# Apache Spark

Шугаепов Ильнур
VK.com
ilnur.shug@gmail.com

январь 2020 г.

Ограничения MapReduce
Spark RDD
Spark Programs
Implementation
Список литературы

Reminder
Проблемы
Частичные решения

| Input files | Map phase | Intermediate files (on local disks) | Reduce phase | Output files |

Ограничения MapReduce
Spark RDD
Spark Programs
Implementation
Список литературы

Reminder
Проблемы
Частичные решения

1. Iterative ML algorithms
2. Ad-hoc analytics
3. Interactive data-mining

Map Reduce is Good Enough?[1]

---

[1] Jimmy Lin. "Mapreduce is good enough? if all you have is a hammer, throw away everything that's not a nail!" B: *Big Data* 1.1 (2013), с. 28—37.

Ограничения MapReduce  
Spark RDD  
Spark Programs  
Implementation  
Список литературы

Reminder  
Проблемы  
Частичные решения

- Hive[2] — data wearehousing solution
- Pig[3] — dataflow system

---

[2]Ashish Thusoo и др. "Hive: a warehousing solution over a map-reduce framework".
В: *Proceedings of the VLDB Endowment* 2.2 (2009), с. 1626—1629.

[3]Alan F Gates и др. "Building a high-level dataflow system on top of Map-Reduce: the
Pig experience". В: *Proceedings of the VLDB Endowment* 2.2 (2009), с. 1414—1425.

Ограничения MapReduce
Spark RDD
Spark Programs
Implementation
Список литературы

Reminder
Проблемы
Частичные решения

## Hive
Main components

- HiveQL — SQL like language
- Metastore — catalog with metadata about tables
- Compiler — converts query to a execution plan (MapReduce jobs)

Ограничения MapReduce
Spark RDD
Spark Programs
Implementation
Список литературы

Reminder
Проблемы
Частичные решения

## Pig

- Like Hive but with different query language and without Metastore

Ограничения MapReduce
**Spark RDD**
Spark Programs
Implementation
Список литературы

RDD abstraction
Пример
Linage graph / Lazy computation

### Определение

*RDD:*

- *Resilient — отказоустойчивый*
- *Distributed — разбитый на партиции*
- *Dataset*

*read-only, partitioned collection of records*

Ограничения MapReduce
Spark RDD
Spark Programs
Implementation
Список литературы

RDD abstraction
Пример
Linage graph / Lazy computation

## Efficient Fault-tolerance

- Запомним граф вычислений (linage)
- Тогда если часть данных будет потеряна, то их легко можно восстановить

Ограничения MapReduce
Spark RDD
Spark Programs
Implementation
Список литературы

RDD abstraction
Пример
Linage graph / Lazy computation

## RDD creation

RDD можно построить одним из следующих способов:

- Из данных находящихся на HDFS или в RAM
- Выполнив операцию над существующим RDD:
  - Transformations
  - Actions

Ограничения MapReduce
Spark RDD
Spark Programs
Implementation
Список литературы

RDD abstraction
Пример
Linage graph / Lazy computation

## Transformations

- map
- filter
- join
- reduceByKey
- ...

Ограничения MapReduce
**Spark RDD**
Spark Programs
Implementation
Список литературы

RDD abstraction
Пример
Linage graph / Lazy computation

## Actions

- count — количество элементов в RDD
- save — сохранение RDD, например, на HDFS
- collect — отправить RDD на driver

Ограничения MapReduce
Spark RDD
Spark Programs
Implementation
Список литературы

RDD abstraction
Пример
Linage graph / Lazy computation

## Persistance and Partitioning

- Users can indicate which RDDs they will reuse and choose a storage strategy for them (e.g., in-memory storage).
- They can also ask that an RDD's elements be partitioned across machines based on a key in each record.

Ограничения MapReduce
Spark RDD
Spark Programs
Implementation
Список литературы

RDD abstraction
Пример
Linage graph / Lazy computation

Операции задают декларативное описание того, что мы хотим
сделать

Ограничения MapReduce
Spark RDD
Spark Programs
Implementation
Список литературы

RDD abstraction
Пример
Linage graph / Lazy computation

```
1 lines = spark.textFile("hdfs://...")
2 errors = lines.filter(_.startsWith("ERROR"))
3 errors.persist()
4
5 errors.count()
6
7 // Count errors mentioning MySQL:
8 errors.filter(_.contains("MySQL")).count()
9
10 // Return the time fields of errors mentioning
11 // HDFS as an array (assuming time is field
12 // number 3 in a tab-separated format):
13 errors.filter(_.contains("HDFS"))
14       .map(_.split('\t')(3))
15       .collect()
```

Ограничения MapReduce
**Spark RDD**
Spark Programs
Implementation
Список литературы

RDD abstraction
**Пример**
Linage graph / Lazy computation

## Lineage graph



Рис.: Boxes represent RDDs and arrows represent transformations

Ограничения MapReduce
Spark RDD
Spark Programs
Implementation
Список литературы

RDD abstraction
Пример
Linage graph / Lazy computation

## Linage graph

RDD has enough information about how it was derived from other datasets (its lineage) to compute its partitions from data in stable storage

Ограничения MapReduce
**Spark RDD**
Spark Programs
Implementation
Список литературы

RDD abstraction
Пример
Linage graph / Lazy computation

## Lazy computation

- Spark computes RDDs lazily the first time they are used in an action, so that it can pipeline transformations.
- Spark keeps persistent RDDs in memory by default, but it can spill them to disk if there is not enough RAM.

Ограничения MapReduce
Spark RDD
**Spark Programs**
Implementation
Список литературы

Operations
Примеры
Representing RDDs

Рис.: The user's *driver* program launches multiple *workers*, which read data blocks from a distributed file system and can persist computed RDD partitions in memory.

Ограничения MapReduce    Operations
Spark RDD    Примеры
**Spark Programs**    Representing RDDs
Implementation
Список литературы

## Transformations
Types

$$
\begin{aligned}
map(f\colon T \Rightarrow U) &: & RDD[T] &\Rightarrow RDD[U]\\
filter(f\colon T \Rightarrow Bool) &: & RDD[T] &\Rightarrow RDD[T]\\
flatMap(f\colon T \Rightarrow Seq[U]) &: & RDD[T] &\Rightarrow RDD[U]\\
sample(fraction\colon Float) &: & RDD[T] &\Rightarrow RDD[T]\\
groupByKey() &: & RDD[(K,V)] &\Rightarrow RDD[(K, Seq[V])]\\
reduceByKey(f\colon (V,V) \Rightarrow V) &: & RDD[(K,V)] &\Rightarrow RDD[(K,V)]\\
union() &: & (RDD[T], RDD[T]) &\Rightarrow RDD[T]\\
join() &: & (RDD[(K,V)], RDD[(K,W)]) &\\
 & & &\Rightarrow RDD[(K, (V,W))]\\
cogroup() &: & (RDD[(K,V)], RDD[(K,W)]) &\\
 & & &\Rightarrow RDD[(K, (Seq[V], Seq[W]))]\\
crossProduct() &: & (RDD[T], RDD[U]) &\Rightarrow RDD[(T, U)]\\
mapValues(f\colon V \Rightarrow W) &: & RDD[(K,V)] &\Rightarrow RDD[(K,W)]\\
sort(c\colon Comparator[K]) &: & RDD[(K,V)] &\Rightarrow RDD[(K,V)]\\
partitionBy(p\colon Partitioner[K]) &: & RDD[(K,V)] &\Rightarrow RDD[(K,V)]
\end{aligned}
$$

Ограничения MapReduce
Spark RDD
**Spark Programs**
Implementation
Список литературы

**Operations**
Примеры
Representing RDDs

## Actions
Types

$$
\begin{array}{rcl}
count() & : & RDD[T] \Rightarrow Long \\
collect() & : & RDD[T] \Rightarrow Seq[T] \\
reduce(f\colon (T,T) \Rightarrow T) & : & RDD[T] \Rightarrow T \\
lookup(k\colon K) & : & RDD[(K,V)] \Rightarrow Seq[V] \\
save(path\colon String) & : & \text{Outputs RDD to a storage system}
\end{array}
$$

Ограничения MapReduce
Spark RDD
**Spark Programs**
Implementation
Список литературы

Operations
Примеры
Representing RDDs

## Logistic Regression

LogReg[4]

---

[4]Trevor Hastie, Robert Tibshirani и Jerome Friedman. *The elements of statistical learning: data mining, inference, and prediction*. Springer Science & Business Media, 2009.

Ограничения MapReduce
Spark RDD
**Spark Programs**
Implementation
Список литературы

Operations
**Примеры**
Representing RDDs

# Logistic Regression
Code

```
1 val points = spark.textFile(...)
2                    .map(parsePoint).persist()
3 var w = // random initial vector
4 for (i <- 1 to ITERATIONS) {
5   val gradient = points.map{ p =>
6     p.x * (1/(1+exp(-p.y*(w dot p.x)))-1)*p.y
7   }.reduce((a,b) => a+b)
8   w -= gradient
9 }
```

Ограничения MapReduce
Spark RDD
**Spark Programs**
Implementation
Список литературы

Operations
Примеры
Representing RDDs

# Logistic Regression
Performance



Рис.: Duration of the first and later iterations in Hadoop, HadoopBinMem and Spark for logistic regression and k-means using 100 GB of data on a 100-node cluster.

Ограничения MapReduce
Spark RDD
Spark Programs
Implementation
Список литературы

Operations
Примеры
Representing RDDs

# Logistic Regression
Performance



(a) Logistic Regression

(b) K-Means

Рис.: Running times for iterations after the first in Hadoop, HadoopBinMem, and Spark

Ограничения MapReduce
Spark RDD
Spark Programs
Implementation
Список литературы

Operations
Примеры
Representing RDDs

## Logistic Regression
Performance

Keeping points in memory across iterations can yield a $20\times$ speedup

Ограничения MapReduce
Spark RDD
**Spark Programs**
Implementation
Список литературы

Operations
Примеры
Representing RDDs

## PageRank

PageRank[56]

---

[5] Lawrence Page и др. *The pagerank citation ranking: Bringing order to the web.*
Тех. отч. Stanford InfoLab, 1999.

[6] Jure Leskovec, Anand Rajaraman и Jeffrey David Ullman. *Mining of massive data sets.* Cambridge university press, 2019.

Ограничения MapReduce
Spark RDD
**Spark Programs**
Implementation
Список литературы

Operations
**Примеры**
Representing RDDs

# PageRank
Code

```scala
1 val links = spark.textFile(...).map(...).persist()
2 var ranks = // RDD of (URL, rank) pairs
3 for (i <- 1 to ITERATIONS) {
4   // Build an RDD of (targetURL, float) pairs
5   // with the contributions sent by each page
6   val contribs = links.join(ranks).flatMap {
7     (url, (links, rank)) =>
8       links.map(dest => (dest, rank/links.size))
9   }
10  // Sum contributions by URL and get new ranks
11  ranks = contribs.reduceByKey((x,y) => x+y)
12             .mapValues(sum => a/N + (1-a)*sum)
13 }
```

Ограничения MapReduce
Spark RDD
**Spark Programs**
Implementation
Список литературы

Operations
**Примеры**
Representing RDDs

# PageRank
Linage graph

Ограничения MapReduce
Spark RDD
**Spark Programs**
Implementation
Список литературы

Operations
Примеры
Representing RDDs

# PageRank
Performance

Preserving partitioning might help

```
1 links = spark.textFile(...).map(...)
2                .partitionBy(myPartFunc).persist()
```

If ranks and links are co-partitioned then join requires no
communication

Ограничения MapReduce
Spark RDD
**Spark Programs**
Implementation
Список литературы

Operations
**Примеры**
Representing RDDs

# PageRank
Performance



Рис.: Performance of PageRank on Hadoop and Spark.

Ограничения MapReduce
Spark RDD
**Spark Programs**
Implementation
Список литературы

Operations
Примеры
**Representing RDDs**

- Partititons — atomic pices of the dataset
- Dependencies — dependencies on parent RDDs

Ограничения MapReduce
Spark RDD
**Spark Programs**
Implementation
Список литературы

Operations
Примеры
**Representing RDDs**

## Dependencies



Рис.: Examples of narrow and wide dependencies. Each box is an RDD, with partitions shown as shaded rectangles.

Ограничения MapReduce
Spark RDD
**Spark Programs**
Implementation
Список литературы

Operations
Примеры
Representing RDDs

Dependencies
Narrow

- Narrow dependencies allow for pipelined execution on one cluster node, which can compute all the parent partitions.
- Recovery after a node failure is more efficient with a narrow dependency, as only the lost parent partitions need to be recomputed, and they can be recomputed in parallel on different nodes.

Ограничения MapReduce
Spark RDD
Spark Programs
**Implementation**
Список литературы

Spark Application
Job Scheduling
Memory Managment
Remarks

Рис.: Starting a Spark application on a distributed system

Ограничения MapReduce
Spark RDD
Spark Programs
Implementation
Список литературы

Spark Application
Job Scheduling
Memory Managment
Remarks

## Spark Application

### Замечание

- *One node can have multiple Spark executors, but an executor cannot span multiple nodes.*
- *An RDD will be evaluated across the executors in partitions (shown as red rectangles).*
- *Each executor can have multiple partitions, but a partition cannot be spread across multiple executors.*

Ограничения MapReduce
Spark RDD
Spark Programs
Implementation
Список литературы

Spark Application
Job Scheduling
Memory Managment
Remarks

## SparkContext

SparkContext - connection between user's program and cluster.
Containes information about requested resources, type of resources
allocation (dynamic/static), etc

Ограничения MapReduce
Spark RDD
Spark Programs
**Implementation**
Список литературы

Spark Application
**Job Scheduling**
Memory Managment
Remarks

## Pipeline

1. User runs an action on RDD
2. Scheduler builds a DAG of *stages* (each stage containes piplened transformations with narrow dependencies. *Boundaries* - wide deps.)
3. Scheduler launches *tasks* to compute missing partitions

Ограничения MapReduce
Spark RDD
Spark Programs
**Implementation**
Список литературы

Spark Application
**Job Scheduling**
Memory Managment
Remarks

Рис.: The Spark application tree

Ограничения MapReduce
Spark RDD
Spark Programs
Implementation
Список литературы

Spark Application
Job Scheduling
Memory Managment
Remarks

Рис.: Boxes with solid outlines are RDDs. Partitions are shaded rectangles, in black if they are already in memory.

Ограничения MapReduce      Spark Application
Spark RDD      Job Scheduling
Spark Programs      Memory Managment
Implementation      Remarks
Список литературы

- User can choose how to persist data:
  - in-memory storage as deserialized Java objects
  - in-memory storage as serialized data
  - on-disk storage.
- LRU eviction policy at the level of RDDs

📄 Gates, Alan F и др. "Building a high-level dataflow system on top of Map-Reduce: the Pig experience". В: *Proceedings of the VLDB Endowment* 2.2 (2009), с. 1414—1425.

📄 Hastie, Trevor, Robert Tibshirani и Jerome Friedman. *The elements of statistical learning: data mining, inference, and prediction*. Springer Science & Business Media, 2009.

📄 Karau, Holden и Rachel Warren. *High performance Spark: best practices for scaling and optimizing Apache Spark*. " O'Reilly Media, Inc.", 2017.

📄 Leskovec, Jure, Anand Rajaraman и Jeffrey David Ullman. *Mining of massive data sets*. Cambridge university press, 2019.

📄 Lin, Jimmy. "Mapreduce is good enough? if all you have is a hammer, throw away everything that's not a nail!" В: *Big Data* 1.1 (2013), с. 28—37.

📄 Page, Lawrence и др. *The pagerank citation ranking: Bringing order to the web.* Тех. отч. Stanford InfoLab, 1999.

📄 *Spark Overview*.
   https://spark.apache.org/docs/latest/index.html.

📄 Thusoo, Ashish и др. "Hive: a warehousing solution over a
   map-reduce framework". В: *Proceedings of the VLDB Endowment*
   2.2 (2009), с. 1626—1629.

📄 Vavilapalli, Vinod Kumar и др. "Apache hadoop yarn: Yet another
   resource negotiator". В: *Proceedings of the 4th annual
   Symposium on Cloud Computing*. 2013, с. 1—16.

📄 Zaharia, Matei и др. "Resilient distributed datasets: A fault-tolerant
   abstraction for in-memory cluster computing". В: *Presented as
   part of the 9th {USENIX} Symposium on Networked Systems
   Design and Implementation ({NSDI} 12)*. 2012, с. 15—28.

📄 Zaharia, Matei и др. "Spark: Cluster computing with working sets.".
   В: *HotCloud* 10.10-10 (2010), с. 95.