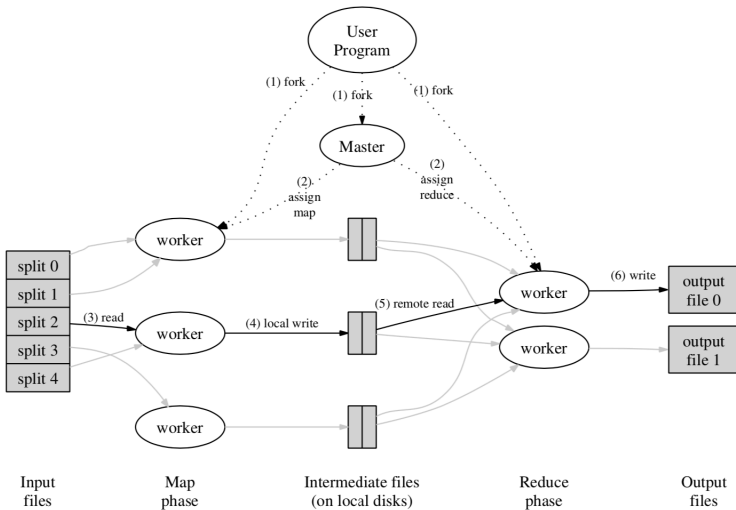




Apache Spark

Ильнур Шугаев

MapReduce



Основные приложения MR¹

- ✓ Iterative ML algorithms
- ✓ Ad-hoc analytics
- ✓ Interactive data-mining

Map Reduce is Good Enough?

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

Ограничения MapReduce

- 1 Сохраняет (временные) результаты всегда на HDFS
- 2 Ничего не знает про структуру данных
- 3 Написание программ из большого числа map,reduce фаз - проблематично

Ограничения MapReduce

Hive

data warehousing solution

Частичные решения

Pig

dataflow system

Hive²

Main components

HiveQL

SQL like query language

Metastore

catalog with metadata
about tables

Compiler

converts query to a ex-
ecution plan

²Ashish Thusoo et al. "Hive: a warehousing solution over a map-reduce framework". In: *Proceedings of the VLDB Endowment* 2.2 (2009), pp. 1626–1629.

Pig³

Like Hive but with different query language and without Metastore



³Alan F Gates et al. "Building a high-level dataflow system on top of Map-Reduce: the Pig experience". In: *Proceedings of the VLDB Endowment* 2.2 (2009), pp. 1414–1425.

Table of Contents

1. Ограничения MapReduce

2. Основные понятия

Пример

RDD abstraction

Lineage graph / Lazy computation

3. Производительность

4. Implementation

Spark Program

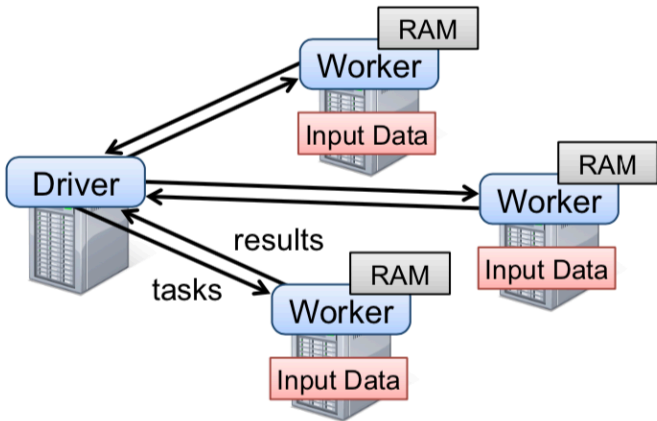


Figure: The user's *driver* program launches multiple *workers*, which read data blocks from a distributed file system

Пример

Поиск по логам

```
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()
```



Lineage graph

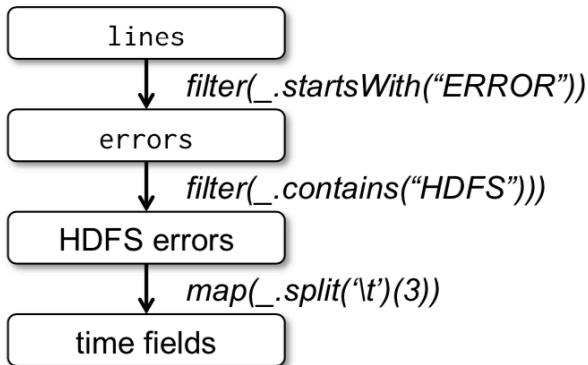


Figure: Boxes represent RDDs and arrows represent transformations

RDD abstraction

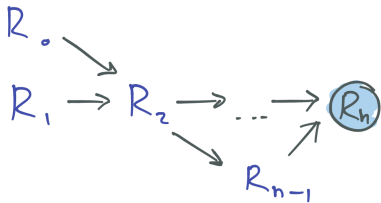
Definition 1 (RDD):

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

read-only, partitioned collection of records

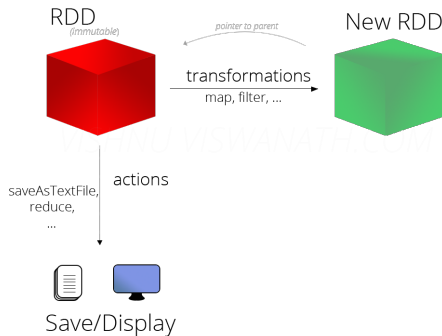
Efficient Fault-tolerance

- Запомним граф вычислений (lineage)
- Тогда если часть данных будет потеряна, то их легко можно восстановить
- RDD знает от каких данных (других RDD) он зависит



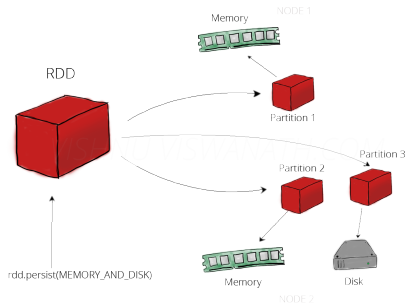
Построение RDD

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



Persistence and Partitioning

- Пользователь может задать каким образом будет храниться RDD
- Пользователь может указать способ партицирования для RDD



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.

ОСНОВНЫЕ ПОНЯТИЯ

Краткий итог

- 1 Программы на спарке — высокоуровневое описание манипуляций над RDDs.
- 2 Все вычисления ленивые
- 3 Пользователь может управлять тем, где будут храниться временные результаты
- 4 Граф зависимостей обеспечивает высокую надежность вычислений

Table of Contents

1. Ограничения MapReduce

2. Основные понятия

3. Производительность

Transformations and Actions

Примеры

4. Implementation

Transformations

Types

<i>map</i> (<i>f</i> : $T \Rightarrow U$)	:	$RDD[T] \Rightarrow RDD[U]$
<i>filter</i> (<i>f</i> : $T \Rightarrow Bool$)	:	$RDD[T] \Rightarrow RDD[T]$
<i>flatMap</i> (<i>f</i> : $T \Rightarrow Seq[U]$)	:	$RDD[T] \Rightarrow RDD[U]$
<i>sample</i> (<i>fraction</i> : <i>Float</i>)	:	$RDD[T] \Rightarrow RDD[T]$
<i>groupByKey</i> ()	:	$RDD[(K, V)] \Rightarrow RDD[(K, Seq[V])]$
<i>reduceByKey</i> (<i>f</i> : $(V, V) \Rightarrow V$)	:	$RDD[(K, V)] \Rightarrow RDD[(K, V)]$
<i>union</i> ()	:	$(RDD[T], RDD[T]) \Rightarrow RDD[T]$
<i>join</i> ()	:	$(RDD[(K, V)], RDD[(K, W)])$ $\Rightarrow RDD[(K, (V, W))]$
<i>cogroup</i> ()	:	$(RDD[(K, V)], RDD[(K, W)])$ $\Rightarrow RDD[(K, (Seq[V], Seq[W]))]$
<i>crossProduct</i> ()	:	$(RDD[T], RDD[U]) \Rightarrow RDD[(T, U)]$
<i>mapValues</i> (<i>f</i> : $V \Rightarrow W$)	:	$RDD[(K, V)] \Rightarrow RDD[(K, W)]$
<i>sort</i> (<i>c</i> : <i>Comparator</i> [<i>K</i>])	:	$RDD[(K, V)] \Rightarrow RDD[(K, V)]$
<i>partitionBy</i> (<i>p</i> : <i>Partitioner</i> [<i>K</i>])	:	$RDD[(K, V)] \Rightarrow RDD[(K, V)]$

Actions

Types

<i>count()</i>	:	$RDD[T] \Rightarrow Long$
<i>collect()</i>	:	$RDD[T] \Rightarrow Seq[T]$
<i>reduce(f: (T, T) \Rightarrow T)</i>	:	$RDD[T] \Rightarrow T$
<i>lookup(k: K)</i>	:	$RDD[(K, V)] \Rightarrow Seq[V]$
<i>save(path: String)</i>	:	Outputs RDD to a storage system

Logistic Regression⁴

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 }
```

- $x \in \mathbb{R}^n$, $\hat{y} = \sigma(\omega \cdot x)$, $\omega \in \mathbb{R}^n$
- $\mathcal{L} = \frac{1}{N} \sum y_i \log \hat{y}_i + (1-y_i) \log(1-\hat{y}_i)$
- $\omega^t \leftarrow \omega^{t-1} - \eta \frac{\partial \mathcal{L}}{\partial \omega}$

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

Logistic Regression

Performance

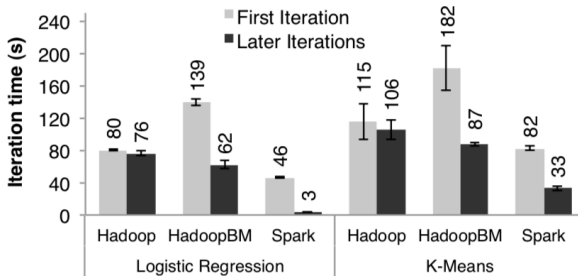


Figure: Duration of the first and later iterations in Hadoop, HadoopBinMem and Spark for logistic regression and k-means using 100 GB of data on a 100-node cluster.

Logistic Regression

Performance

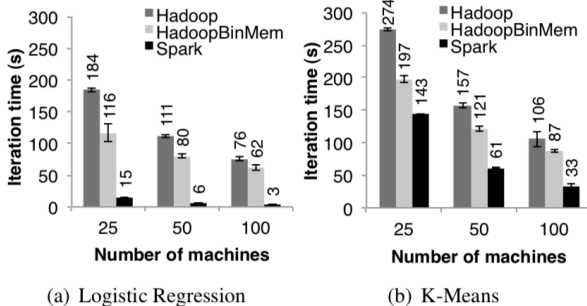


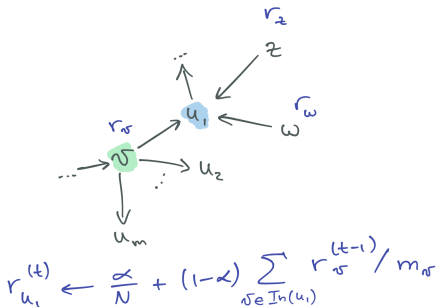
Figure: Running times for iterations after the first in Hadoop, HadoopBinMem, and Spark

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

PageRank⁵⁶

Code

```
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 }
```

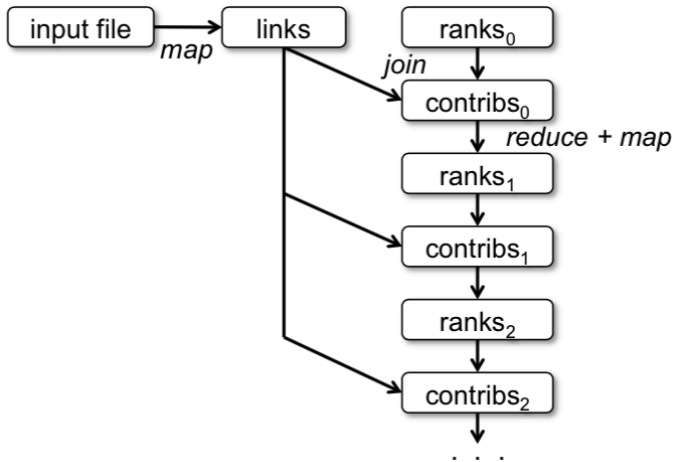


⁵Lawrence Page et al. *The pagerank citation ranking: Bringing order to the web*. Tech. rep. Stanford InfoLab, 1999.

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

PageRank

Lineage graph



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

PageRank

Performance

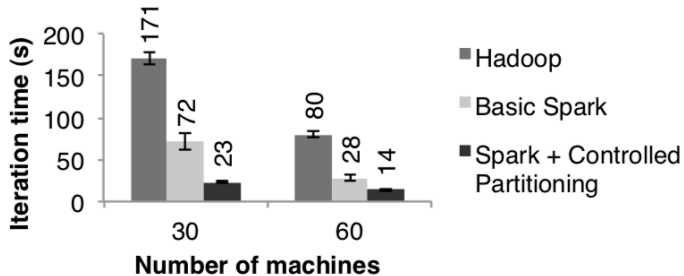


Figure: Performance of PageRank on Hadoop and Spark.

Производительность

Краткий итог

- 1 Spark работает в 20-100 раз быстрее чем MapReduce
- 2 Писать программы можно сильно быстрее

Table of Contents

1. Ограничения MapReduce

2. Основные понятия

3. Производительность

4. Implementation

Representing RDDs

Job Scheduling

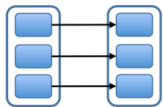
Representing RDDs

Граф

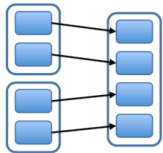
- ✓ Partititons — atomic pices of the dataset
- ✓ Dependencies — dependencies on parent RDDs

Dependencies

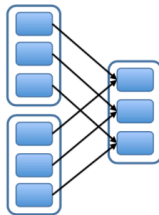
Narrow Dependencies:



map, filter

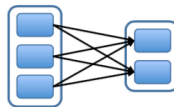


union

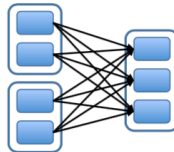


join with inputs
co-partitioned

Wide Dependencies:



groupByKey



join with inputs not
co-partitioned

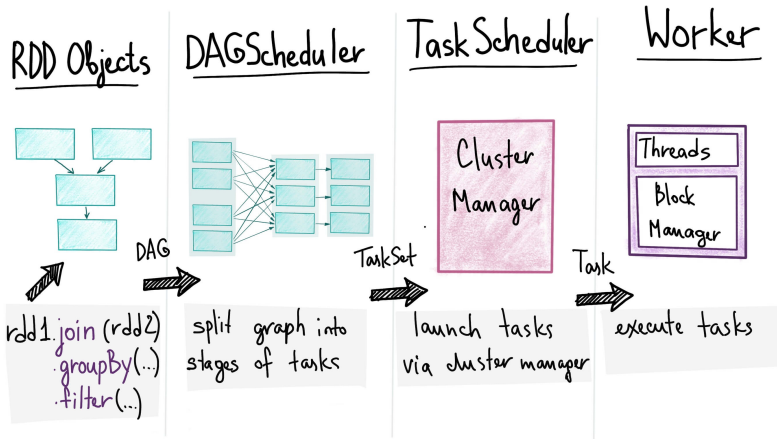
Figure: Examples of narrow and wide dependencies. Each box is an RDD, with partitions shown as shaded rectangles.

Dependencies

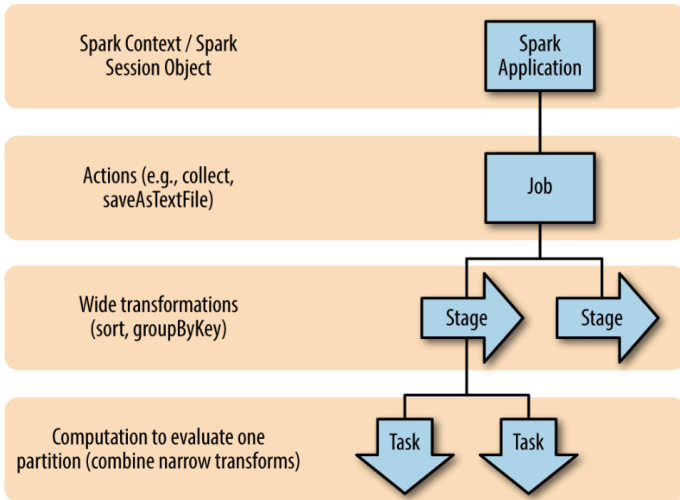
Narrow

- Narrow dependencies allow for pipelined execution on one cluster node
- Recovery after a node failure is more efficient with a narrow dependency

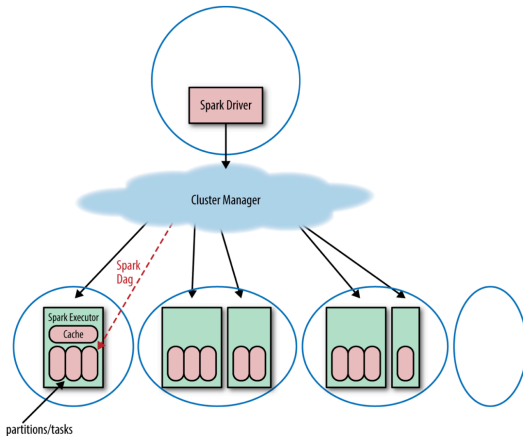
Spark Application Pipeline



Spark Application Tree



Замечания



- 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.

Замечания

SparkContext

Definition 2 (SparkContext): *Connection between user's program and cluster. Contains information about requested resources, type of resources allocation (dynamic/static), etc*

Замечания

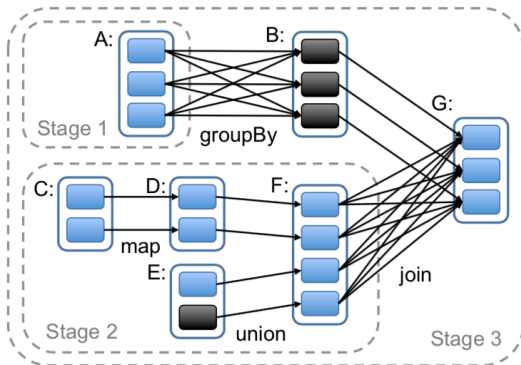













Figure: Boxes with solid outlines are RDDs. Partitions are shaded rectangles, in black if they are already in memory.

Итог

-  Gates, Alan F et al. "Building a high-level dataflow system on top of Map-Reduce: the Pig experience". In: *Proceedings of the VLDB Endowment* 2.2 (2009), pp. 1414–1425.
-  Hastie, Trevor, Robert Tibshirani, and Jerome Friedman. *The elements of statistical learning: data mining, inference, and prediction*. Springer Science & Business Media, 2009.
-  Karau, Holden and Rachel Warren. *High performance Spark: best practices for scaling and optimizing Apache Spark*. " O'Reilly Media, Inc.", 2017.
-  Leskovec, Jure, Anand Rajaraman, and 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!" In: *Big Data* 1.1 (2013), pp. 28–37.
-  Page, Lawrence et al. *The pagerank citation ranking: Bringing order to the web*. Tech. rep. Stanford InfoLab, 1999.
-  Spark Overview. <https://spark.apache.org/docs/latest/index.html>.
-  Thusoo, Ashish et al. "Hive: a warehousing solution over a map-reduce framework". In: *Proceedings of the VLDB Endowment* 2.2 (2009), pp. 1626–1629.
-  Vavilapalli, Vinod Kumar et al. "Apache hadoop yarn: Yet another resource negotiator". In: *Proceedings of the 4th annual Symposium on Cloud Computing*. 2013, pp. 1–16.

-  Zaharia, Matei et al. "Resilient distributed datasets: A fault-tolerant abstraction for in-memory cluster computing". In: *Presented as part of the 9th {USENIX} Symposium on Networked Systems Design and Implementation ({NSDI} 12)*. 2012, pp. 15–28.
-  Zaharia, Matei et al. "Spark: Cluster computing with working sets.". In: *HotCloud 10.10-10* (2010), p. 95.



Вопросы?

Ильнур Шугаев