive-streams.org>

# Конвейер



# Reactive Stream



```
1 interface Publisher<T> {
2     void subscribe(Subscriber<? super T> s);
3 }
4 interface Subscriber<T> {
5     void onSubscribe(Subscription s);
6     void onNext(T t);
7     void onError(Throwable t);
8     void onComplete();
9 }
10 interface Subscription {
11     void request(long n);
12     void cancel();
```



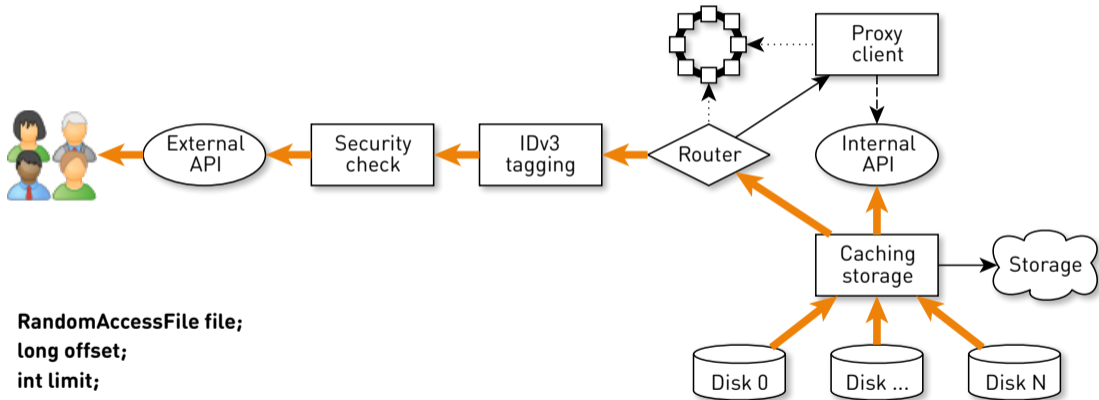
# Поток Chunkов

```
1 interface Chunk {  
2     int read(ByteBuffer dst);  
3     int write(Socket socket);  
4     void write(FileChannel channel, long offset);  
5 }
```

- **Ссылка** на данные
- Ограниченный интерфейс
- Множество реализаций



# Chunk over RandomAccessFile

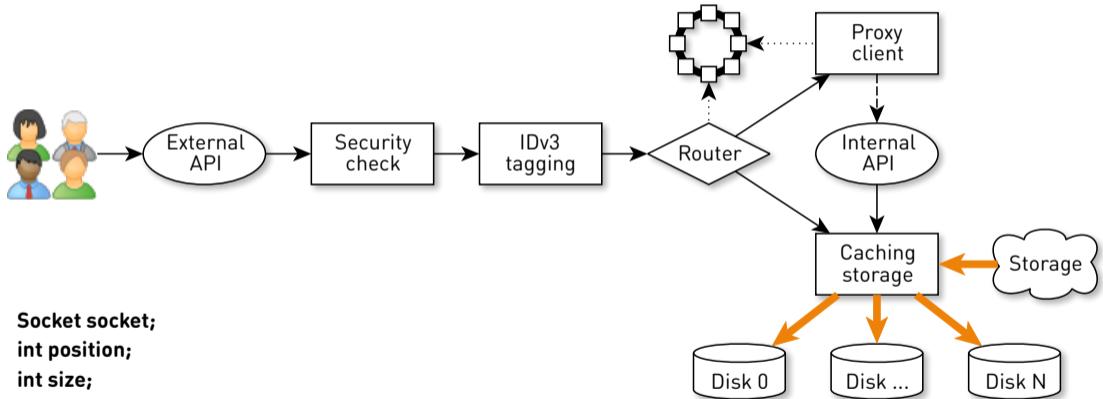


**RandomAccessFile file;**  
**long offset;**  
**int limit;**





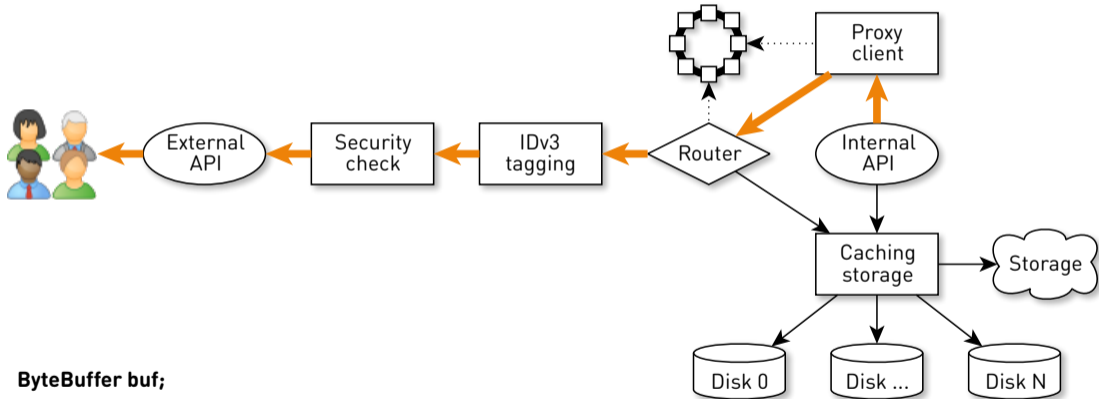
# Chunk over Socket



Socket socket;  
int position;  
int size;



# Chunk over ByteBuffer



**ByteBuffer buf;**



# Неблокирующий API<sup>16</sup>

```
1 interface Subscriber<T> {
2     void onSubscribe(Subscription s);
3     void onNext(T t);
4     void onError(Throwable t);
5     void onComplete();
6 }
7 interface Subscription {
8     void request(long n);
9     void cancel();
```

---

<sup>16</sup><https://jokerconf.com/2018/talks/6wmp33pmjgism04u4wckgy/>



# Чтение на ночь

## [concurrency-interest] Synchronization primitives in Reactive Streams implementations

Pavel Rappo [pavel.rappo@gmail.com](mailto:pavel.rappo@gmail.com)

*Fri Sep 28 08:51:58 EDT 2018*

- Previous message (by thread): [\[concurrency-interest\] JLS 17.5 Examples](#)
- Next message (by thread): [\[concurrency-interest\] Synchronization primitives in Reactive Streams implementations](#)
- **Messages sorted by:** [\[\\_date\\_\]](#) [\[\\_thread\\_\]](#) [\[\\_subject\\_\]](#) [\[\\_author\\_\]](#)

---

Hello,

Let me start this discussion as a branch of the "Reactive Streams utility API" thread here:

<http://mail.openjdk.java.net/pipermail/core-libs-dev/2018-September/055671.html>

I would like to talk about a particular synchronization primitive I find to be repeatedly reinvented in Reactive Streams implementations. I hope that we could discuss different incarnations of this primitive and maybe extract it into a reusable component.

# В духе Typed Actor Model<sup>17</sup>

---

<sup>17</sup><https://github.com/reactive-streams/reactive-streams-jvm/tree/master/examples>



# В духе Typed Actor Model<sup>17</sup>

- Вызов метода → сообщение
- Сообщение — в Queue
- Шедулимся на Executor

---

<sup>17</sup><https://github.com/reactive-streams/reactive-streams-jvm/tree/master/examples>



# В духе Typed Actor Model<sup>17</sup>

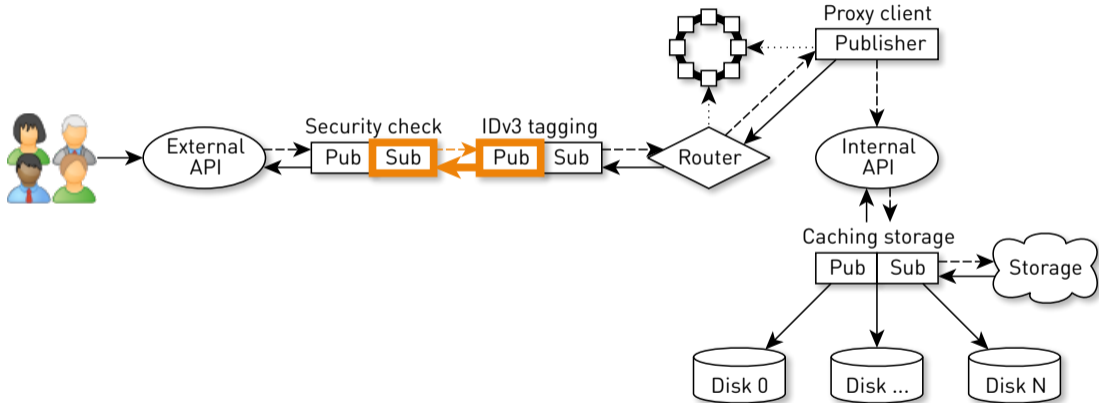
- Вызов метода → сообщение
- Сообщение — в Queue
- Шедулимся на Executor
- Обрабатываем **последовательно**

---

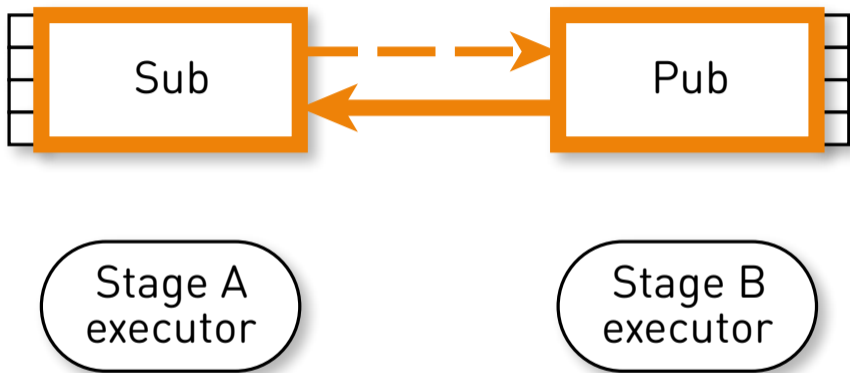
<sup>17</sup><https://github.com/reactive-streams/reactive-streams-jvm/tree/master/examples>

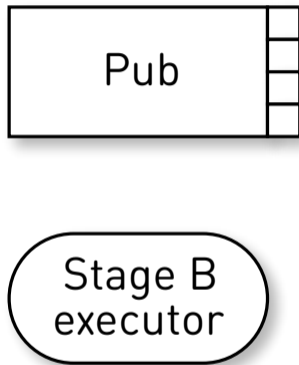
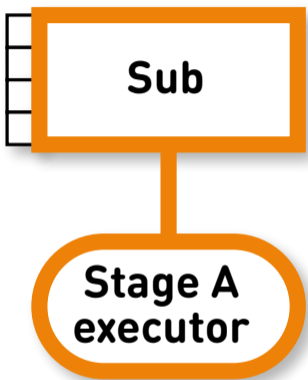


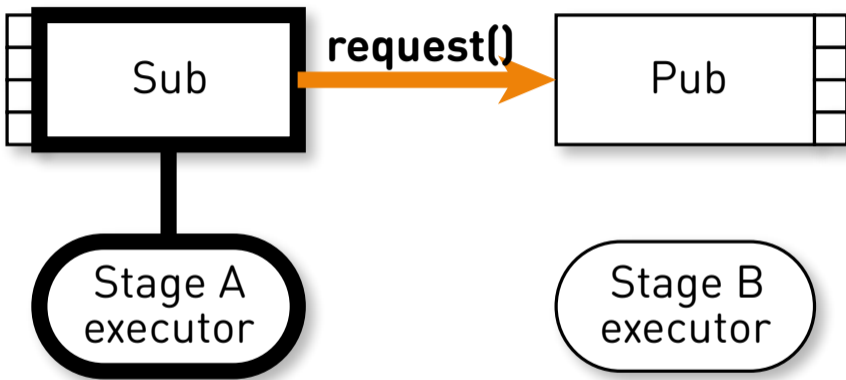
# Reactive Stream

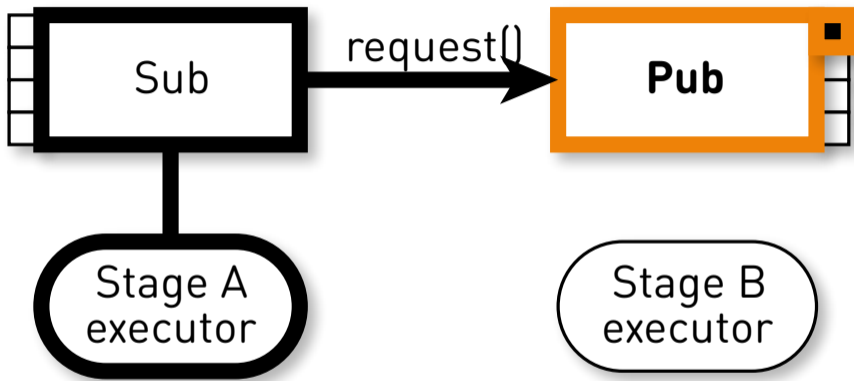


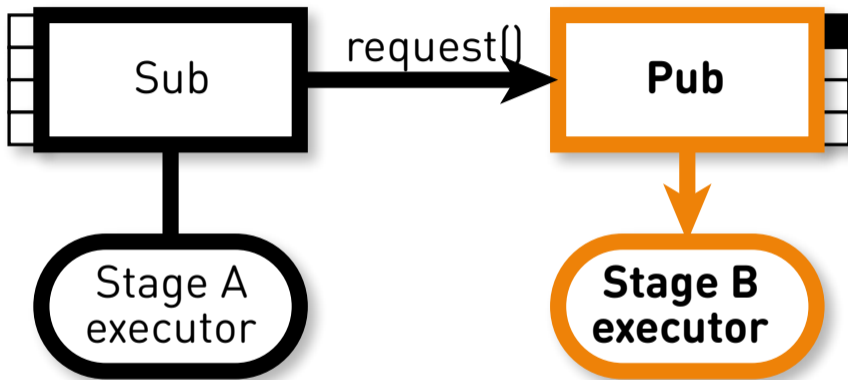


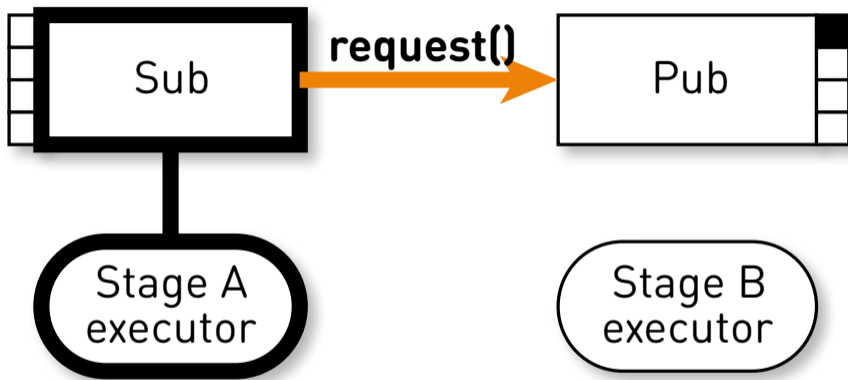


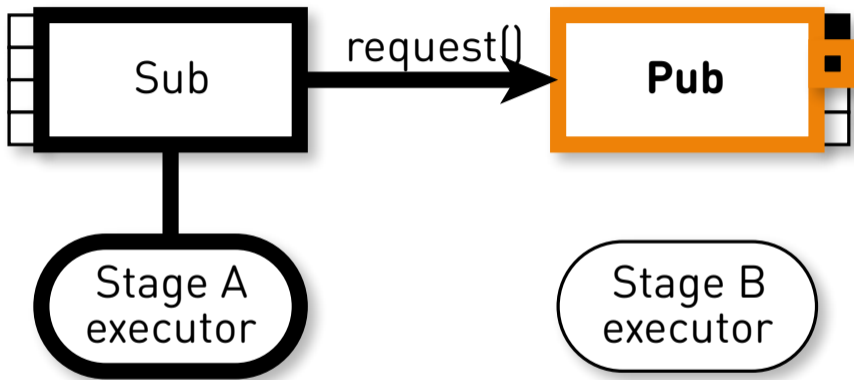


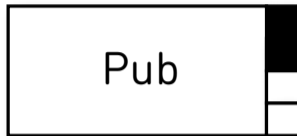
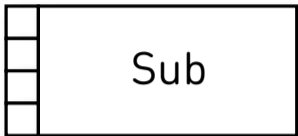




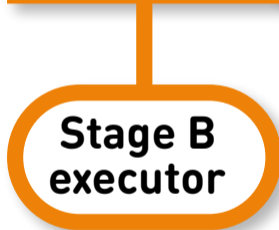
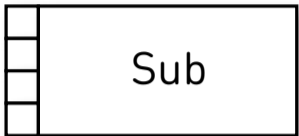


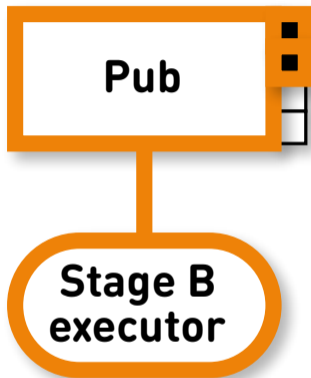
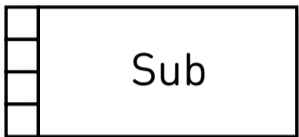


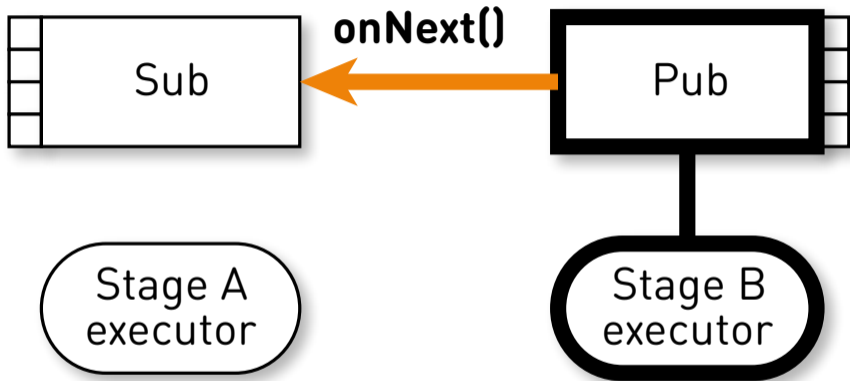


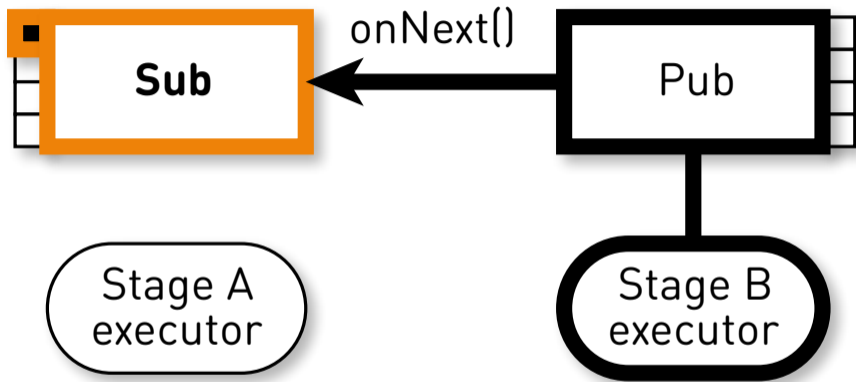


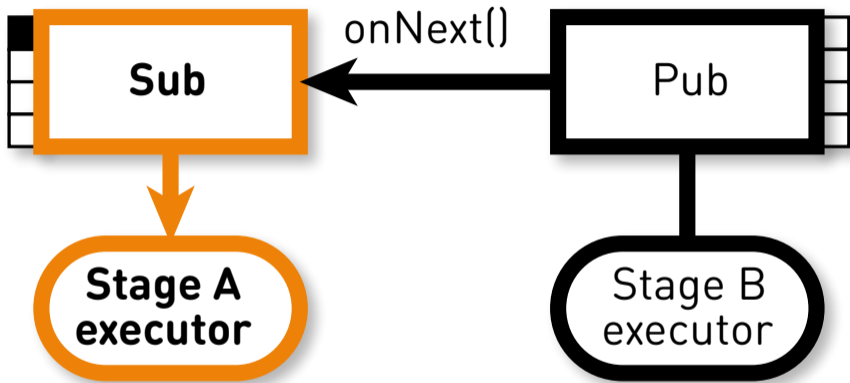


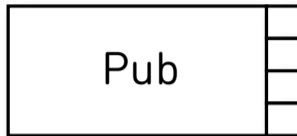
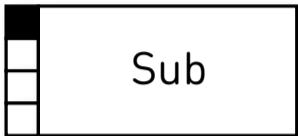












# Состояние

```
1 // Incoming messages
2 final Queue<M> mailbox;
3
4 // Message processing works here
5 final Executor executor;
6
7 // To ensure HB relationship between runs
8 final AtomicBoolean on = new AtomicBoolean();
```



# Асинхронизация

```
1 @Override
2 void request(final long n) {
3     enqueue(new Request(n));
4 }
5
6 void enqueue(M message) {
7     mailbox.offer(message);
8     tryScheduleToExecute();
9 }
```





# tryScheduleToExecute()

```
1 if (on.compareAndSet(false, true)) {  
2   try {  
3     executor.execute(this);  
4   } catch (Exception e) {  
5     ...  
6   }  
7 }
```



# run()

```
1 if (on.get())
2   try {
3     dequeueAndProcess();
4   } finally {
5     on.set(false);
6     if (!messages.isEmpty()) {
7       tryScheduleToExecute();
8     }
9   }
10 }
```



# dequeueAndProcess()

```
1 M message;
2 while ((message = mailbox.poll()) != null) {
3     // Pattern match
4     if (message instanceof Request) {
5         doRequest(((Request) message).n);
6     } else {
7         ...
8     }
9 }
```



# В итоге

- **Неблокирующая** реализация



# В итоге

- **Неблокирующая** реализация
- Простой **последовательный** код



# В итоге

- **Неблокирующая** реализация
- Простой **последовательный** код
- Никакого contention



# В итоге

- **Неблокирующая** реализация
- Простой **последовательный** код
- Никакого contention
- **Ограниченное** количество потоков



# Production

- **12 машин (2x запас)**





# Production

- **12 машин** (2x запас)
- До **20 Гб/с** с машины через 100К соединений



# Production

- **12 машин (2x запас)**
- До **20 Гб/с** с машины через 100К соединений
- Масштабируемость
  - Пропускная способность
  - Ёмкость



# Production

- **12 машин (2x запас)**
- До **20 Гб/с** с машины через 100К соединений
- Масштабируемость
  - Пропускная способность
  - Ёмкость
- Отказоустойчивость
  - ДЦ
  - Машины
  - Диски

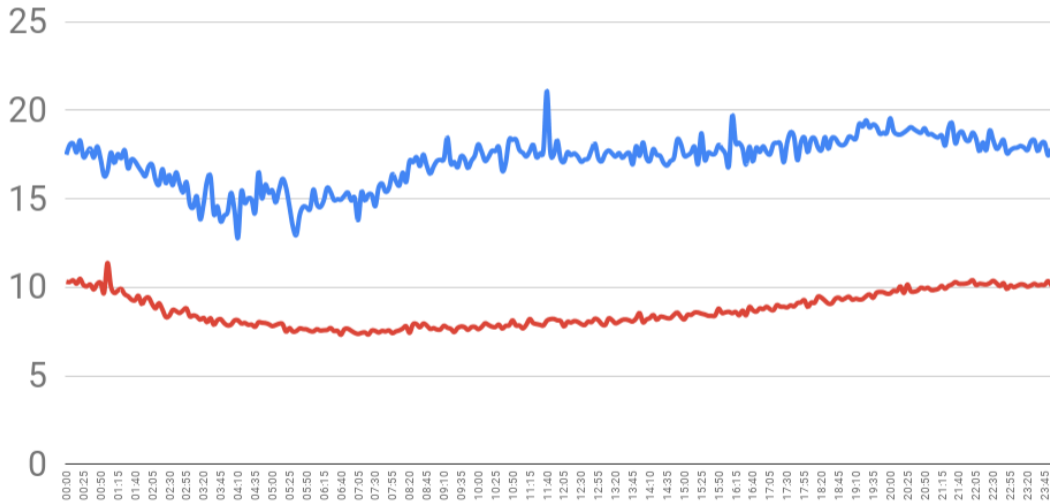


# Production

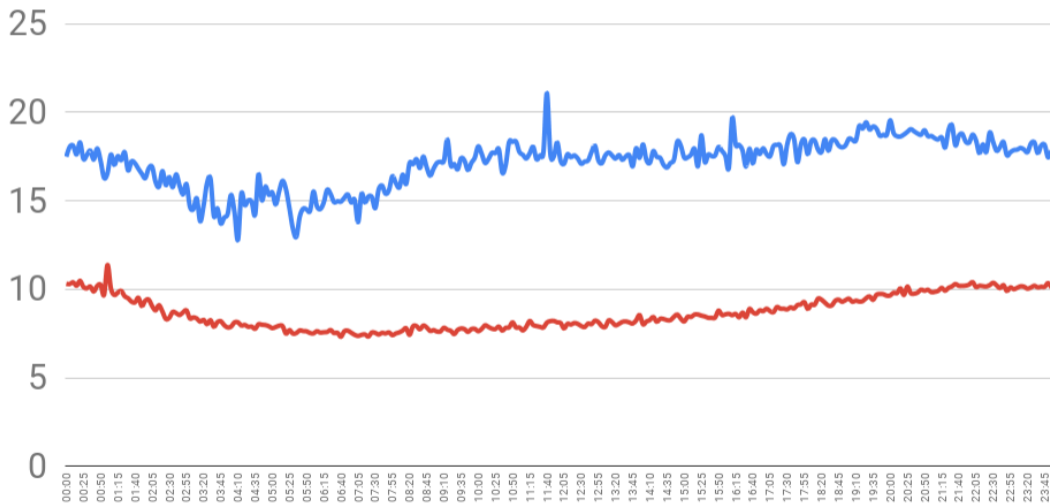
- **12 машин (2x запас)**
- До **20 Гб/с** с машины через 100К соединений
- Масштабируемость
  - Пропускная способность
  - Ёмкость
- Отказоустойчивость
  - ДЦ
  - Машины
  - Диски
- Java + one-nio



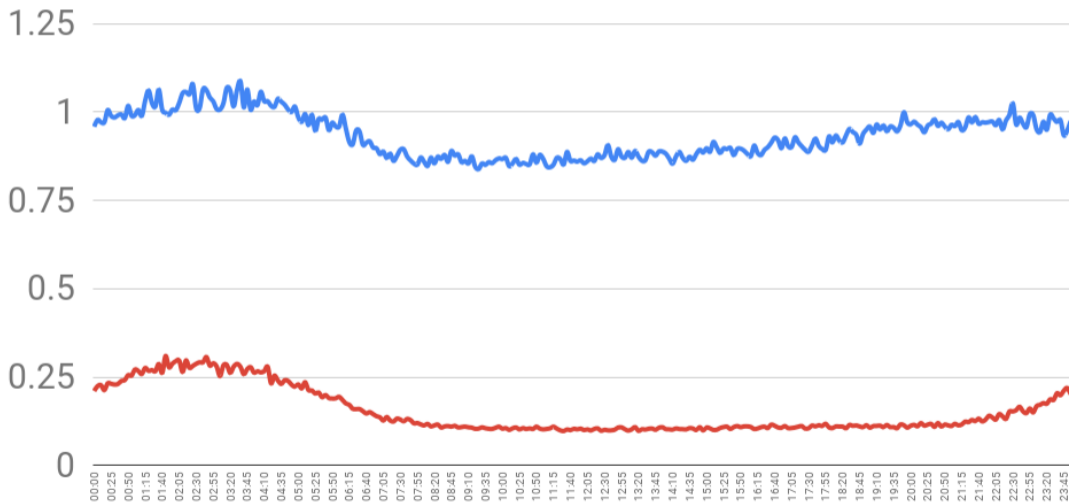
99%, MC



# 99%, mc (upstream)



# 75%, мс (с дисков)



**High Load**  
+  
**Reactive Streams**  
=  




# Вопросы?



**Вадим Цесько**

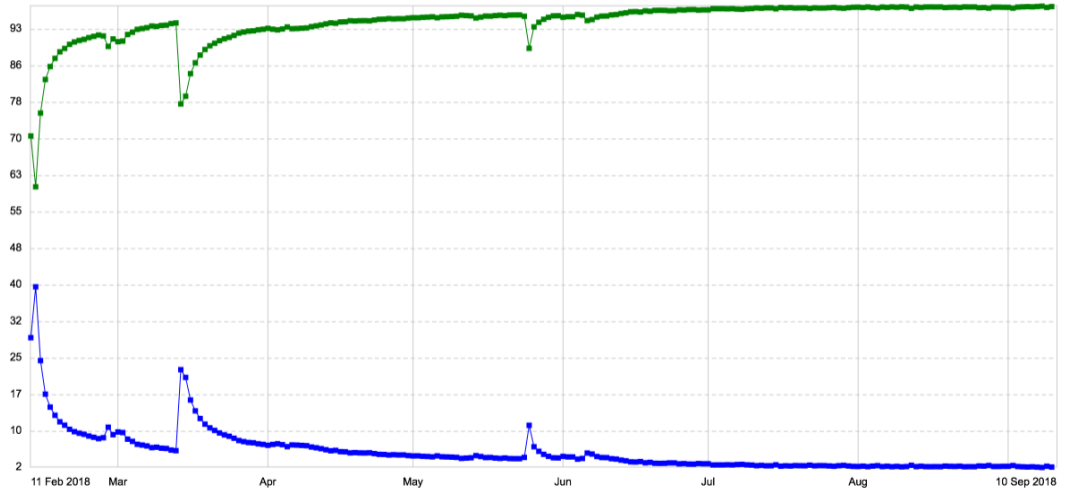
**Одноклассники**

<https://v.ok.ru>

<https://incubos.org>



# CDN cache hit



# Backlog

- Адаптивный фактор репликации

---

<sup>18</sup>[https://en.wikipedia.org/wiki/Rendezvous\\_hashing](https://en.wikipedia.org/wiki/Rendezvous_hashing)



# Backlog

- Адаптивный фактор репликации
- Read repair

---

<sup>18</sup>[https://en.wikipedia.org/wiki/Rendezvous\\_hashing](https://en.wikipedia.org/wiki/Rendezvous_hashing)



# Backlog

- Адаптивный фактор репликации
- Read repair
- Touch/prefetch

---

<sup>18</sup>[https://en.wikipedia.org/wiki/Rendezvous\\_hashing](https://en.wikipedia.org/wiki/Rendezvous_hashing)



# Backlog

- Адаптивный фактор репликации
- Read repair
- Touch/prefetch
- Hash ring поверх дисков

---

<sup>18</sup>[https://en.wikipedia.org/wiki/Rendezvous\\_hashing](https://en.wikipedia.org/wiki/Rendezvous_hashing)



# Backlog

- Адаптивный фактор репликации
- Read repair
- Touch/prefetch
- Hash ring поверх дисков
- Rendezvous hashing<sup>18</sup>

---

<sup>18</sup>[https://en.wikipedia.org/wiki/Rendezvous\\_hashing](https://en.wikipedia.org/wiki/Rendezvous_hashing)



# Backlog

- Адаптивный фактор репликации
- Read repair
- Touch/prefetch
- Hash ring поверх дисков
- Rendezvous hashing<sup>18</sup>
- Асинхронный HTTP клиент

---

<sup>18</sup>[https://en.wikipedia.org/wiki/Rendezvous\\_hashing](https://en.wikipedia.org/wiki/Rendezvous_hashing)





# Backlog

- Адаптивный фактор репликации
- Read repair
- Touch/prefetch
- Hash ring поверх дисков
- Rendezvous hashing<sup>18</sup>
- Асинхронный HTTP клиент
- Open-source one-download

---

<sup>18</sup>[https://en.wikipedia.org/wiki/Rendezvous\\_hashing](https://en.wikipedia.org/wiki/Rendezvous_hashing)

