



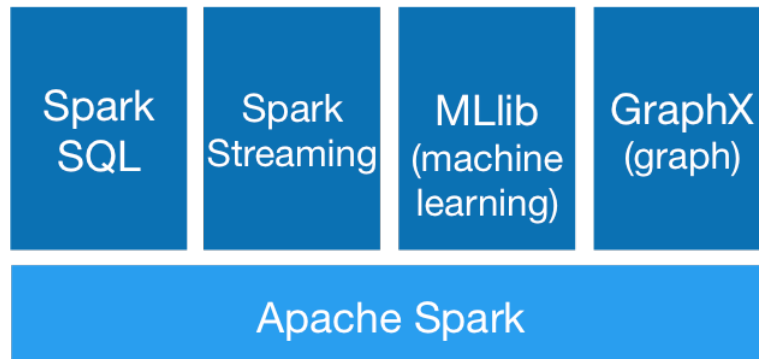
Apache Spark™

February 2015

Grigor Vladimir
vladimir.grigor@kiosked.com

What is Apache Spark™?

- Spark is a fast and general engine for large-scale data processing
- High-level APIs in Java, Scala and Python
- Higher-level tools:
 - Spark SQL for SQL and structured data processing
 - MLlib for machine learning
 - GraphX for graph processing
 - Spark Streaming.



Who uses Apache Spark™?

 amazon.com

 Bai du 百度

 淘宝网 全球
Taobao.com

YAHOO!

 ebay™

 shopify

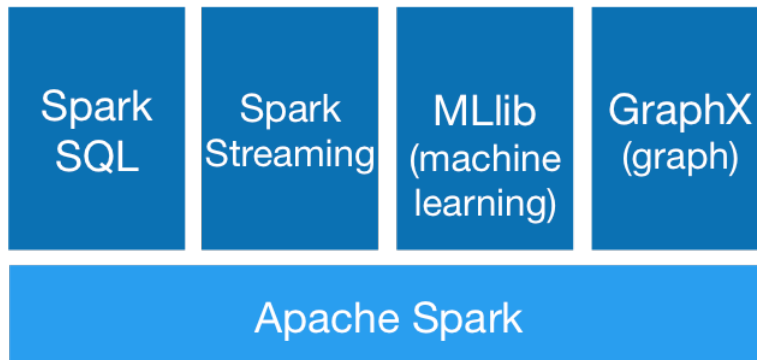
 NOKIA

 Yandex

 GROUPON®

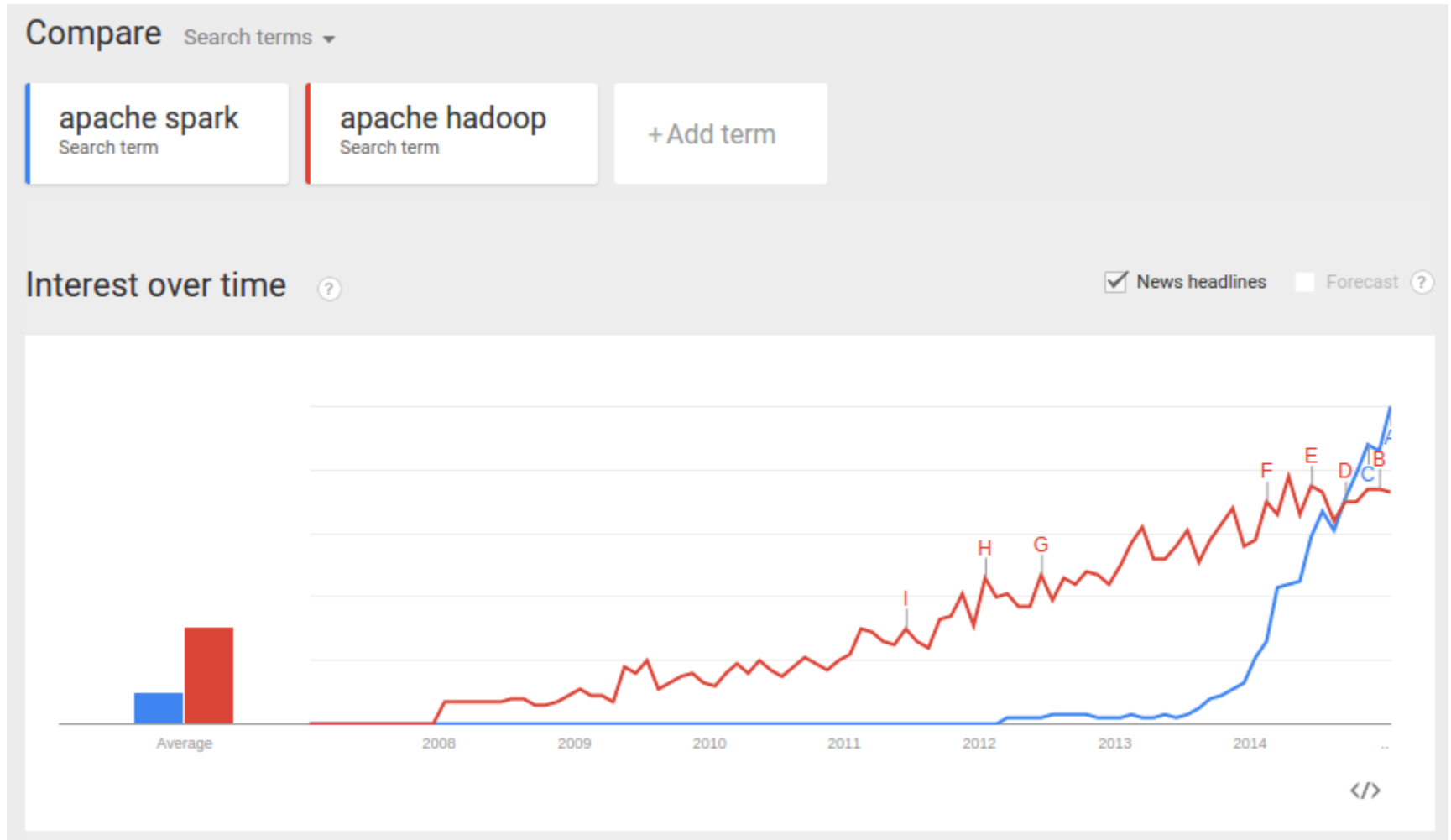
What Apache Spark can be used for?

- Exploratory big data analysis
- Real-time analytics
- Batch processing - EMR jobs



WARNING:
TECHNICAL
STUFF AHEAD

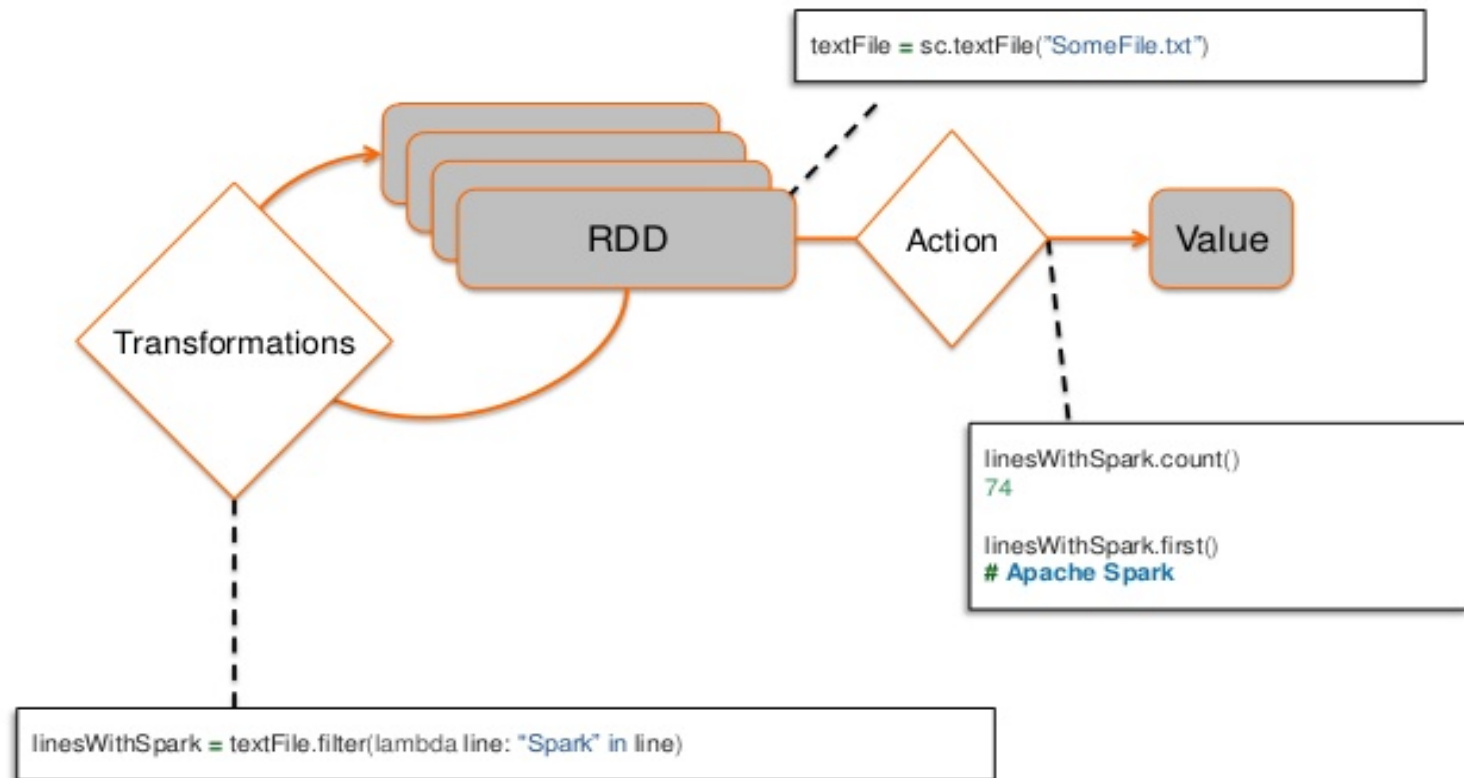
Google Trends



Demo

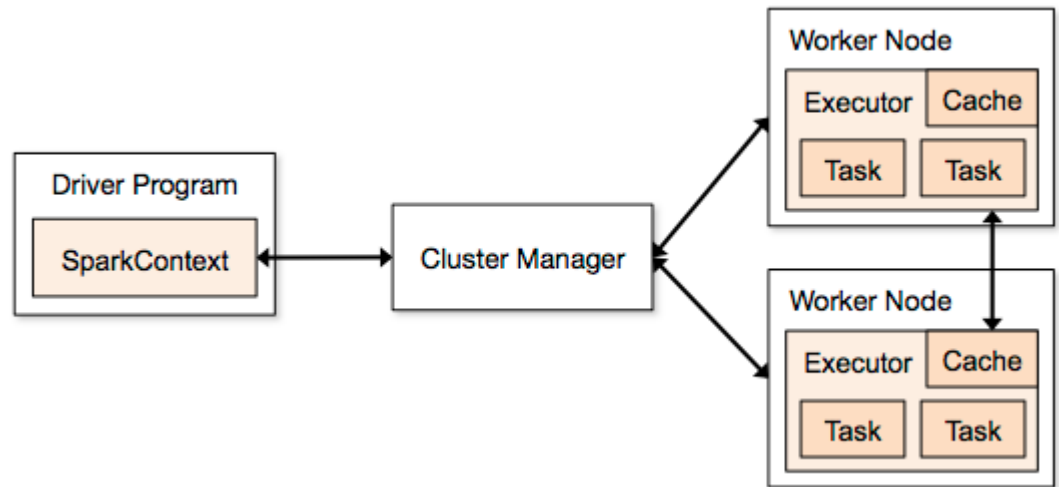


How it works: Resilient Distributed Dataset (RDD)



How it works: cluster mode

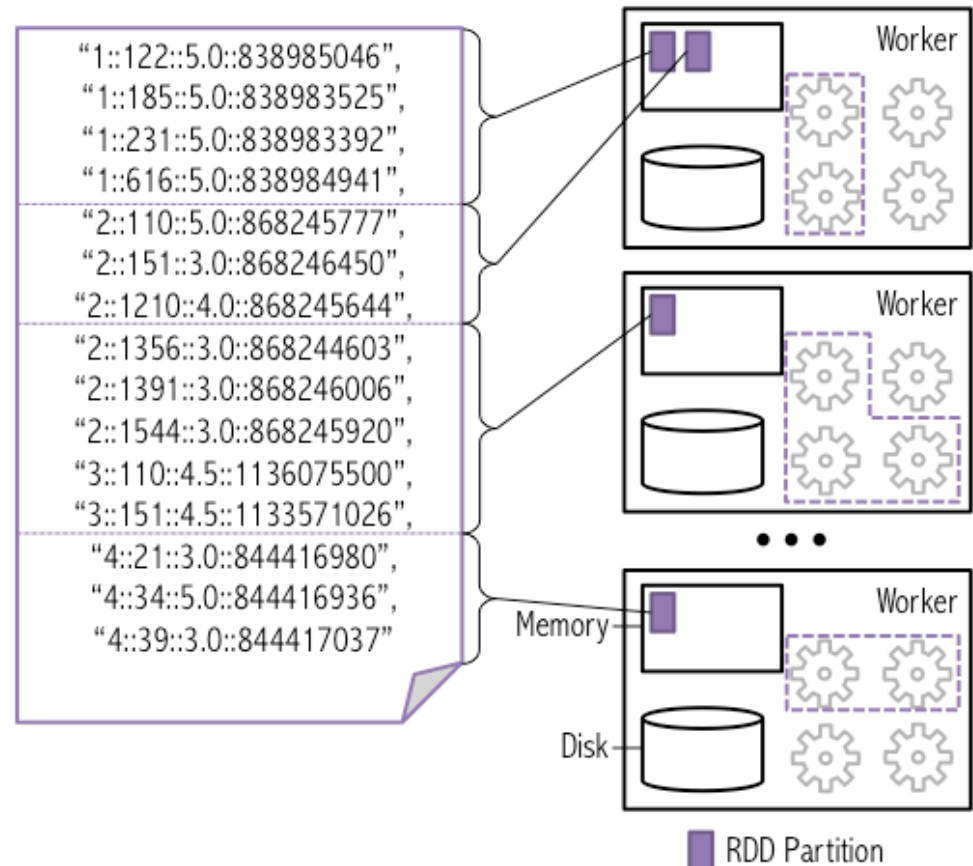
- Application is coordinated by SparkContext run within Driver
- Cluster Manager allocate resources - Executors
- SparkContext sends application code to Executors
- SparkContext sends tasks to Executors



How it works: RDD partitions

Dataset is broken into
partitions

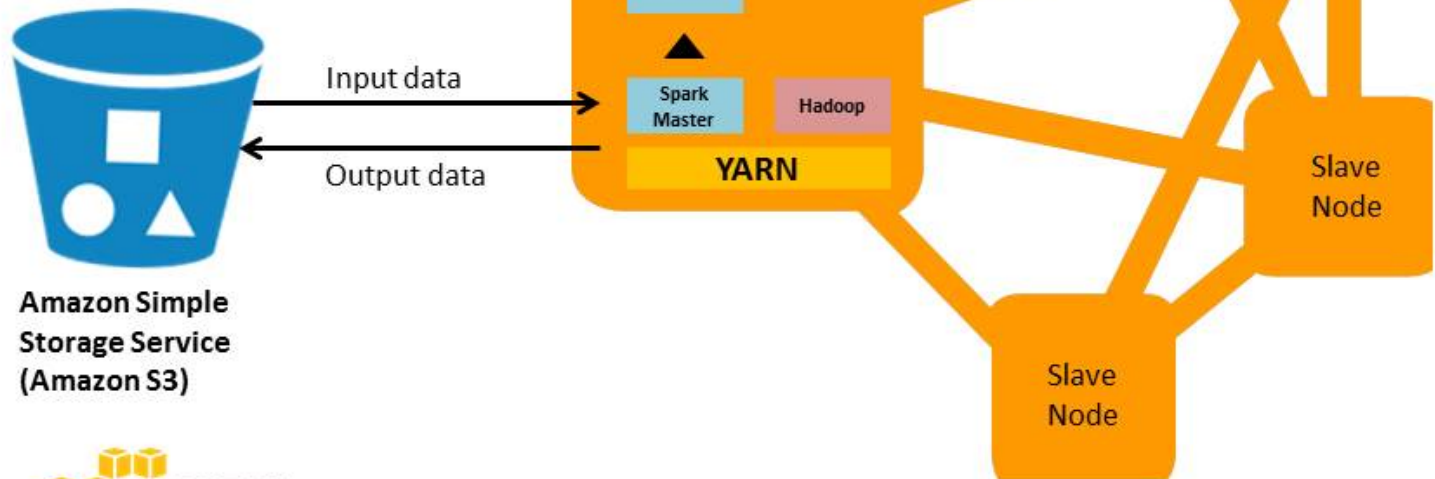
Partitions are each stored
in a worker's memory



How it works: Yarn on EMR

Spark is managed by YARN,
and runs alongside Hadoop on
your Amazon EMR cluster.
Each node is a Amazon EC2
instance

Spark and Spark SQL are
installed on the master node
of your cluster.



Amazon Elastic MapReduce (Amazon EMR) Cluster

Pitfalls

- “cluster mode is currently not supported for standalone clusters, Mesos clusters, or python applications”
 - <http://spark.apache.org/docs/latest/submitting-applications.html>
- Pickle serialization issues with Python code
- No clear deployment
- JVM
- Myriad of configuration possibilities

Application	User program built on Spark. Consists of a <i>driver program</i> and <i>executors</i> on the cluster.
Application jar	A jar containing the user's Spark application.
Driver program	The process running the <code>main()</code> function of the application and creating the <code>SparkContext</code>
Cluster manager	An external service for acquiring resources on the cluster (e.g. standalone manager, Mesos, YARN)
Deploy mode	Distinguishes where the driver process runs. In "cluster" mode, the framework launches the driver inside of the cluster. In "client" mode, the submitter launches the driver outside of the cluster.
Worker node	Any node that can run application code in the cluster
Executor	A process launched for an application on a worker node, that runs tasks and keeps data in memory or disk storage across them. Each application has its own executors.
Task	A unit of work that will be sent to one executor
Job	A parallel computation consisting of multiple tasks that gets spawned in response to action (e.g. <code>save</code> , <code>collect</code>)
Stage	Each job gets divided into smaller sets of tasks called <i>stages</i> that depend on each other

Spark WordCount vs 50+ lines of Java MR

Scala:

```
val f = sc.textFile("README.md")
val wc = f.flatMap(l => l.split(" ")).map(word => (word, 1)).reduceByKey(_ + _)
wc.saveAsTextFile("wc_out")
```

Python:

```
from operator import add
f = sc.textFile("README.md")
wc = f.flatMap(lambda x: x.split(' ')).map(lambda x: (x, 1)).reduceByKey(add)
wc.saveAsTextFile("wc_out")
```



... is hiring!

Vladimir Grigor
vladimir.grigor@kiosked.com

[Facebook.com/kiosked](https://www.facebook.com/kiosked)
[@Kiosked](https://twitter.com/Kiosked)