

# Введение в Big Data и распределенные файловые системы

**Алексей Драль**, [aadral@bigdatateam.org](mailto:aadral@bigdatateam.org)

CEO at BigData Team, <http://bigdatateam.org>

<https://www.facebook.com/bigdatateam>

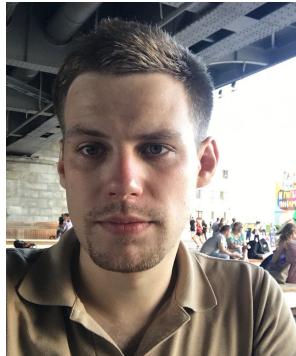


18.02.2021, Moscow, Russia

# Разработчики курса



Алексей Драль



Артем Выборнов



Андрей Титов



Антон Горохов



Павел Клеменков



**RAMBLER  
GROUP**



**NVIDIA.**



**Yandex**



**NVIDIA.**

и многие другие!



BIGDATA  
TEAM



Канал Telegram

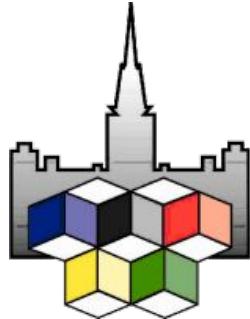


Виктория Брынза



**BIGDATA  
TEAM**

- ▶ СУНЦ МГУ
- ▶ МГУ
- ▶ ШАД Яндекса
- ▶ Рамблер
- ▶ Яндекс
- ▶ Amazon AWS
- ▶ МФТИ
- ▶ Сбербанк
- ▶ BigData Team



**RAMBLER  
GROUP**



**Yandex**



**Обо мне**



- 1.** Введение в Big Data. Распределенные файловые системы, HDFS
- 2.** Hadoop MapReduce
- 3.** Оптимизация MapReduce вычислений
- 4.** Hive: SQL поверх больших данных
- 5.** Модель вычислений Spark: RDD
- 6.** Spark DataFrames, Spark SQL
- 7.** Оптимизация Spark вычислений
- 8.** Потоковая обработка данных (Kafka, Spark Streaming)
- 9.** Cassandra + Spark
- 10.** Оптимизация хранилища (Data Layout)



*The 70-20-10 rule was developed by Morgan McCall, Robert W. Eichinger and Michael M. Lombardo at the Center for Creative Leadership.*

- ▶ 10 домашних заданий
- ▶ 2 тестирования



до  
20-30%



после  
70-80%

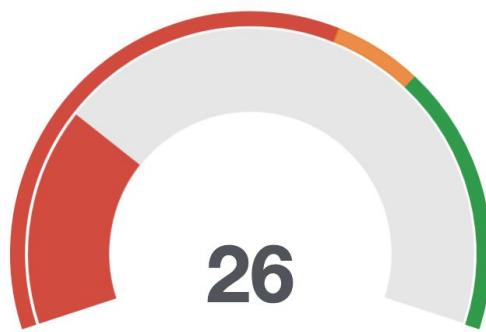
HDFS	MR	MR	MR	MR	Hive	Hive	Hive	Spark	Spark	Spark	Spark	Spark	RT	RT	RT	RT	NoSQL	NoSQL
0	0	0.75	0	0	0	0	0	1	1	0	0	1	0	0	0.66	0	0	0
0.5	0	1	0	0.6	0	1	0	1	0.75	0	0	1	0	0	0.66	0.66	1	1
0	0	1	0	0	0	0	0	1	0.5	0	0	0	0.5	0	0	0	0	0
0	1	0.75	0	0.4	0	0	0	1	0.75	1	0	0	0	1	0	0	0	0
0.75	0	1	0	0.8	0	0	0	1	0.75	1	1	1	1	0	0	0	0	0
0.25	0	1	0	0	0	1	0	1	0.5	1	0	0	0.5	0	0.66	0.66	0	0
0.75	0	1	0	0.8	0	0	0	1	0.75	0	0	0	0	0	0.99	0.66	1	0
0.75	1	1	1	0.8	0	0	1	1	0	1	0	0	0	0	0	0	0	0
0	0	0.75	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0.75	0	0.5	0	0	0	0	0	1	1	0	1	0	0.75	0	0.33	0.66	1	1
0.5	0	0.75	0	1	0	0	1	1	0.75	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0
0	0	1	0	1	0	0	0	0	1	1	1	1	0	0.75	0	0.99	0.33	1
0.75	0	1	0	0	0	0	0	0	0	0	0	0	1	1	0	0	0	1
0.25	0	0.75	1	0	0	0.75	0	1	0.75	1	1	1	1	1	0	0	0	1
1	0	0.5	0	0	0	0	0	0	0.75	0	0	0	0	0	0	0	0	1
0.75	0	1	0	0	0	0	0	1	0.75	0	0	0	1	1	0.66	0	0	0
0.5	0	1	0	1	1	1	0	1	0	1	1	0	1	0	0	0	0	0

HDFS	MR	MR	MR	MR	Hive	Hive	Hive	Spark	Spark	Spark	Spark	Spark	RT	RT	RT	RT	NoSQL	NoSQL		
1	0	1	1	0.8	1	1	1	1	1	1	1	1	0.6	1	1.00	1.00	1	0.25	0.75	
0.75	1	1	1	0.6	1	1	1	0.75	0.5	1	0.5	1	0.6	1	1.00	1.00	1	1	1	
1	0	1	1	1	1	0	1	0.5	0.5	0	0	1	1	0	0.67	1.00	0	0.5	0.75	
1	0	1	1	1	1	0	1	0.5	0.5	0	0	1	1	0	0.67	1.00	0	0.5	0.75	
0.5	1	1	1	1	1	1	1	0.75	0.25	1	0.75	1	0.8	1	1.00	1.00	0	0.75	1	
1	1	1	1	0	0.2	1	0	1	0.75	1	0	1	0	1	1.00	0.50	0	0.75	0.75	
0.5	1	0	0	0	0.8	0	0	0	0.5	0.75	0	0.75	1	0.6	0	0.67	0.50	0	0.5	0.75
0.5	1	0	1	0.8	1	1	0	0.75	0.75	1	0.5	1	0.8	1	1.00	0.50	1	0.5	1	
0.75	1	1	1	1	1	0	0	1	0.5	0	0.5	1	0.8	1	0.33	0.83	0	0.75	0.5	
1	1	1	1	0.8	1	1	1	0.75	1	0	1	1	0.6	1	1.00	1.00	0	0.75	1	
1	0	1	1	1	1	0	0	1	0.5	0	1	1	0.6	1	1.00	0.50	1	0.75	1	
0.75	0	1	1	1	1	1	1	0.75	1	0	0.5	0	0.8	1	1.00	1.00	1	0.5	1	
0.75	1	1	1	0.6	1	1	1	0.75	0	0	0.5	1	0	0	0.33	0.00	1	1	1	
1	0	1	1	0.8	1	0	1	0.75	0.25	0	0.5	1	1	1	0.67	1.00	1	1	1	
0.5	1	1	1	0.8	1	0	0	1	1	0	0.5	1	0.8	0	1.00	0.67	0	0.75	0.25	

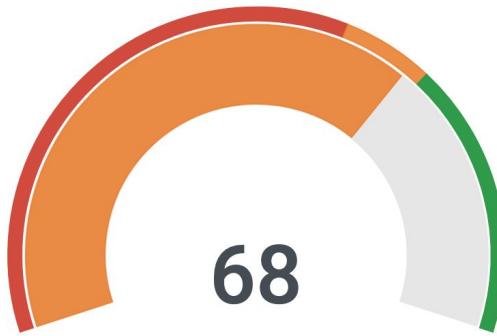
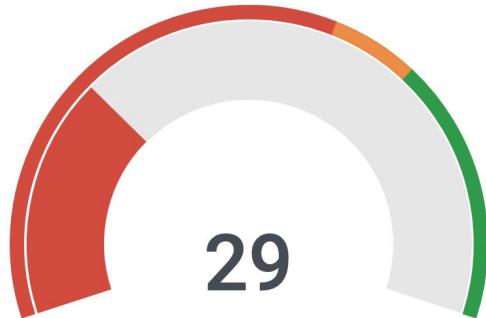


# Точка отсчета 21-Q1

**20-Q1**



**20-Q3**



Intro quiz: [https://rebrand.ly/x5bd2021q1\\_quiz\\_intro](https://rebrand.ly/x5bd2021q1_quiz_intro)

- ▶ Мотивация работы с Big Data
- ▶ Многопроцессорные вычислительные системы
- ▶ Распределенные файловые системы (GFS / HDFS)
- ▶ Hadoop Sizing и архитектура Namenode
- ▶ HDFS Workshop



# Мотивация работы с Big Data



# Data Scientist: The Sexiest Job of the 21st Century

by Thomas H. Davenport and D.J. Patil

FROM THE OCTOBER 2012 ISSUE





## Q&A

Откуда берутся (большие) **данные**?



# Internet of Things (IoT)

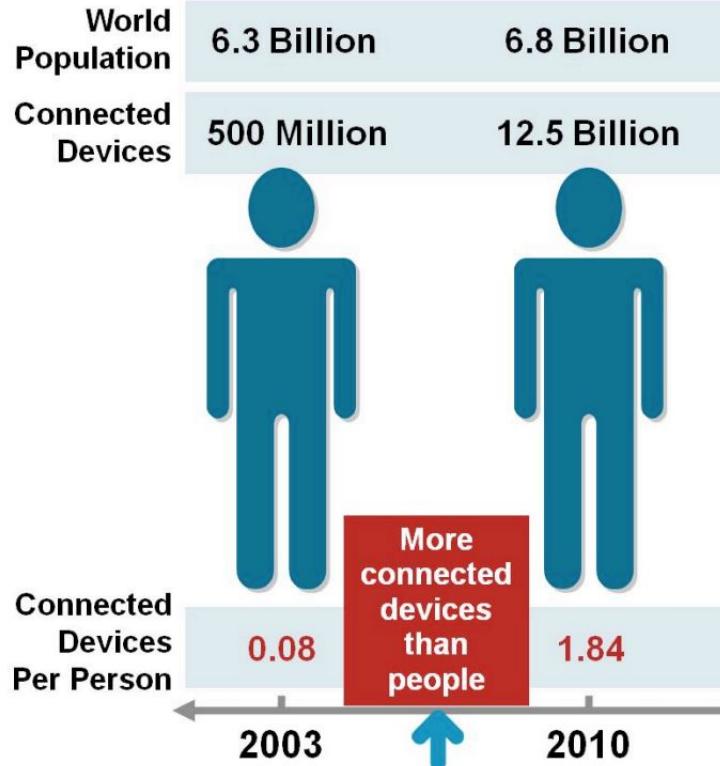




# **Q&A: IoT vs Human Beings?**

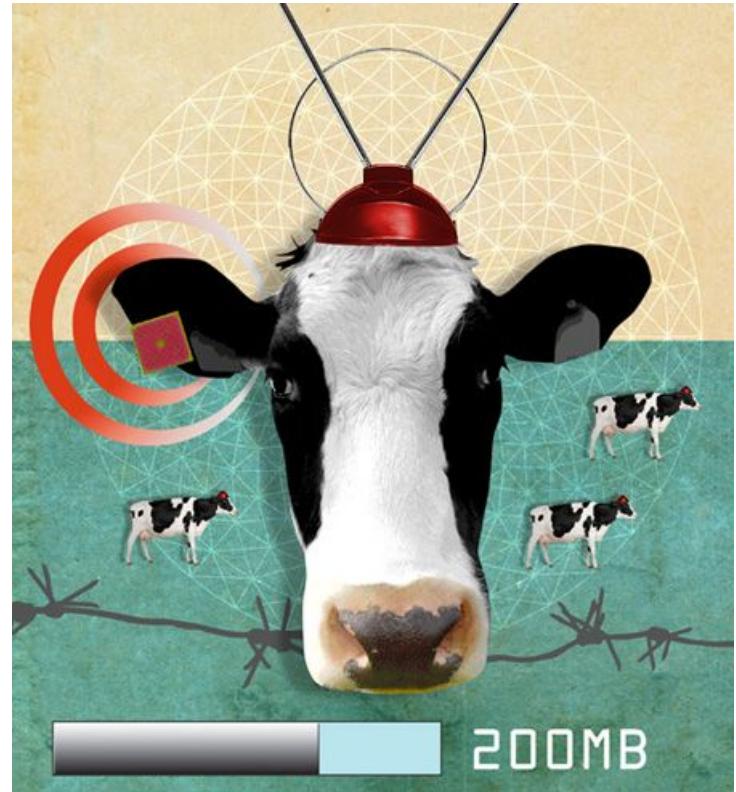


# BIGDATA TEAM

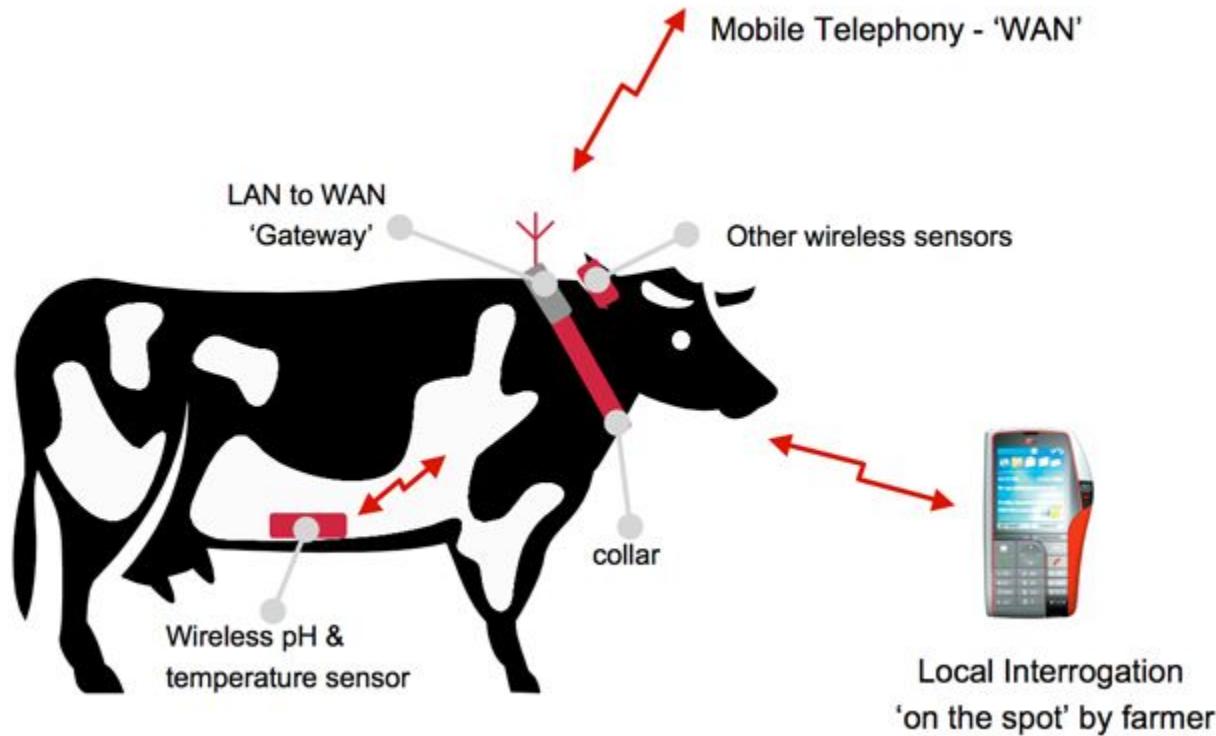


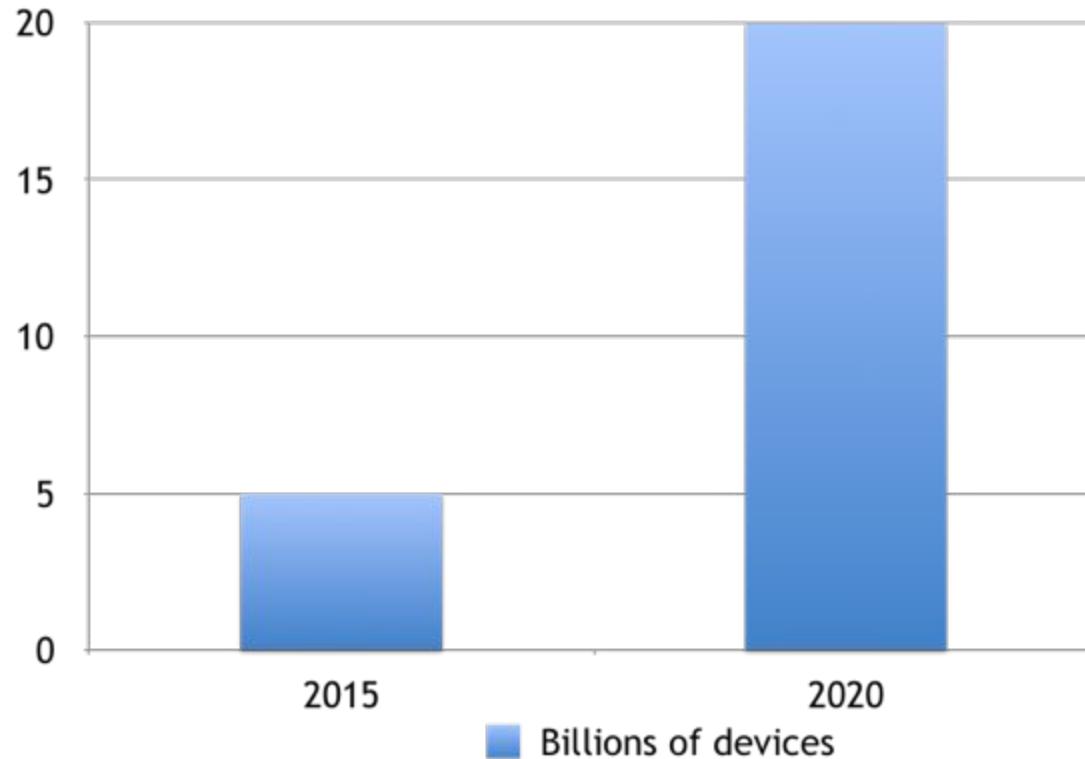
Source: Cisco IBSG, April 2011

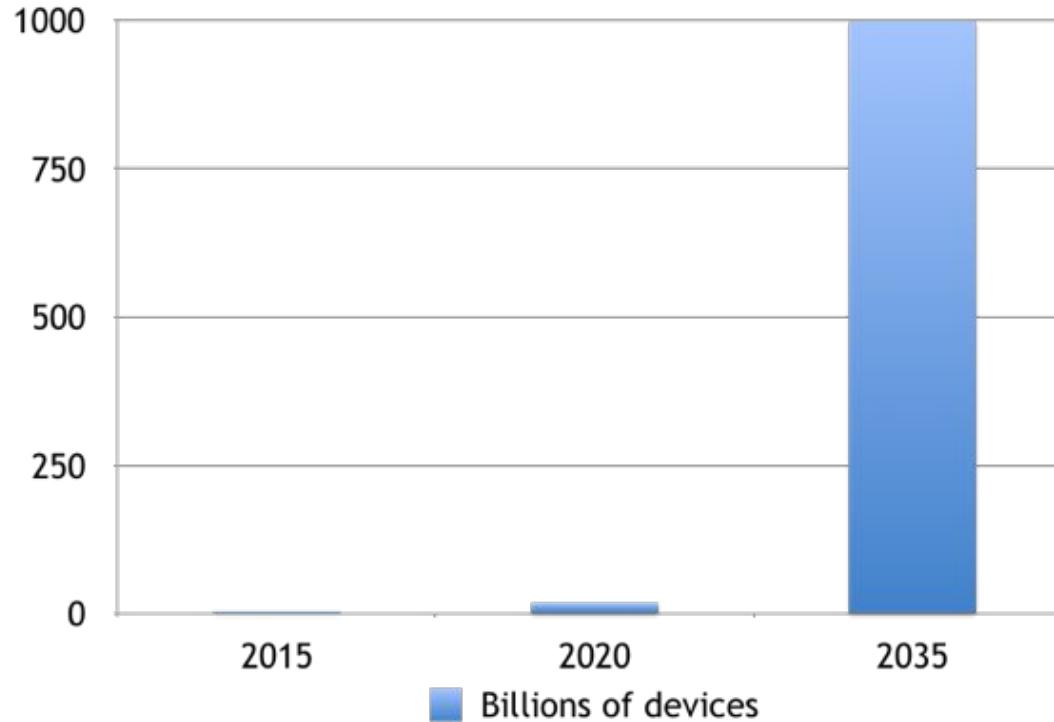
# Internet of Things (IoT)



Source: The Economist, 2010







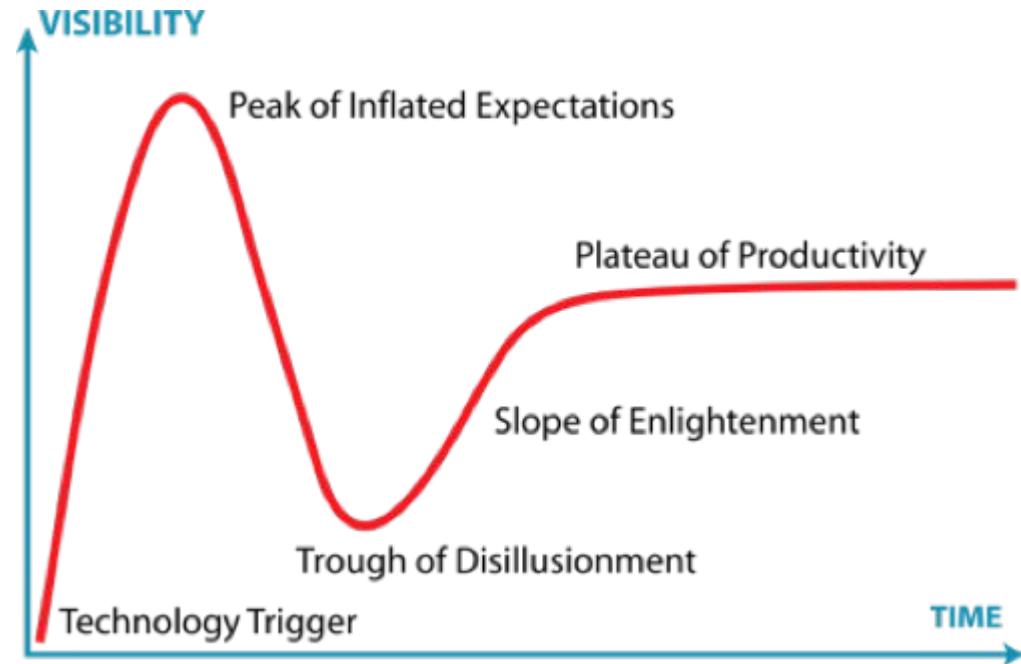
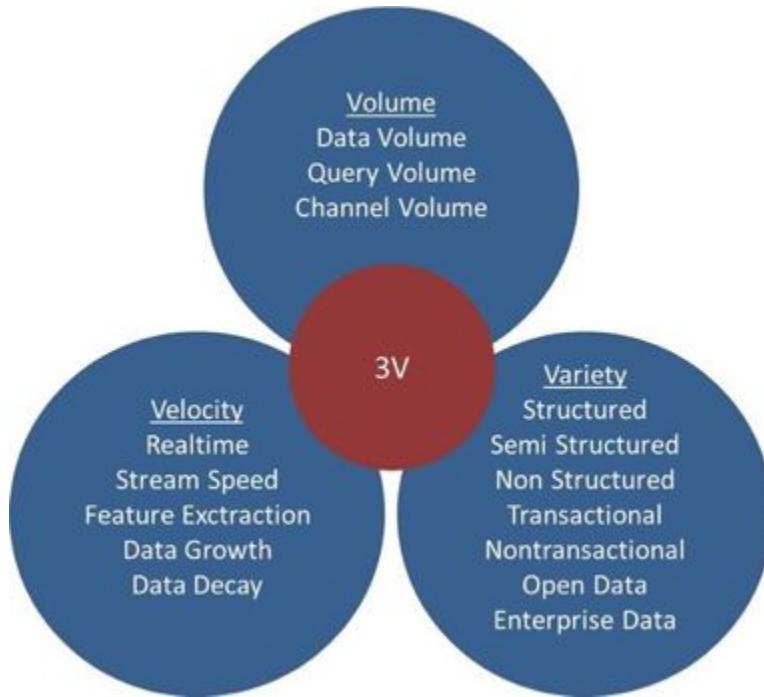


## **Q&A**

**Что такое Big Data?**

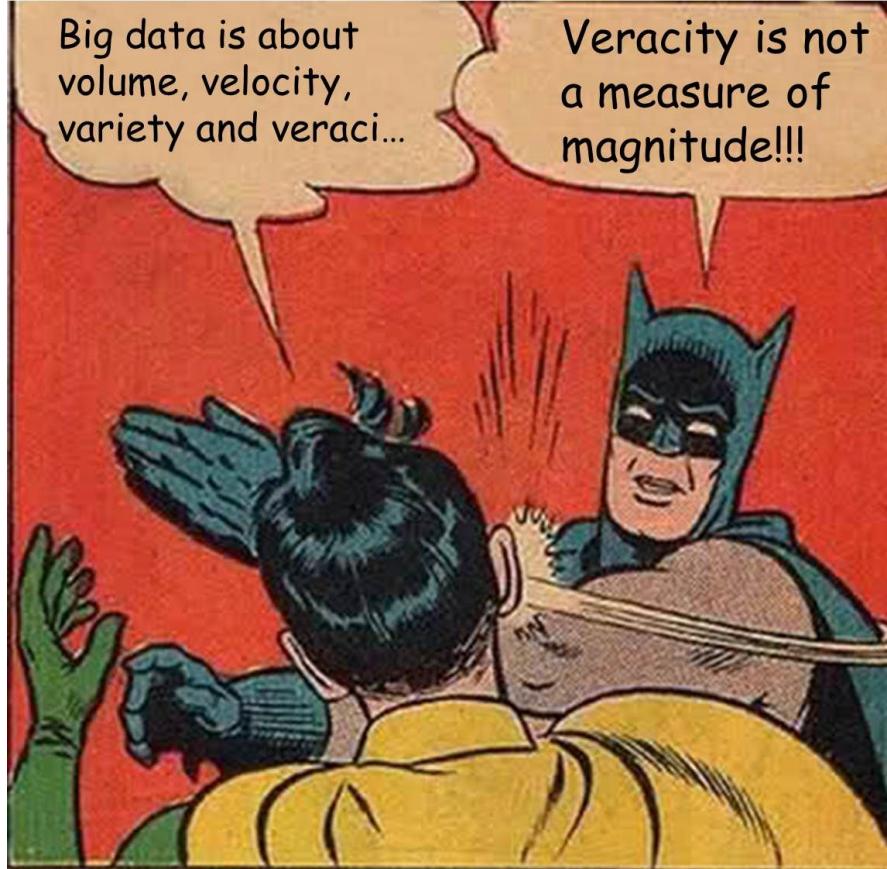


# Big and Small Data





# 4+ Vs of Big Data





## Q&A

Где применяется (data) science?



# **Q&A: Why Big Data?**



BIGDATA  
TEAM

# Data Scientist? Data Engineer? Data ...?

# indeed

## Index jobs

My recent searches

data engineer - 34,228 new

data scientist - 9,279 new

# 4X

glassdoor

data engineer United States Jobs

Job Type Date Posted Salary Range Distance More

Data Engineer Jobs in United States 101,151 Jobs



glassdoor

data scientist United States Jobs

Job Type Date Posted Salary Range Distance More

Data Scientist Jobs in United States 27,202 Jobs





**BIGDATA  
TEAM**

# Coursera Specialization on Big Data

## Capstone Project



Hadoop, Spark

**Yandex**

Hive, Spark

<https://bigdatateam.org/big-data-engineering>

Spark ML

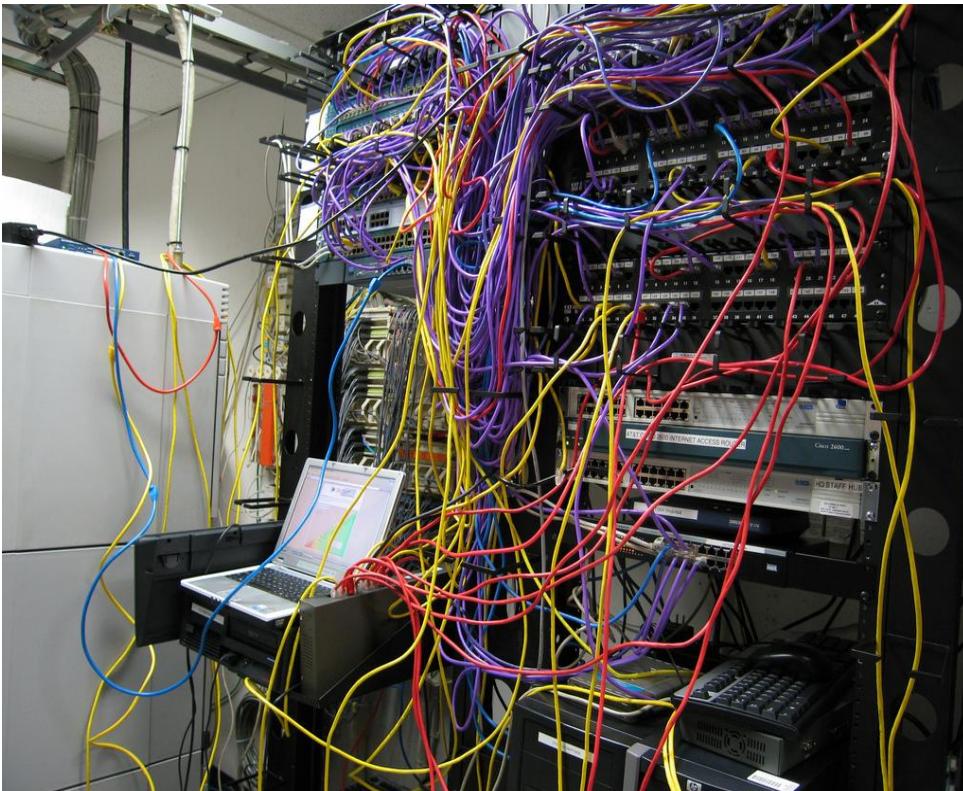
Real Time



# Многопроцессорные Вычислительные Системы (МВС)

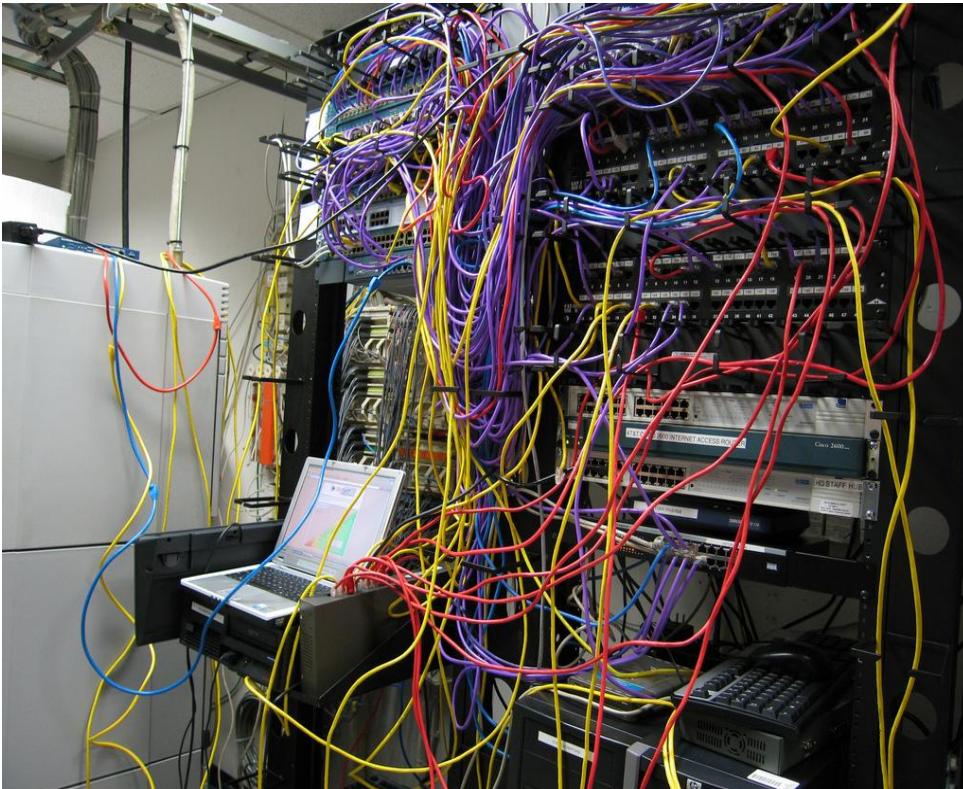


# Q&A: Что здесь может сломаться?





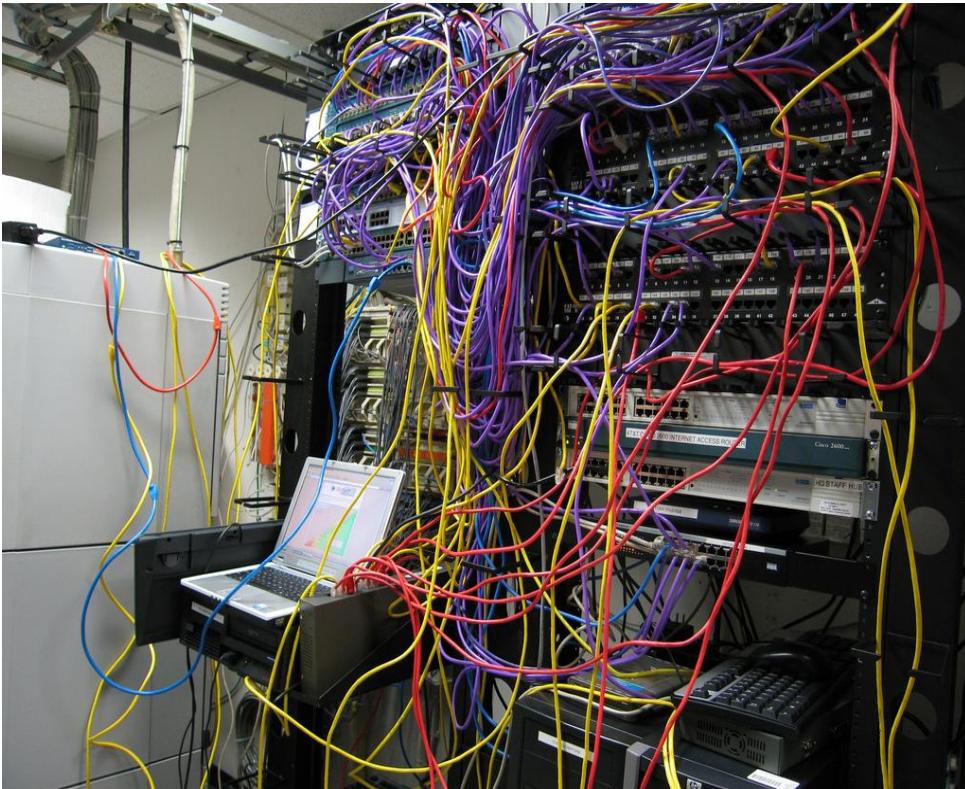
# Компоненты МВС



- ▶ Сервера (узлы)
- ▶ Сеть (соединения)
- ▶ Модель синхронности



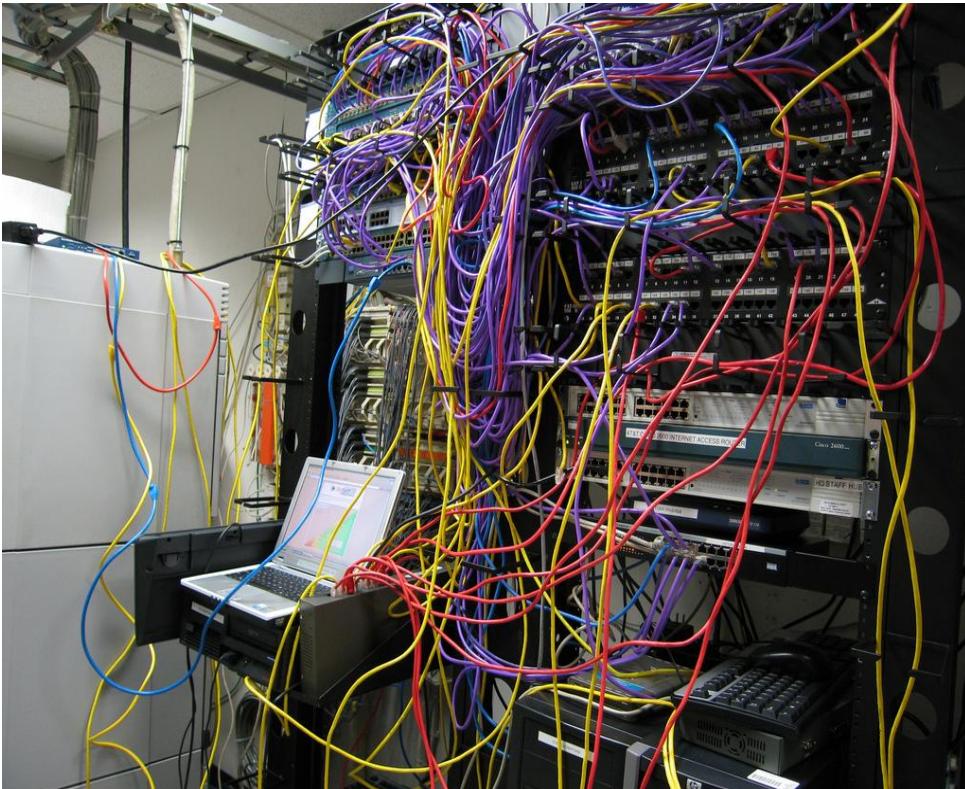
# Гарантии относительно поломки узлов (Node Failures)



- ▶ Fail-Stop
- ▶ Fail-Recovery
- ▶ Byzantine



# Гарантии относительно поломок сети (Link Failures)



- ▶ Perfect Link
- ▶ Fair-Loss Link
- ▶ Byzantine



## Модель синхронности



- ▶ Clock Skew
- ▶ Clock Drift



- ▶ Параллельные вычисления
- ▶ Распределенные вычисления
- ▶ Грид вычисления



BIGDATA  
TEAM

# Параллельные вычисления

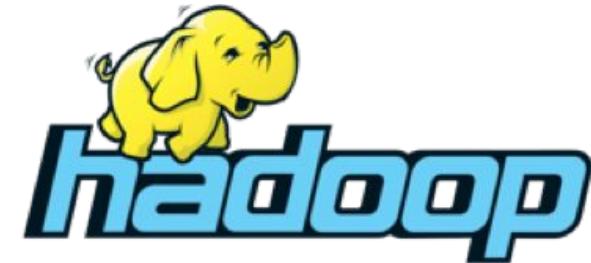
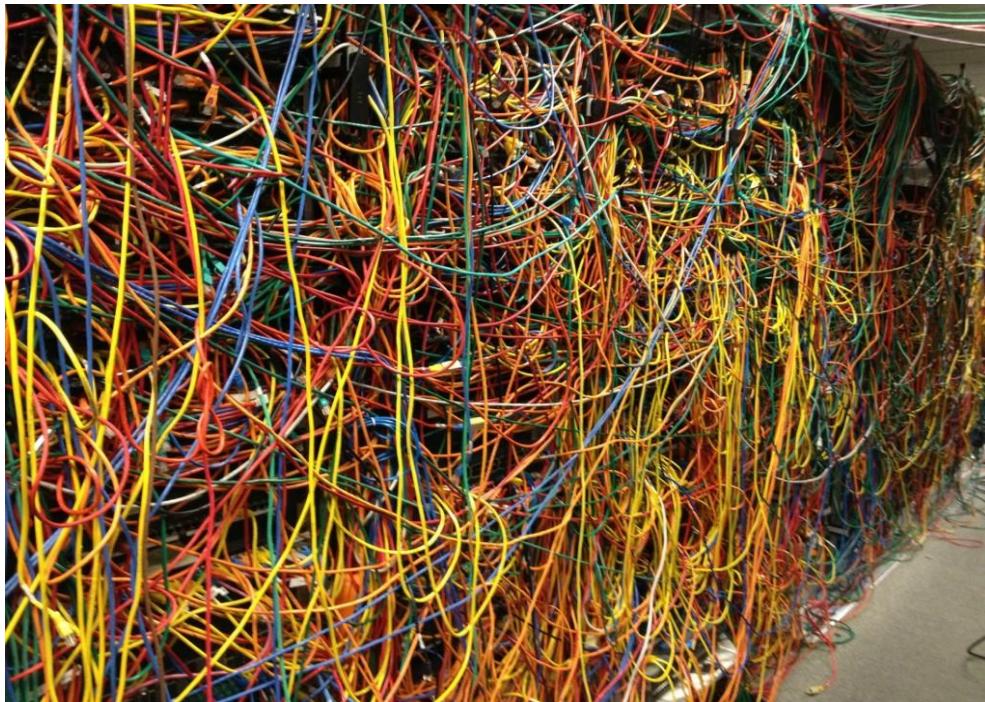


Fail-Stop + Perfect Link + Synchronous



BIGDATA  
TEAM

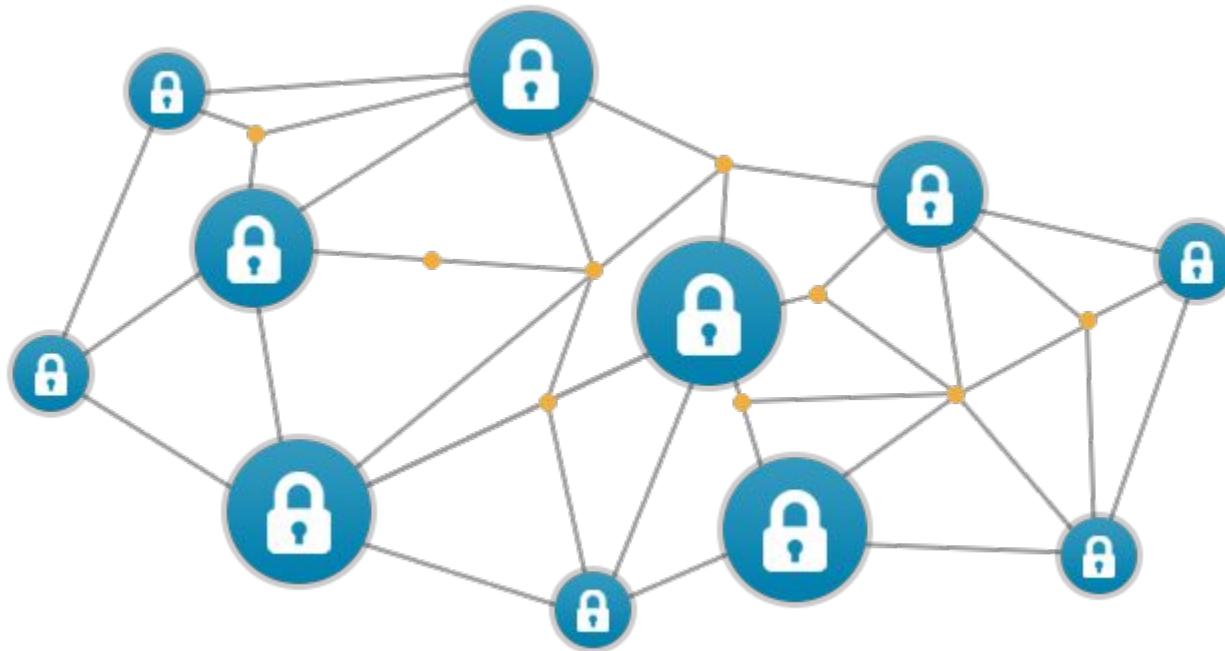
# Распределенные вычисления



Fail-Recovery + Fair-Loss Link + Asynchronous

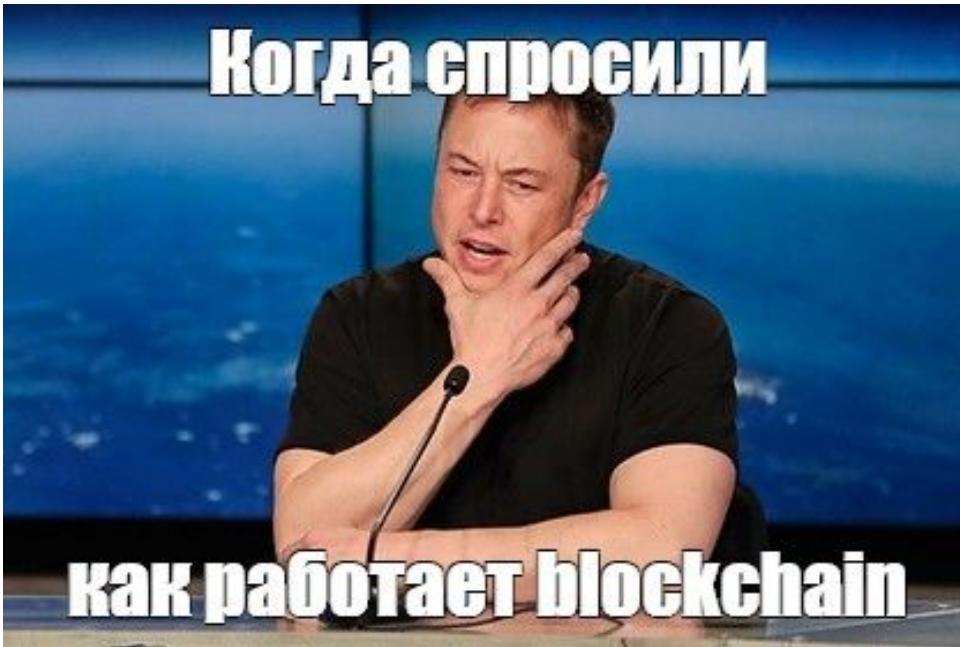


Byzantine-Failure + Byzantine Link + Asynchronous





# Грид vs Распределенные системы



Visa:

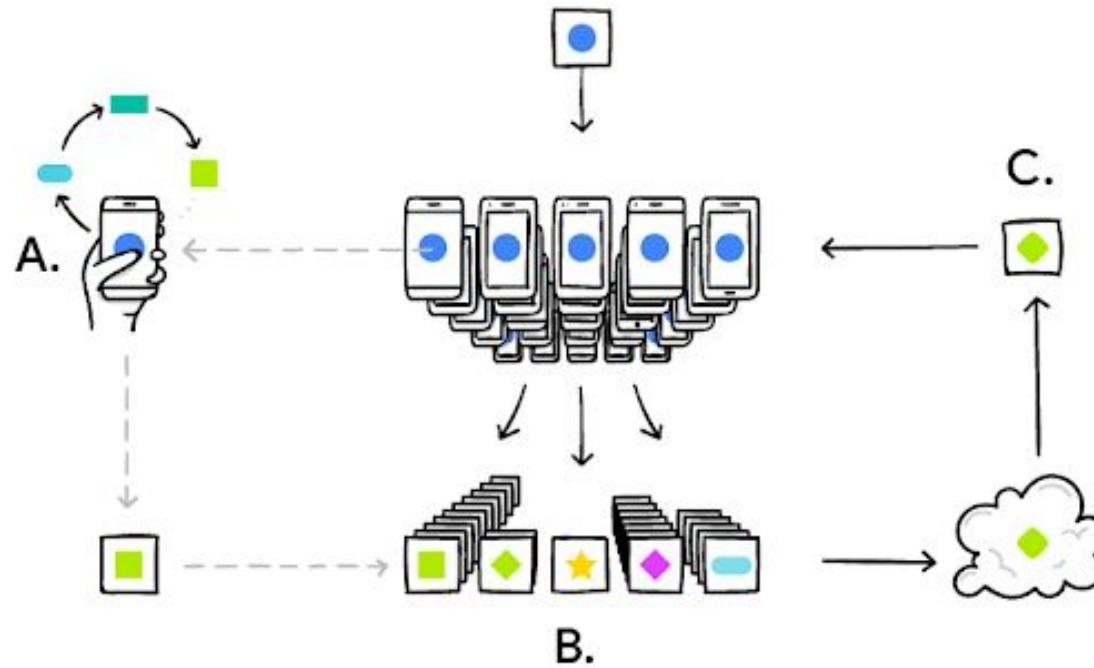
- ▶ avg: 1,700 TPS
- ▶ Scale up to: 24,000 TPS

Bitcoin:

- ▶ avg: 5-6 TPS
- ▶ Scale up to: 150-160 TPS



# Federated Machine Learning

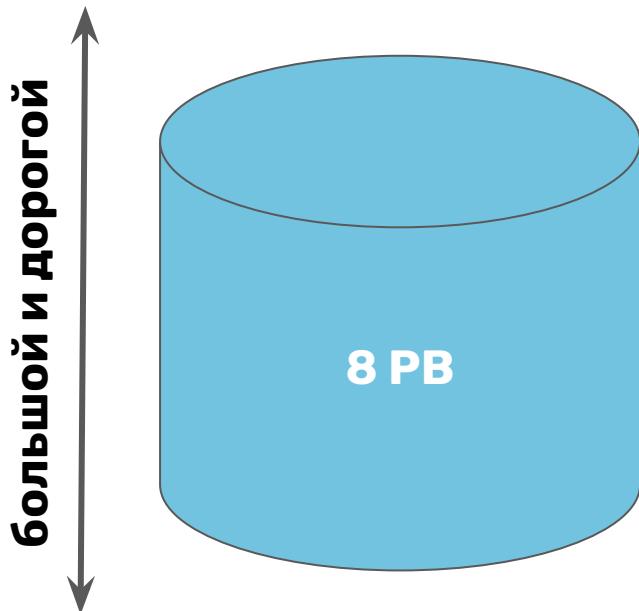




# Распределенные файловые системы (GFS / HDFS)



# Распределенные файловые системы



**вертикальное масштабирование  
(scale-up)**



**горизонтальное масштабирование  
(scale-out)**



## The Google File System

Sanjay Ghemawat, Howard Gobioff, and Shun-Tak Leung  
Google\*

### ABSTRACT

We have designed and implemented the Google File System, a scalable distributed file system for large distributed data-intensive applications. It provides fault tolerance while running on inexpensive commodity hardware, and it delivers high aggregate performance to a large number of clients.

While sharing many of the same goals as previous distributed file systems, our design has been driven by observations of our application workloads and technological environment, both current and anticipated, that reflect a marked departure from some earlier file system assumptions. This has led us to reexamine traditional choices and explore rad-

### 1. INTRODUCTION

We have designed and implemented the Google File System (GFS) to meet the rapidly growing demands of Google's data processing needs. GFS shares many of the same goals as previous distributed file systems such as performance, scalability, reliability, and availability. However, its design has been driven by key observations of our application workloads and technological environment, both current and anticipated, that reflect a marked departure from some earlier file system design assumptions. We have reexamined traditional choices and explored radically different points in the design space.



## Q&A

Как обеспечить гарантии хранения  
данных в рамках fail-recovery?

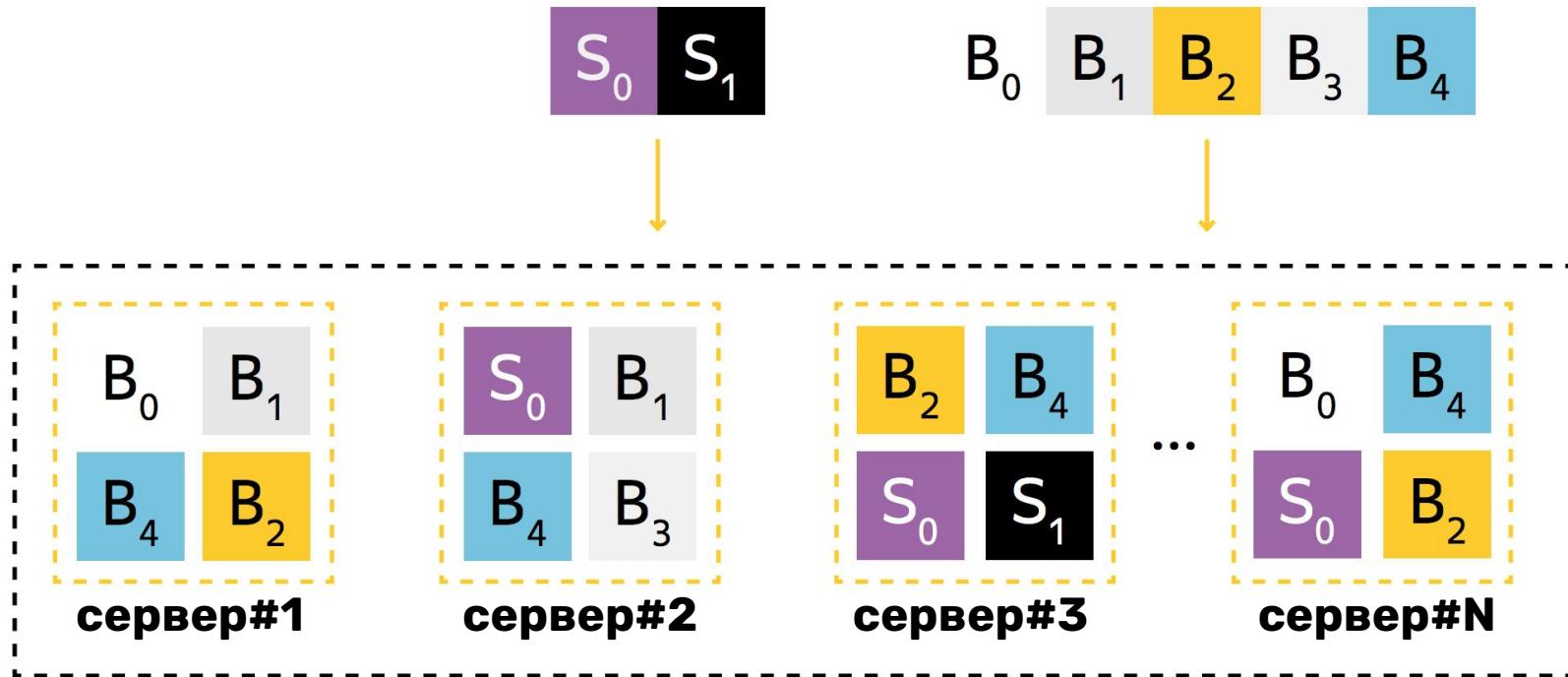


## Q&A

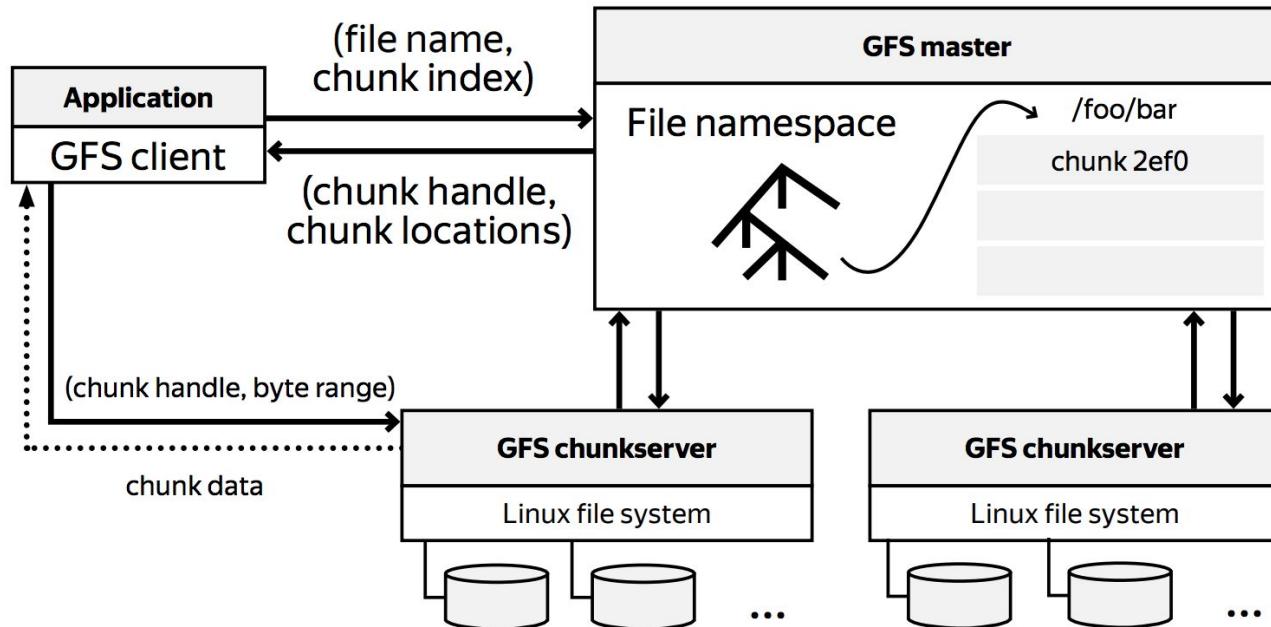
Как храним файлы размером 2 ТВ и  
10 GB в кластере?



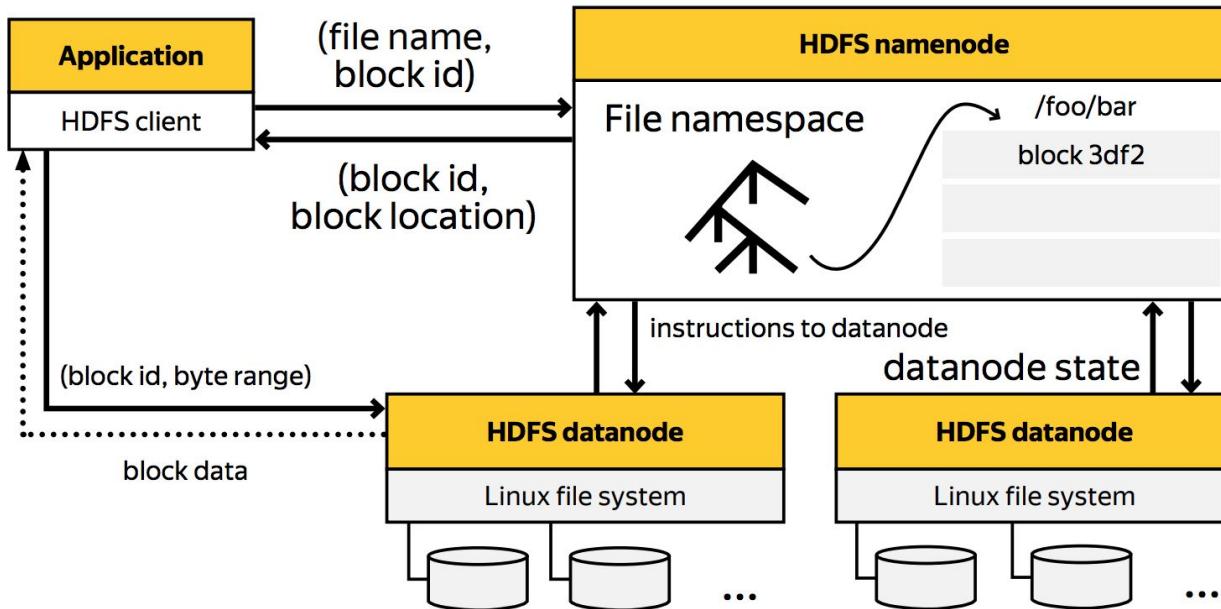
# Распределенная ФС



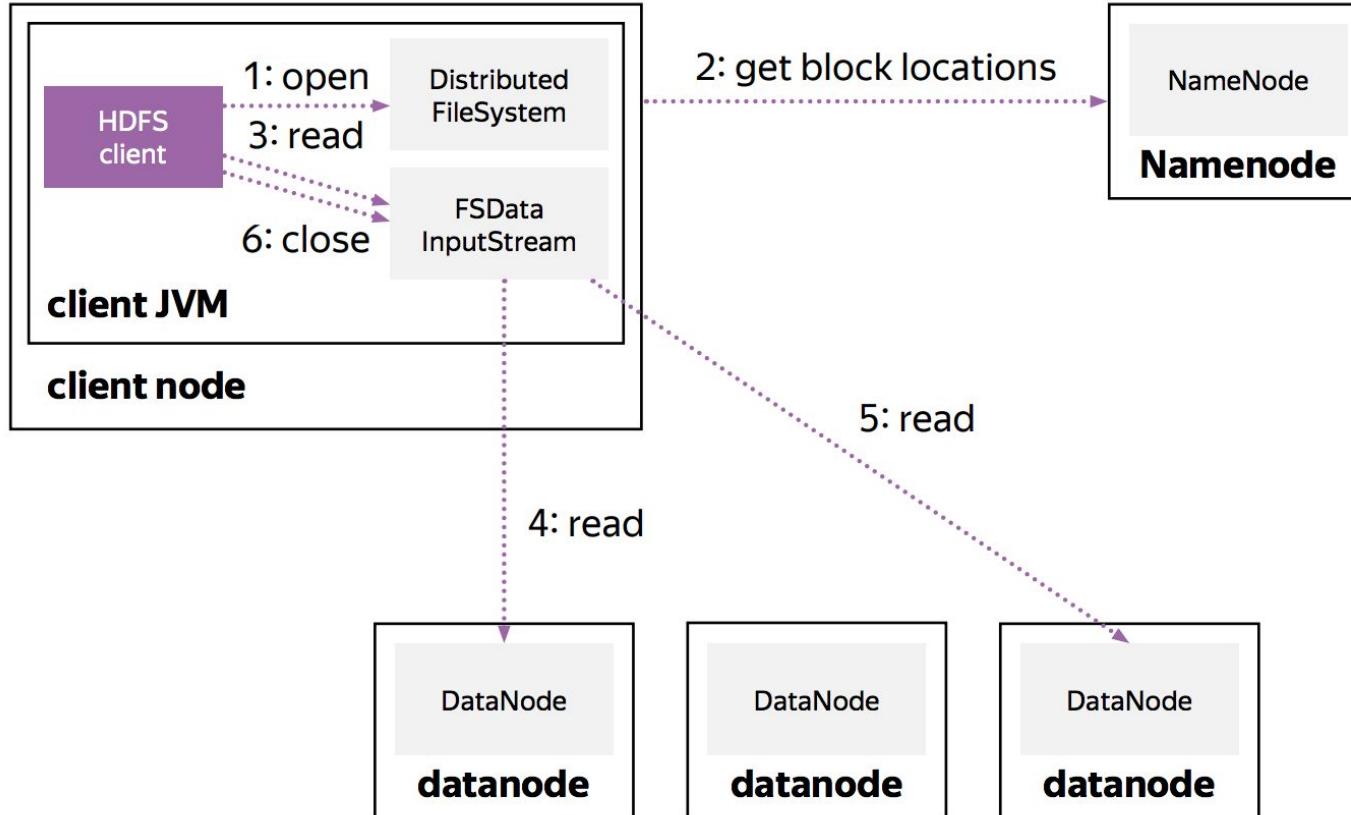
- ▶ Поломки компонент - это норма (используем репликацию)
- ▶ Равномерная утилизация компонент кластера
- ▶ Семантика write-once-read-many

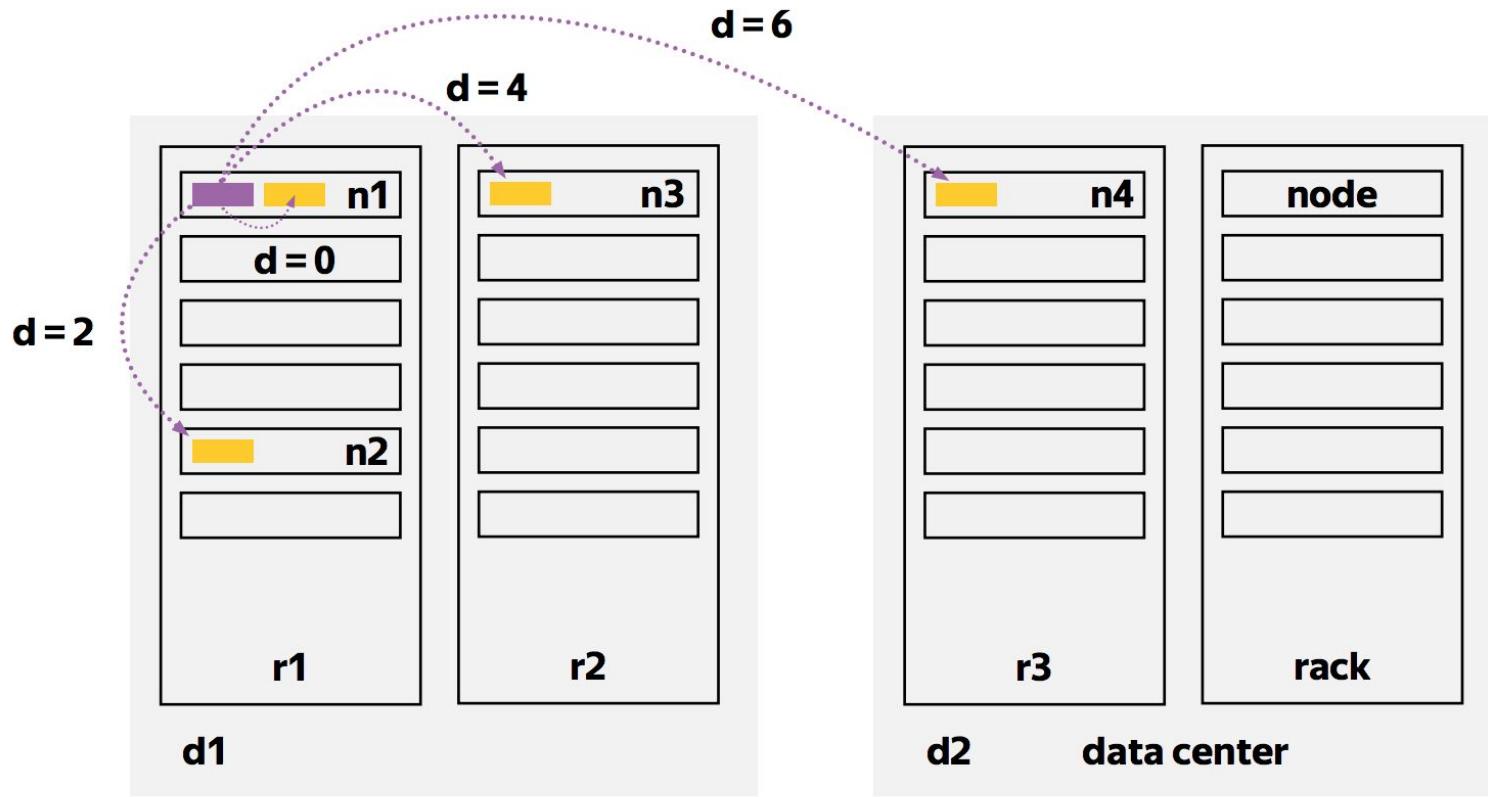


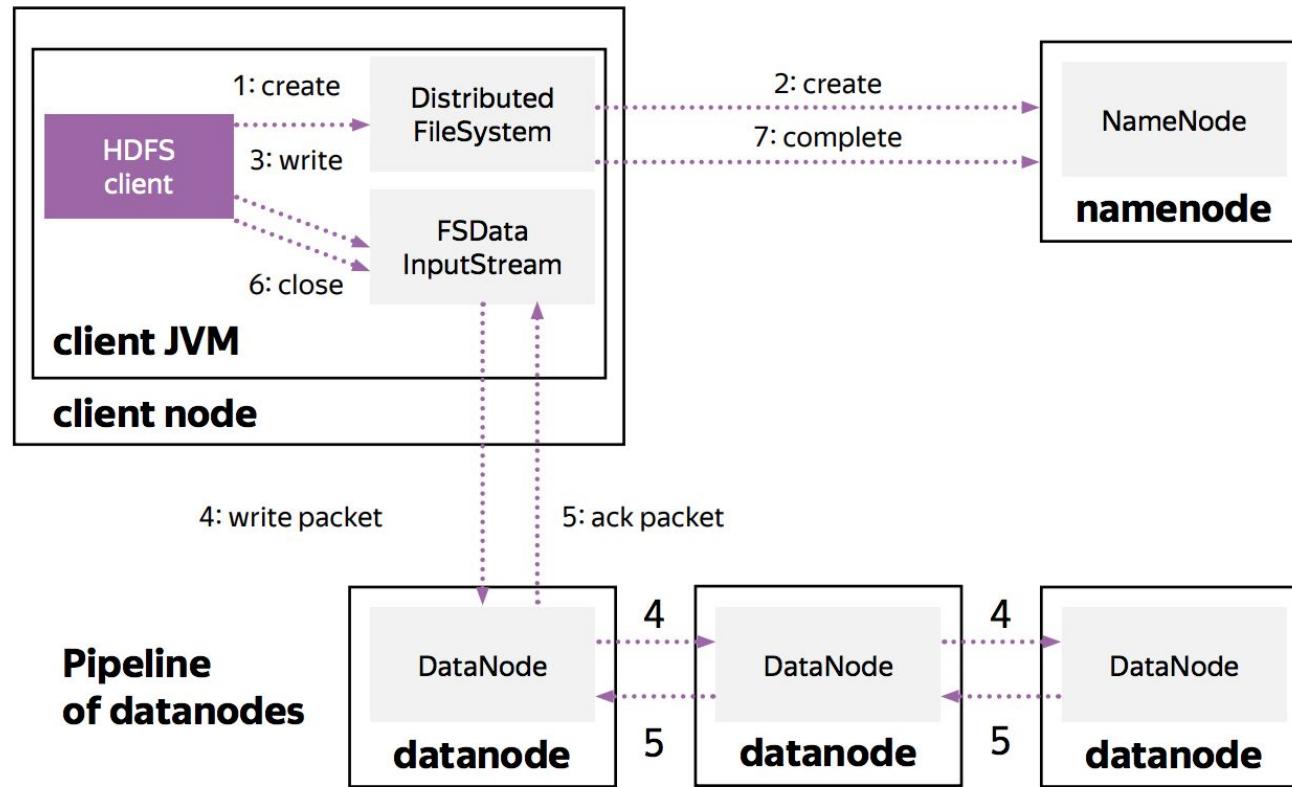
# GFS



# HDFS

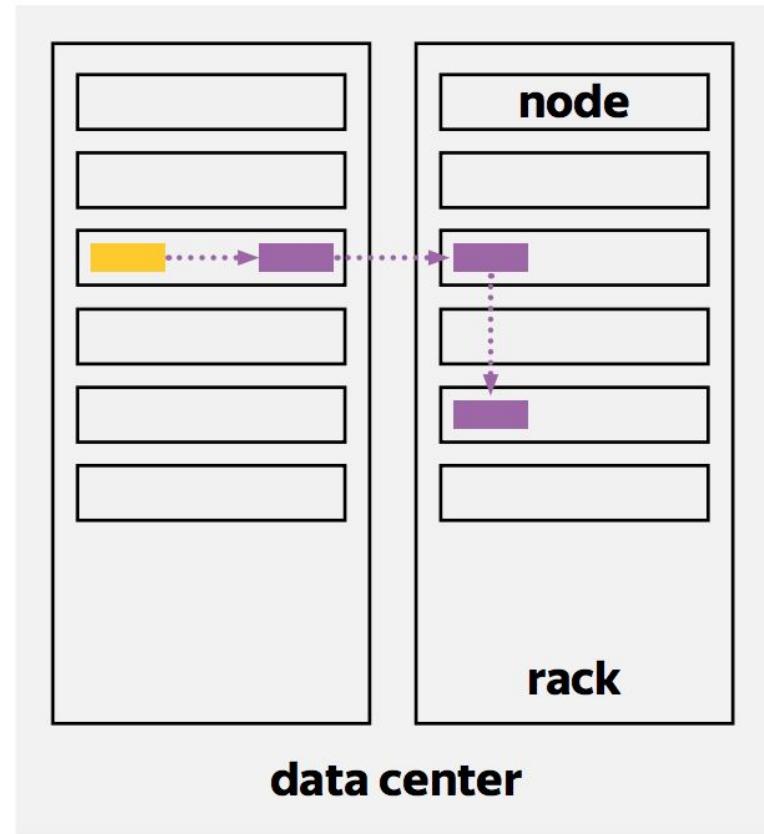






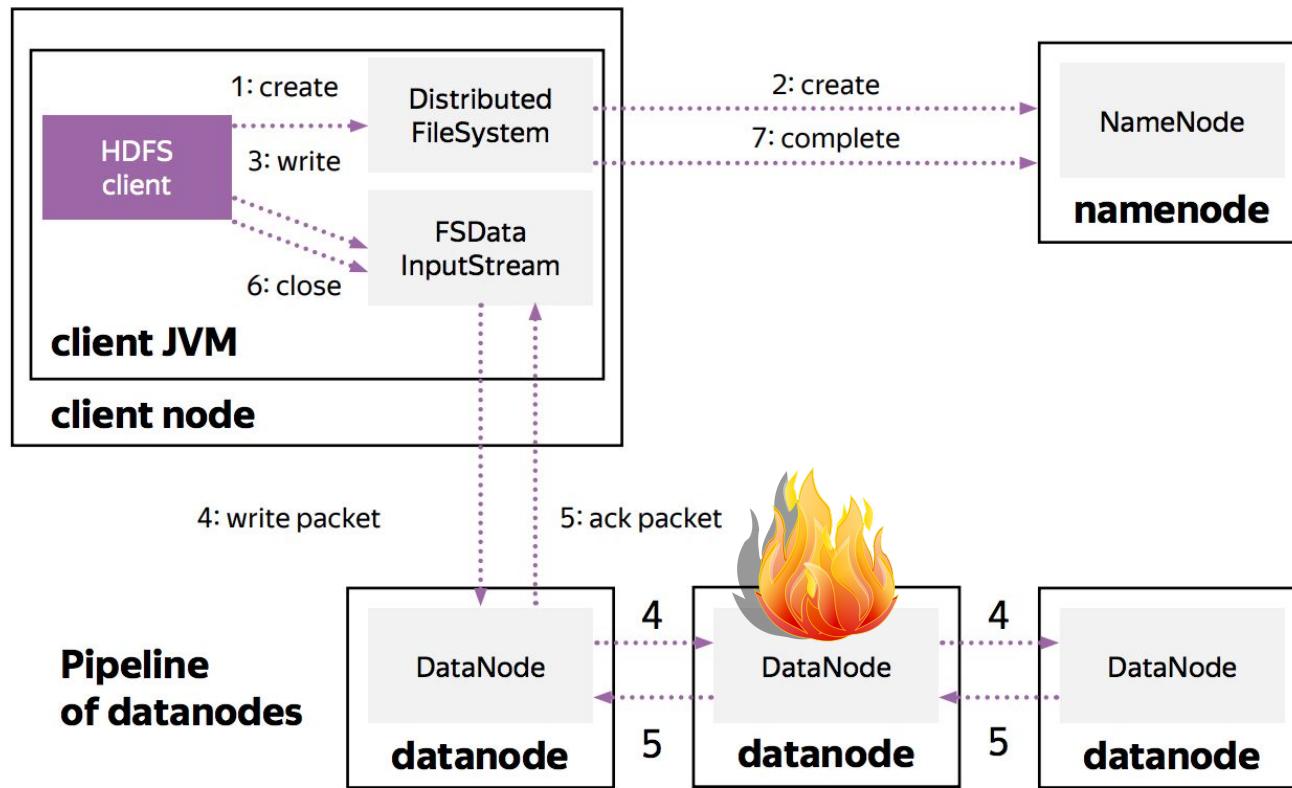


# HDFS Replica Placement



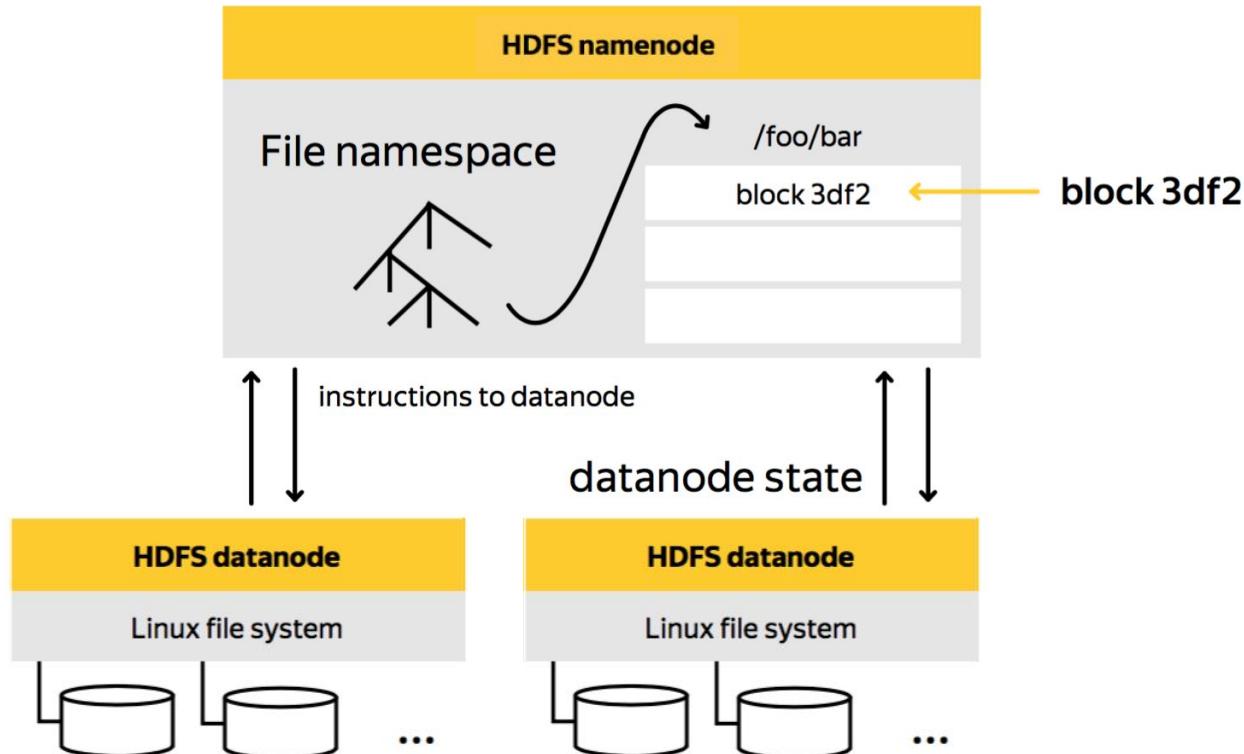


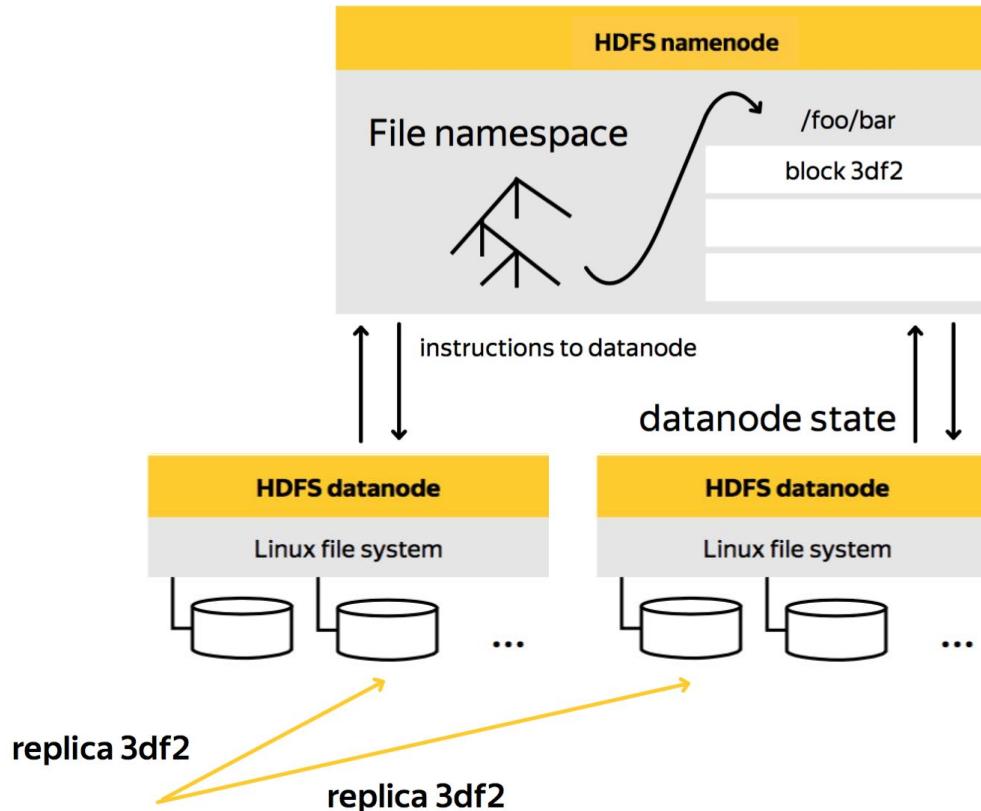
# Запись данных в HDFS v.2

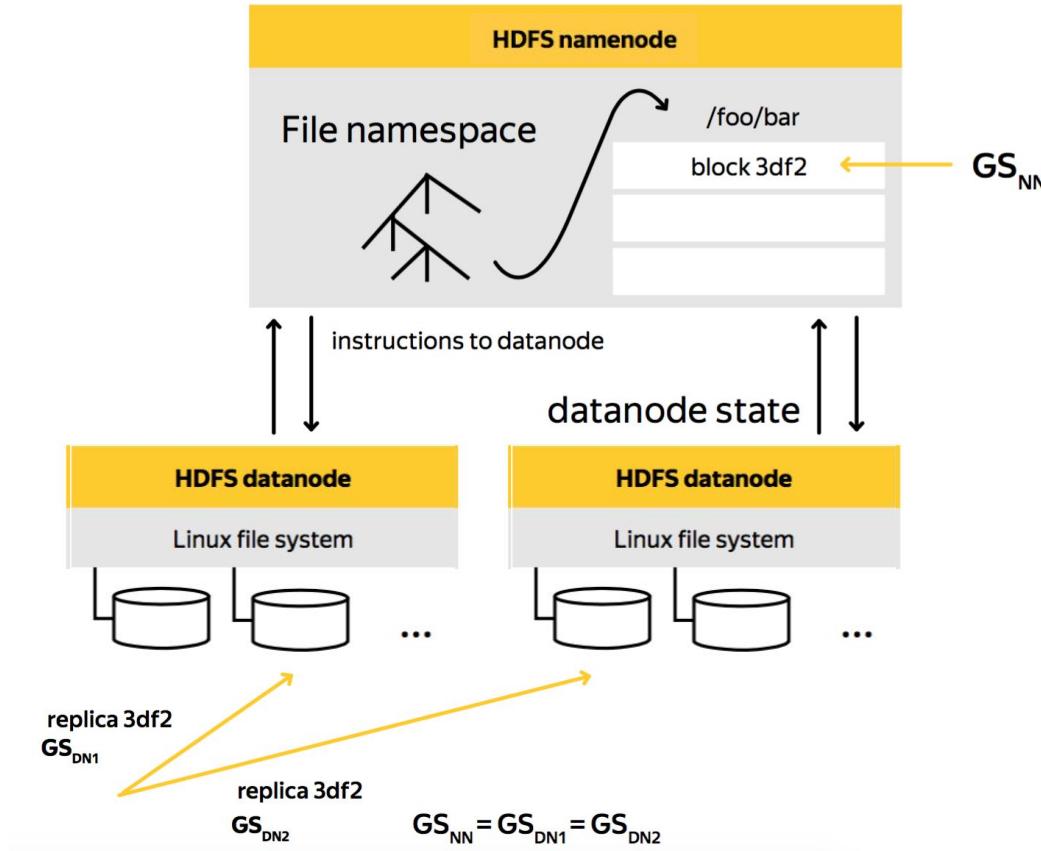




# чанки, блоки и реплики



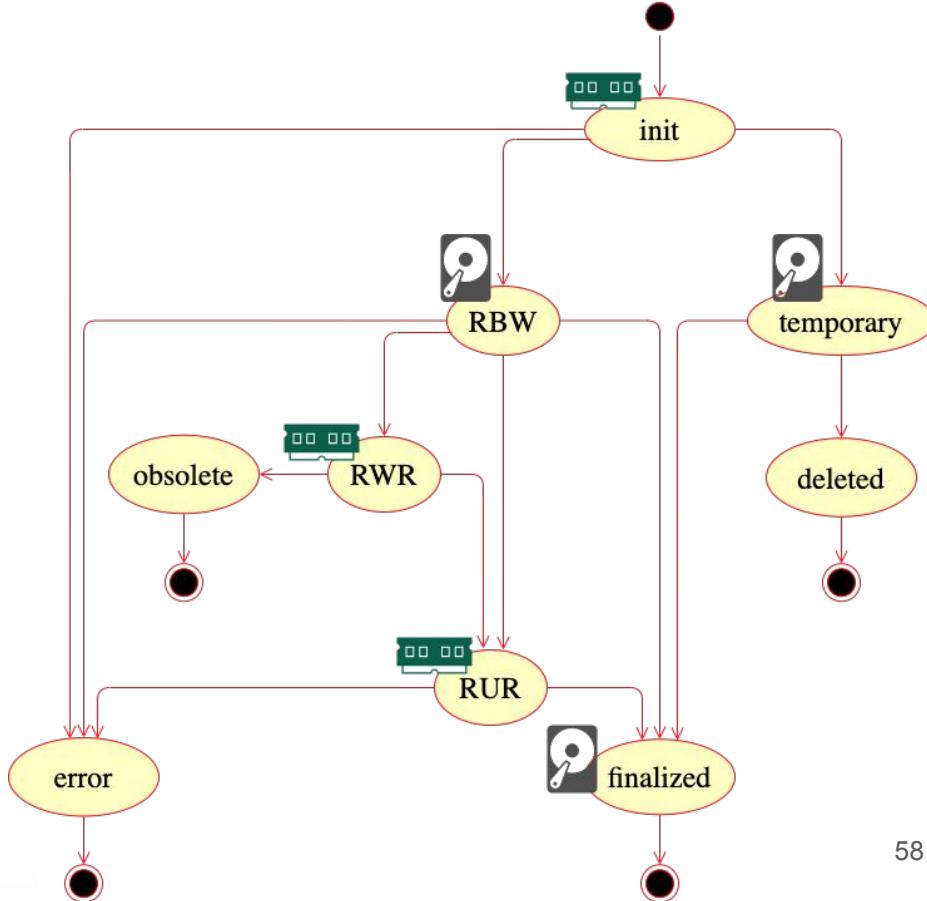
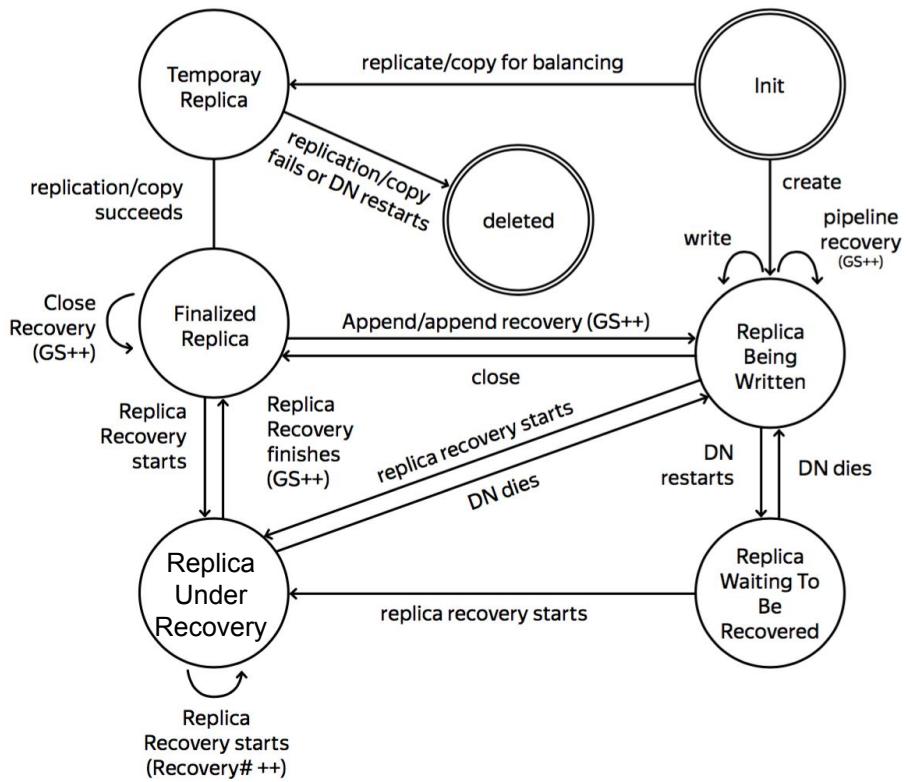






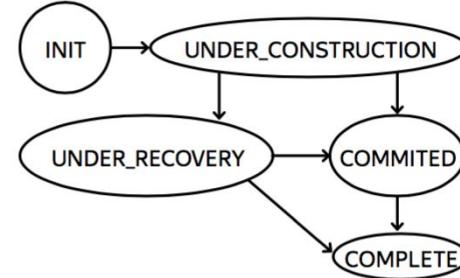
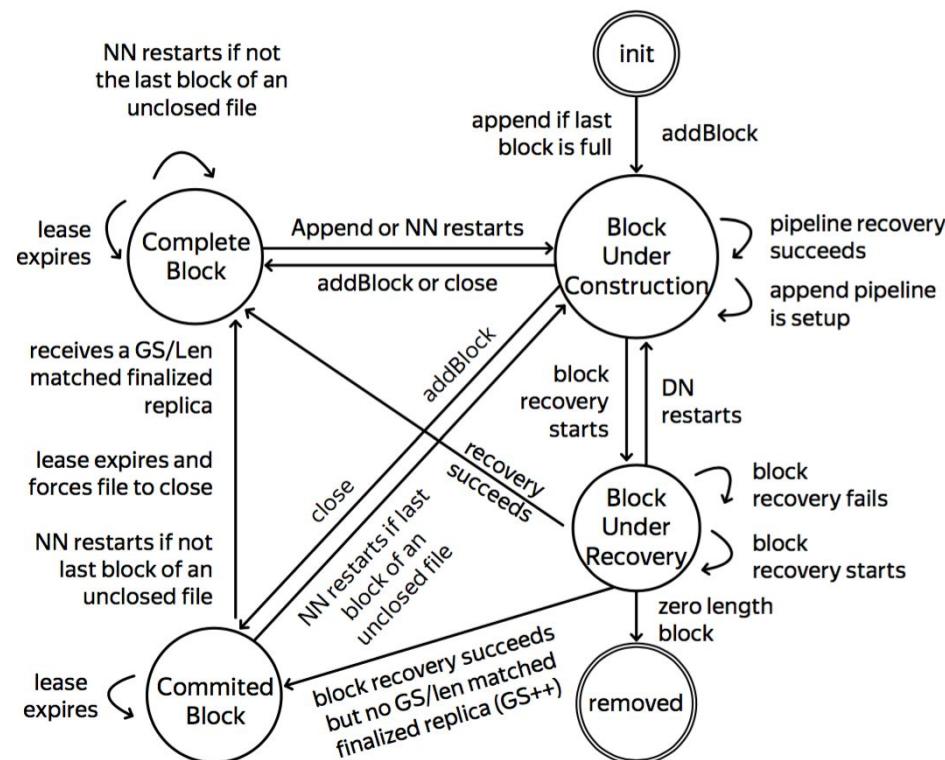
# Процедуры восстановления HDFS

- ▶ Replica Recovery
- ▶ Block Recovery
- ▶ Lease Recovery
- ▶ Pipeline Recovery





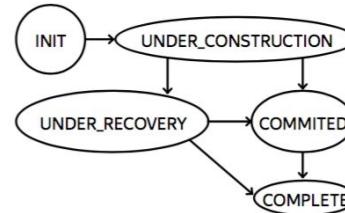
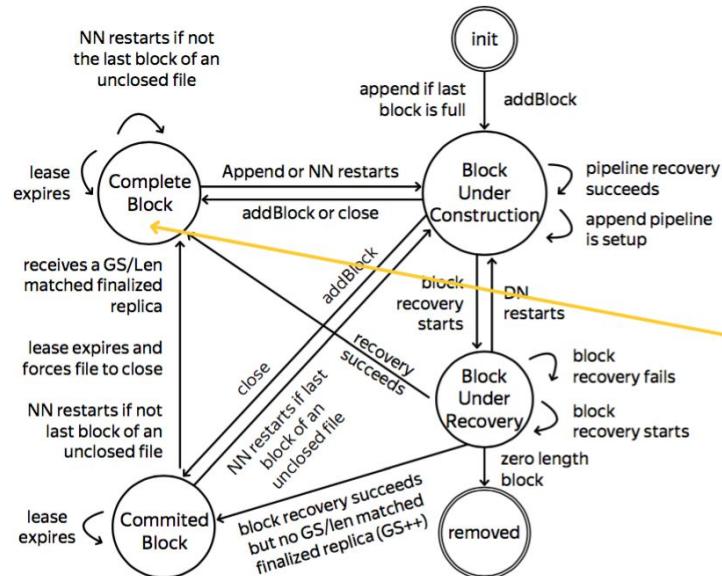
# Simplified Block State Transition



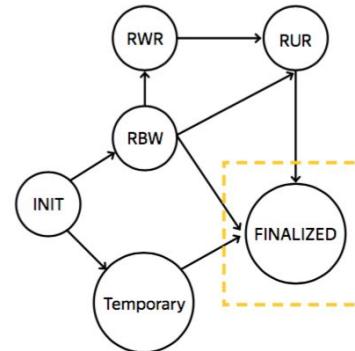
Simplified  
Block State  
Transition



# Block & Replica States Matching



Simplified  
Block State  
Transition  
**complete**



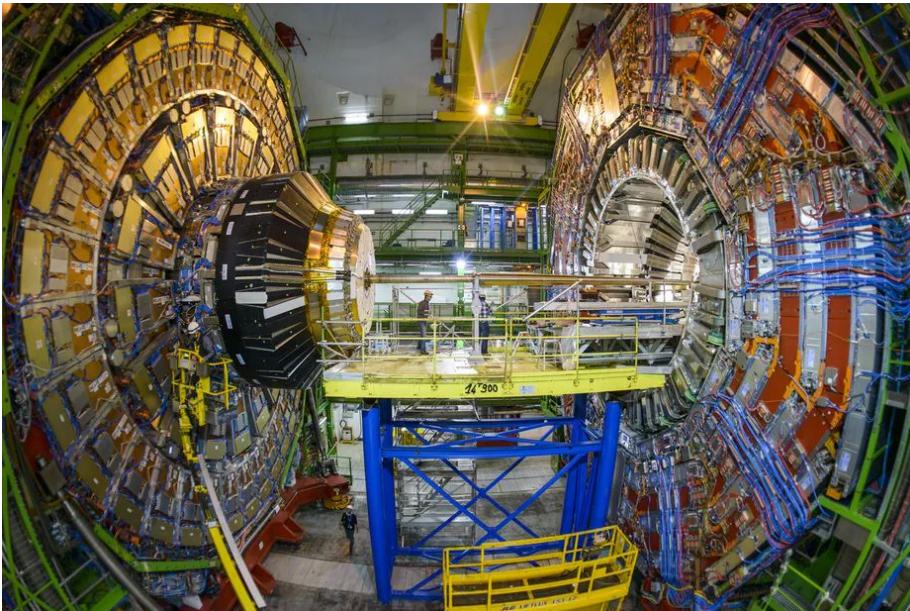




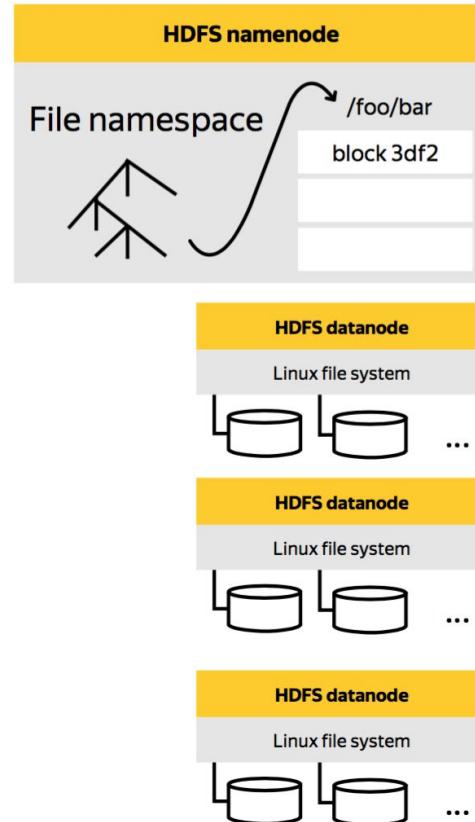
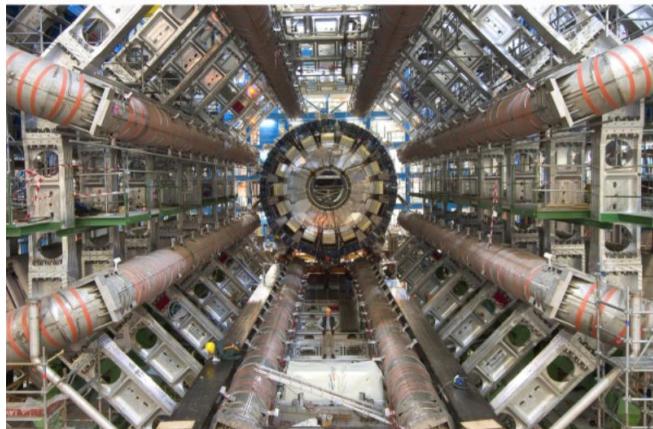
# Hadoop Sizing

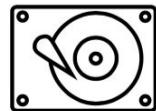
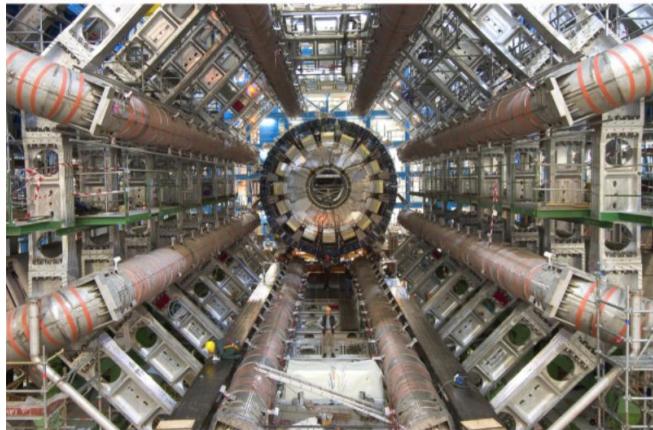


# Запрос на консультацию

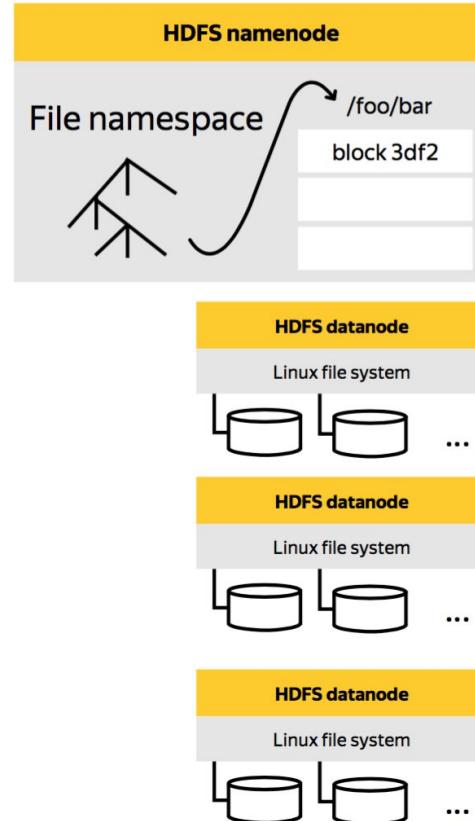


The CERN Data Centre passed a major milestone on 29 June 2017 with more than **200 petabytes** of data now archived on tape





10 PB / 2 TB \* 3 ~ 15 k





BIGDATA  
TEAM

Считаем сбои



# DRIVE STATS 2020





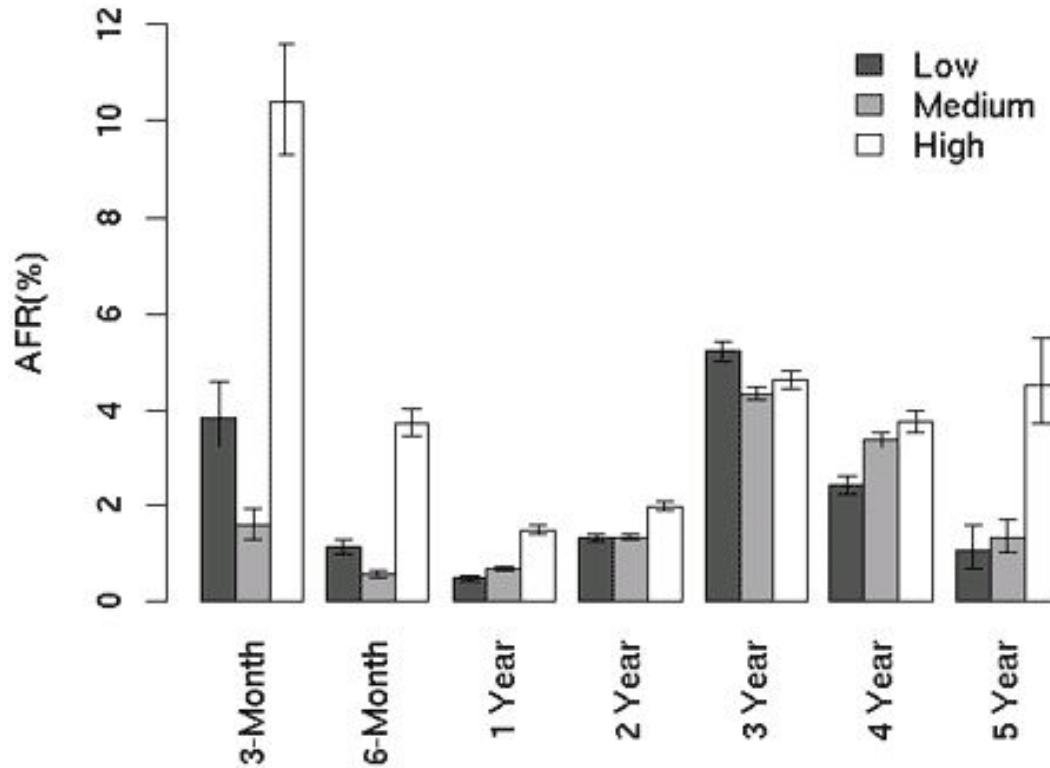
## Backblaze Hard Drive Failure Rates for 2020

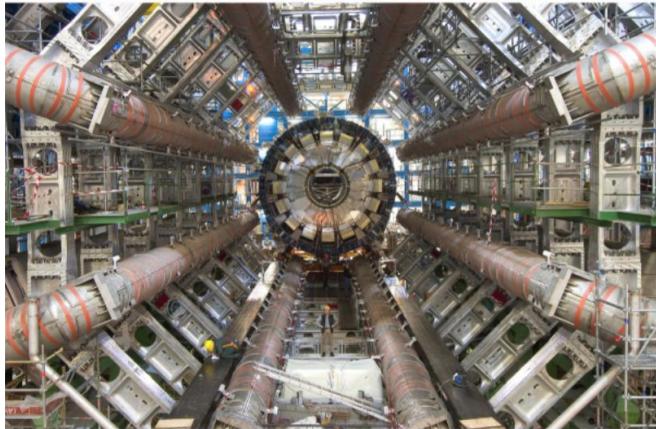
Reporting period 1/1/2020 - 12/31/2020 inclusive

MFG	Model	Drive Size	Drive Count	Avg Age (months)	Drive Days	Drive Failures	AFR
HGST	HMS5C4040ALE640	4TB	3,100	56.65	1,083,774	8	0.27%
HGST	HMS5C4040BLE640	4TB	12,744	50.43	4,663,049	34	0.27%
HGST	HUH728080ALE600	8TB	1,075	34.85	372,000	3	0.29%
Seagate	ST12000NM0010	12TB	7,100	0.00	1,200,170	0	0.00%
Seagate	ST14000NM001G	14TB	5,987	2.89	454,090	13	1.04%
Seagate	ST14000NM0138	14TB	360	1.56	5,784	0	0.00%
Seagate	ST16000NM001G	16TB	59	12.93	21,323	1	1.71%
Seagate	ST18000NM000J	18TB	60	3.27	5,820	2	12.54%
Total	All Drives	ATD	00	07.00	22,224	0	0.00%



# Вероятность сбоя от перегрузки





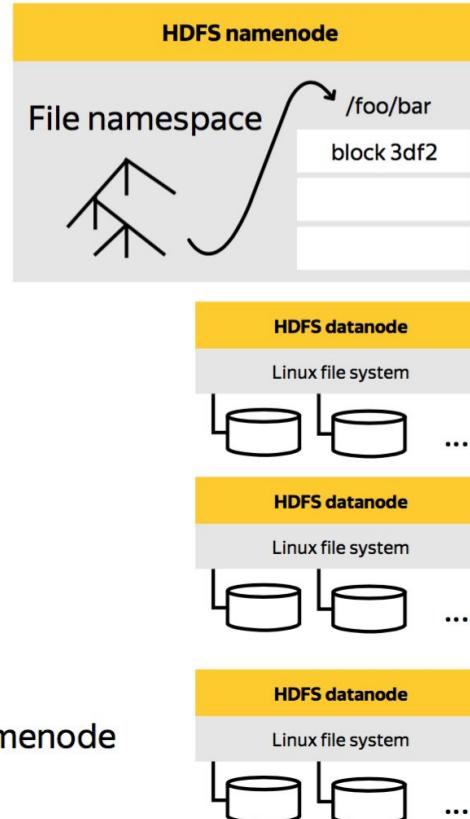
10 PB / 2 TB \* 3 ~ 15 k

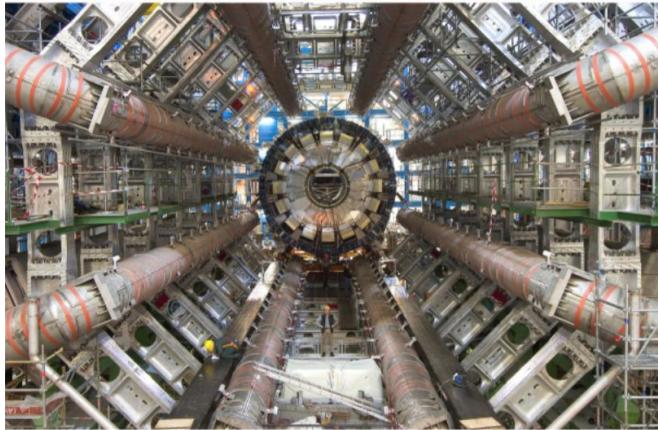


150 B - average block size on Namenode

<https://issues.apache.org/jira/browse/HADOOP-1687>

# Hadoop Sizing



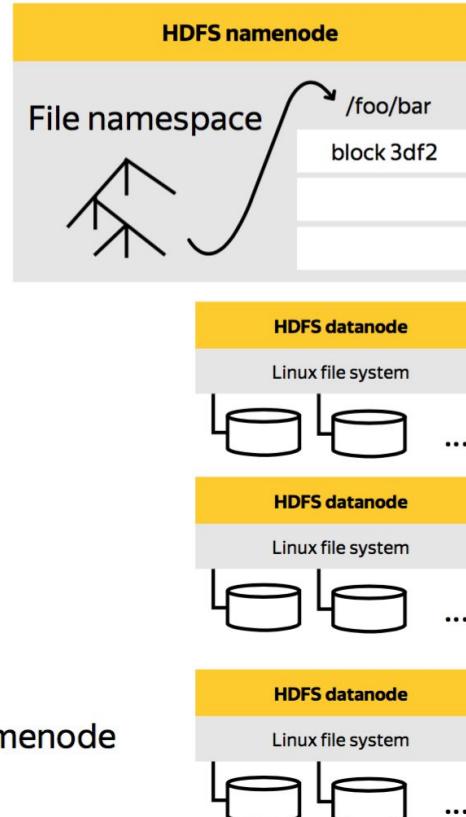


10 PB / 2 TB \* 3 ~ 15 k



150 B - average block size on Namenode

<https://issues.apache.org/jira/browse/HADOOP-1687>



# Hadoop Sizing

Input interpretation:

$10 \text{ PB}$  (petabytes)  $\times 150 \text{ bytes}$   
 $128 \text{ MB}$  (megabytes)

Result:

$1.172 \times 10^{10} \text{ bytes}$

Unit conversions:

$11.72 \text{ GB}$  (gigabytes)



# Hadoop “Small Files Problem”



- ▶ reading speed - 600 MB/sec
- ▶ 10 PB → **207** дней

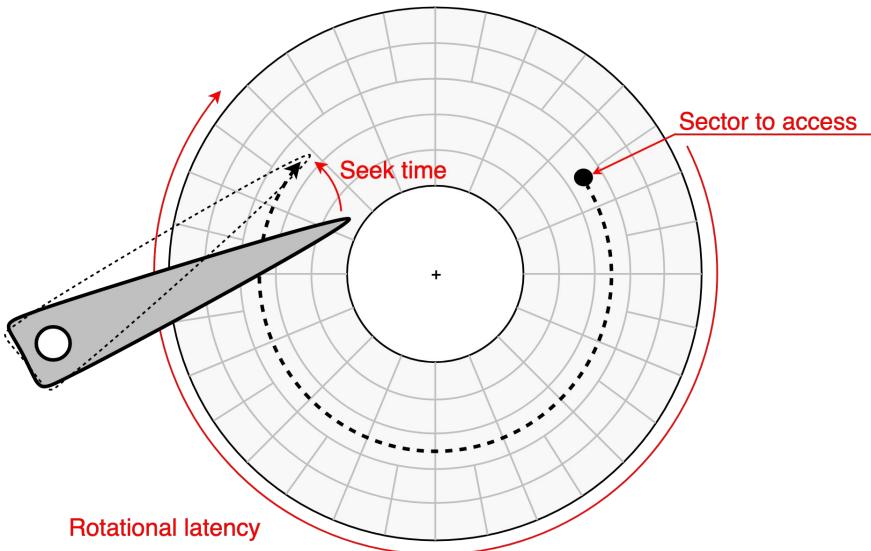


- ▶ reading speed - 600 MB/sec
- ▶ 10 PB → **103.5** дней



RAM: 16GB, block size: **32MB**, replication factor: 3

- максимальный размер хранилища:  $16\text{GB} / 150\text{B} \times \text{32MB} / 3 = 1.138\text{PB}$



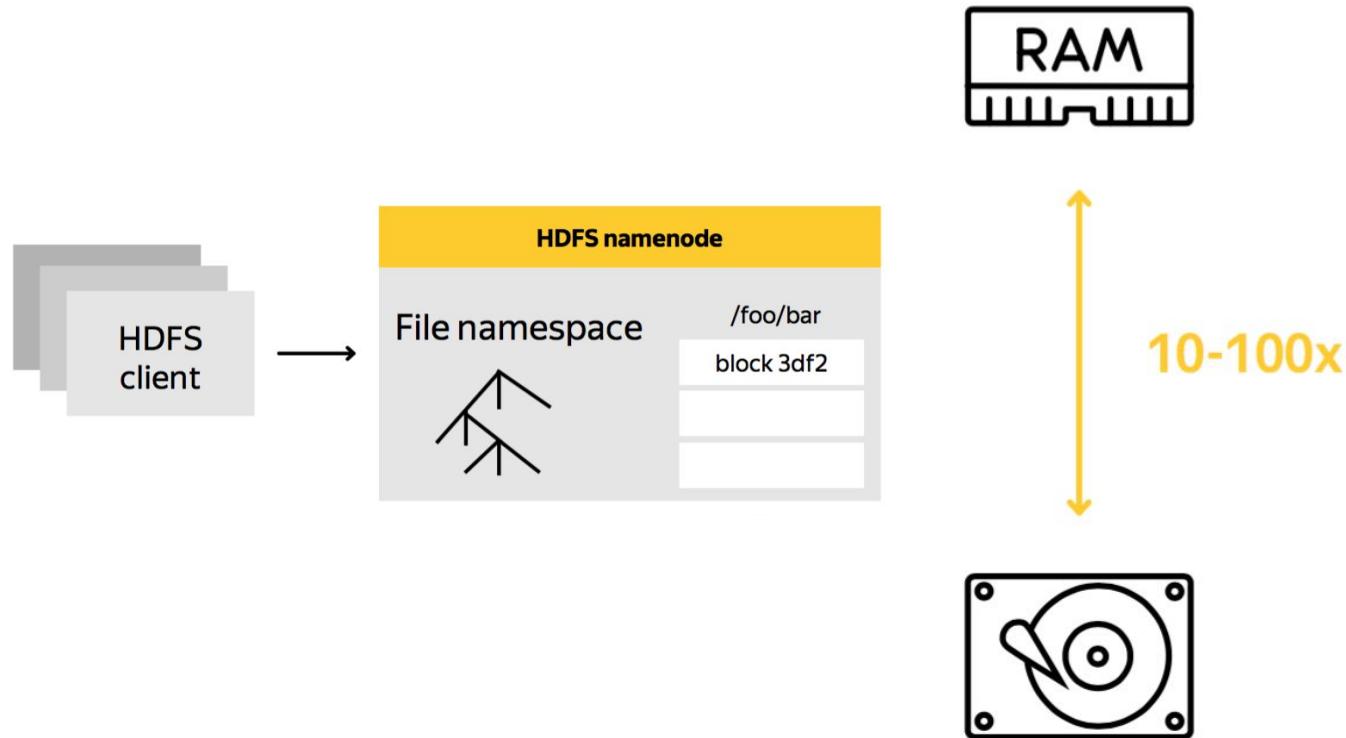
- HDD seek time: 0.2-0.8 мс
- SSD ~seek time: 0.08-0.16 мс
- reading speed - 600 MB/сек
- 32 MB ~ 50 мс



# Архитектура Namenode

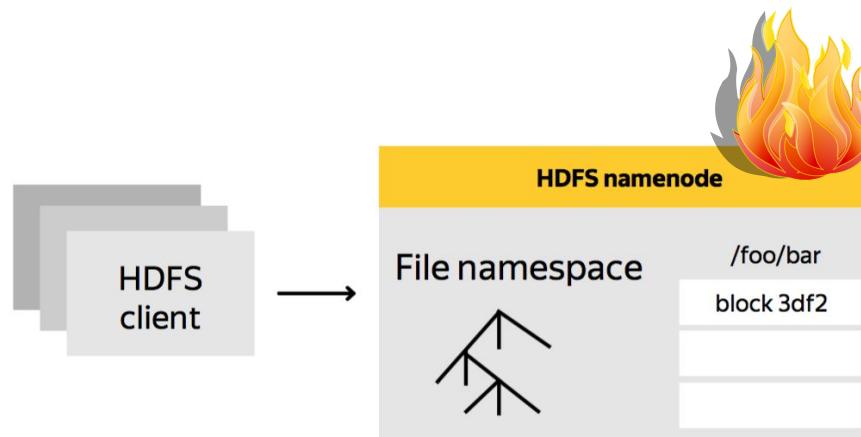


# Быстрый доступ





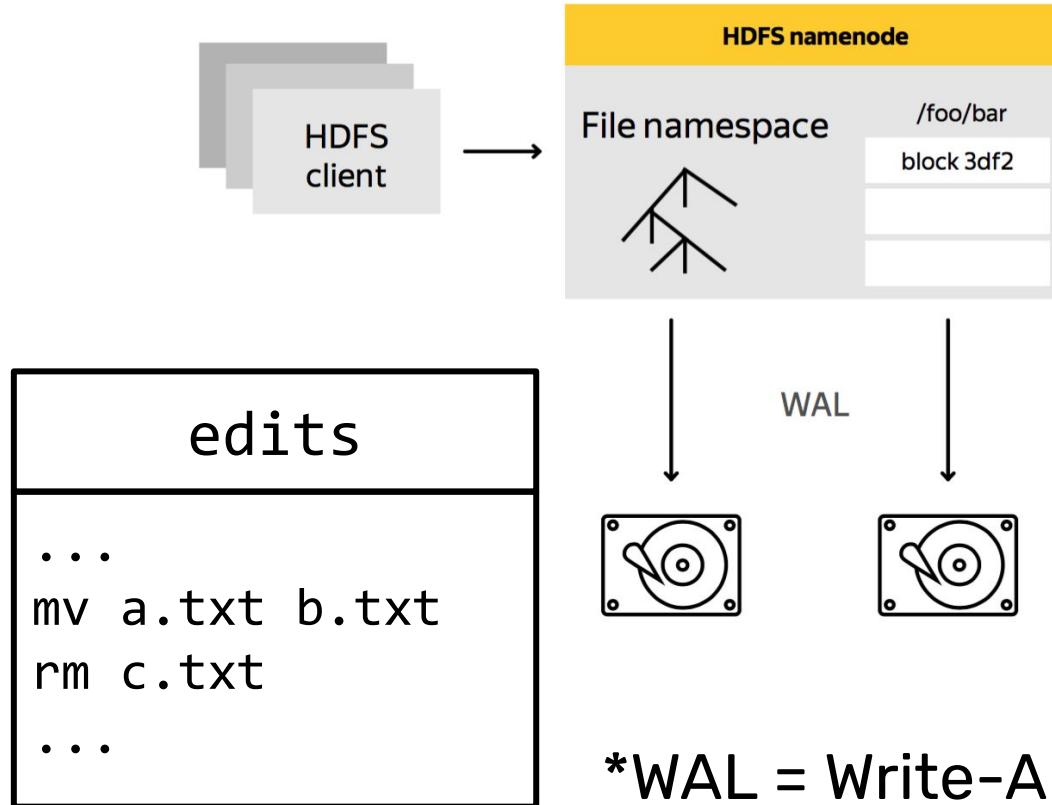
# Namenode - это SPoF\*



\*SPoF = Single Point of Failure

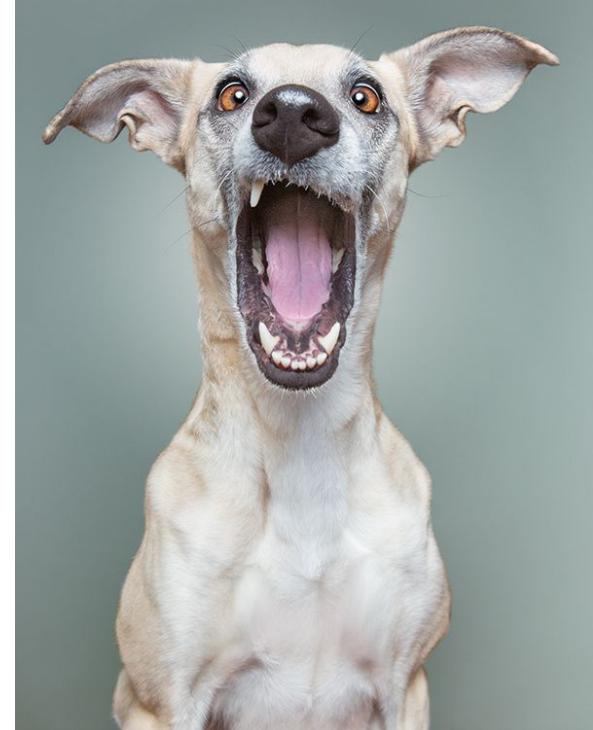
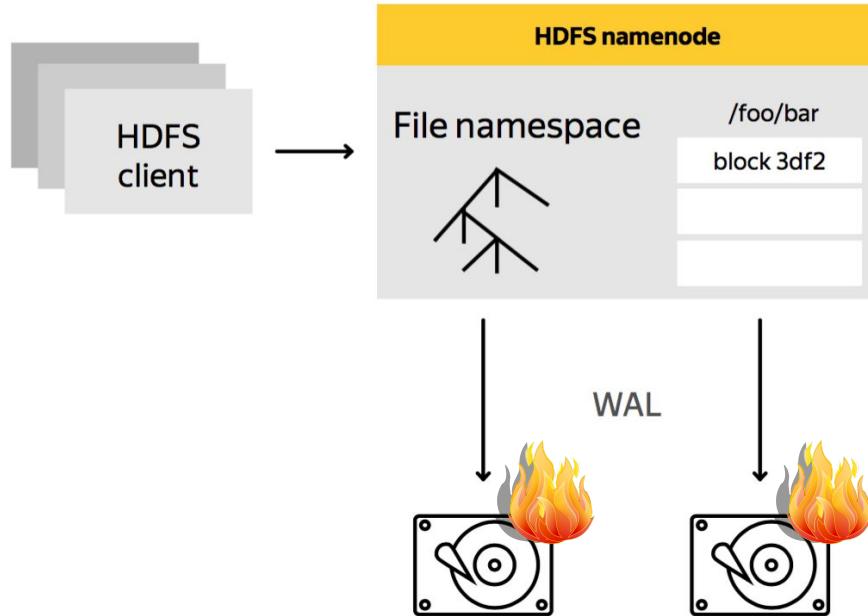


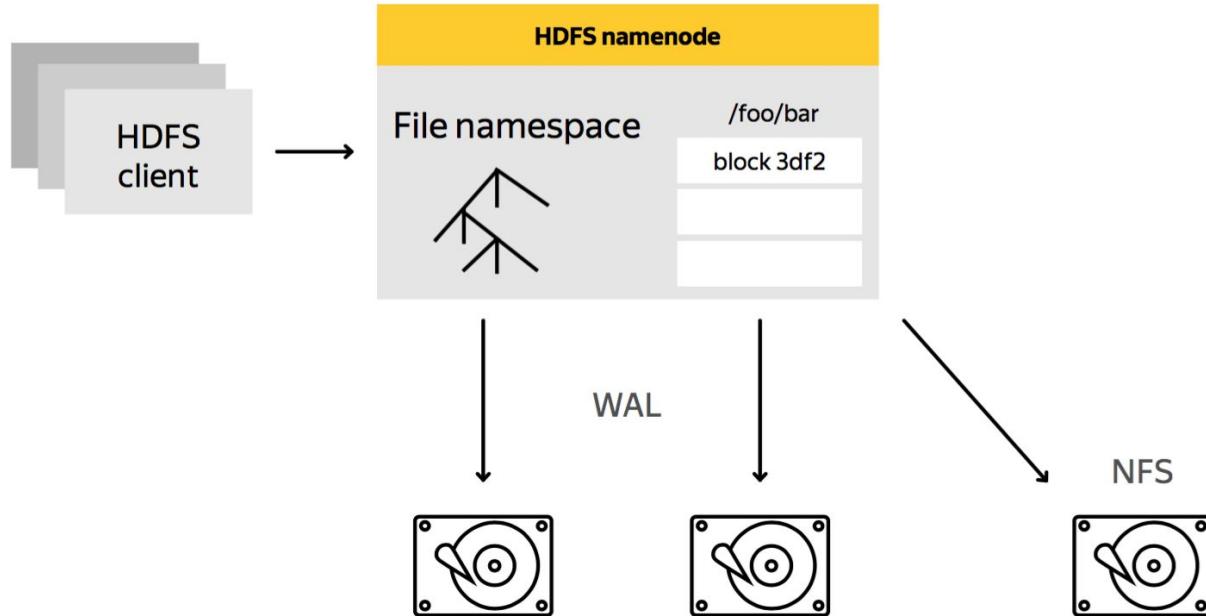
# Борьба со сбоями Namenode

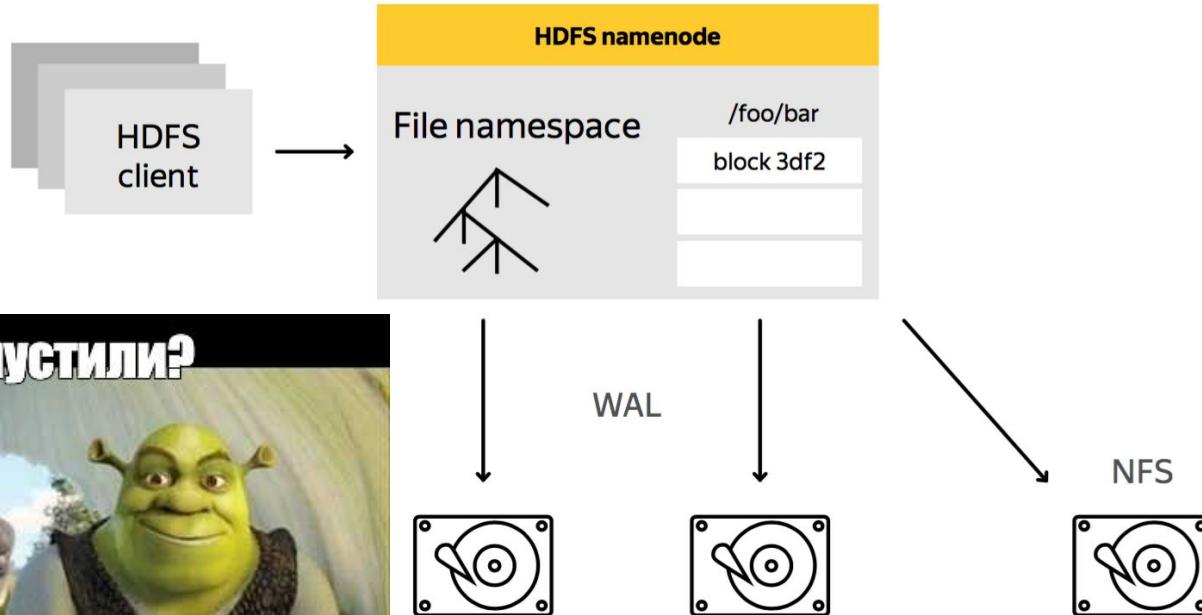




А что если...

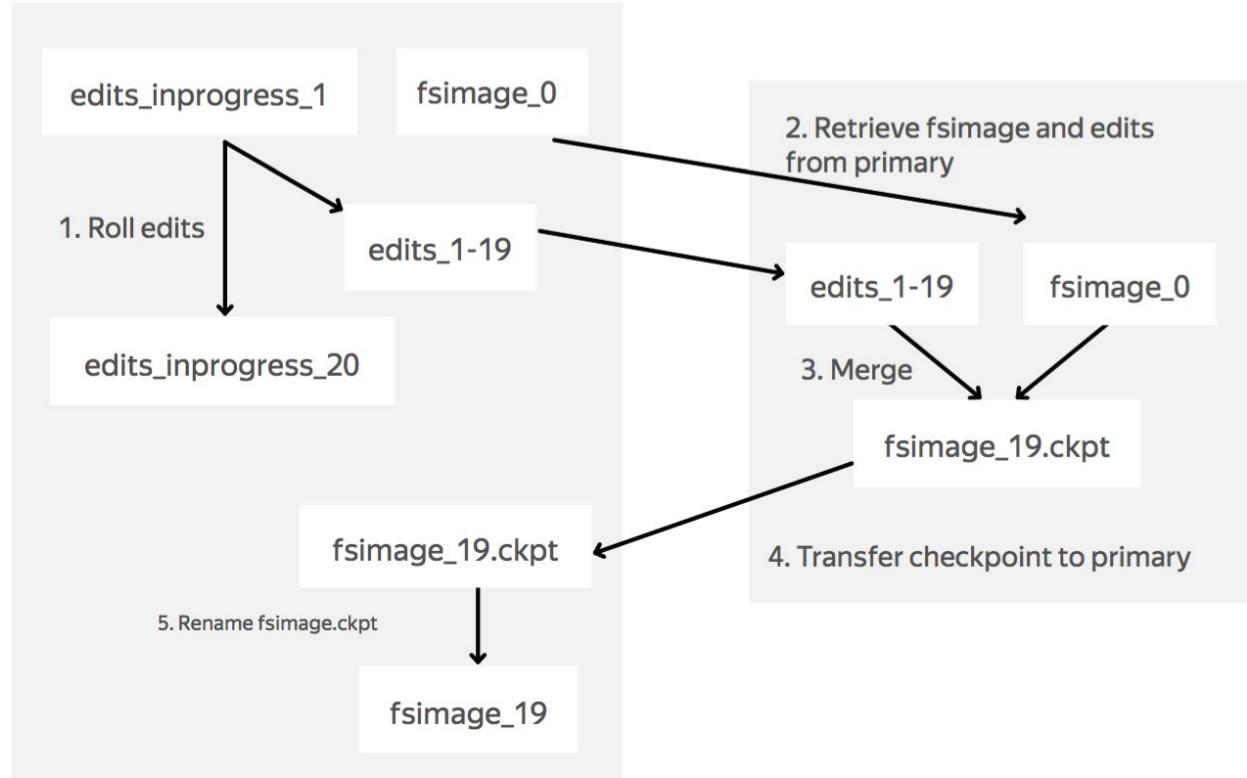








# Checkpoint Namenode





## Hadoop Sizing:

- ▶ **2x** ресурсов на Namenodes
- ▶ NFS

## Избегать:

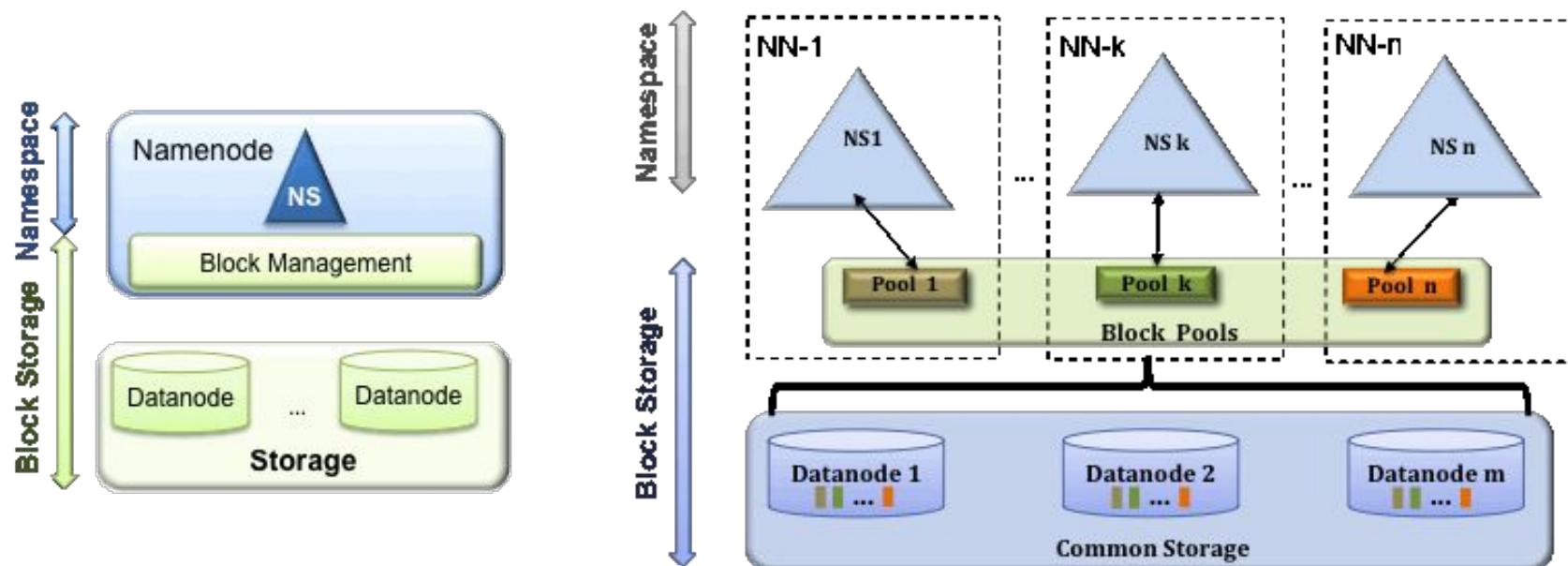
- ▶ Secondary Namenode
- ▶ Backup Namenode

## Использовать:

- ▶ Checkpoint Namenode



# HDFS 2.0: HDFS Federation

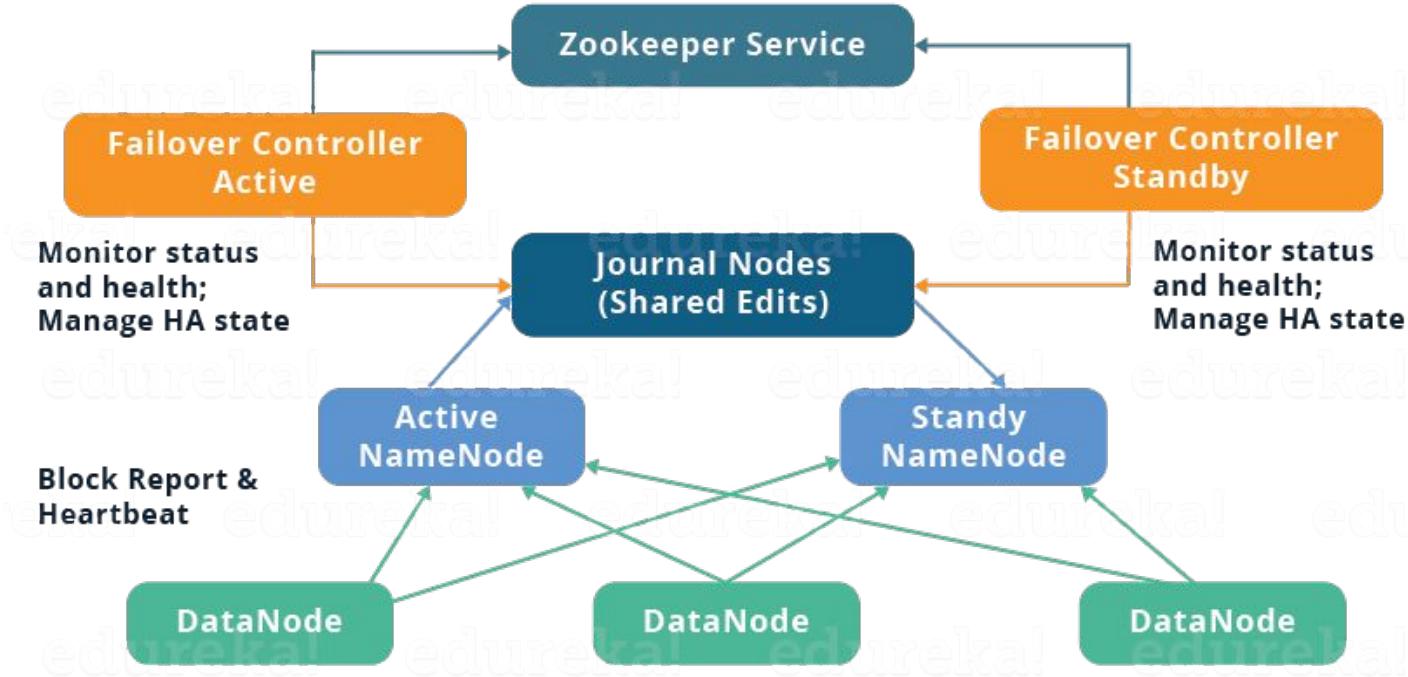


**было**

**стало**



# Синхронизация через Journal Nodes





Hadoop and NoSQL Downfall Parody

<https://www.youtube.com/watch?v=hEqQMLSXQIY>



- ▶ Вы можете перечислить 3 типа Многопроцессорных Вычислительных Систем, а также их типовые приложения
- ▶ Вы можете нарисовать диаграмму переходов состояний для блока и реплики
- ▶ Вы можете объяснить “design choice” архитектуры Namenode, а также указать на разницу между различными типами Namenode: Primary / Secondary / Checkpoint / Backup.
- ▶ Вы можете оценить ресурсы, необходимые для Hadoop-кластера (Hadoop sizing) для решения вашей задачи, а также можете объяснить что такое small files problem

# Спасибо! Вопросы?

Feedback: [http://rebrand.ly/x5bd2021q1\\_feedback\\_01\\_hdfs](http://rebrand.ly/x5bd2021q1_feedback_01_hdfs)

**Алексей Драль**, [aadral@bigdatateam.org](mailto:aadral@bigdatateam.org)

CEO at BigData Team, <http://bigdatateam.org>

<https://www.linkedin.com/in/alexey-dral>

<https://www.facebook.com/bigdatateam>