

Big Data Introduction

Distributed File Systems

Dral Alexey (aadral@bigdatateam.org)

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

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

Course Team

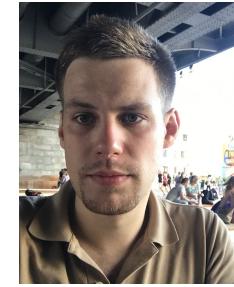
Big Data Instructors



Alexey Dral



Anton Gorokhov

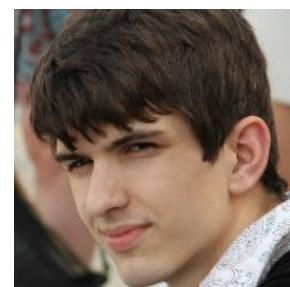


Artym Vybornov

TA / Administration



Victoria



Nikolay

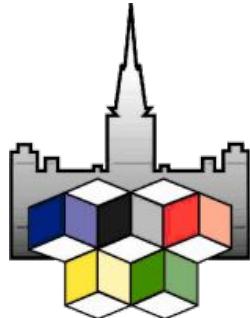


Oleg



**BIGDATA
TEAM**

- ▶ AESC MSU
- ▶ MSU
- ▶ YSDA
- ▶ Rambler
- ▶ Yandex
- ▶ Amazon AWS
- ▶ MIPT
- ▶ Sberbank
- ▶ BigData Team



Yandex



SBERBANK



SCHOOL OF DATA ANALYSIS



About Myself



- ▶ Distributed File Systems, HDFS
- ▶ MapReduce and Optimization Techniques
- ▶ Hive and Optimization Techniques
- ▶ Spark RDD, DataFrames and its Optimizations
- ▶ RealTime: Kafka and Spark Streaming
- ▶ NoSQL with Cassandra
- ▶ Data Ingestion and Layout Optimization
- ▶ Spark ML



IMPACT

Development is an ongoing process driven largely by on-the-job experiences

10%

Formal
Training

20%

Feedback
and Coaching

70%

On-the-Job Experiences/
Development in Role

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



- 12 домашних заданий
- 2 тестирования

Материалы: http://rebrand.ly/mf2019q2_gdrive_public



HDFS	HDFS	MR	MR	MR	MR	Hive	Hive	Hive	Spark	Spark	Spark	Spark	RT	RT	RT	RT	NoSQL	NoSQL
0.75	0.4	0	1	0	0.8	0	0	0	1	1	1	0	0.75	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
1	0.6	0	1	0	0.8	0	0	0	1	1	0	1	0	1	1	0	0	0
0	0	0	0.75	0	0	0	0	0	0	1	1	0	0	0	0	0	0	0
0.75	0.4	0	1	0	0.8	0	0.5	0	1	0	0	0	1	0.67	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0.75	0.6	0	0.75	1	0.8	0	0	0	1	0	1	1	0.75	1	0	1	0.33	0
0	0	0	0	0	0	0	0	0	0	0	0	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	0	0.75	0	1	1	0	1	1	1	0	0	0.75	0	0	0	0	0
1	0	0	1	0	0	0	0	1	0	1	1	0	0	0	1	1	0	0
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
0.75	0.4	0	1	0	0.6	1	1	0	1	1	1	0	1	0.67	0	0	0	0
1	0	0	0	1	0	1	0	0.75	1	1	0	1	0	0	0	0	0	0



MF-2019-Q2	before	after	gain	support	<hidden>	before	after	gain	support
HDFS	28%			2	HDFS	42%	80%	38%	1
MR	28%			4	MR	35%	79%	44%	4
Hive	27%			4	Hive	30%	71%	41%	4
Spark	30%			4	Spark	38%	63%	25%	5
RT	13%			4	RT	21%	68%	47%	4
NoSQL	4%			2	NoSQL	17%	76%	59%	2
avg-score max: 20	4.6			20	avg-score max: 20	6.1	14.2	8.1	20
	23%	0%	0%			31%	71%	40%	



-  [Big Data: Motivation](#)
-  [Distributed File System](#)
-  HDFS Workshop

Big Data: Motivation



**BIGDATA
TEAM**

Big Data: Motivation

Data Scientist: The Sexiest Job of the 21st Century

by Thomas H. Davenport and D.J. Patil

FROM THE OCTOBER 2012 ISSUE





Откуда берутся (большие) **данные**?



Google search results for "data science".

data science

All Images News Videos Books More Settings Tools

About 55,900,000 results (0.58 seconds)

Data Scientist Masters Program - 12+ industry-based projects
Ad www.simplilearn.com/Data_Scientist/certification ▾
Mentorship from industry experts, industry-recommended learning path. Start Now!

Data science, also known as **data-driven science**, is an interdisciplinary field about **scientific methods**, processes, and systems to extract knowledge or insights from **data** in various forms, either structured or unstructured, similar to **data mining**.

Data science - Wikipedia
https://en.wikipedia.org/wiki/Data_science

www.edx.org

About this result Feedback

Data science - Wikipedia

https://en.wikipedia.org/wiki/Data_science ▾

Data science, also known as data-driven science, is an interdisciplinary field about scientific methods, processes, and systems to extract knowledge or insights from data in various forms, either structured or unstructured, similar to data mining.

[Data mining](#) · [List of statistical packages](#) · [Business analytics](#) · [Information science](#)



Twitter, Inc.



VKontakte



Instagram

...







PSTS1_MYCT : MKIRLHTLLAVLTAAPLLLAAAGCGSKPPSGSPETGAGAGTVATTPASSP---VTIAETGSTILLYPLFNLNGPAEHERYPN : 78
M 6 gaG 3f yP6 W 5 tgn

PSTS_ECOLI : K-VNYCGIGSSGGVVKQIIANTVDFGASDAPISDEKLAQEGIFQFPPTVIGGVVLAV-NIPGILKSGELVIDGKTLGDIYLGKIK : 134
PSTS_PASMU : K-VNYCSIGSGGGQQCIIAKTIDFGASDDPMKAELLAQHQILQCFPAIIGGTVPVV-NLPEITAGQLKLSGEVLAIDIFLGKIK : 129
PSTS_RHILo : VGLNYCSIGSGGGIRQVIAKTVTFGATDKPMSDADIEKNGIVQFPMVMGGIVPIV-NLTGVKPGELVIDGKTLAQIYLGAIT : 135
PSTS_HAEIN : K-VNYCSIGSGGGQQCIIAKTIDFGASDDPMKSELLQQHQIVQFPNAVIGGIVPVV-NLPEIKPGKLKISCKLLAEIFLGKIK : 127
PSTS_XYLFA : K-INYCSIGSGGGIAQIKAATIDFGSSDKPLDSSEITQAGIGQFPSSAIGGVVPPV-NLDNIEPGKIRLITCPPLLADIFLGKIS : 158
PSTS1_MYCT : VTITAQGTGSGAGIAQAAAAGTVNIGASDAYLSEGDMAAHKGImnialaisaqqvnyNIPGV-SEHLKLNGKVLAAMYQGTIK : 159
6nyQ iGSggG Q A T6 fGa3D p6 6 1 qfp gg v v N6 6 g L L G La 65IG I

PSTS_ECOLI : KWDDDEAIAKLNPGLKLESQNIAVVRRADGSQTSEVFTSYLAKVNNEE-WKNNVGTGSTVVRNPIGLGGKNDGIAAFVQRLPGA : 215
PSTS_PASMU : KWNDPAIAKLNQGANLPDKAIIVVHRSRDGSQTTEFWGWTNYLSKVSTE-WKETVQGKSVWKPQGGKNEGVAAYVSKIKYS : 210
PSTS_RHILo : TWDDAAIKALNPSLTLPESTAIAVVHRSRDGSQTTFNFTNYLVKLSPD-WKDVKVGSDDTAWEWPVGAKGSEGVANTVKQTDGG : 216
PSTS_HAEIN : KWNPDPDVALNPTLPLPNKNIIVIHRSDGSQTTFEGFTNYLSKISND-WKNQVGECKSVWNLTGQGGKNEGVAASYVRQMKYS : 208
PSTS_XYLFA : KWNDAAIIISANPGLHLPTKINIVIHRSDGSQTTFNFSNYLSKVSAAE-WKQKVGEGETSVQWNPQGVGGKNEGVAASYVQQIKGS : 239
PSTS1_MYCT : TWDEPQIAALNPGVNVLEGTAVVPLHESDGSGDTELETQYLSKQDPEGWKGSPGFRTVDEPAVPGALGENNGGMVTGCAET : 241
W1D 6 1Np LP 6 6hRsDGSGt3F 53 YL K Wk vG g V 5p g G eG G a v

PSTS_ECOLI : IGYVEYAY----AKQNNLAYTKLISADGKPVSPTEENFANAAGKADW--SK--TFAQDLTNQKGEDAIPITSTTFILIHK : 287
PSTS_PASMU : IGYVEYAY----AKQNQLAWSILQNKAGQFVQPSAESFMAAAANAQWESAV--GMGVILITNEEGDTSWPVTAASFILLHK : 284
PSTS_RHILo : IGYVEYAY----AKQNNLNSYSKMLNAAAGKVVVERPSESFGAAAASNADFKGAK--NFNVIITNEPGDTTWPIAASTWVLIHK : 290
PSTS_HAEIN : IGYVEYAY----AKQNQLAWSIQNQAGQFVQPSNESFMAAAASHAKW--HekaGMGVILITNETGEKSWPITAASFILLNK : 282
PSTS_XYLFA : IGYVELAY----ALQNKMSYTALQNAAGQWVQPSAESFAAAASNADWSNAK--DFNLVITNATGEAAWPITATNFILMRK : 313
PSTS1_MYCT : PGCVayigisfldqasRCIGEAQIGNSSGNFILEDAOSICAAAAGFASKTPA--NqaismidgPAPDCYPIINYEYAIVNN : 321
iGyVeyay A Qn 6 6 n G 6 P 2sf aAA a 6t1 g 5P6 5 66 k

* 340 * 360 * 380 *



Internet of Things (IoT)

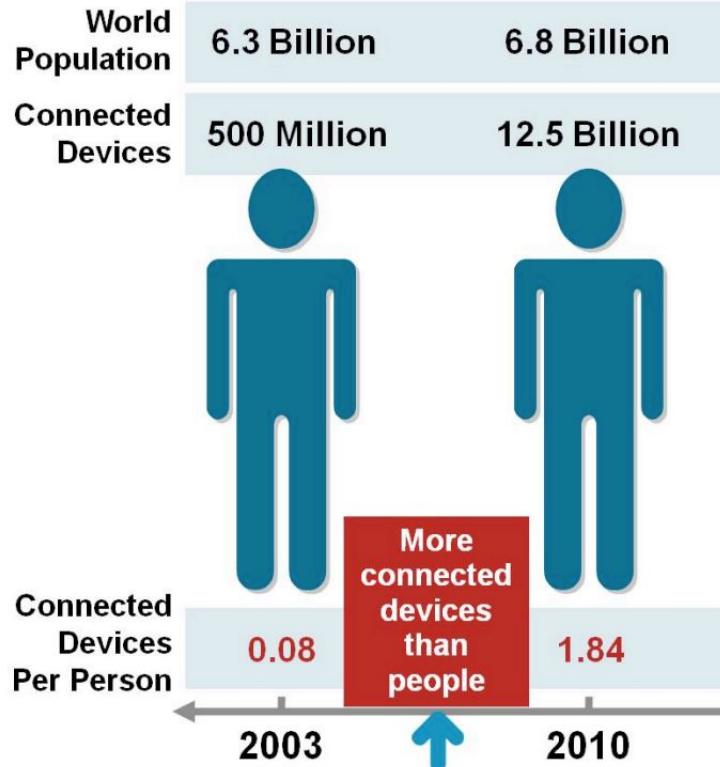




Q&A: IoT vs Human Beings?

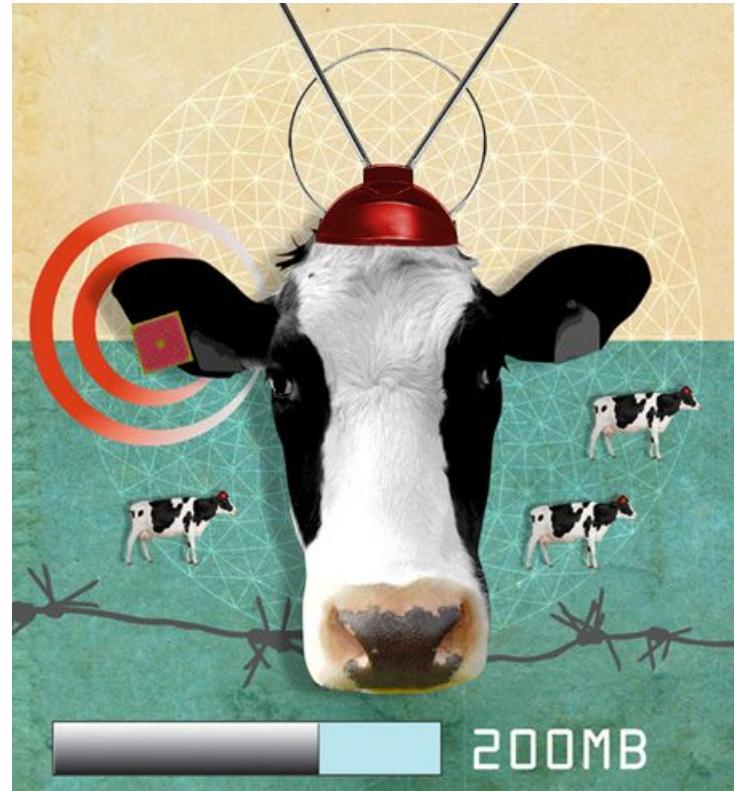


**BIGDATA
TEAM**

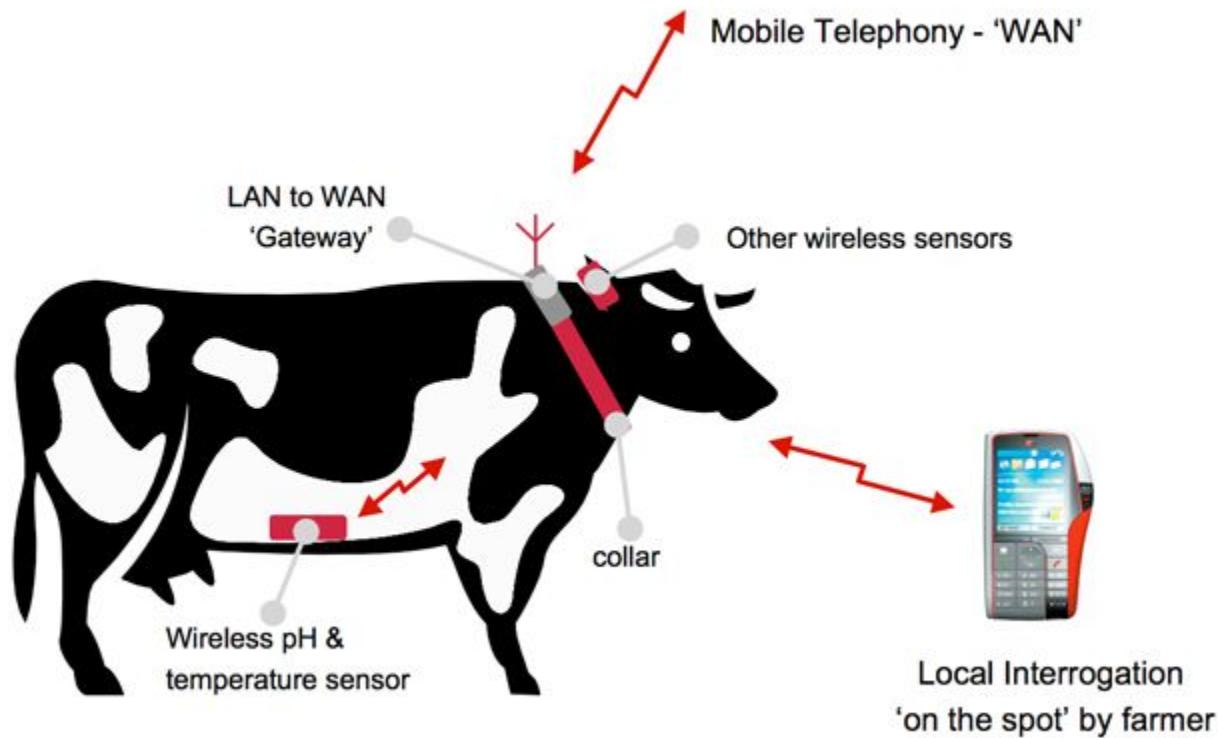


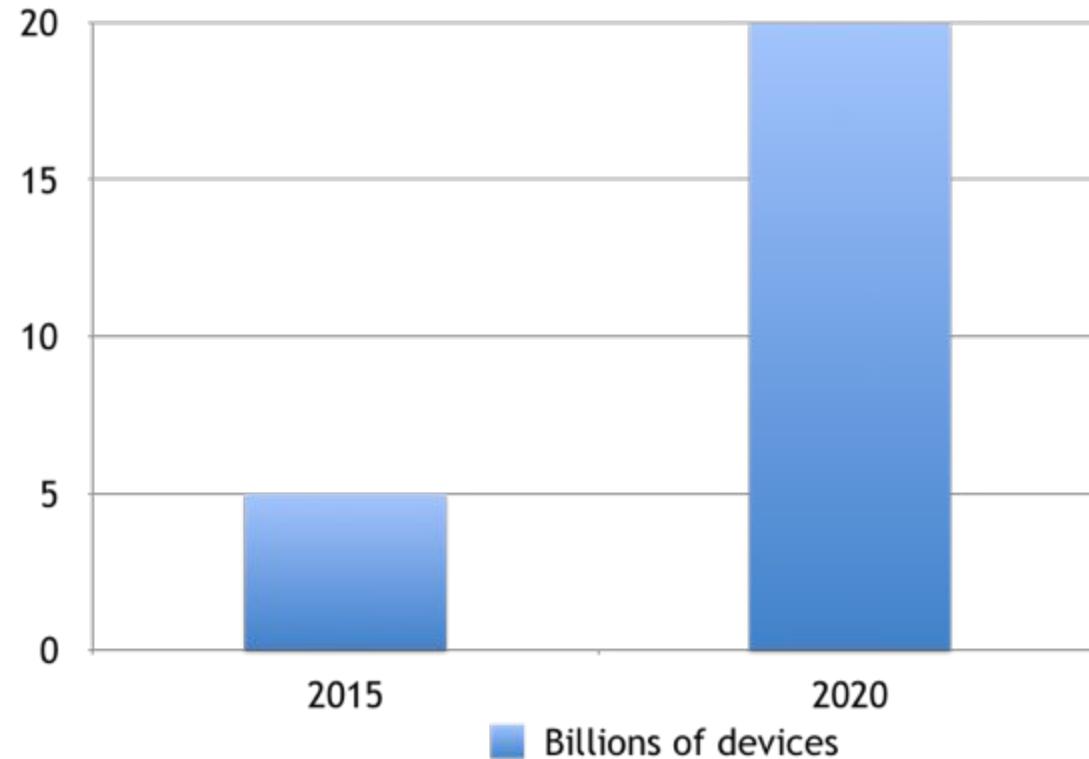
Source: Cisco IBSG, April 2011

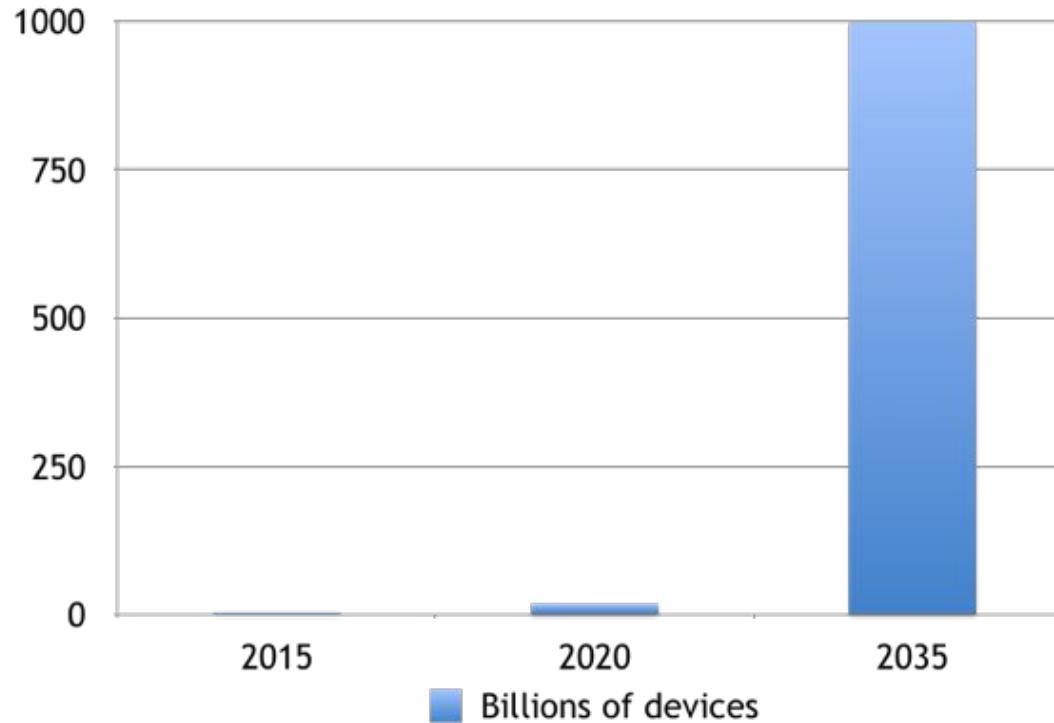
Internet of Things (IoT)



Source: The Economist, 2010



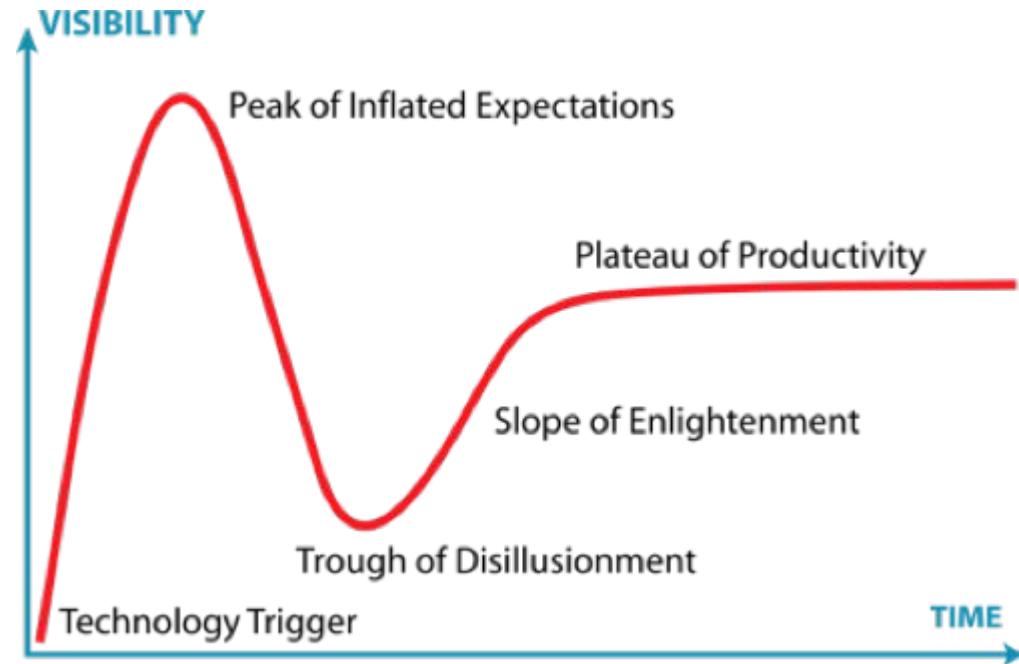
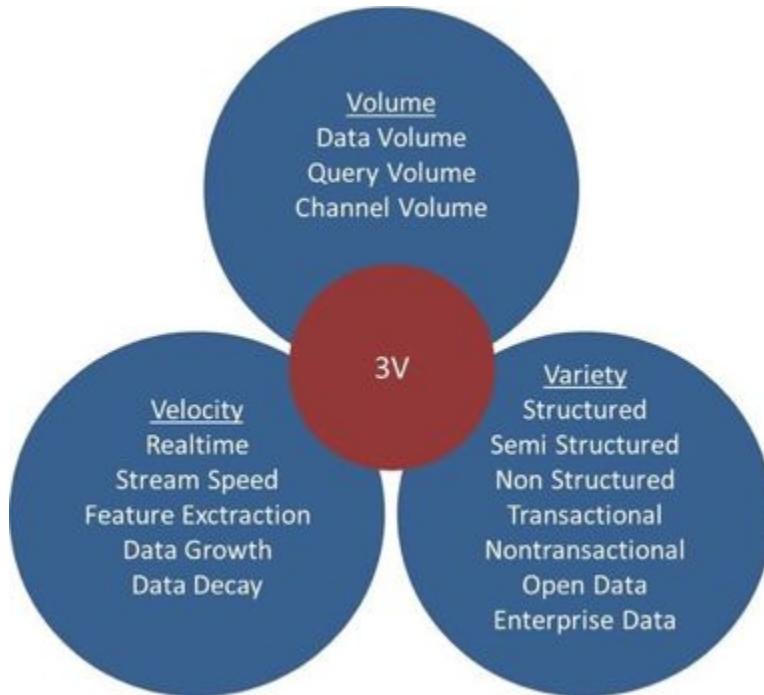




Q&A: Что такое Big Data?

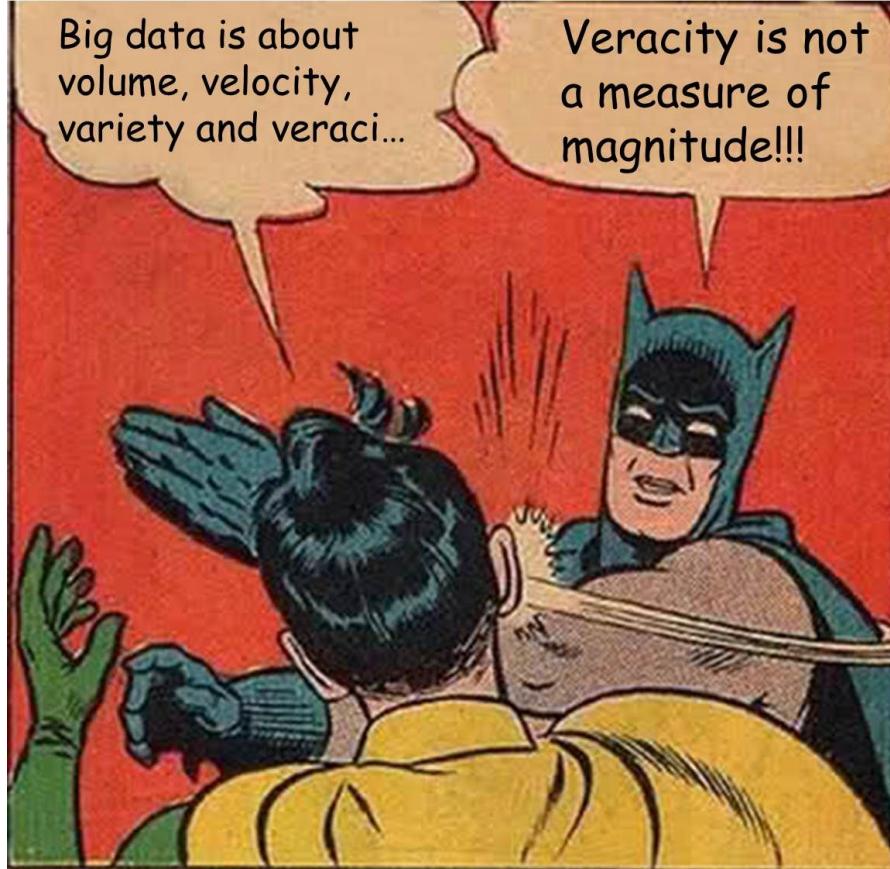


Big and Small Data





4+ Vs of Big Data





Где применяется (data) science?



Why Big Data?



Google data science

All Images News Videos Books More Settings Tools

About 55,900,000 results (0.58 seconds)

Data Scientist Masters Program - 12+ industry-based projects
Ad www.simplilearn.com/Data_Scientist/certification
Mentorship from industry experts, industry-recommended learning path. Start Now!

Data science, also known as **data-driven science**, is an interdisciplinary field about **scientific** methods, processes, and systems to extract knowledge or insights from **data** in various forms, either structured or unstructured, similar to **data mining**.

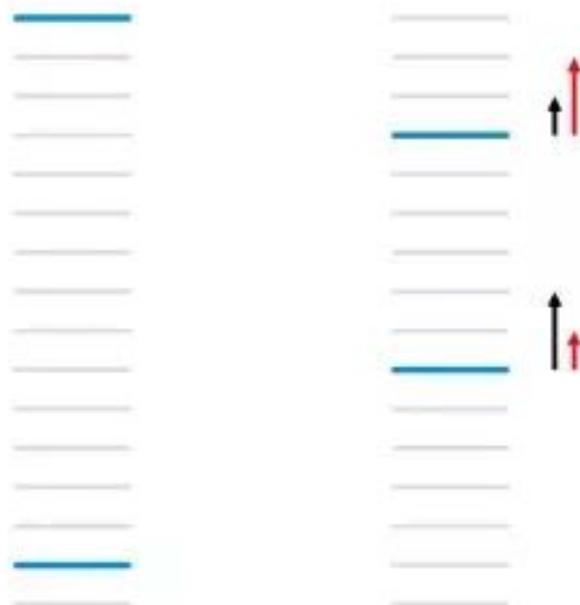
Data science - Wikipedia
https://en.wikipedia.org/wiki/Data_science

Data science - Wikipedia
https://en.wikipedia.org/wiki/Data_science ▾
Data science, also known as data-driven science, is an interdisciplinary field about scientific methods, processes, and systems to extract knowledge or insights from data in various forms, either structured or unstructured, similar to data mining.
Data mining · List of statistical packages · Business analytics · Information science

About this result Feedback



**BIGDATA
TEAM**



Learning to Rank



University of Illinois at Urbana-Champaign

Text Retrieval and Search Engines

★★★★★ 358 ratings

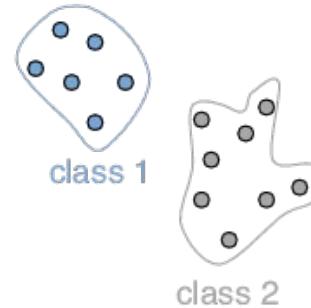
Course 2 of 6 in the Specialization [Data Mining](#)





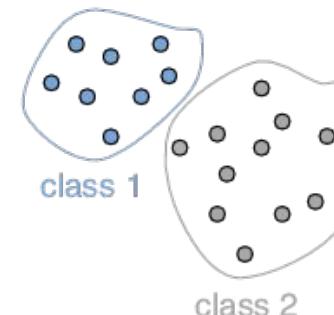
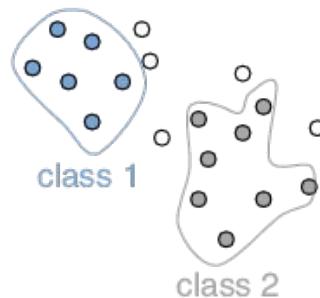
Classification

TRAINING



Classifier trained on labeled data

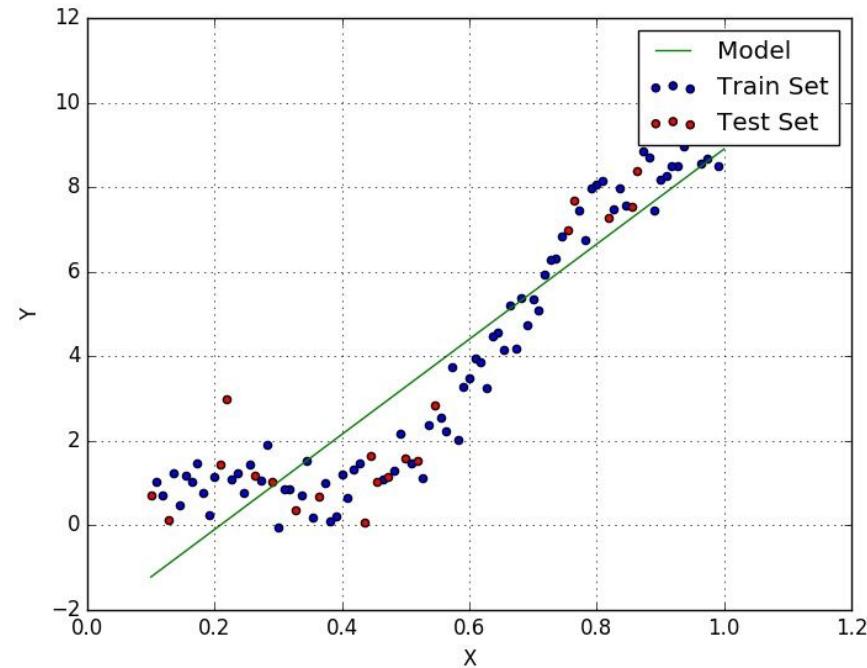
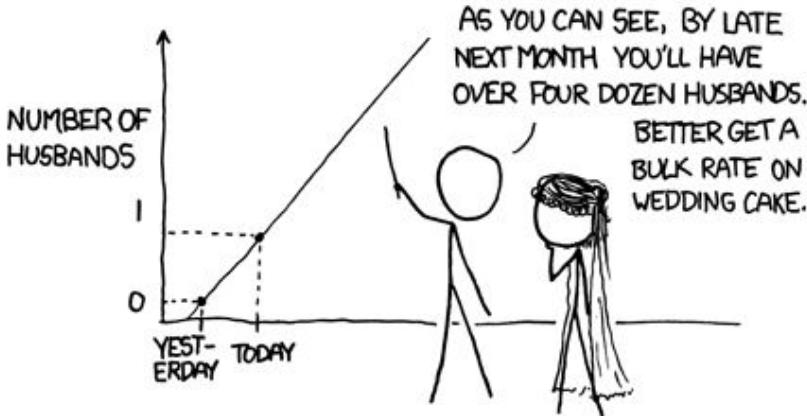
TESTING



Incoming unlabeled data → Classification



MY HOBBY: EXTRAPOLATING





**BIGDATA
TEAM**

Churn Prediction

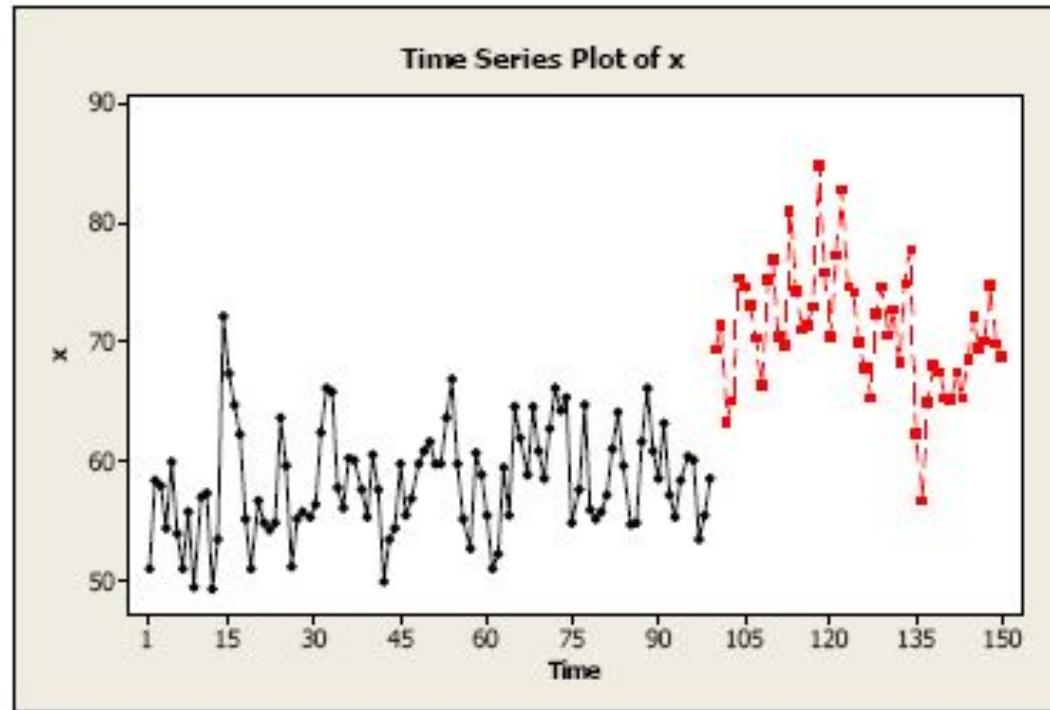




**BIGDATA
TEAM**

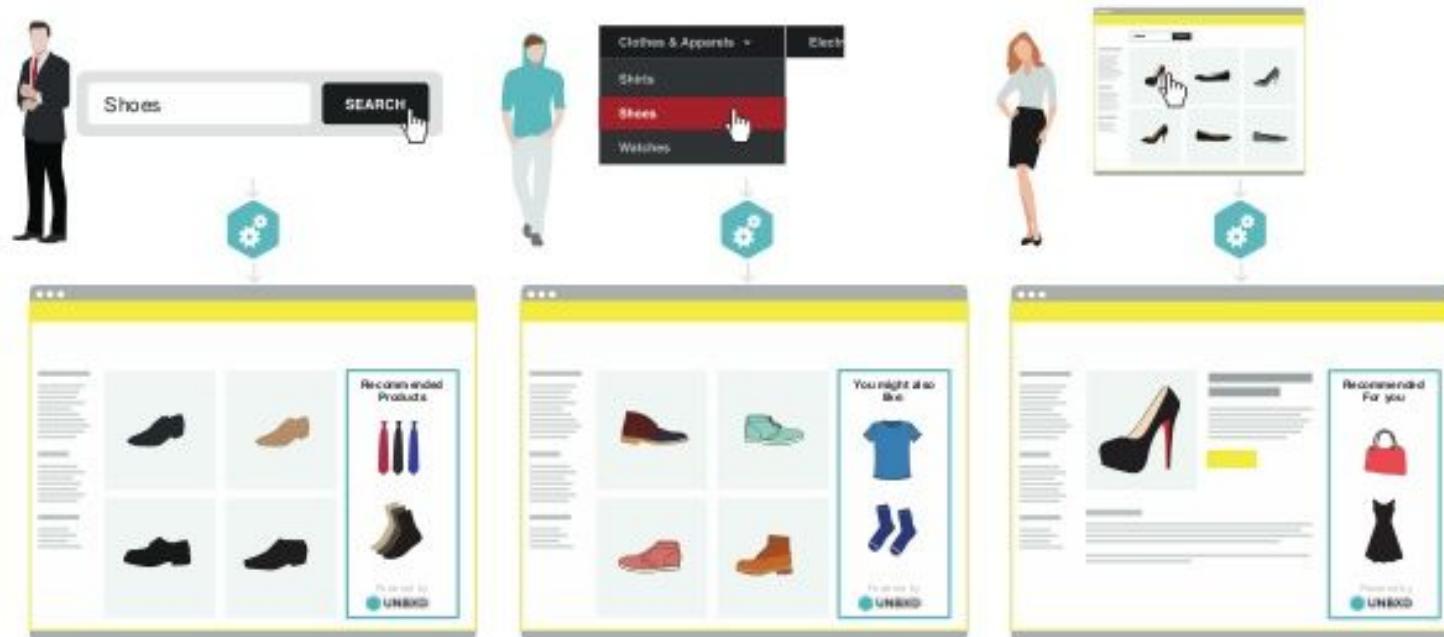
Demand Forecast







Personal Recommendations





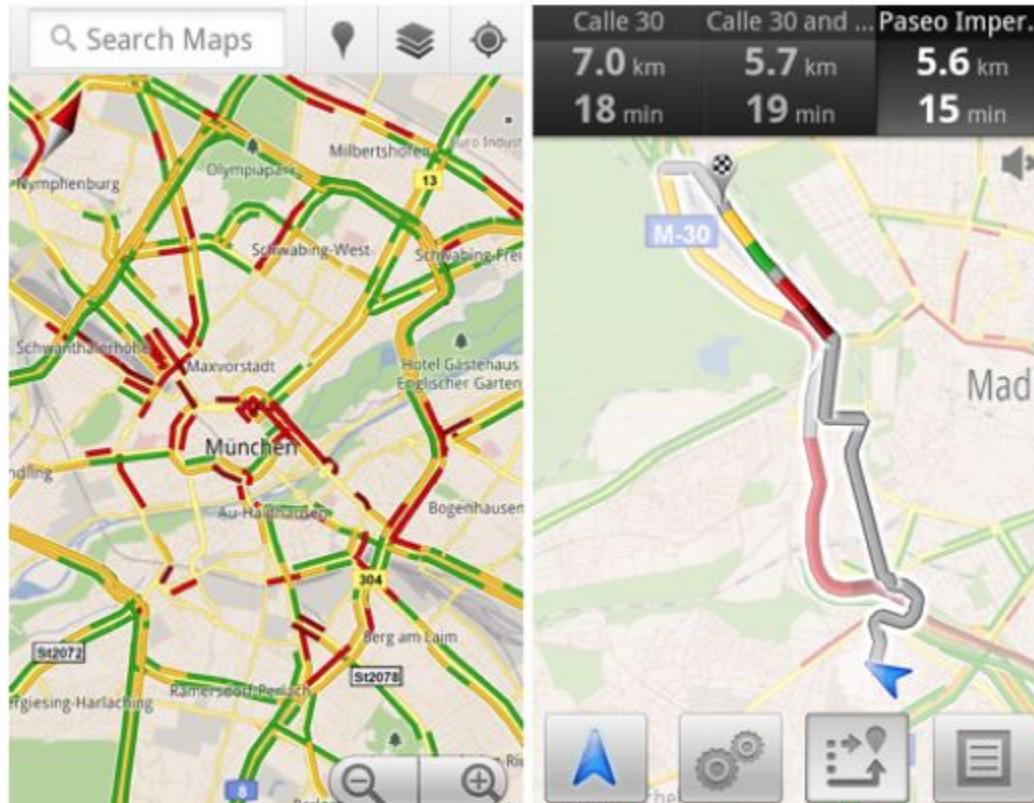


[https://yandex.ru/blog/company/stal-neft-i-iskusstvennyy-intellekt-yandex-data-factor
y-o-novoy-promyshlennoy-revolutsii](https://yandex.ru/blog/company/stal-neft-i-iskusstvennyy-intellekt-yandex-data-factor-y-o-novoy-promyshlennoy-revolutsii)



**BIGDATA
TEAM**

Traffic Jam Prediction & Routing







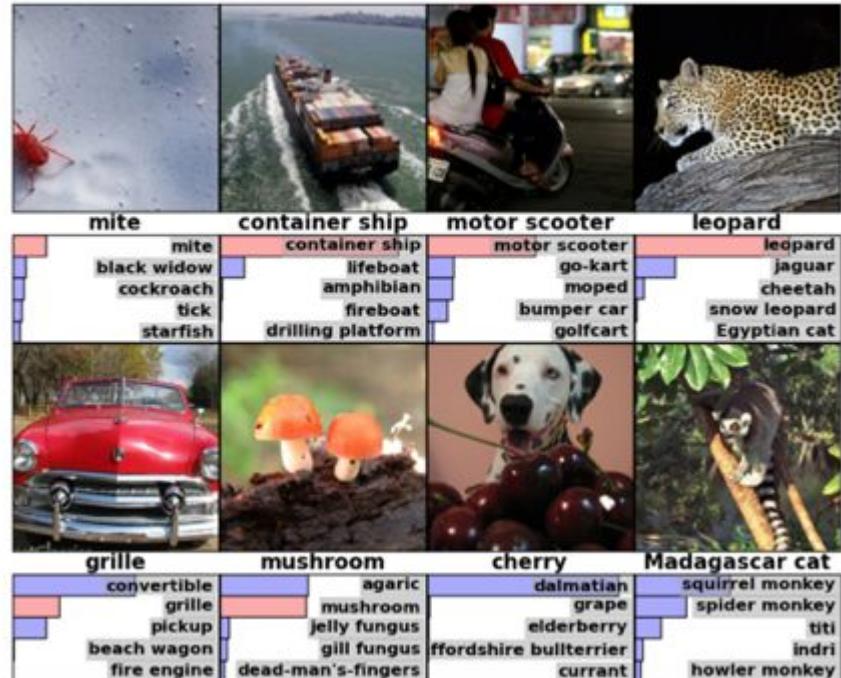
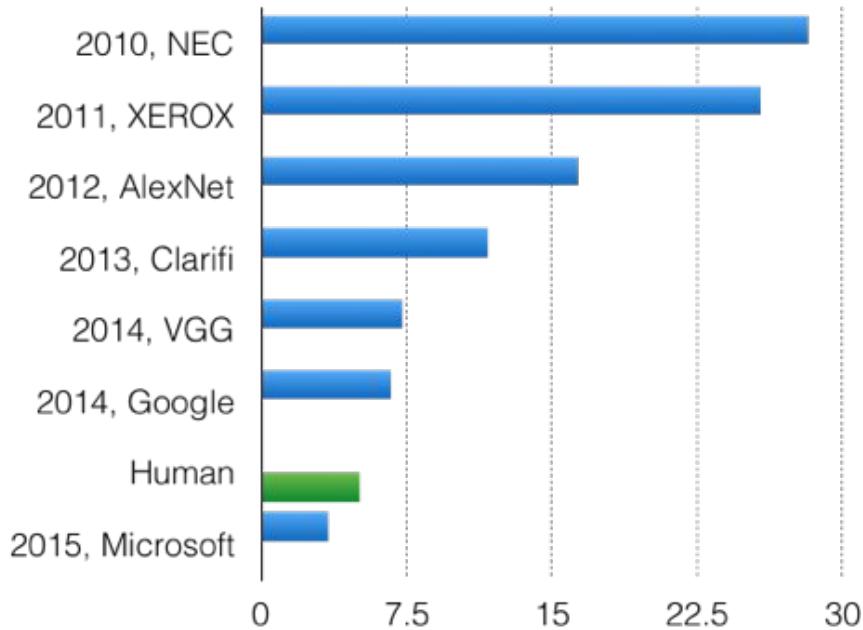
**BIGDATA
TEAM**

Urban Air-Transport



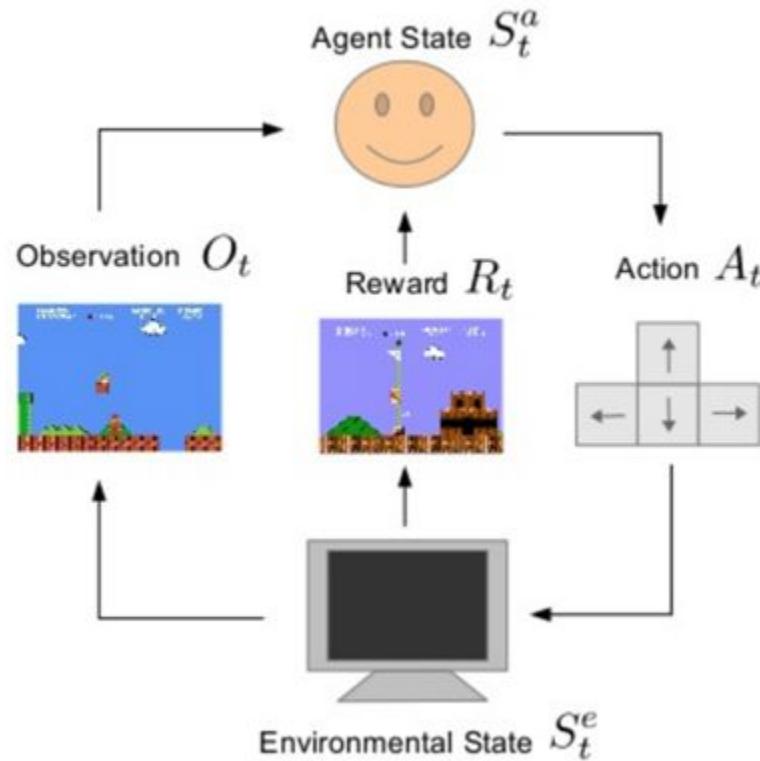


Image Recognition: ImageNet





Reinforcement Learning





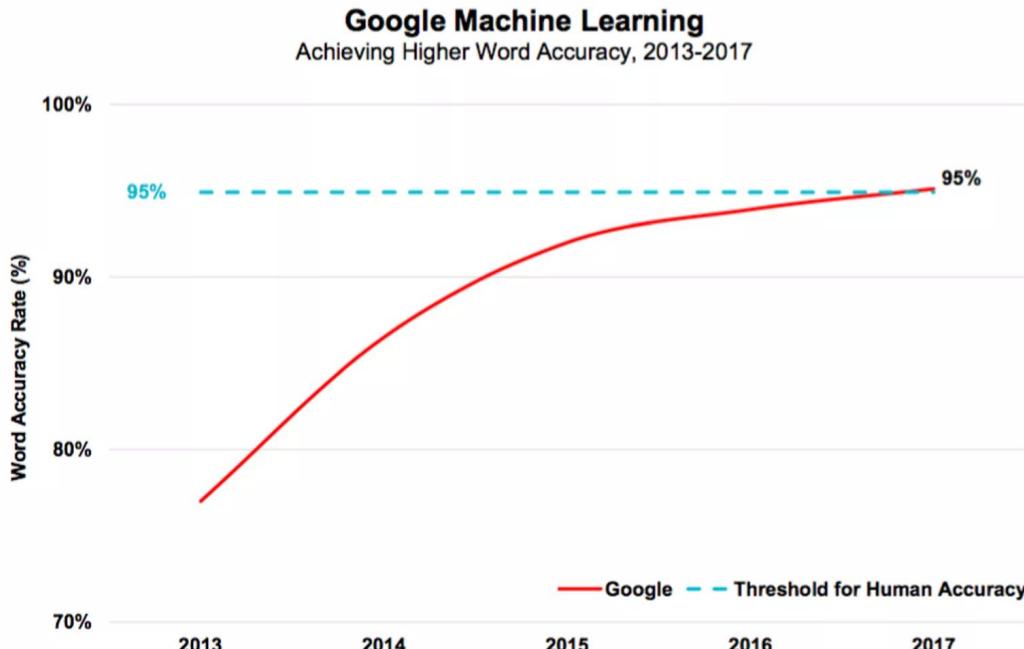
BIGDATA
TEAM

Data Science for People





...Voice-Based Platform *Back-Ends* =
Voice Recognition Accuracy Continues to Improve



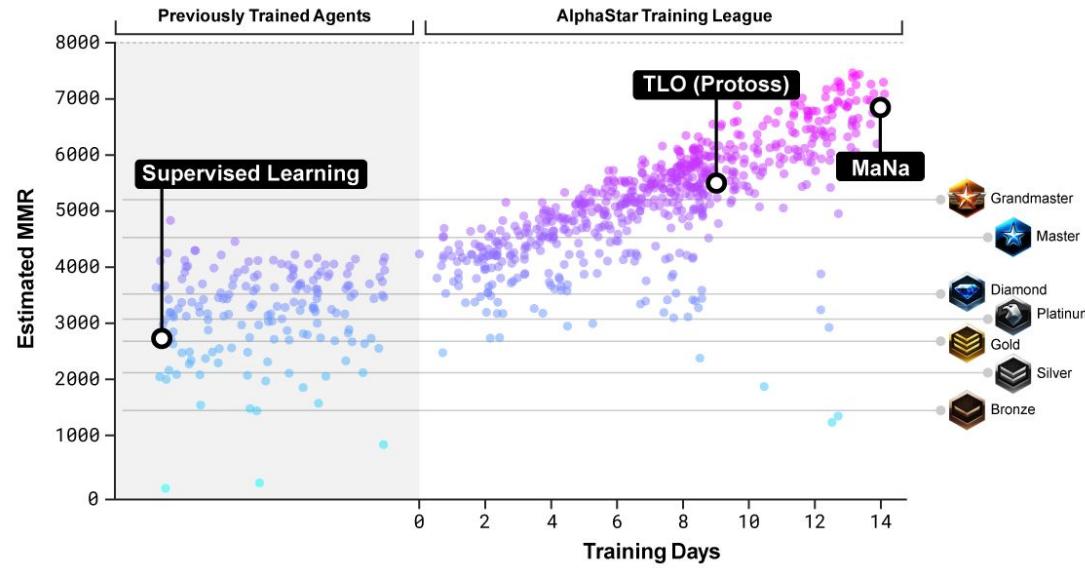


ПОМОЩНИК GOOGLE НАУЧИЛСЯ РАЗГОВАРИВАТЬ ПО ТЕЛЕФОНУ



"Mm-hmm."





- ▶ Training - 14 days, 16 TPUs / agent
- ▶ 200 years of real-time StarCraft play / agent
- ▶ Technical description is being prepared for publication

<https://deepmind.com/blog/alphastar-mastering-real-time-strategy-game-starcraft-ii/>



<http://aggeek.net/ru/technology/id/top-5-naibolee-mnogoobeschajuschih-tehnologij-v-smart-agriculture-431/>



**BIGDATA
TEAM**

Data Science for Business





BIGDATA
TEAM

Data Science: Outsourcing

The Home of Data Science & Machine Learning

Kaggle helps you learn, work, and play

Create an account

or

Host a competition

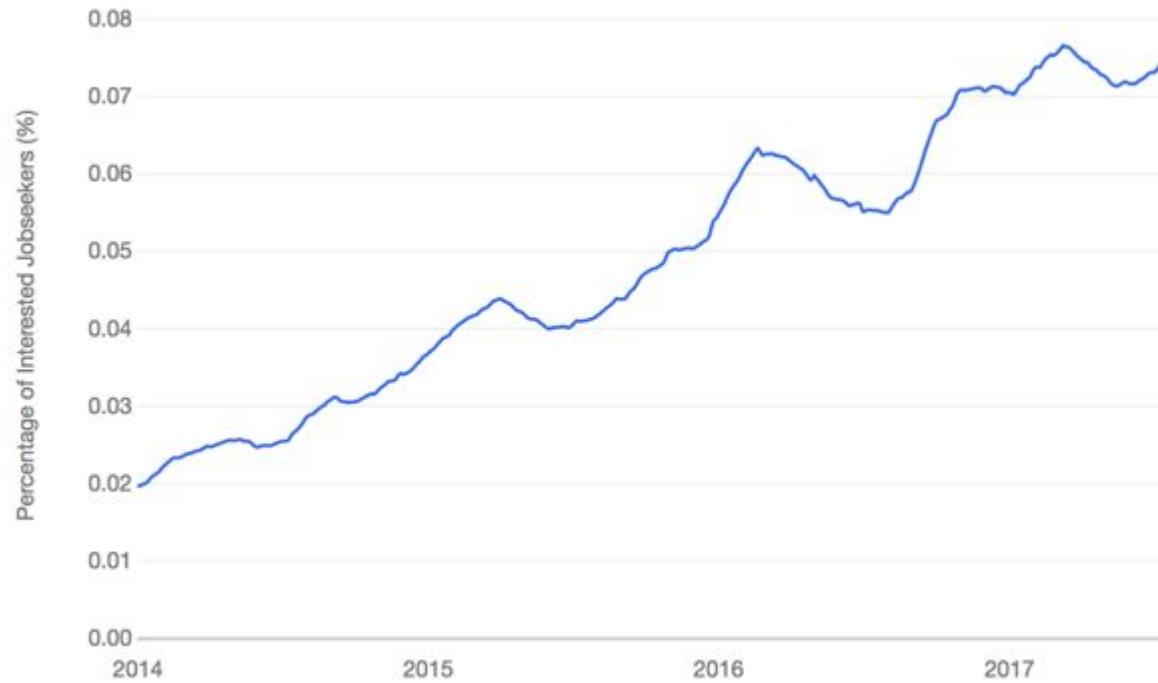
kaggle



Big Data: Motivation, v2



Data Science Job Index



indeed®



**BIGDATA
TEAM**

indeed

Index jobs

My recent searches

[data engineer - 34,228 new](#)

[data scientist - 9,279 new](#)

4X

Data Scientist? Data Engineer? Data ...?

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

Hive, Spark

Spark ML

Real Time

Yandex

<https://www.coursera.org/specializations/big-data-engineering>



**BIGDATA
TEAM**

Telegram Channel

http://rebrand.ly/mf2019q2_telegram_join



**BIGDATA
TEAM**

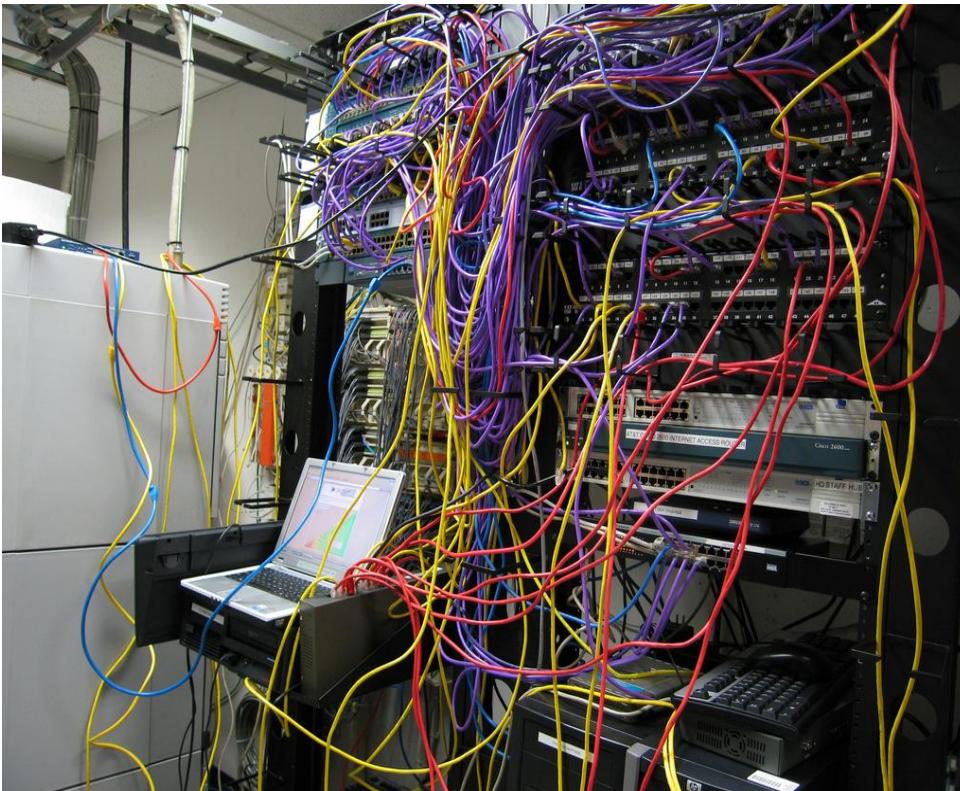
Tea / Coffee Break

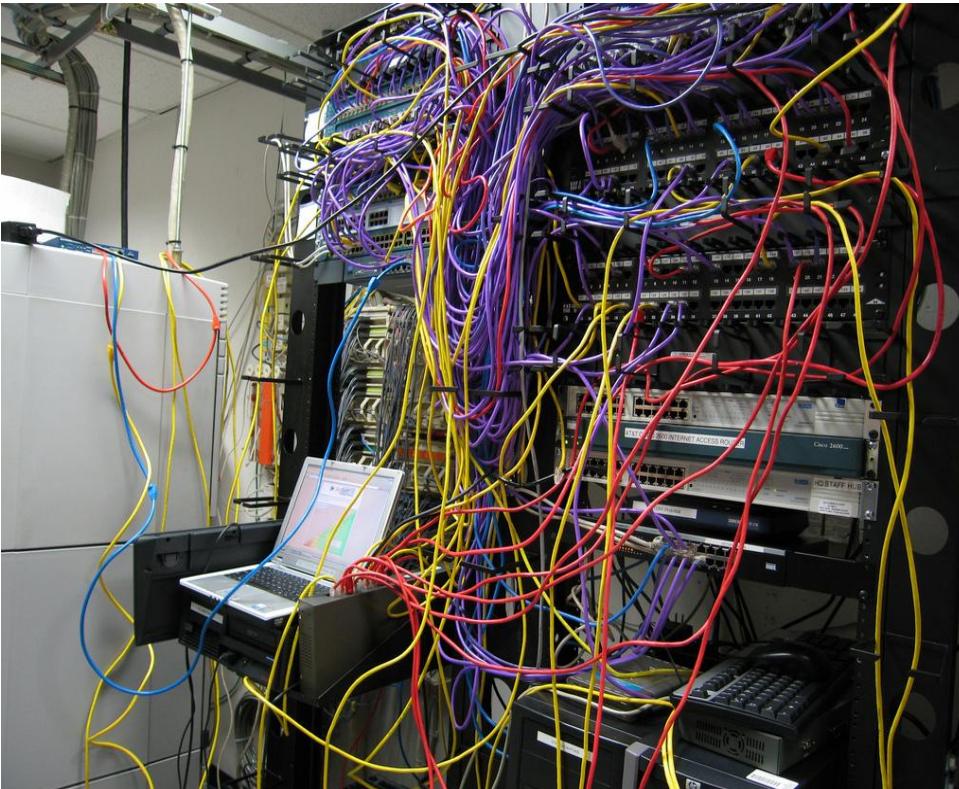


Distributed File Systems



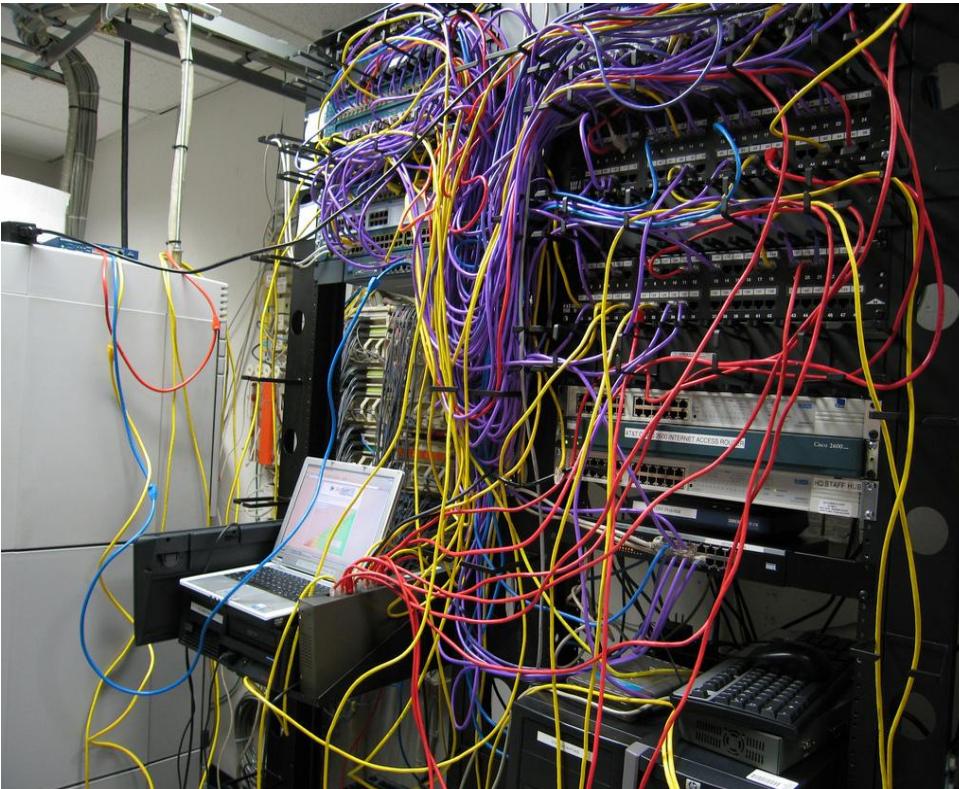
Q: what can break?





Node Failures

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



Link Failures

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



(a) Synchronous Model

- ▶ Clock Drift
- ▶ Clock Skew



- ▶ Parallel Computing
- ▶ Distributed Computing
- ▶ Grid Computing



**BIGDATA
TEAM**



Parallel Computing

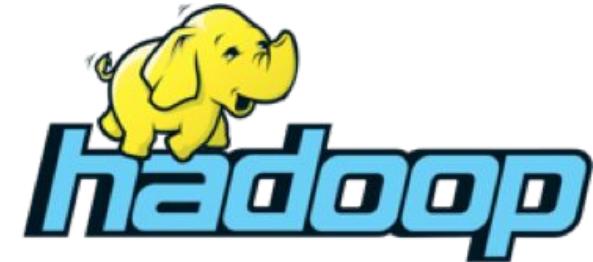
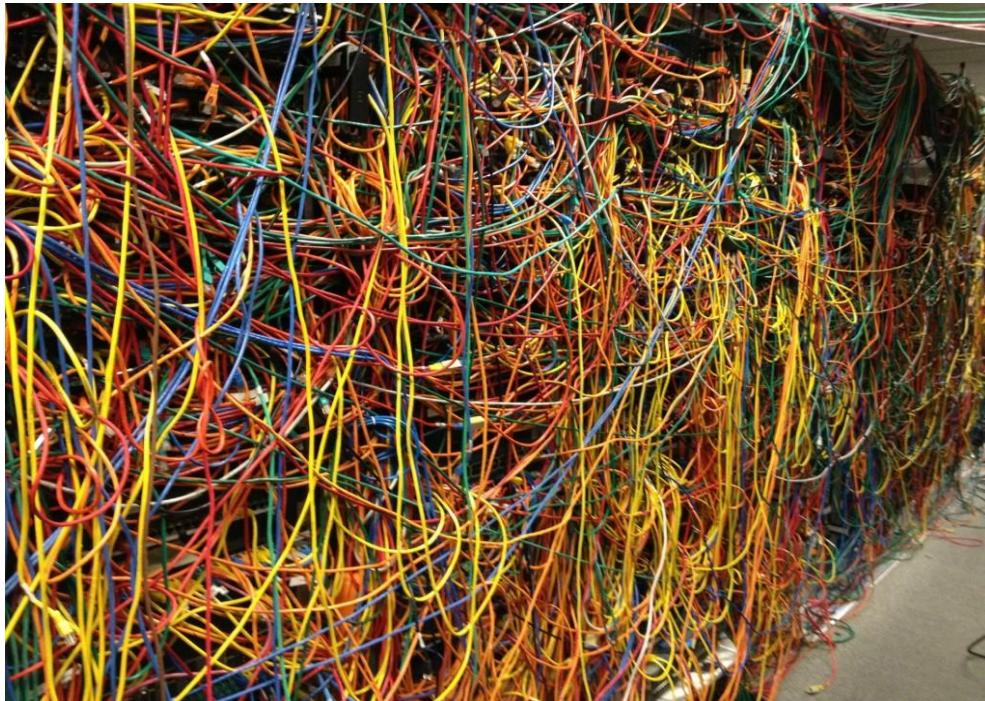


Fail-Stop + Perfect Link + Synchronous



BIGDATA
TEAM

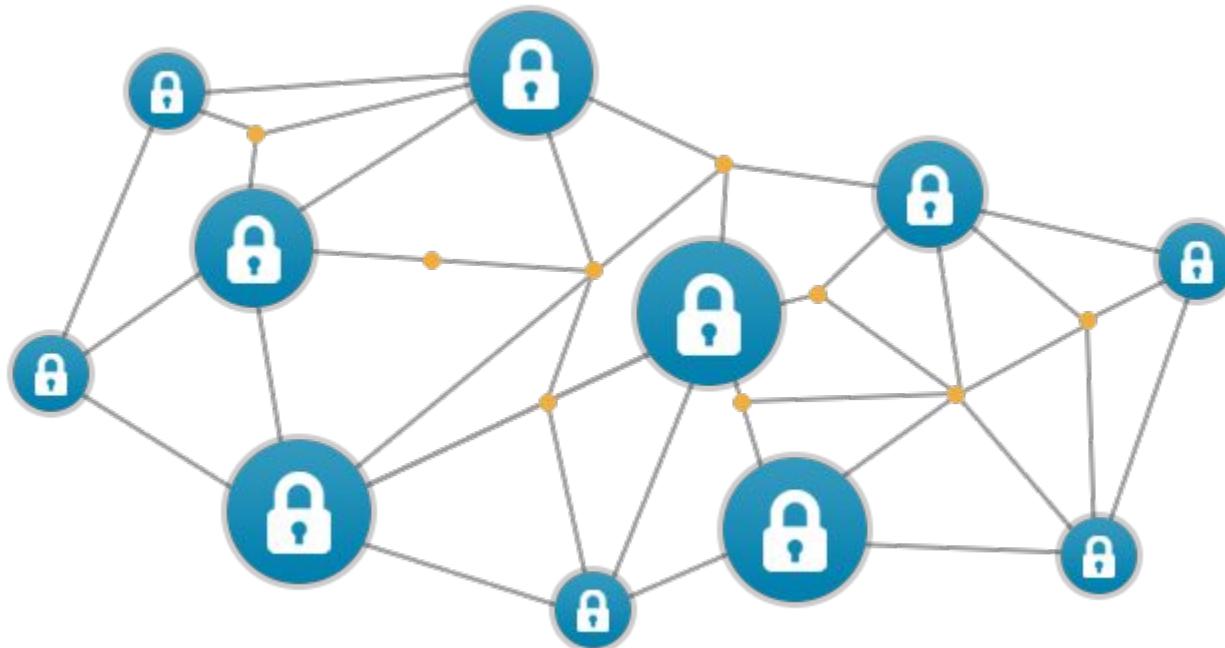
Distributed Computing



Fail-Recovery + Fair-Loss Link + Asynchronous

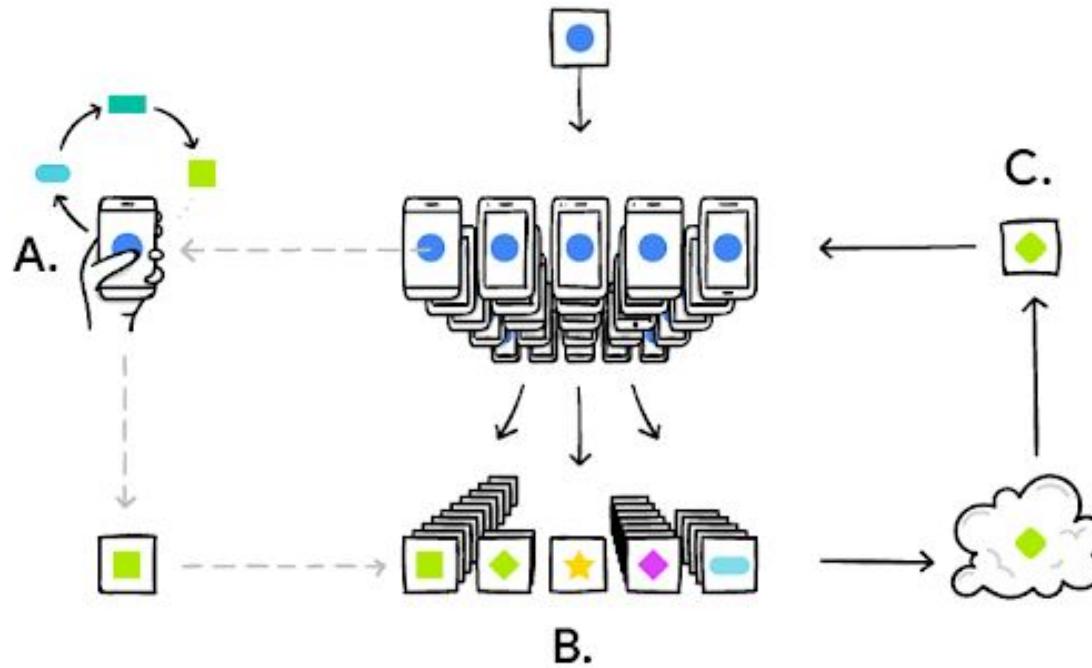


Byzantine-Failure + Byzantine Link + Asynchronous





Federated Machine Learning





**BIGDATA
TEAM**

GFS / HDFS



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

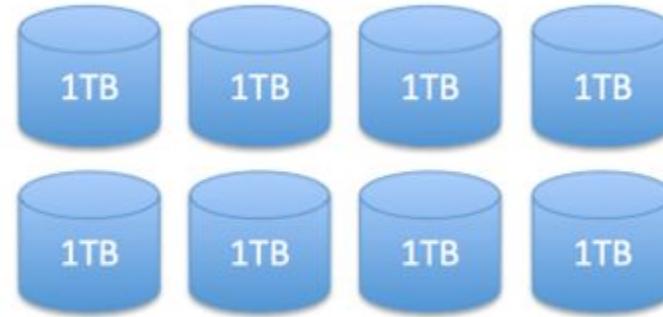
We have designed and implemented the Google File System (GFS) to meet the rapidly growing demands of Google's data processing needs. GFS shares many of the same goals as previous distributed file systems such as performance, scalability, reliability, and availability. However, its design has been driven 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.



Distributed File Systems



Scale-up



Scale-out

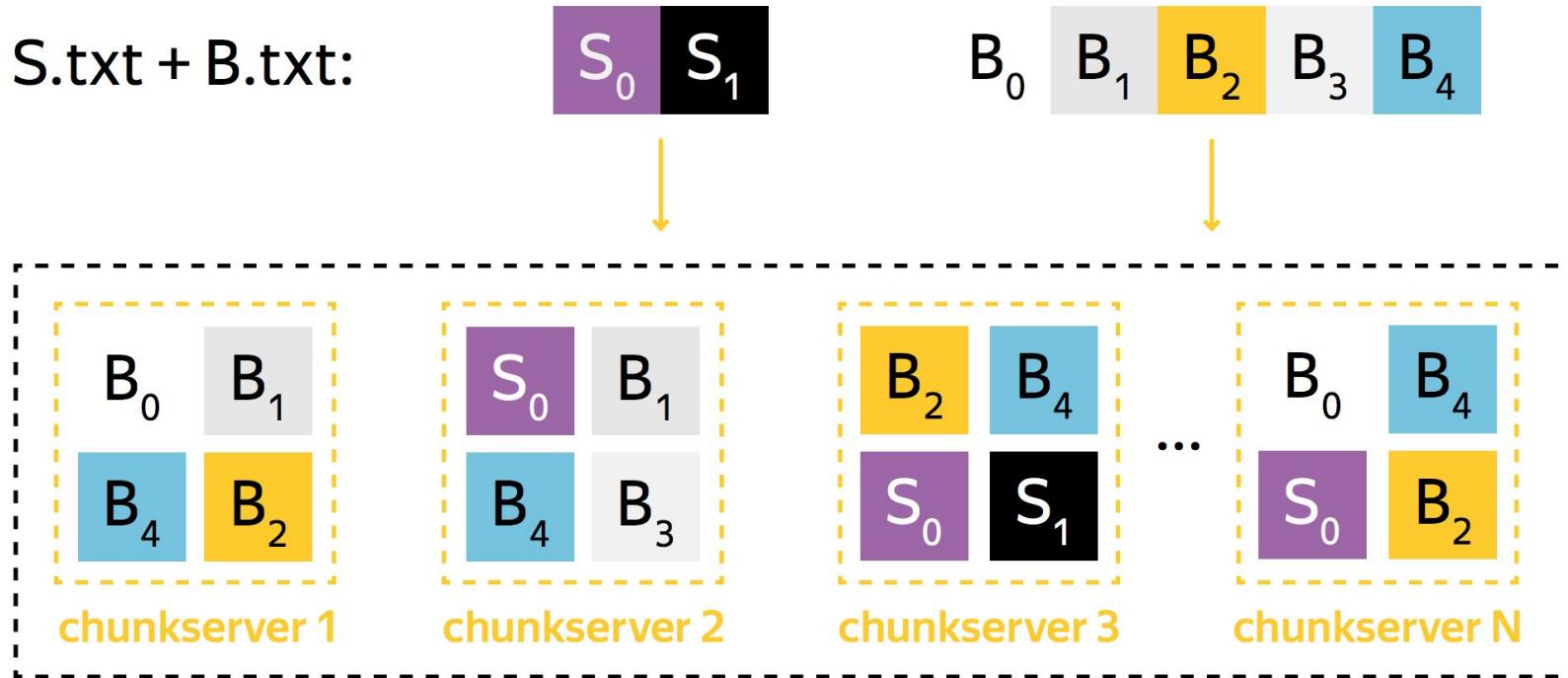


Q: Как выживать в случае поломок серверов?

Q: Как решать проблему равномерной нагрузки серверов в случае наличия файлов разных размеров (2TB vs 10GB)?

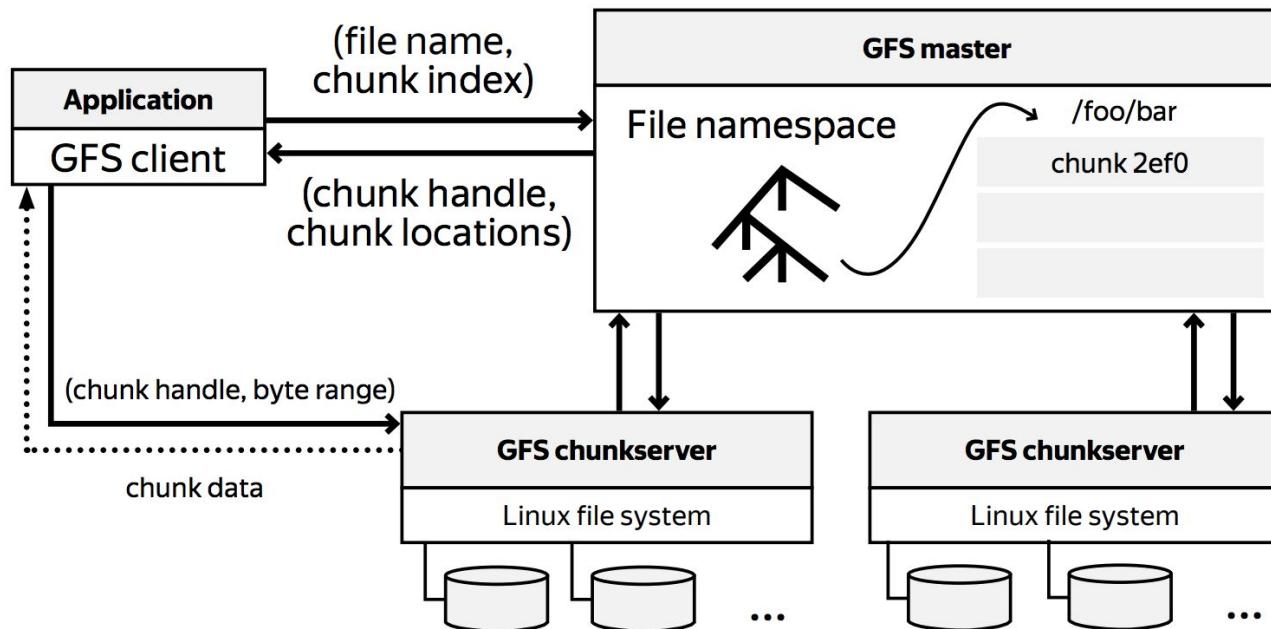


S.txt + B.txt:

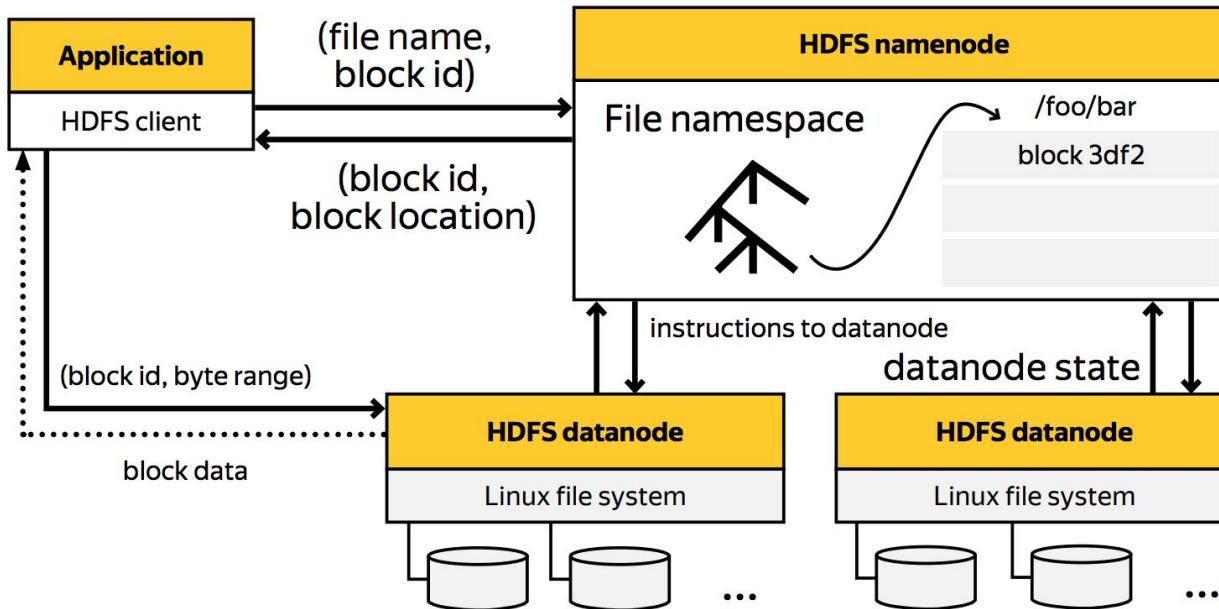




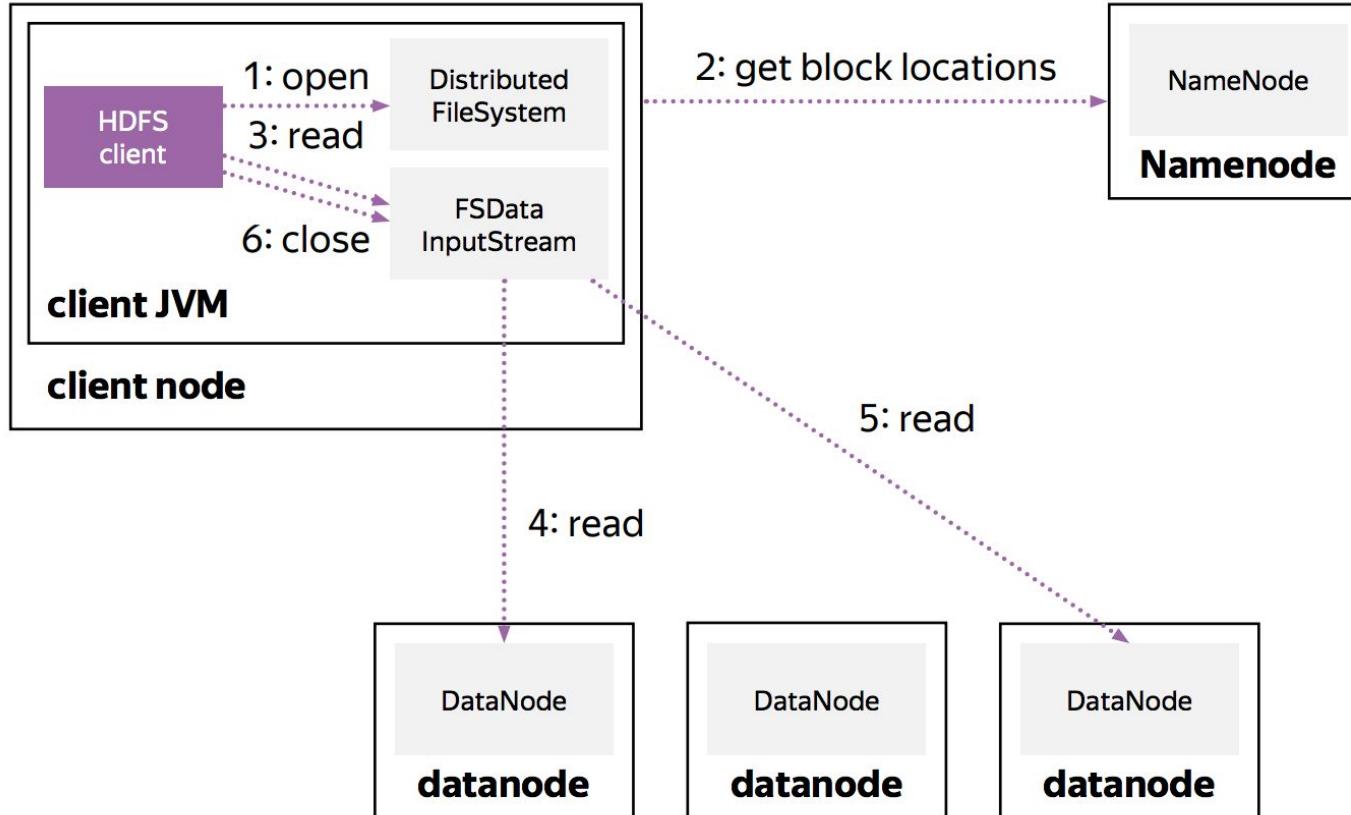
- ▶ Components failures are a norm → replication
- ▶ Even space utilization
- ▶ Write-once-read-many semantic

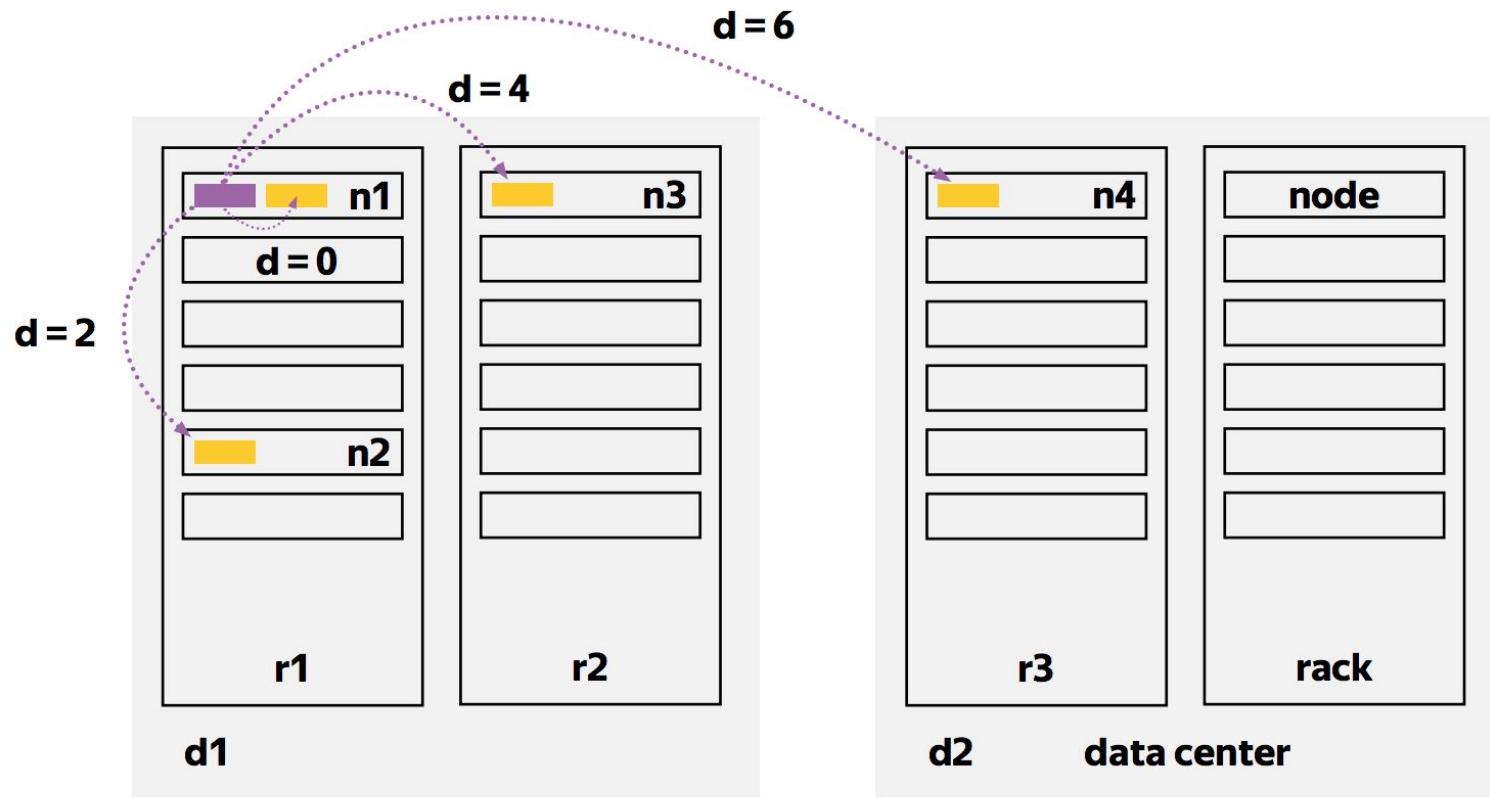


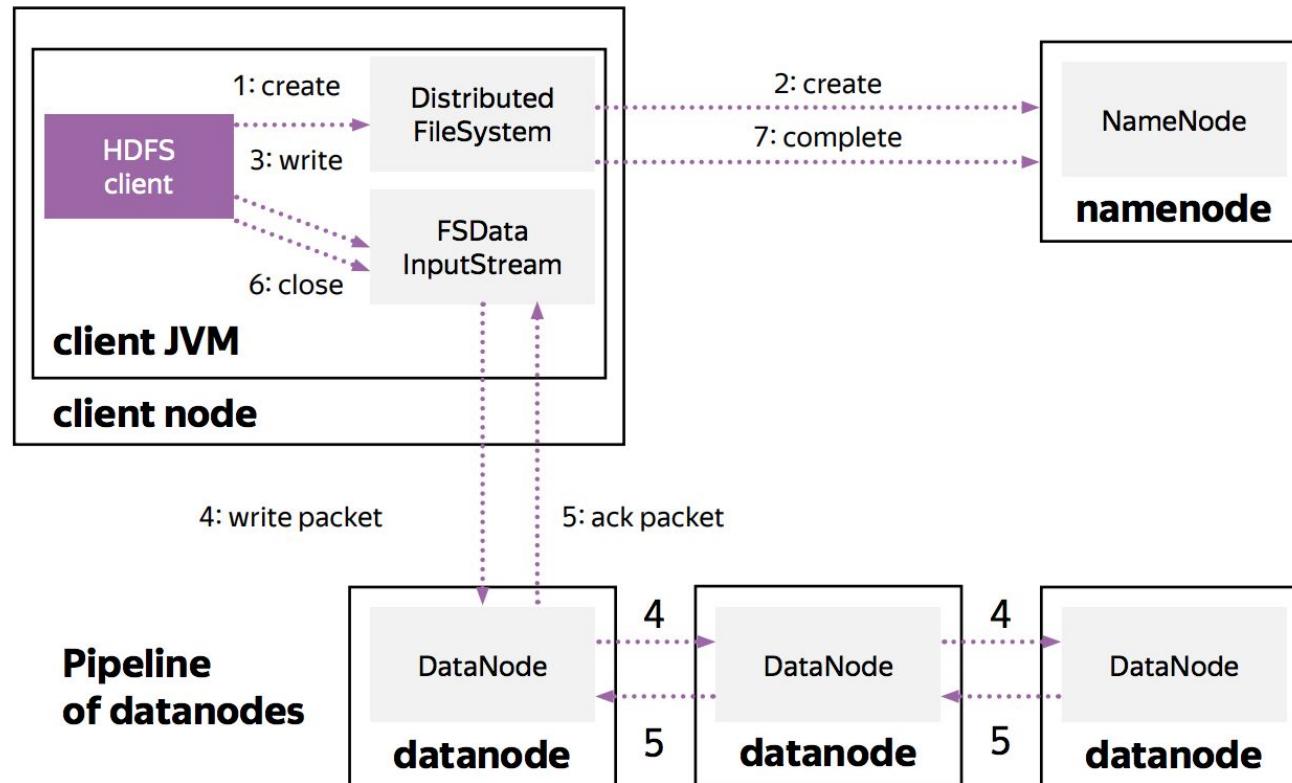
GFS



HDFS

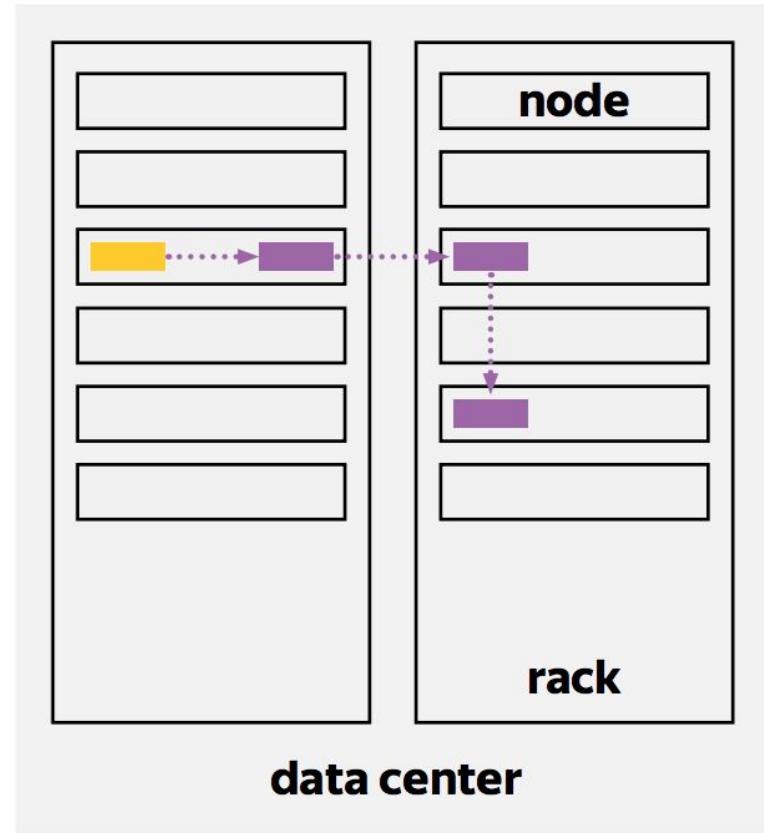


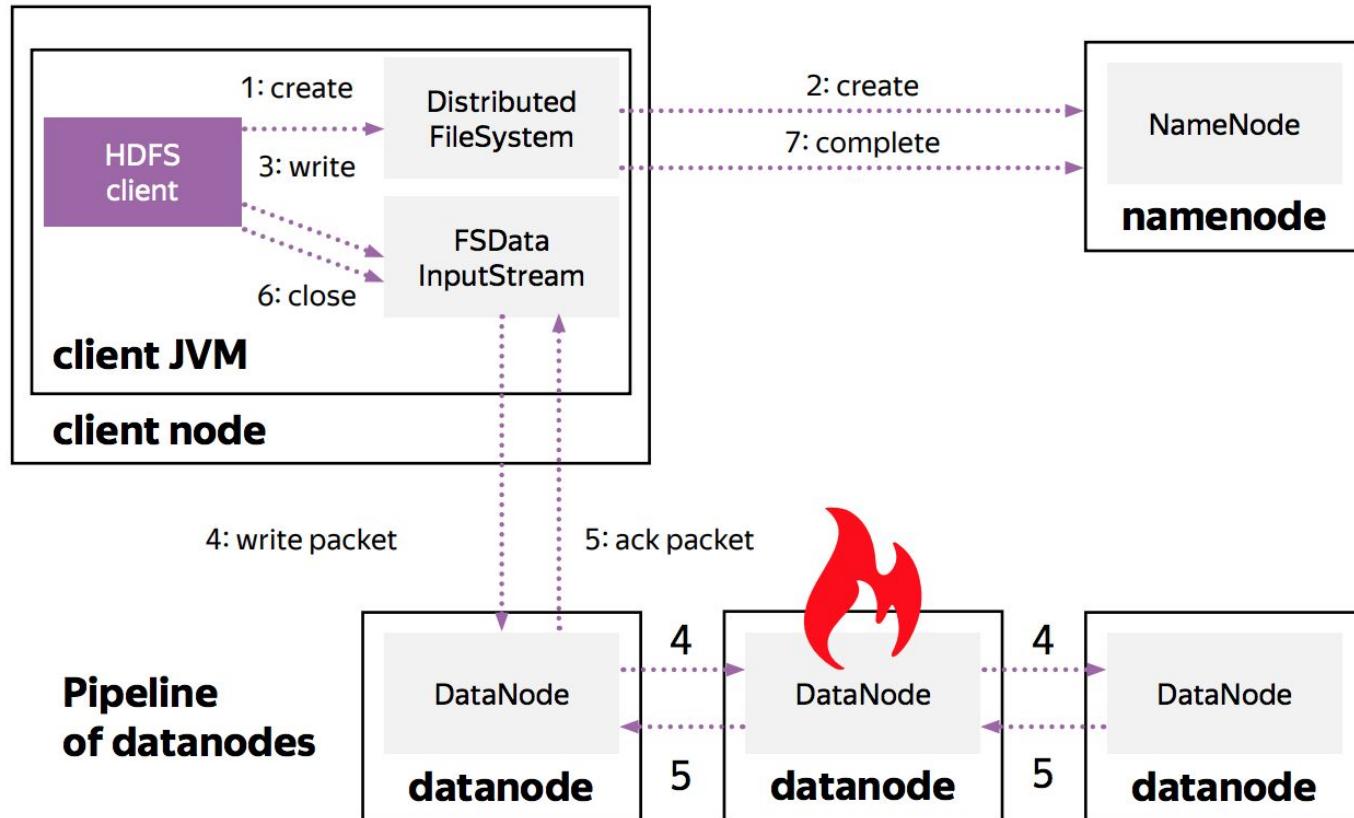




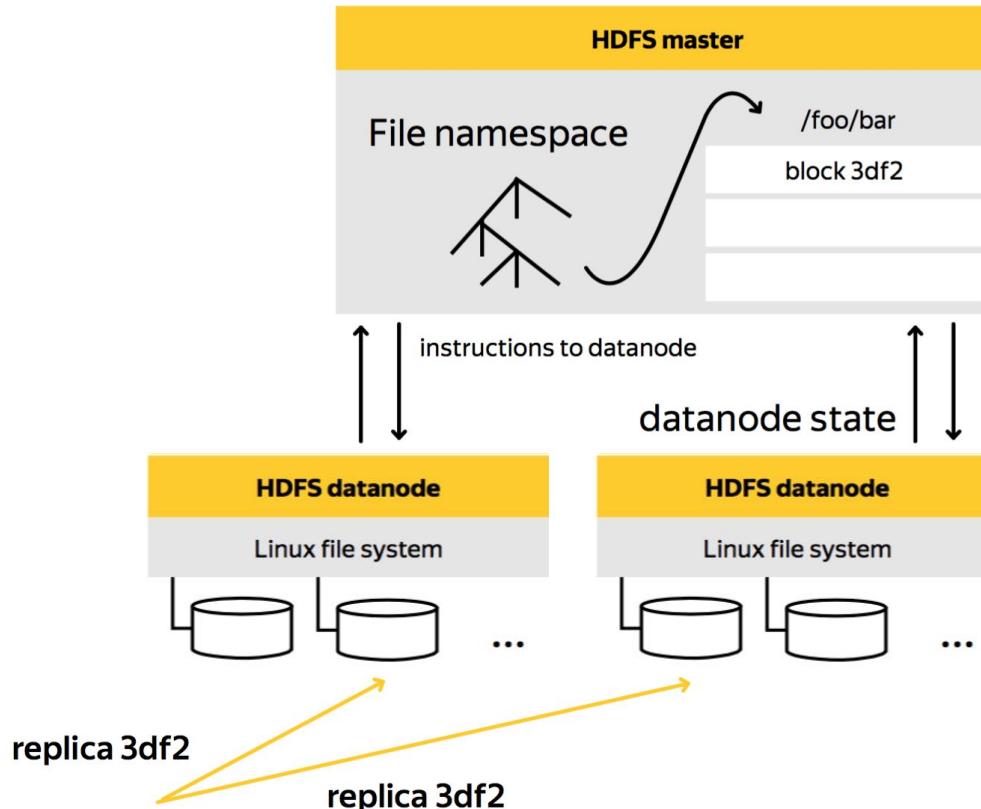


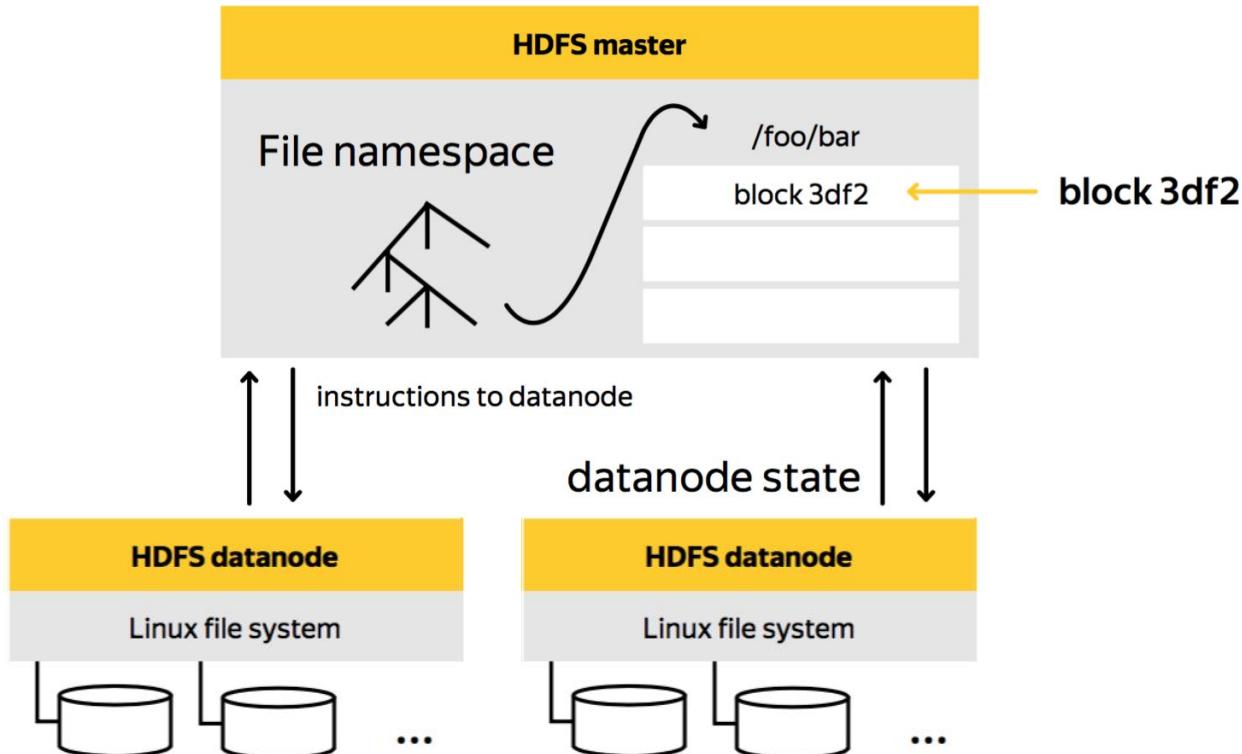
HDFS Replica Placement

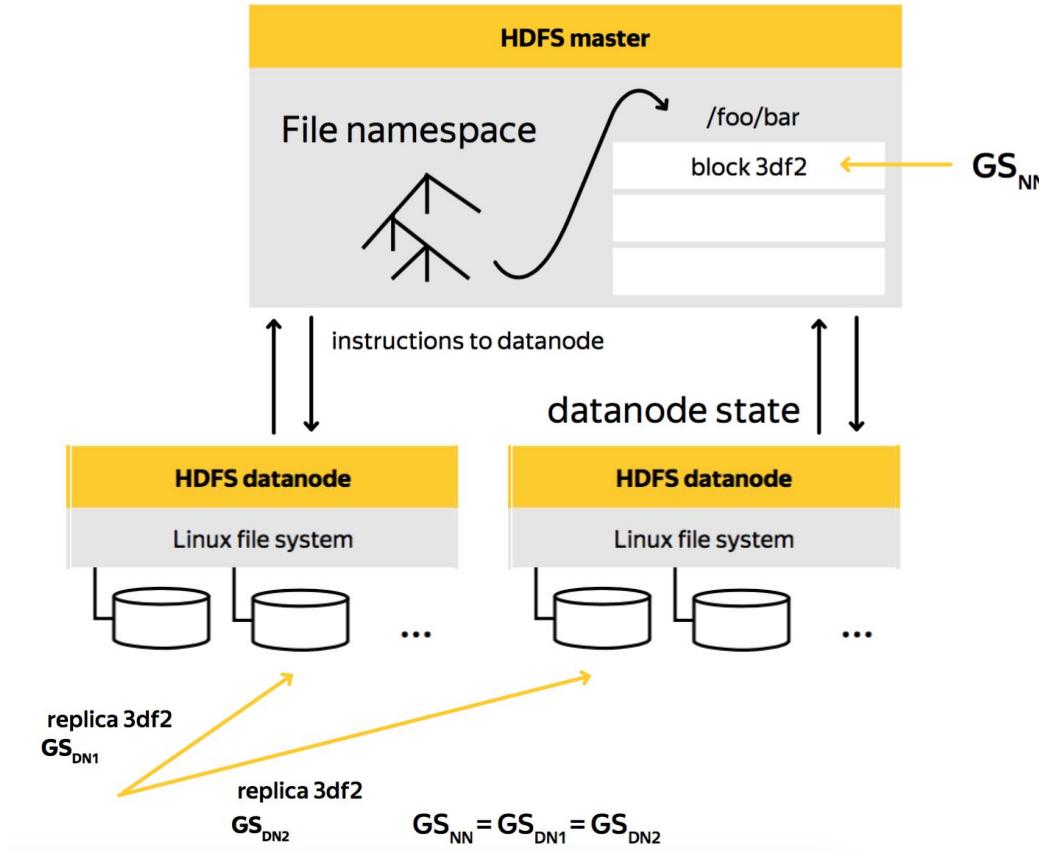




Blocks & Replicas

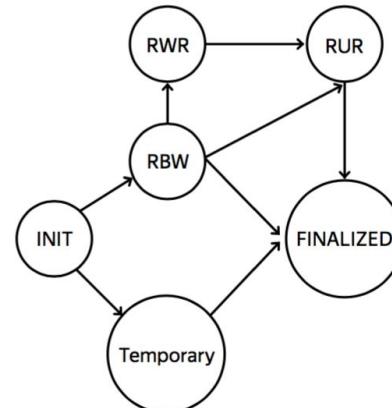
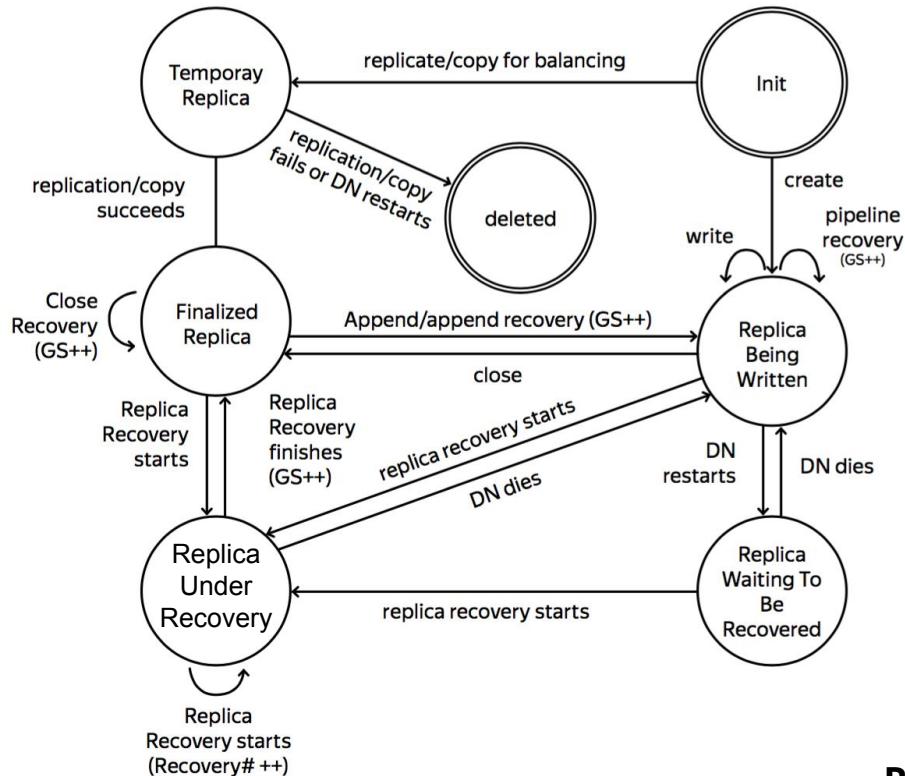








Replica Recovery Process

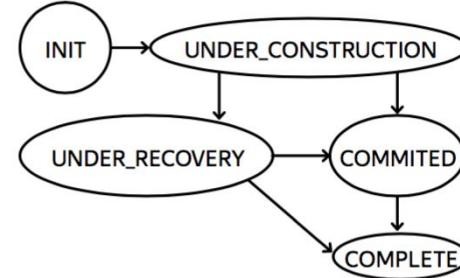
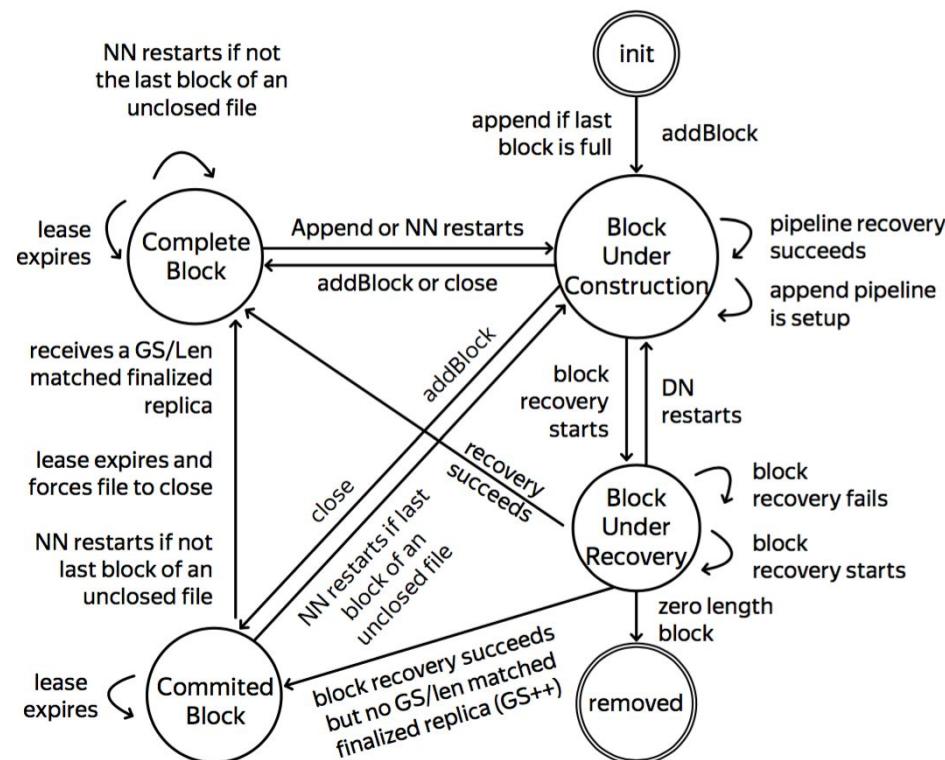


Simplified
Replica State
Transition

RBW / RWR / RUR / Temporary / ...



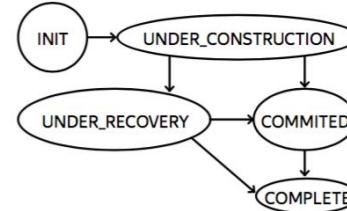
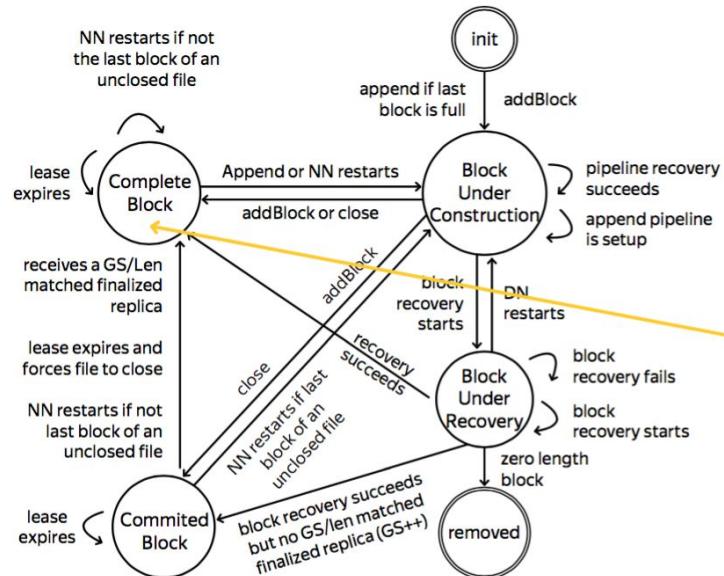
Simplified Block State Transition



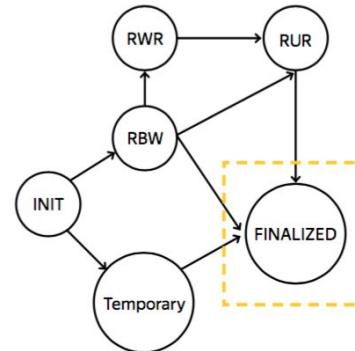
Simplified
Block State
Transition



Block & Replica States Matching



Simplified
Block State
Transition
complete

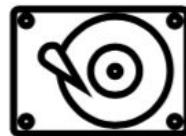
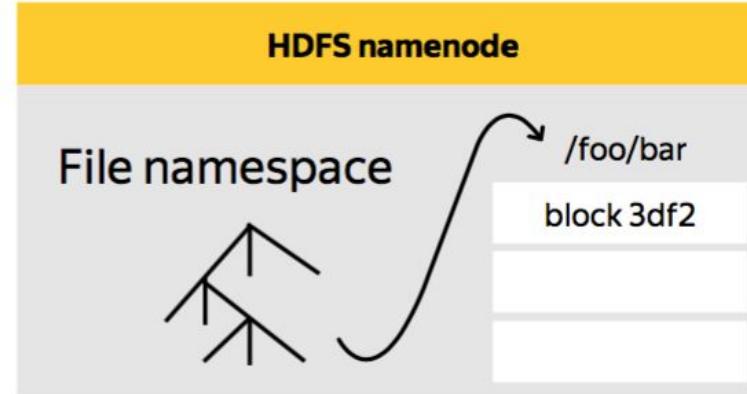




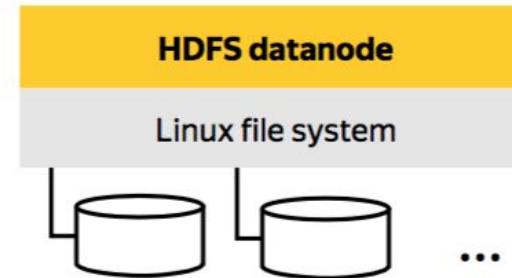
Memo: Block vs Replica



**block
state**



**replica
state**

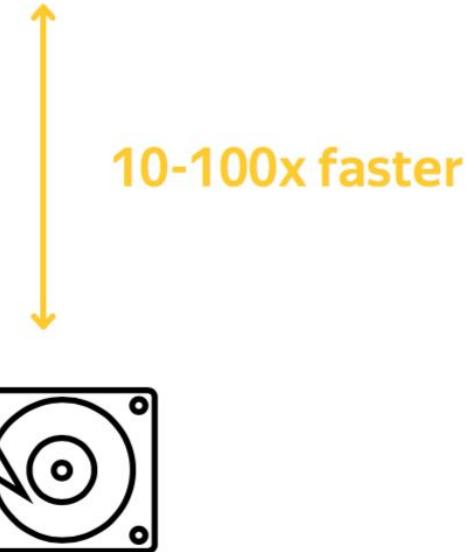
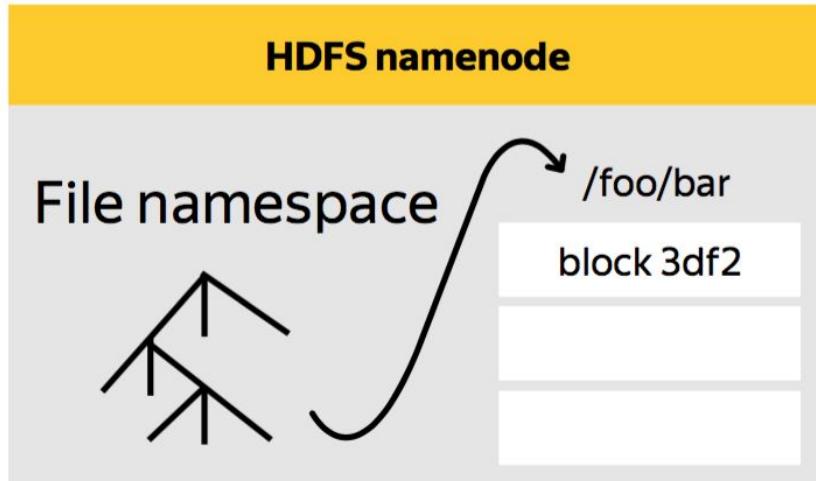


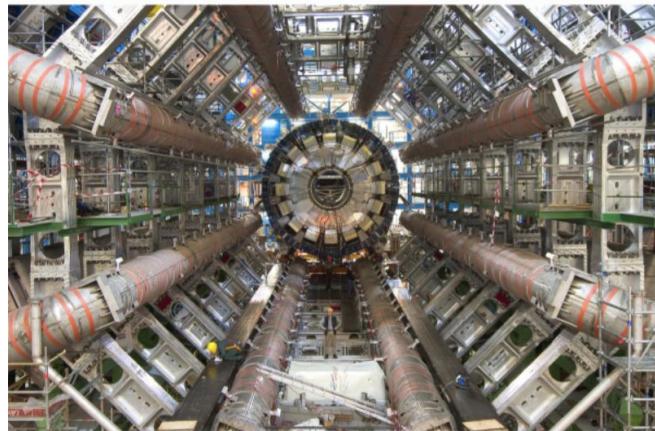


HDFS Namenode

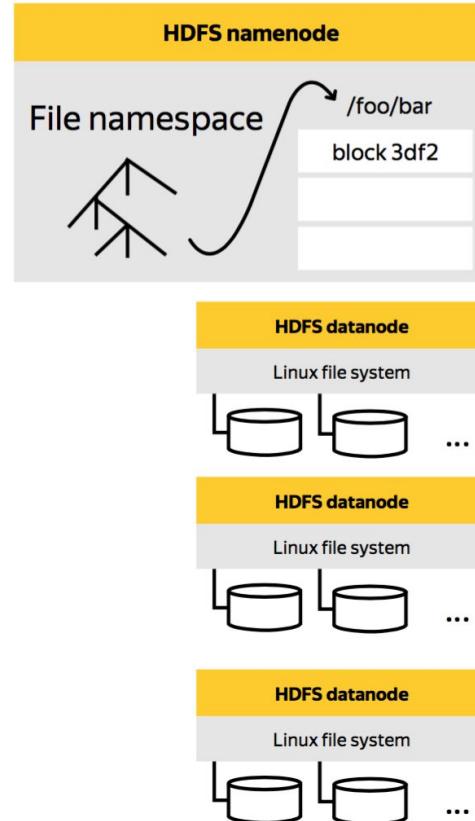


Namenode Design Choices



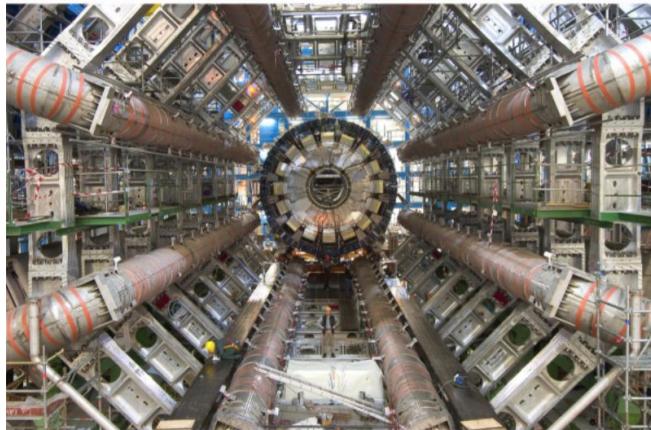


1 year ~ 10 PB

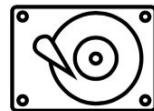




Hadoop Sizing: Small Files Problem



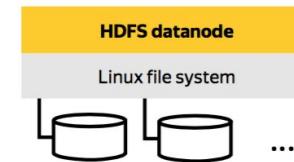
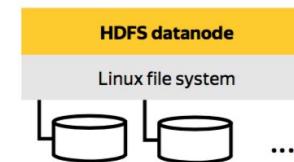
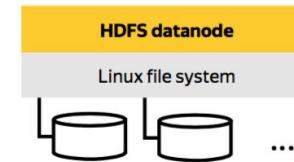
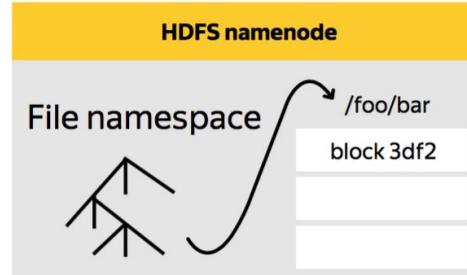
1 year ~ 10 PB



10 PB / 2 TB * 3 ~ 15 k



150 B - average block size on Namenode





Speeding Up the Cluster



2 TB

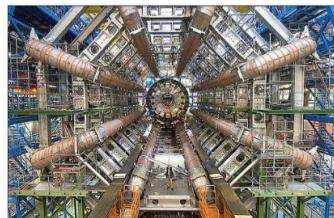
vs



1 TB



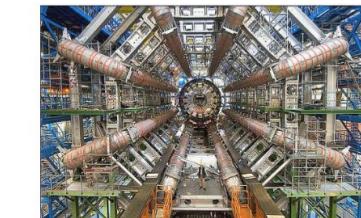
1 TB



1 year ~ 10 PB



35 days



1 year ~ 5 PB



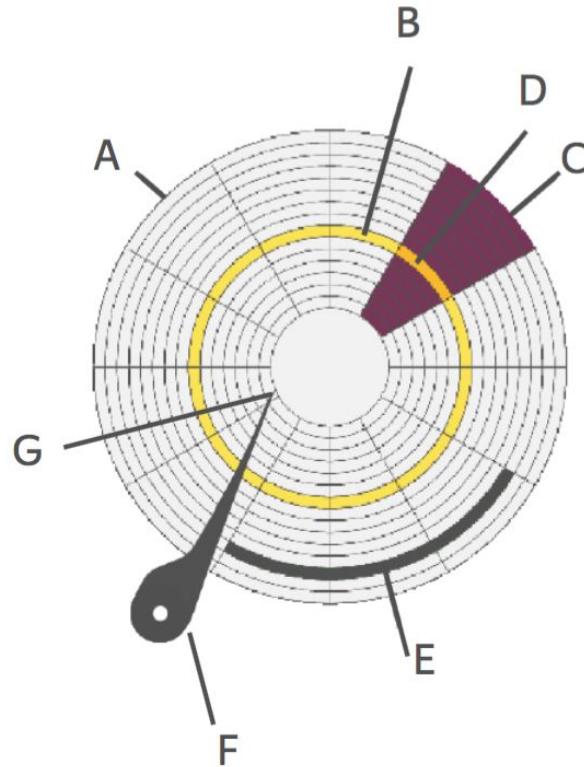
+ 5 PB



17.5 days



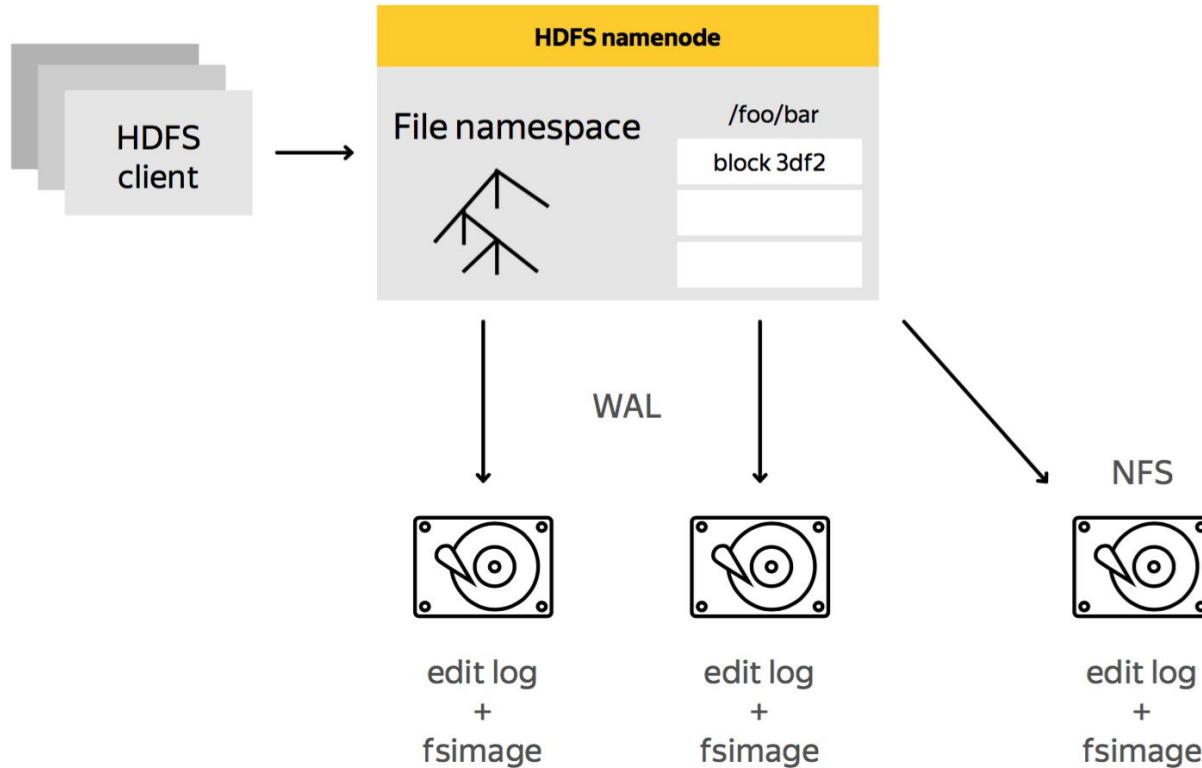
Default Block Size



- reading speed - 3.5 GB/sec
- 128 MB - 30-40 ms
- seek time: 0.2-0.8 ms



Namenode - SPoF





Checkpoint Namenode

Primary Namenode

edits_inprogress_1

1. Roll edits

edits_inprogress_20

fsimage_0

edits_1-19

fsimage_19.ckpt

fsimage_19

Secondary Namenode

= Checkpoint Namenode

≠ Backup Node

2. Retrieve fsimage and edits from primary

edits_1-19

fsimage_0

3. Merge

fsimage_19.ckpt

4. Transfer checkpoint to primary

fsimage.ckpt

fsimage



Hadoop and NoSQL Downfall Parody

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



- ▶ Вы можете привести примеры приложения Data Science, а также где можно наработать ML-навыки
- ▶ 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

Thank you! Questions?

Feedback: http://rebrand.ly/mf2019q2_feedback_01_hdfs

Dral Alexey (aadral@bigdatateam.org)

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

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

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



Appendix

[http://rebrand.ly/mf2019q2 user guides](http://rebrand.ly/mf2019q2_user_guides)



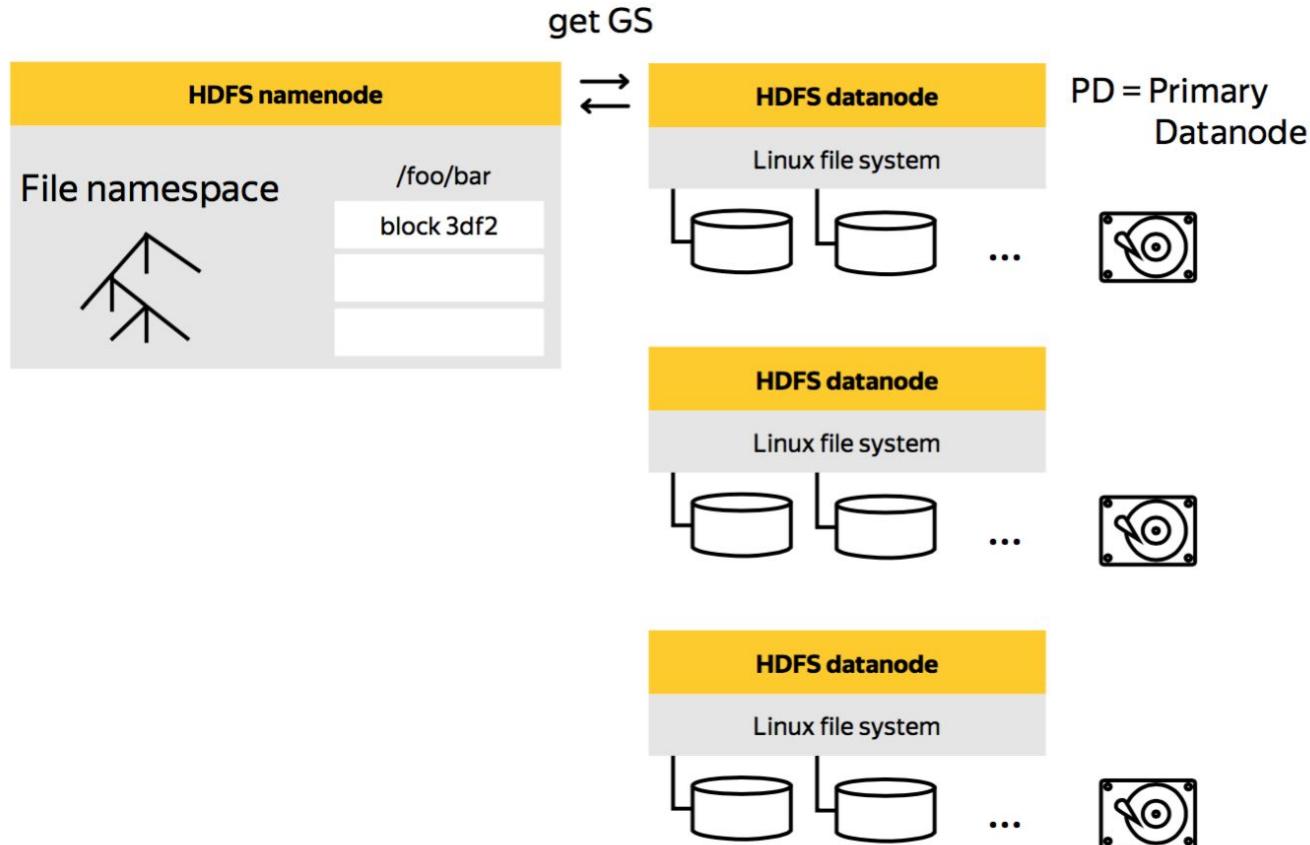
Lease Recovery

- ▷ Block Recovery

- Replica Recovery

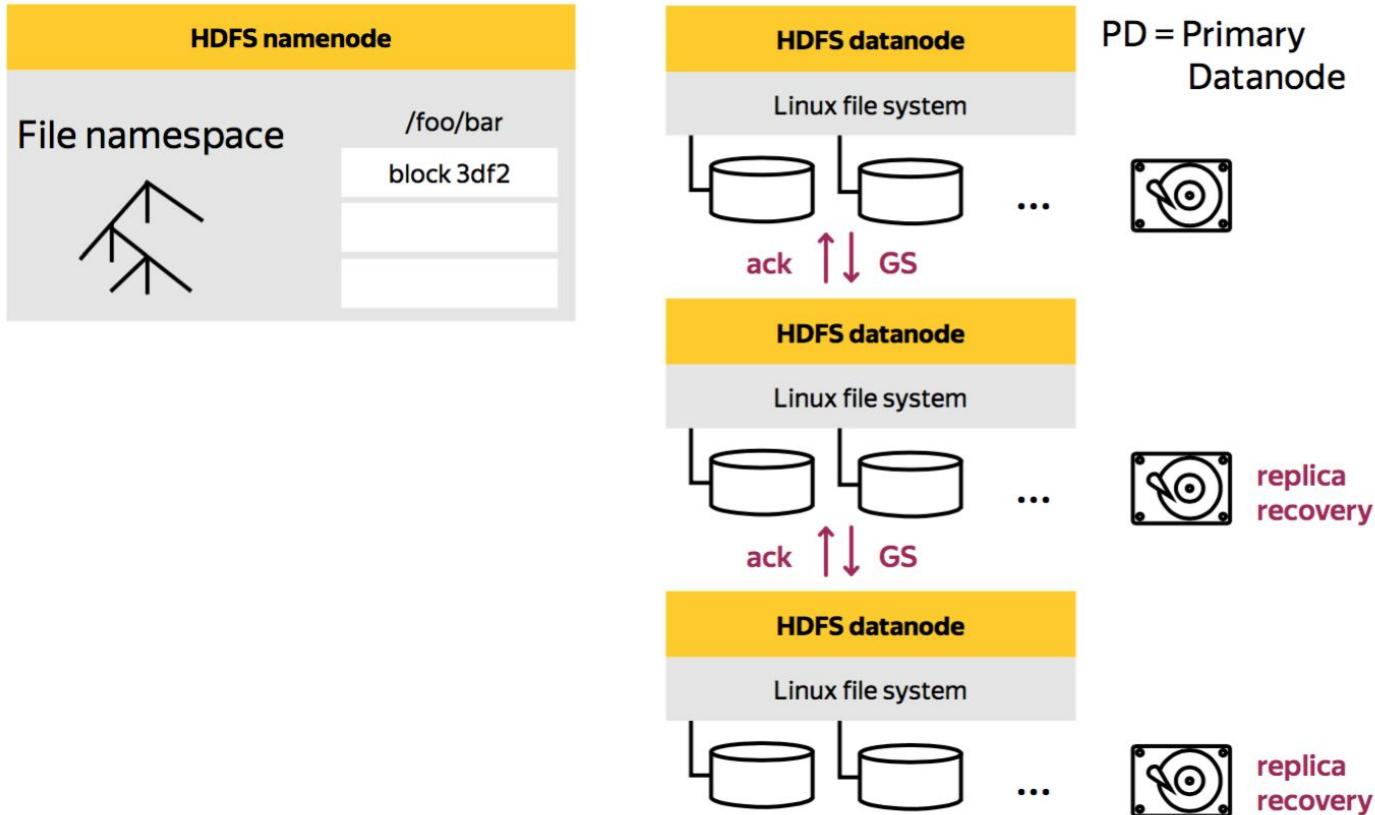


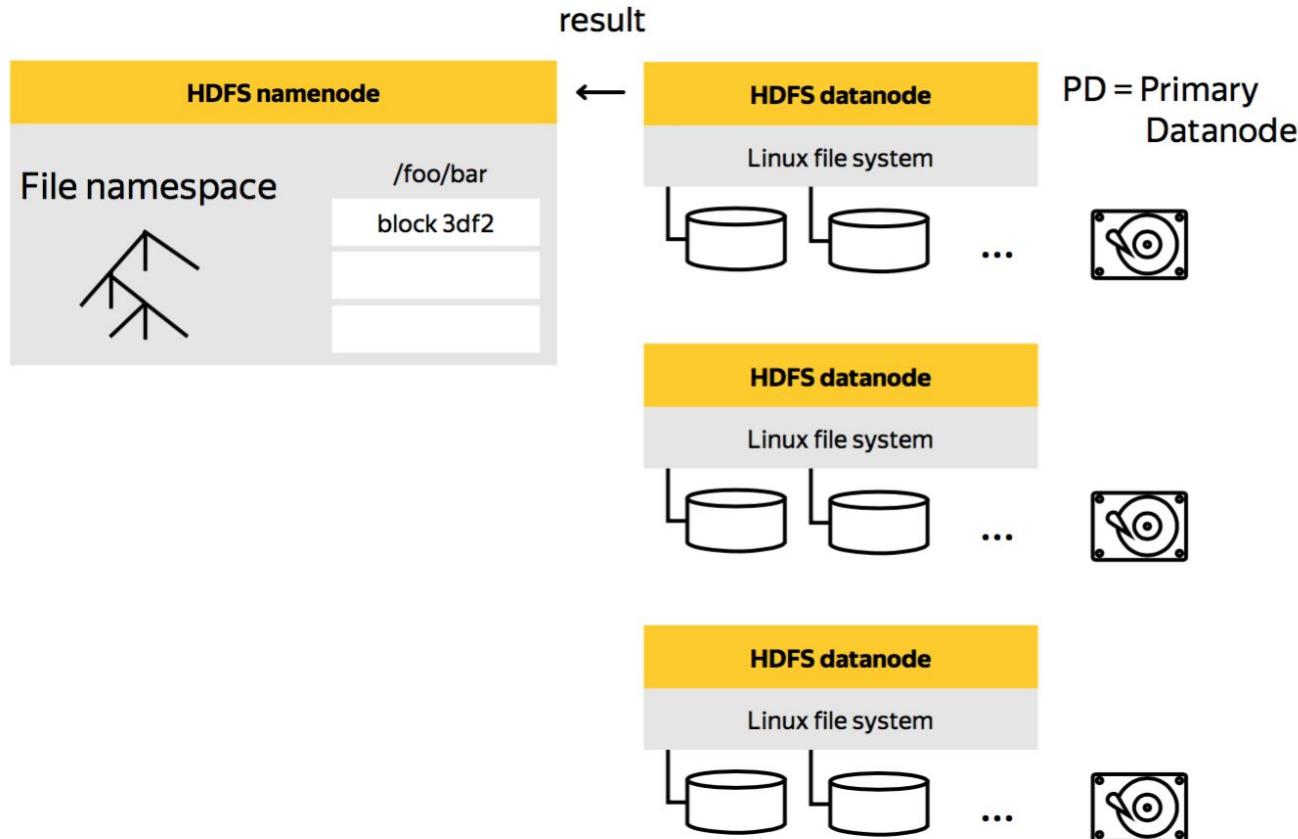
Pipeline Recovery

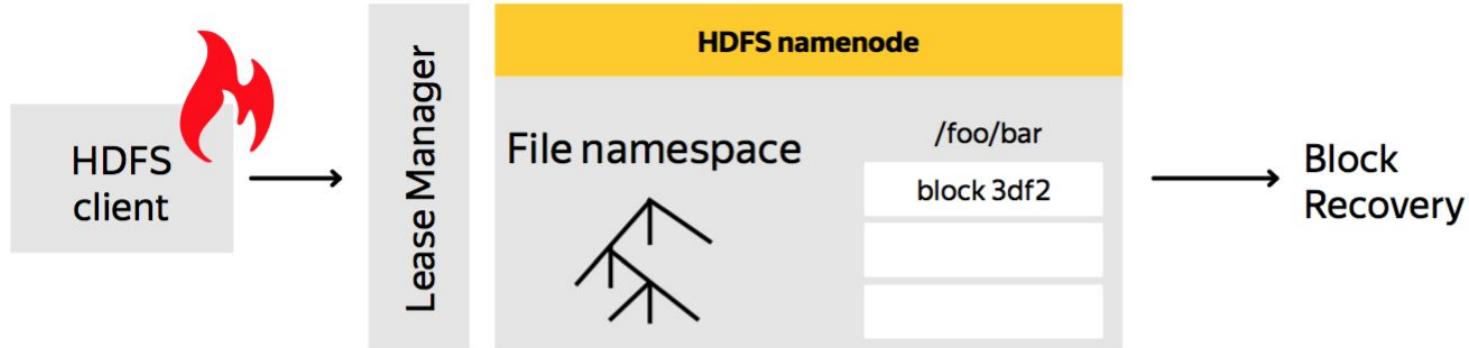




Block Recovery

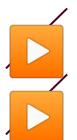
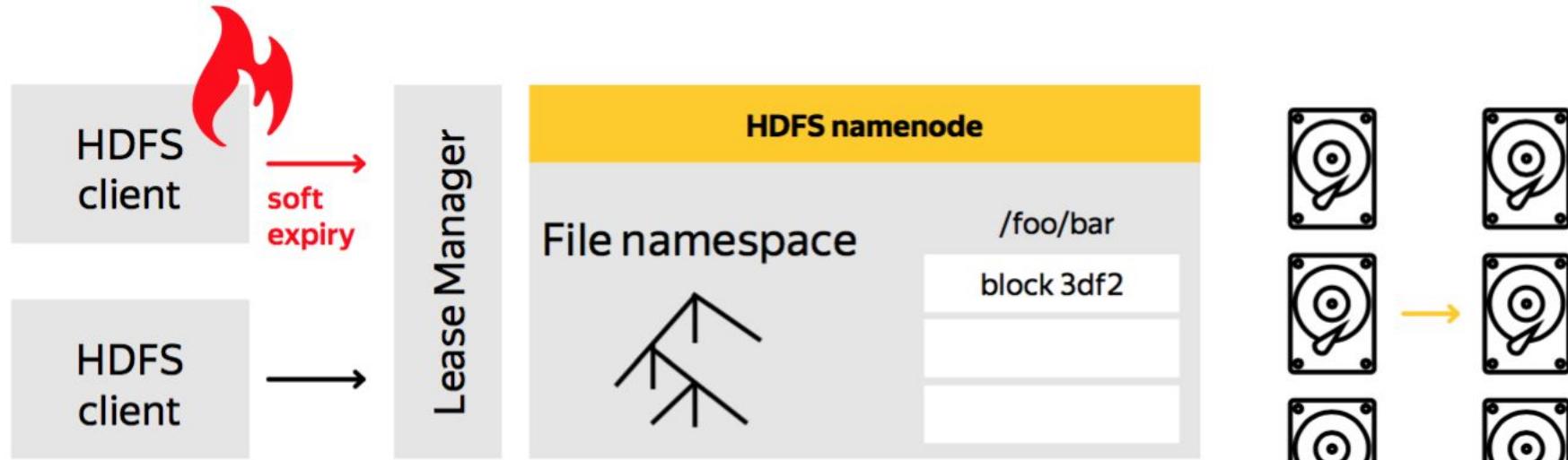






user₁, /path/1, lease (soft + hard)
user₁, /path/2, lease (soft + hard)
user₂, /path/3, lease (soft + hard)
user₃, /path/4, lease (soft + hard)

...

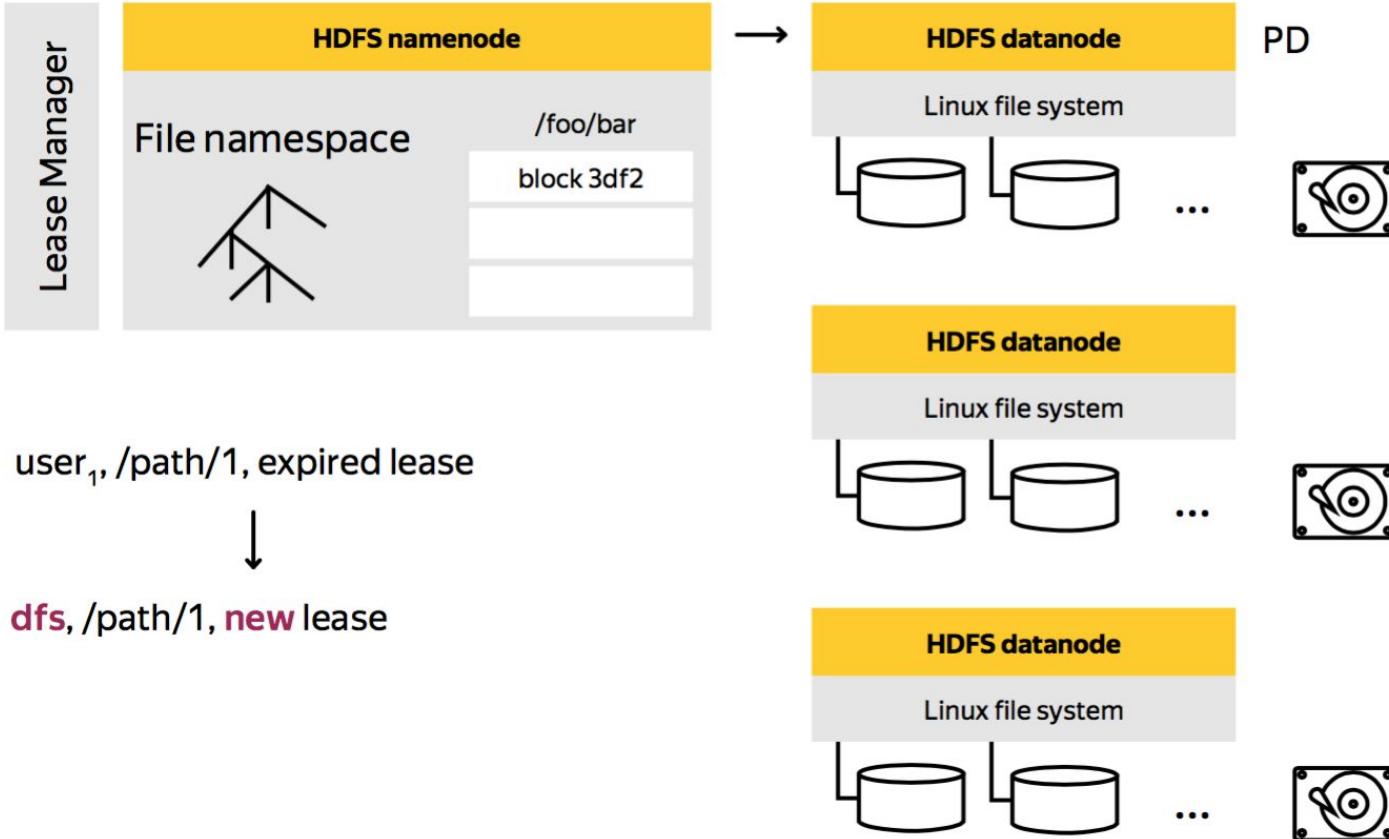


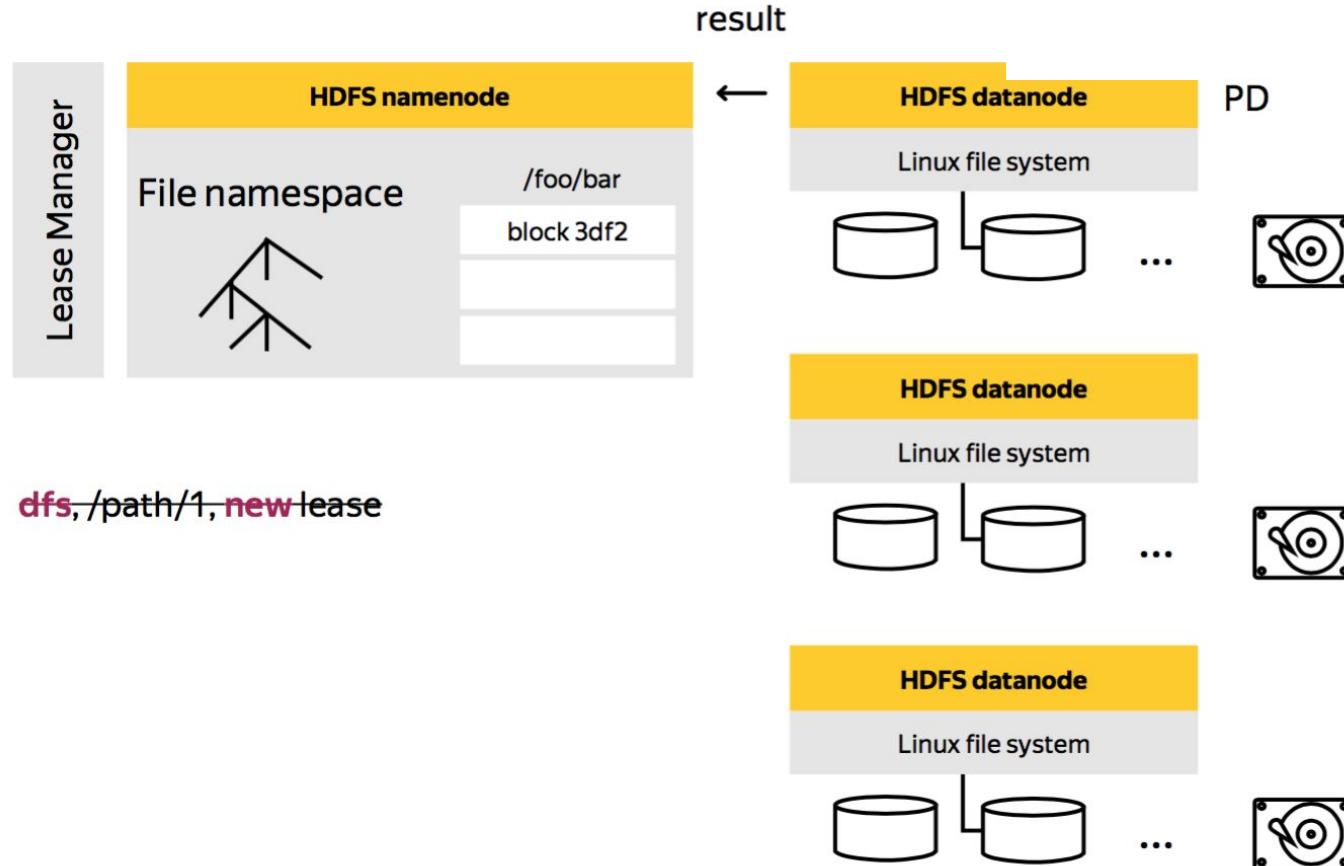
- Concurrency control
- Consistency guarantee

**consistency
guarantee**



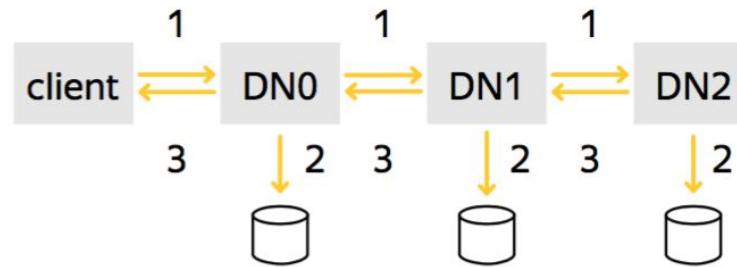
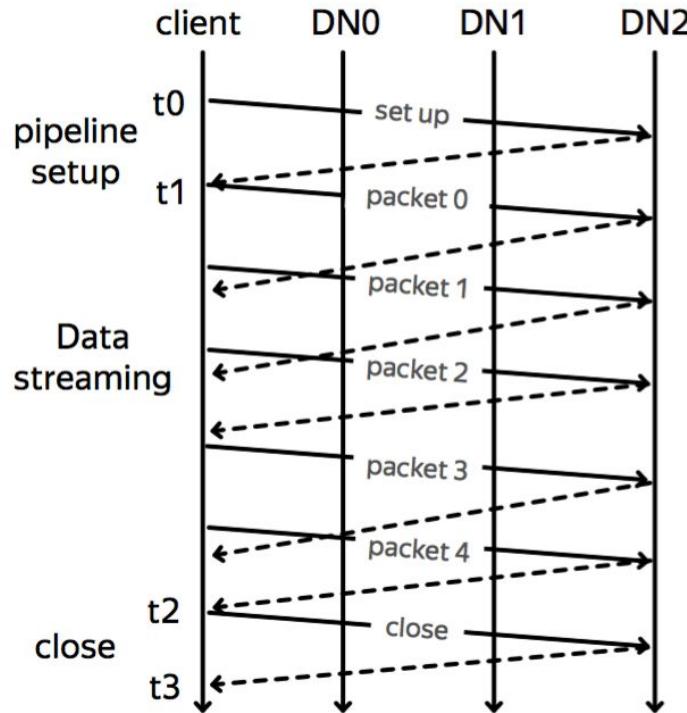
Lease Recovery





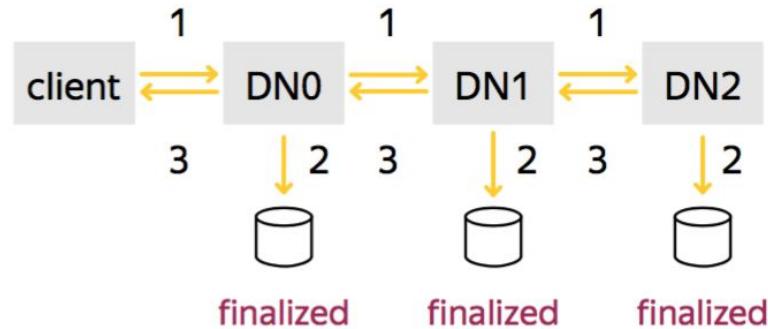
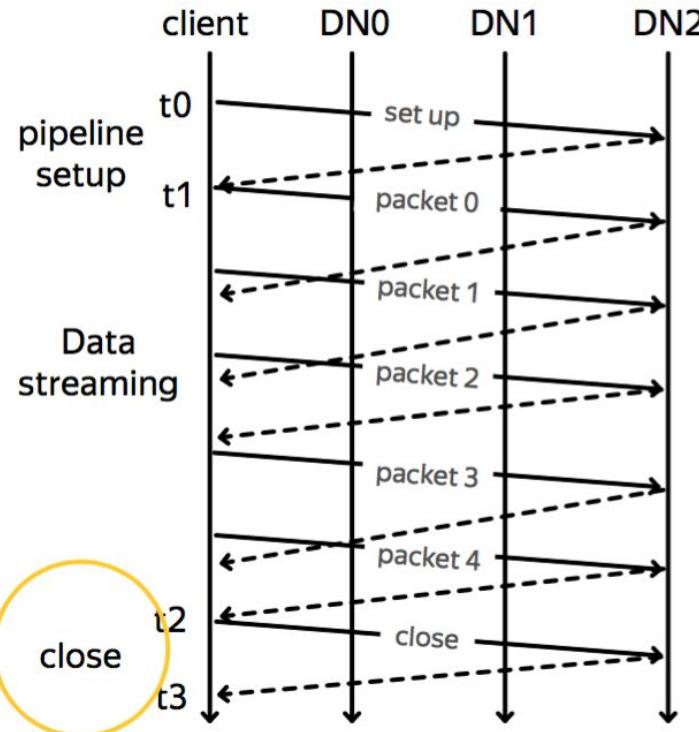


Pipeline Recovery



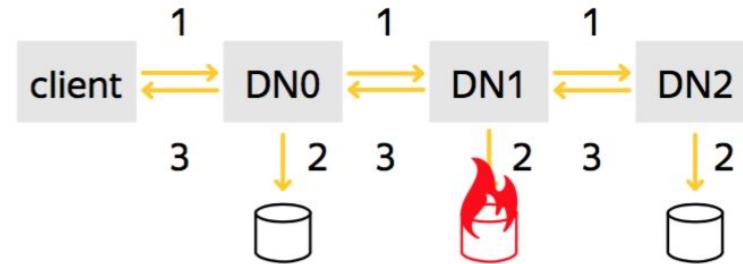
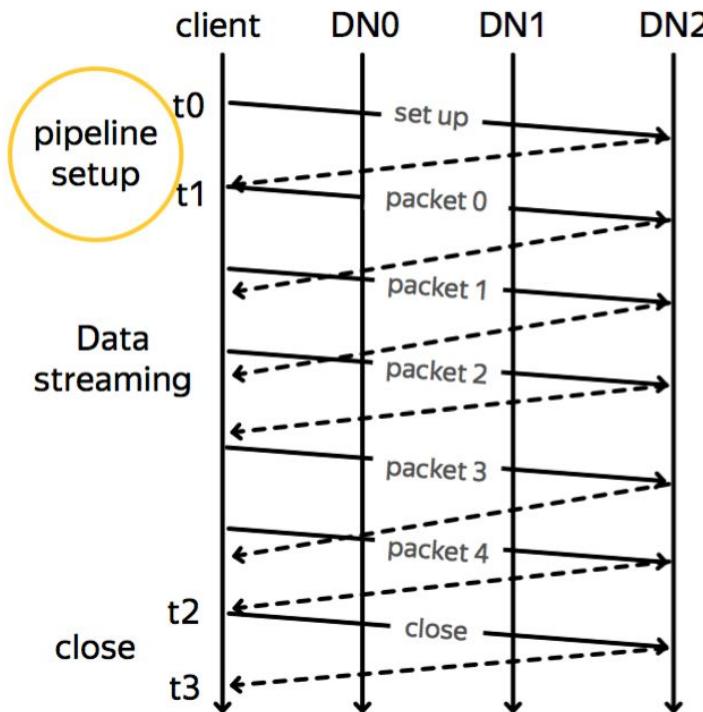


Pipeline Stages





Q&A: Что будем делать?

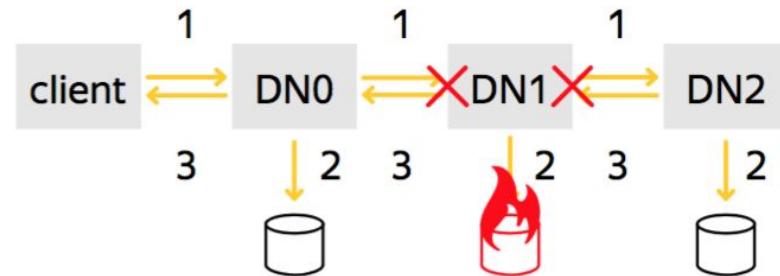
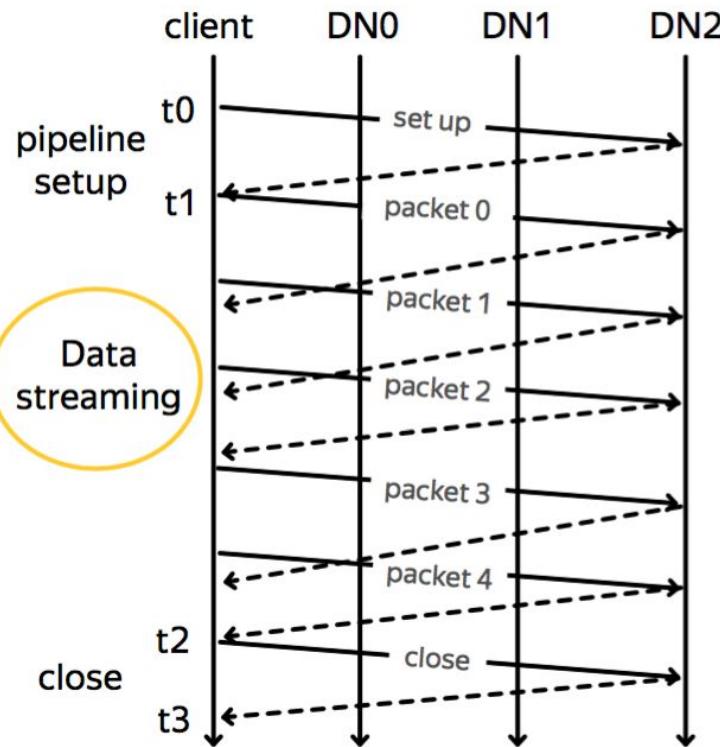


cases:

1. new file
2. append mode

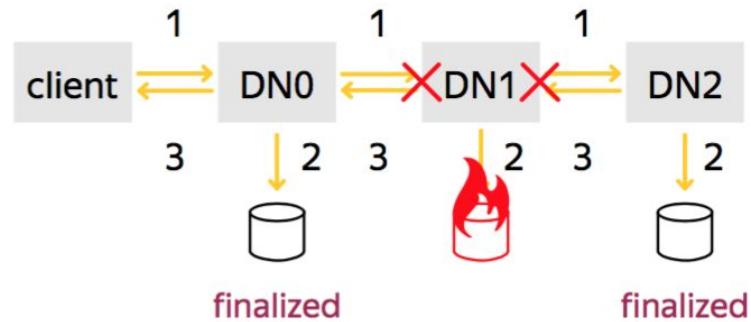
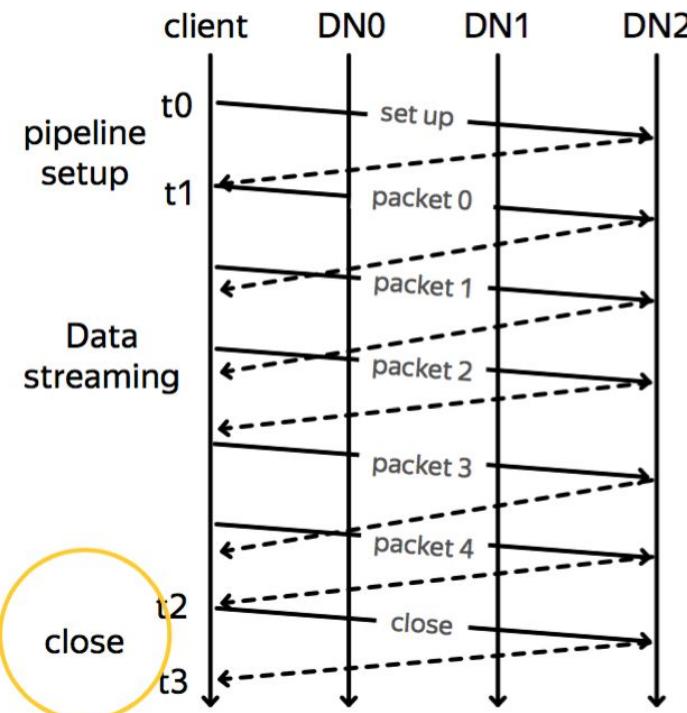


Q&A: Что будем делать?





Q&A: Что будем делать?





Fallacies of Distributed Computing

1. The network is reliable
2. Latency is zero
3. Bandwidth is infinite
4. The network is secure
5. Topology doesn't change
6. There is one administrator
7. Transport cost is zero
8. The network is homogeneous