Universitatea POLITEHNICA din București

FACULTATEA DE
**AUTOMATICĂ** ȘI
**CALCULATOARE**

# Apache Spark -  Fast unified analytics engine
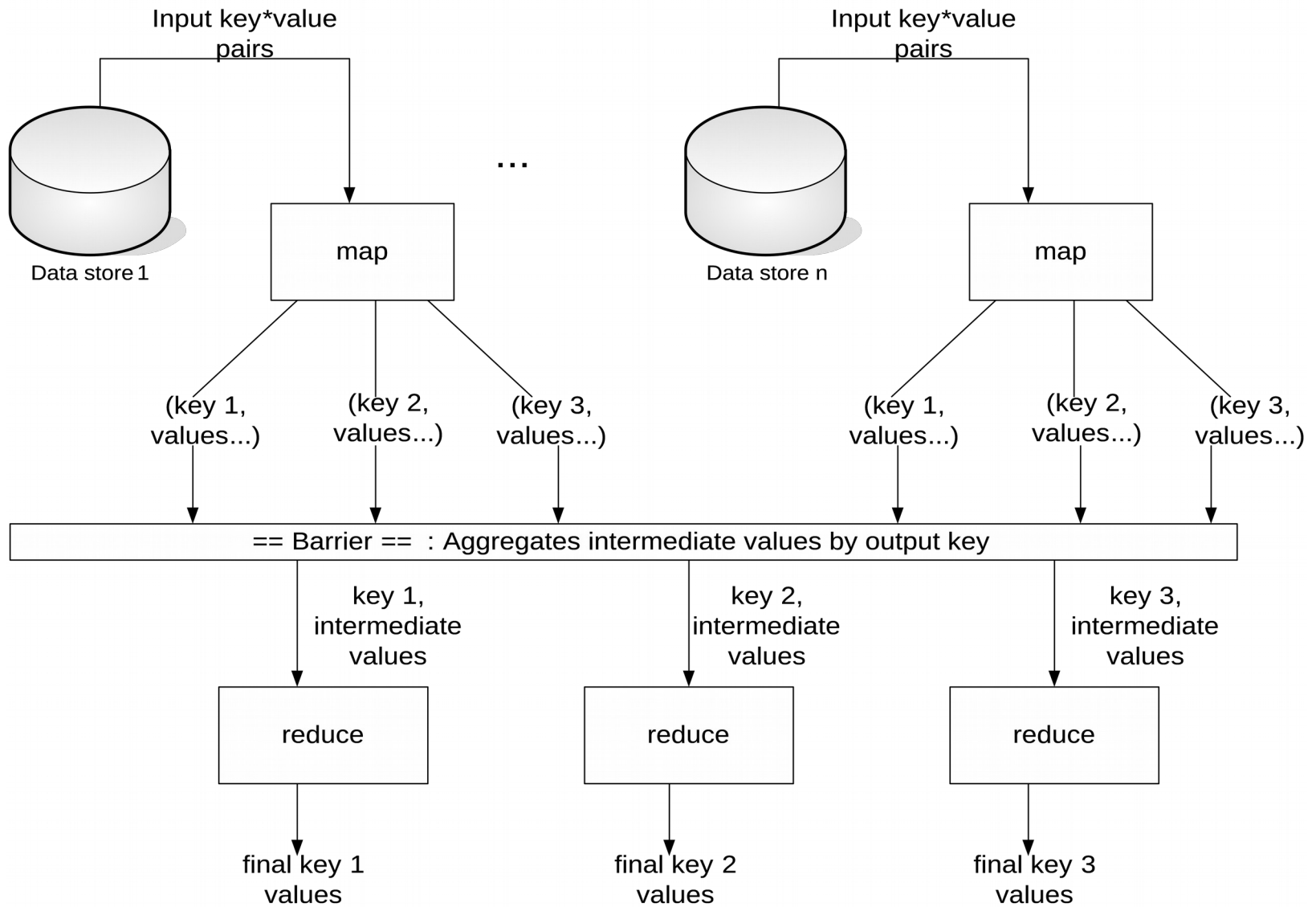
## Authors

Ioana – Laura Popescu

# Motivation App

- ## Web Search
  - 12 PB Web data
  - Must search quickly
  - How?

- ## Web Search Primer
  - Users supply words in query --> find all documents with a specified word
  - Read 12PB of web pages, find keywords in a single machine:
    - 100MB/s = 1GB/10s → 12PB in 120000000s -> ~ 4 years!

# Motivation App (2)

- We need parallelism!
    - Run previous task in one day?
    - Need a cluster (at least 1400 machines)

- Functionality of the computation system?
    - Move data around
    - Check liveness
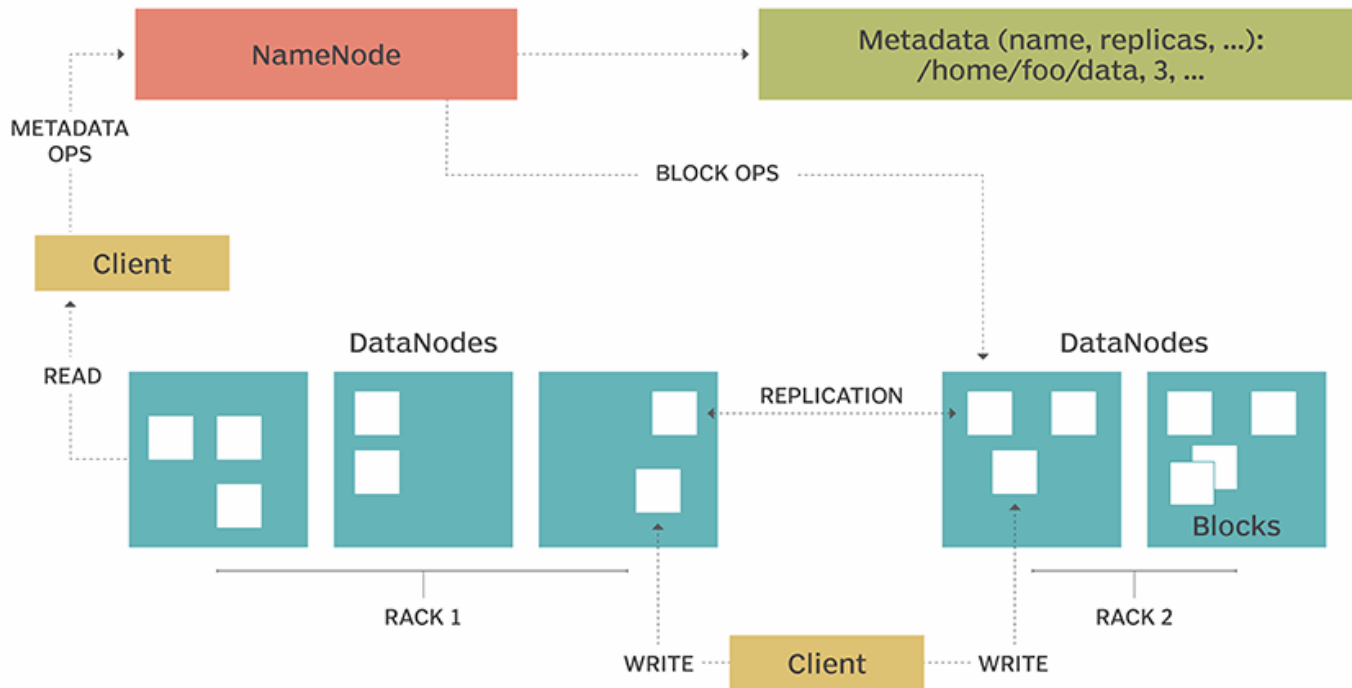    - **Deal with failures**
    - Process stuff

Input key*value pairs

Input key*value pairs

Data store 1

Data store n

...

map

map

(key 1, values...)

(key 2, values...)

(key 3, values...)

(key 1, values...)

(key 2, values...)

(key 3, values...)

== Barrier == : Aggregates intermediate values by output key

key 1, intermediate values

key 2, intermediate values

key 3, intermediate values

reduce

reduce

reduce

final key 1 values

final key 2 values

final key 3 values

# MapReduce

- Automatic parallelization & distribution

  - map() and reduce() functions run on parallel

- Fault-tolerant

  - Master detects worker failures and re-executes in-progess reduce tasks

- Clean abstraction for programmers

# Hadoop

- Uses MapReduce
- Integrates with HDFS
- Commonly used
- Writes all intermediary results to disk!
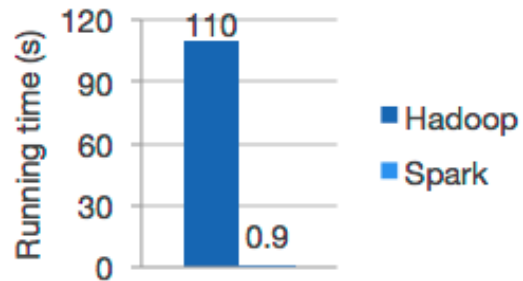  - It slows down significantly even the smallest of jobs
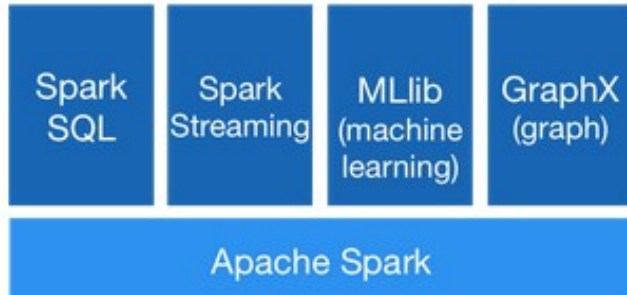
# HDFS

# Apache Spark

- Tries to use RAM whenever available

- Uses RDD's
  - Collections of objects spread over the datanodes
  - Built with parallel transformations
  - Rebuilt on failure detection

Logistic regression in Hadoop and Spark

| | | | |
|---|---|---|---|
| Spark SQL | Spark Streaming | MLib (machine learning) | GraphX (graph) |
| Apache Spark | | | |

# Spark Architecture