With IT since 2007
With Java since 2009
With Hadoop since 2012
With Spark since 2014
With EPAM since 2015

**About**

# Contacts

E-mail : Alexey_Zinovyev@epam.com

Twitter : @zaleslaw @BigDataRussia

vk.com/big_data_russia **Big Data Russia**

**+ Telegram** @bigdatarussia

vk.com/java_jvm **Java & JVM langs**

**+ Telegram** @javajvmlangs

# Main parts

- What is BIG DATA?

- Intro in Hadoop

- HDFS & YARN

- MapReduce Java API

- JOINs techniques*

- JVM Settings*

- File formats*

# WHAT IS BIG DATA?

# Joke about Excel

**DevOps Borat**
@DEVOPS_BORAT

Following

Big Data is any thing which is crash Excel.

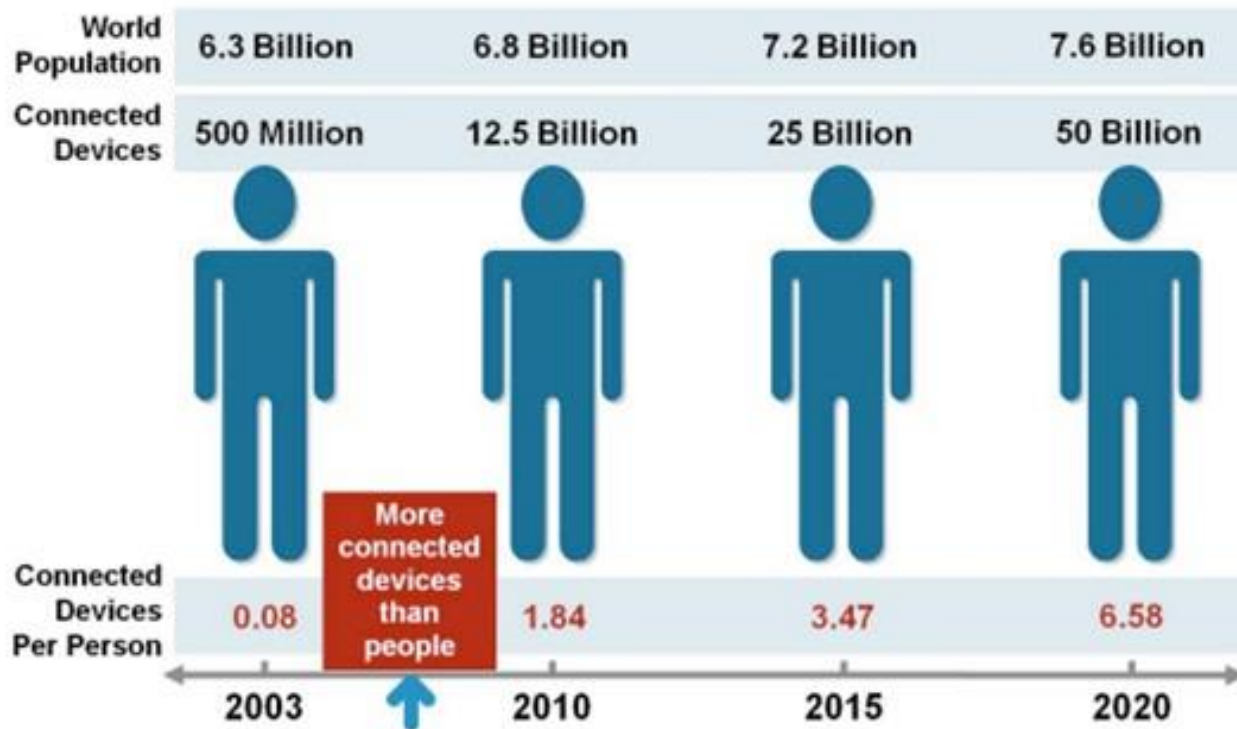← Reply    ⇄ Retweet    ★ Favorite    ••• More
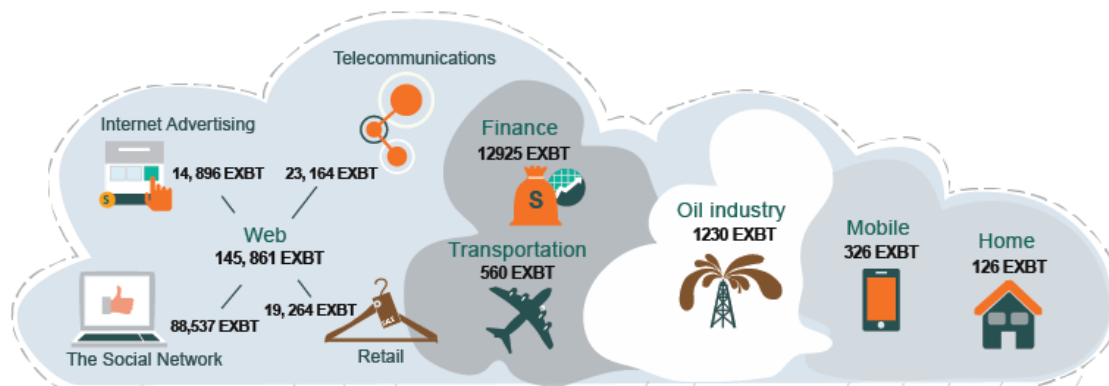
**1,879** RETWEETS    **384** FAVORITES

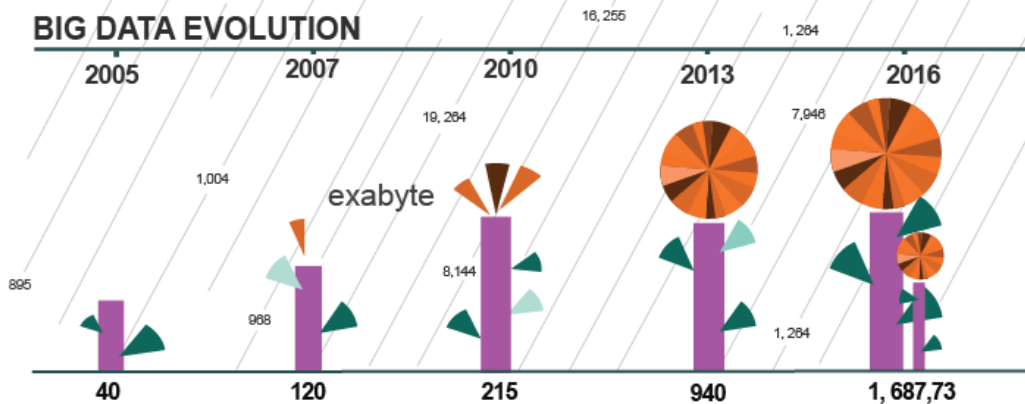11:25 AM - 8 Jan 13

# Every 60 seconds...

# From Mobile Devices



| | | | |
|---|---|---|---|
| **World Population** | 6.3 Billion | 6.8 Billion | 7.2 Billion | 7.6 Billion |
| **Connected Devices** | 500 Million | 12.5 Billion | 25 Billion | 50 Billion |

More connected devices than people

| | | | |
|---|---|---|---|
| **Connected Devices Per Person** | 0.08 | 1.84 | 3.47 | 6.58 |

2003    2010    2015    2020

# From Industry



Telecommunications

Internet Advertising
14, 896 EXBT    23, 164 EXBT

Finance
12925 EXBT

Oil industry
1230 EXBT

Web
145, 861 EXBT

Transportation
560 EXBT

Mobile
326 EXBT

Home
126 EXBT

88,537 EXBT    19, 264 EXBT

The Social Network    Retail

**BIG DATA EVOLUTION**

16, 255

2005    2007    2010    2013    2016

1, 264

19, 264

1,004

exabyte

7,946

895

8,144

968

1, 264

40    120    215    940    1, 687,73

# We started to keep and handle ~~stupid~~ new things!



Traditional systems under pressure

**1 Challenges**
- Constrains data to app
- Can't manage new data
- Costly to Scale

**2 New Data**

Business Value

ERP   CRM   SCM

2012
2.8 Zettabytes

INDUSTRY LEADERS

New

LAGGARDS

Traditional

2020
40 Zettabytes
- Clickstream
- Geolocation
- Web Data
- Internet of Things
- Docs, emails
- Server logs

# Big Data Landscape

Crazy Zoo 2012

Copyright © 2012 Dave Feinleib       dave@vcdave.com       blogs.forbes.com/davefeinleib

BIG DATA LANDSCAPE 2017

**10^6 rows in MySQL**

# GB->TB->PB->?



MEGABYTES
1000 MEGABYTES = 1GB

GIGABYTES
1000 GIGABYTES = 1TB

TERABYTES
1000 TERABYTES = 1PB

PETABYTES
1000 PETABYTES = 1EB

EXABYTES
1000 EXABYTES = 1 ZETTABYTE

1ZB

# Is BigData about PBs?

# Is BigData about PBs?

# It's hard to ...

- .. store

- .. handle

- .. search in

- .. visualize

- .. send in network

# Just do it ... in parallel

# Parallel Computing vs Distributed Computing



Distributed Computing

Parallel Computing

# You need to develop

- .. distributed on-disk storage

- .. in-memory storage (or shared memory buffer)

- .. thread pool to run hundreds of threads

- .. synchronize all components

- .. provide API for reusing by other developers

# All we love reinvent bicycles, but...

# HADOOP

Hadoop

# Disks Performance

# The main concept

Let's read data in parallel

# "Cheap" cluster

# Das Ist Musst surviven!

# Main components

- Hadoop Commons

- Hadoop Clients

- HDFS

- YARN

- MapReduce

# Hadoop frameworks

- Universal (MapReduce, Tez, RDD in Spark)

- Abstract (Pig, Pipeline Spark)

- SQL - like (Hive, Impala, Spark SQL)

- Processing graph (Giraph, GraphX)

- Machine Learning (Mahout, MLib)

- Stream processing (Spark Streaming, Storm)

# Hadoop Architecture

# Key features

- Automatic parallelization and distribution

- Fault-tolerance

- Data Locality

- Writing the Map and Reduce functions only

- Single-threaded model

# HDFS DAEMONS

# The main idea

'Time to transfer' > 'Time to seek'

# Main idea

# HDFS node types

- NameNode

- DataNode

- SecondaryNode

- StandbyNode

- Checkpoint Node

- Backup Node

# The main thought about HDFS

HDFS node is JVM daemon

# You can do it with HDFS node

- monitor with JMX

- use jmap, jps and so on..

- configure NameNode Heap Size

- use power of JVM flags

# YARN

# From Hadoop 1 to Hadoop 2

# Daemons in YARN

Memory Capacity: 4G

FIFO

Job 1: 4 tasks, 1G demand
Job 2: 4 tasks, 3G demand

Memory Capacity: 4G

FAIR

Memory Capacity: 4G

FFD

Different scheduling algorithms

# MAPREDUCE THEORY

# MapReduce in different languages

| Language | Code sample |
|----------|-------------|
| Java 8 | ```Integer totalAge = persons\n    .stream()\n    .map(Person::getAge)\n    .reduce( 0, (a, b) -> a + b);``` |
| Scala | ```val totalAge = persons\n    .map( (p: Person) => p.getAge )\n    .reduce( _ + _ )``` |
| Python | ```totalAge = reduce(\n    (lambda a, b: a + b),\n    list( map(lambda p: p.getAge, persons) )\n)``` |

# Think in Key-Value style



map (k1, v1) → list(k2, v2)

reduce (k2, list(v2*) )→ list(k3, v3)

# MR Typical Tasks

- WordCount

- Log handling

- Filtering

- Reporting Preparation

# Why should we use MapReduce?

# We try to reduce 'time to act' but keep BigData

# Classic Batch



Input    MapReduce    Output

# Do you like batches?

# Fixed Windows

# Filter all elements

# MAPREDUCE FRAMEWORK

# Pay attention!

# Hadoop != MapReduce

# Yet One YARN's lover

| GOVERNANCE INTEGRATION | DATA ACCESS | | | | | | | | SECURITY | OPERATIONS |
|---|---|---|---|---|---|---|---|---|---|---|

**Data Lifecycle & Governance**

Falcon;
Atlas

**Data Workflow**

Sqoop
Flume
Kafka
NFS
WebHDFS

| Batch | Script | SQL | NoSQL | Stream | Search | In-Mem | Others... |
|---|---|---|---|---|---|---|---|
| MapReduce | Pig | Hive | HBase Accumulo Phoenix | Storm | Solr | Spark | ISV Engines |
| Tez | Tez | Tez | Slider | Slider | | | Hortonworks Certified HDP YARN READY |
| | | | | | | | S / T |

**YARN: Data Operating System**

HDFS Hadoop Distributed File System

DATA MANAGEMENT

**SECURITY**

**Administration Authentication Authorization Auditing Data Protection**

Ranger
Knox
Atlas
HDFS Encryption

**OPERATIONS**

**Provisioning, Managing, & Monitoring**

Ambari
Cloudbreak
ZooKeeper

**Scheduling**

Oozie

# How MapReduce Works?



**Map()**          **Shuffle**          **Reduce()**          http://blog.sqlauthority.com

# Main steps

- Map

- Shuffle

- Reduce

## Minimal Runner

```java
public static void main(String[] args) throws Exception {

  int exitCode = ToolRunner.run(new MinimalMapReduce(), args);
  System.exit(exitCode);

 }
}
```

```java
public class MinimalMapReduce extends Configured implements Tool {
 @Override
 public int run(String[] args) throws Exception {



 }

 public static void main(String[] args) throws Exception {

  int exitCode = ToolRunner.run(new MinimalMapReduce(), args);
  System.exit(exitCode);

 }
}
```

**Minimal Runner**

```java
public class MinimalMapReduce extends Configured implements Tool {
 @Override
 public int run(String[] args) throws Exception {


  Job job = new Job(getConf());
  job.setJarByClass(getClass());
  FileInputFormat.addInputPath(job, new Path(args[0]));
  FileOutputFormat.setOutputPath(job, new Path(args[1]));
  return job.waitForCompletion(true) ? 0 : 1;

 }

 public static void main(String[] args) throws Exception {

  int exitCode = ToolRunner.run(new MinimalMapReduce(), args);
  System.exit(exitCode);

 }
}
```

**Minimal Runner**

```
job.setInputFormatClass(TextInputFormat.class);

job.setMapperClass(Mapper.class);

job.setMapOutputKeyClass(LongWritable.class);

job.setMapOutputValueClass(Text.class);

job.setPartitionerClass(HashPartitioner.class);

job.setNumReduceTasks(1);

job.setReducerClass(Reducer.class);

job.setOutputKeyClass(LongWritable.class);

job.setOutputValueClass(Text.class);

job.setOutputFormatClass(TextOutputFormat.class);
```

Job Config

# Would you like to config in Java?

WordCount

# TESTING & DEVELOPMENT

# MR Unit idea

# MRUnit – Testing Mapper



MR Unit

Unit Test

(1) Set up and execute test

MapDriver

(4) Compare the expected outputs

Mock Output Collector

(2) Call Map method with key / value

Mapper

(3) Map output is captured

Do it separatly

## Simple Test

```java
public class MRUnitHelloWorld {
    MapDriver<LongWritable, Text, Text, IntWritable> mapDriver;

    @Before
    public void setUp() {
        WordMapper mapper = new WordMapper();
        mapDriver = new MapDriver<LongWritable, Text, Text,
IntWritable>();
        mapDriver.setMapper(mapper);
    }

    @Test
    public void testMapper() {
        mapDriver.withInput(new LongWritable(1), new Text("cat
dog"));
        mapDriver.withOutput(new Text("cat"), new IntWritable(1));
        mapDriver.withOutput(new Text("dog"), new IntWritable(1));
        mapDriver.runTest();
    }
}
```

# Testing strategies

- First develop/test in local mode using small amount of data

- Test in pseudo-distributed mode and more data

- Test on fully distributed mode and even more data

- Final execution: fully distributed mode & all data

MRUnit

# FILE FORMATS

# Input/Output Formats

- Text based (CSV, TSV, JSON, XML)

- Sequence Files

- Column based (Parquet, RCFile, ORC)

- Avro

- HBase formats

- Custom formats

# ORC vs PARQUET



**File Size Comparison Across Encoding Methods**
Dataset: TPC-DS Scale 500 Dataset

**585 GB**
(Original Size)

**505 GB**
(14% Smaller)

Impala
**221 GB**
(62% Smaller)

Hive 12
**131 GB**
(78% Smaller)

- Larger Block Sizes
- Columnar format arranges columns adjacent within the file for compression & fast access

Encoded with
**Text**

Encoded with
**RCFile**

Encoded with
**Parquet**

Encoded with
**ORCFile**

# Codec Performance on the Wikipedia Text Corpus

# * LEVEL

# The main performance idea

Reduce shuffle time & resources

# MAPREDUCE ADVANCED

# Customize MapReduce!

# MapReduce for WordCount

The overall MapReduce word count process



| Input | Splitting | Mapping | Shuffling | Reducing | Final result |
|-------|-----------|---------|-----------|----------|--------------|

# WordCount Combiner

# MapReduce – Word Count Example Flow



**Combine**

**Combiner**

# Can we make something around map(), reduce() calls?

```java
/**
 * Called once at the start of the task.
 */
protected void setup(Context context
                    ) throws IOException, InterruptedException {

    // Prepare something for each Mapper or Reducer
    // Validate external sources
}
```

**Setup**

```
/**
 * Called once at the end of the task.
 */
protected void cleanup(Context context
                        ) throws IOException, InterruptedException
{
  // Finish something after each Mapper or Reducer
  // Handle specific exceptions
}
```

# Full control



WE NEED MORE POWER

memecrunch.com

```java
/**
 * Expert users can override this method for more complete control
over the
 * execution of the Mapper.
 * @param context
 * @throws IOException
 */
public void run(Context context) throws IOException,
InterruptedException {
  setup(context);
  try {
    while (context.nextKeyValue()) {
      map(context.getCurrentKey(), context.getCurrentValue(),
context);
    }
  } finally {
    cleanup(context);
  }
}
```

**Run Mapper**

```
/**
 * Advanced application writers can use the
 * {@link #run(*.Reducer.Context)} method to
 * control how the reduce task works.
 */
public void run(Context context) throws IOException,
InterruptedException {
  setup(context);
  try {
    while (context.nextKey()) {
      reduce(context.getCurrentKey(), context.getValues(),
context);
      // If a back up store is used, reset it
      Iterator<VALUEIN> iter = context.getValues().iterator();
      if(iter instanceof ReduceContext.ValueIterator) {

((ReduceContext.ValueIterator<VALUEIN>)iter).resetBackupStore();
      }
    }
  } finally {
    cleanup(context);
  }
}
```

Run Reducer

# REDUCER IS A PARTICIO OF PROBLEMO

# Could we skip shuffle step?

# Map-Only 'MapReduce' Jobs

May I customize Data Flow before shuffling?

# Hash Partitioner just do it..

```
public class HashPartitioner<K2, V2> extends Partitioner<K2, V2> {

public int getPartition(K2 key, V2 value, int numReduceTasks) {

        return (key.hashCode() & Integer.MAX_VALUE) % numReduceTasks;

    }
}
```

# Partitioner's Role in Shuffle and Sort



MapReduce with Partitioner and Combiner

# Full power

The partitioner's job is to logically funnel map outputs to the reducers.

**Input**

InputFormat.getSplits

RecordReader.
nextKeyValue

Mapper.map

Partitioner.getPartition

Reducer.reduce

RecordWriter.write

**Output**

Map phase

Reduce phase

Create split → Read split → Map → k,v → Partition → k,list(v) → Reduce → Write output

The InputFormat and RecordReader are responsible for determining what data to feed into the map function.

The map and reduce functions are typically written by the user to address a specific use case.

The RecordWriter writes the reduce output to the destination data sink, which is the final resting place of this MapReduce data flow.

Custom Partitioner

# JOINS

Employees

| Name | Age | Dept_Id |
|------|-----|---------|
| Alex | 26  | 2       |
| Ben  | 24  | 2       |
| Sara | 34  | 5       |

Department

| Dept_Id | Name  |
|---------|-------|
| 5       | Mkt   |
| 2       | Eng   |
| 3       | Sales |

Map Tasks

{ Key : 2 } { Value : Tag : Employees Record : [Alex, 26, 2] }

........
........

{ Key : 5 } { Value : Tag : Department Record : [5, Mkt] }

Shuffle & Sort

Reduce Tasks

{ Key : 2 } { Value : Tag : Employees Record : [Alex, 26, 2] , Tag : Employees Record : [Ben, 24, 2] , Tag : Department Record : [2, Eng] }

{ Key : 5 } { Value : Tag : Employees Record : [Sara, 34, 5] , Tag : Department Record : [5, Mkt] }

Output to HDFS

{ Key : 2 }  { [Alex, 26, Eng], [Ben, 24, Eng] }
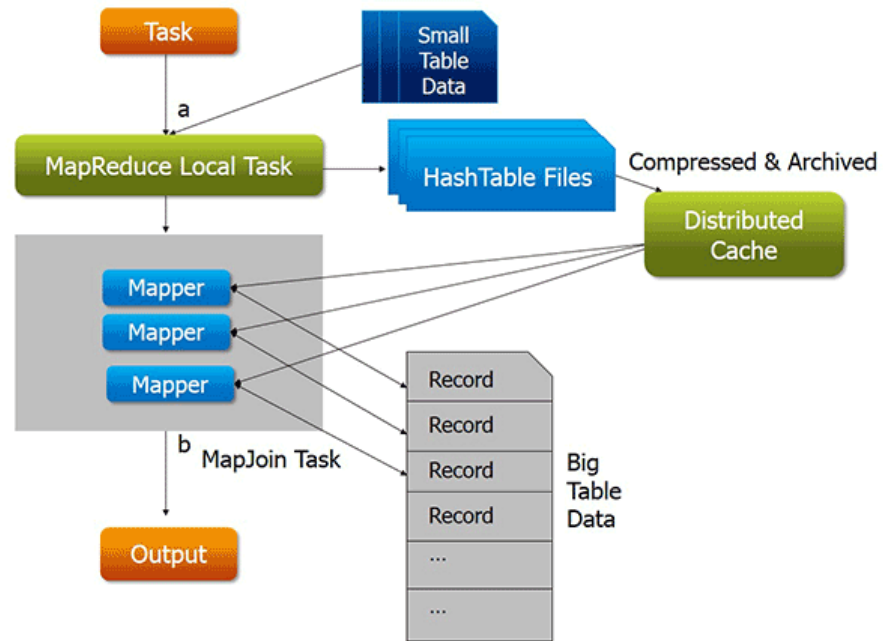
{ Key : 5 }  { [Sara, 34, Mkt] }

Reduce JOIN

# First idea

Let's skip Reduce + Shuffle

# Second idea

Let's copy small table on nodes

# Map-side join

# Map-side join for large datasets



Table b

Table a

Table c

Bucket a1

Bucket a2

Bucket c1

Mapper 1
Bucket b1
a1
c1

Mapper 2
Bucket b1
a1
c1

Mapper 3
Bucket b2
a2
c1

1. Spawn mapper based on the big table
2. Only matching buckets of all small tables are replicated onto each mapper

Normally in production, there will be thousands of buckets!

Table A
| 1, val_1 |
| 3, val_3 |
| 4, val_4 |
| 5, val_5 |

Table B
| 4, val_4 |
| 20, val_20 |
| 23, val_23 |

Table C
| 20, val_20 |
| 25, val_25 |

# What about Really Large Tables?

Employees

| Name | Age | Dept_Id |
|------|-----|---------|
| Alex | 26  | 2       |
| Ben  | 24  | 2       |
| Sara | 34  | 5       |

Department

| Dept_Id | Name  |
|---------|-------|
| 5       | Mkt   |
| 2       | Eng   |
| 3       | Sales |

**SELECT Employees.Name, Employees.Age, Department.Name  FROM Employees INNER JOIN Department ON Employees.Dept_Id=Department.Dept_Id**

# The main JOIN idea for large tables

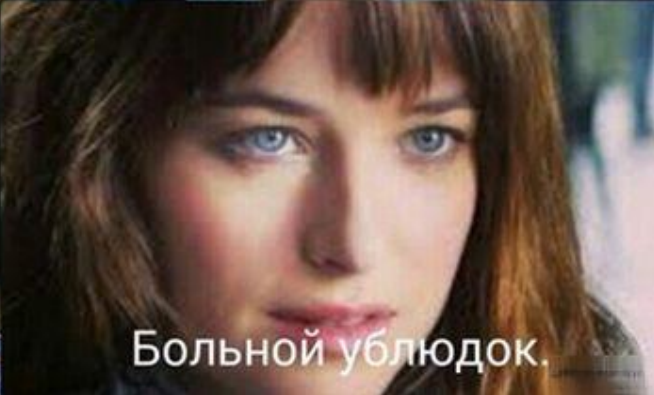Redis or Memcache cluster as Distributed Cache

# PERFORMANCE

Let's run on JVM!

# Typical mistakes

- Collections are stored and sorted in memory

- Logging each input key-value pairs

- JARs hell

- Skew input: all records go to one reducer

- Forget that mapper/reduce is run on different JVM
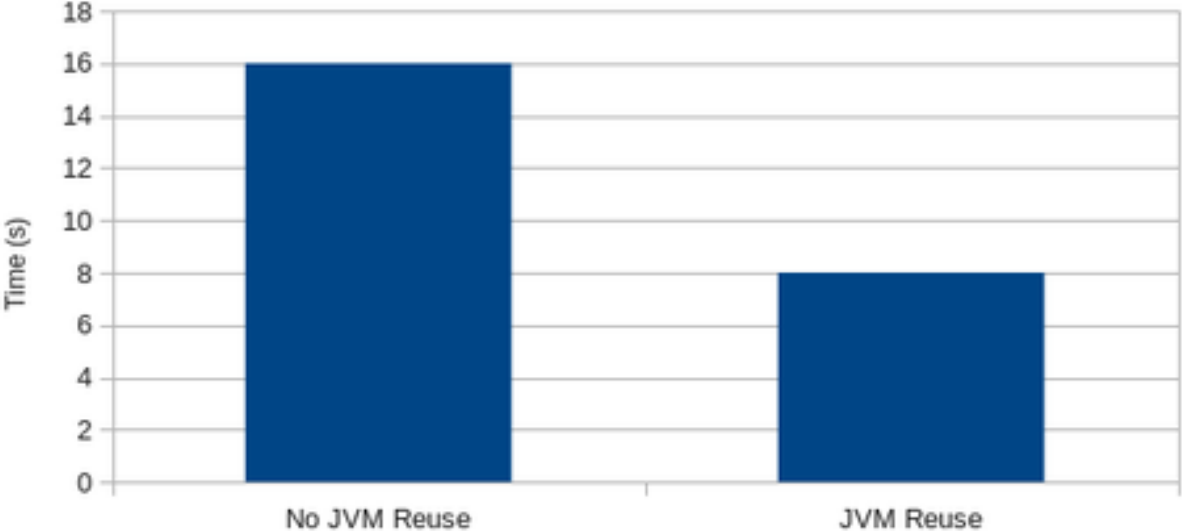
# Performance tips

- Correct data storage (on JVM ☺)

- Don't forget about combiner

- Use appropriate Writable type

- Min required replication factor

- Tune your JVM

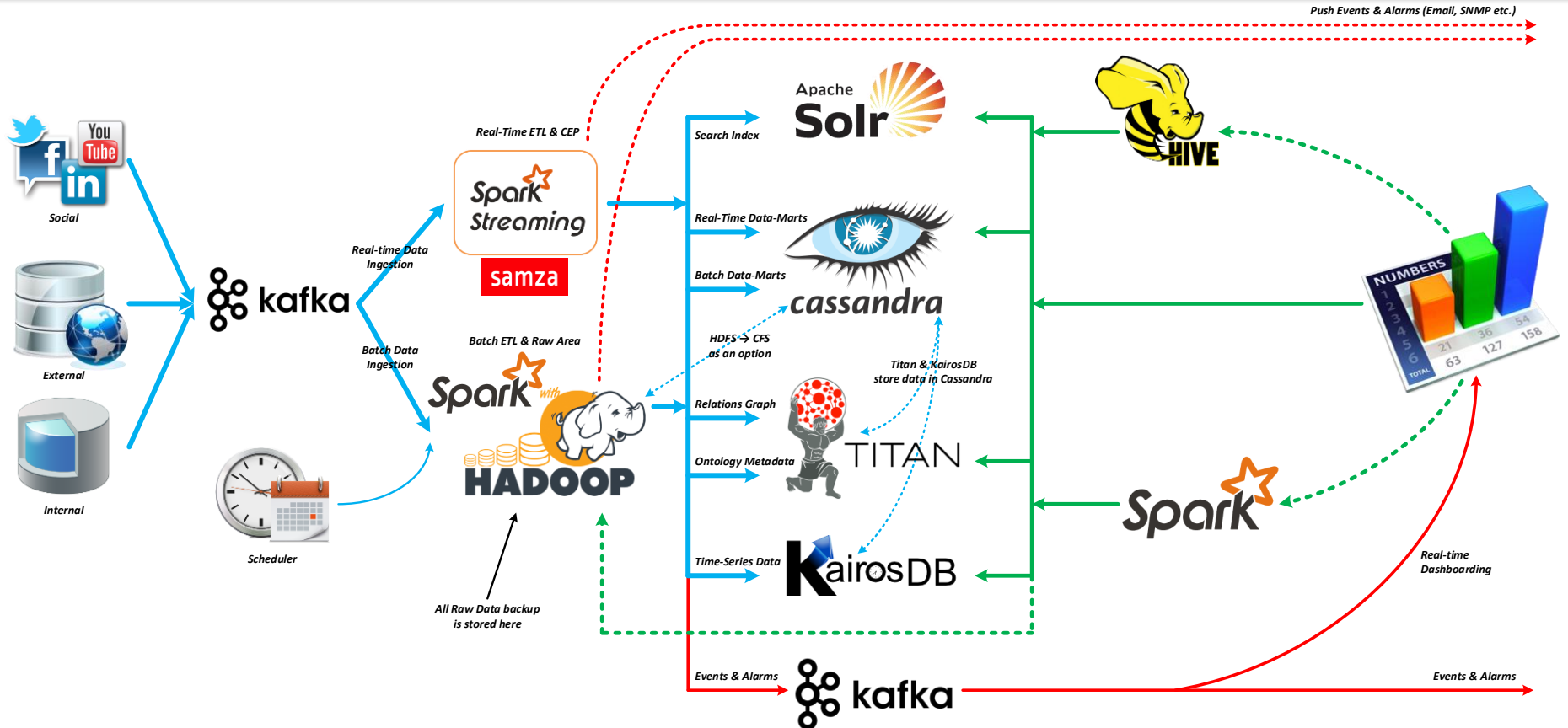- Think in terms Big-O

# JVM tuning

- *mapred.child.java.opts* (heap for tasks)

- *-XX:+PrintGCDetails -XX:+PrintGCTimeStamps*

- Low-latency GC collector *-XX:+UseConcMarkSweepGC,*

*-XX:ParallelGCThreads*

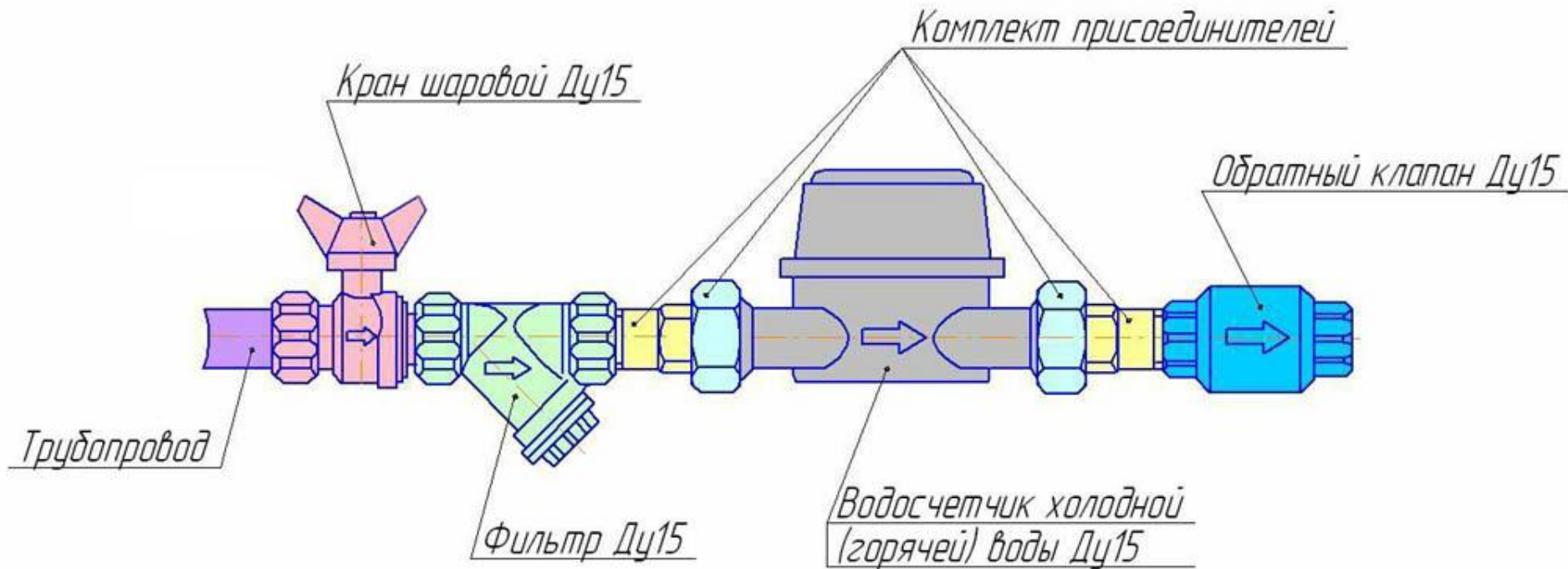- Xmx == Xms (max and starting heap size)

# JVM Reusing: Uber Task



Re-using JVMs across jobs

# And we can DO IT!

# It reminds me ...



Комплект присоединителей

Кран шаровой Ду15

Обратный клапан Ду15

Трубопровод

Фильтр Ду15

Водосчетчик холодной (горячей) воды Ду15
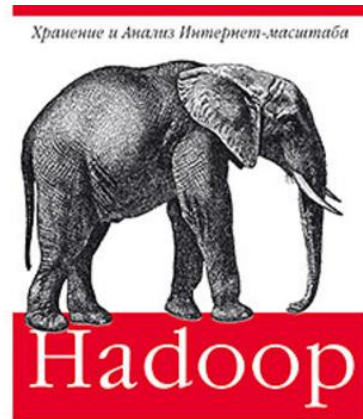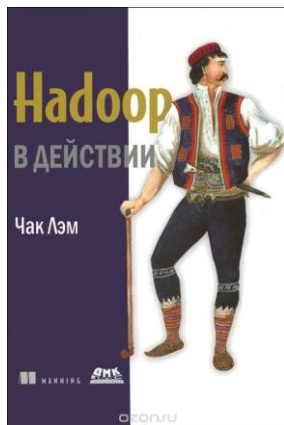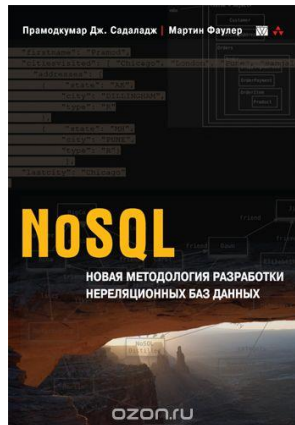
# MapReduce is not a ideal approach! But it works!

# Hadoop 3: Roadmap

- Move to Java 8

- Support more than 2 NameNodes (multiple standby NameNodes)

- Derive heap size or mapreduce.*.memory.mb automatically

- Work with SSD, RAM, HDD, CPU as resources for YARN

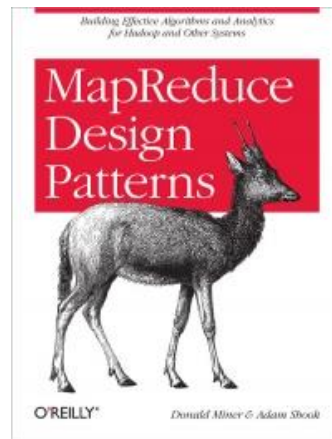- Support Docker containers

# Would you like to know more?

# Recommended Books

# Contacts

E-mail : Alexey_Zinovyev@epam.com

Twitter : @zaleslaw @BigDataRussia

vk.com/big_data_russia **Big Data Russia**

**+ Telegram** @bigdatarussia

vk.com/java_jvm **Java & JVM langs**

**+ Telegram** @javajvmlangs

# Github

Spark Tutorial: Core, Streaming, Machine Learning

https://github.com/zaleslaw/Spark-Tutorial

# Gitbook

Обработка данных на Spark 2.2 и Kafka 0.10

www.gitbook.com/book/zaleslaw/data-processing-book

Any questions?