# pySpark

Fast and Expressive Big Data Analytics with Python

Matei Zaharia

UC Berkeley / MIT

amplab
UC BERKELEY

spark-project.org

# What is Spark?

Fast and expressive cluster computing system interoperable with Apache Hadoop

Improves efficiency through:
- » In-memory computing primitives
- » General computation graphs

→ Up to 100× faster
(2-10× on disk)

Improves usability through:
- » Rich APIs in Scala, Java, Python
- » Interactive shell

→ Often 5× less code

# Project History

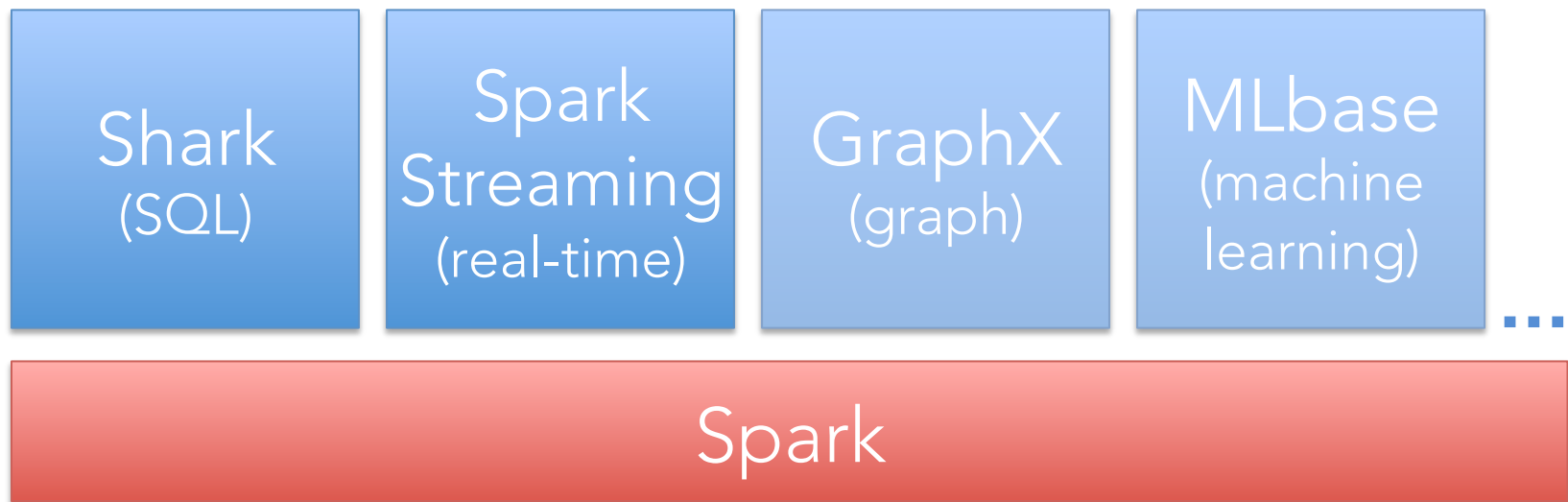Started in 2009, open sourced 2010

17 companies now contributing code
» Yahoo!, Intel, Adobe, Quantifind, Conviva, Bizo, …

Entered Apache incubator in June

Python API added in February

# An Expanding Stack

Spark is the basis for a wide set of projects in the Berkeley Data Analytics Stack (BDAS)

| Shark (SQL) | Spark Streaming (real-time) | GraphX (graph) | MLbase (machine learning) | ... |
|---|---|---|---|---|

**Spark**

More details: amplab.berkeley.edu

# This Talk

Spark programming model

Examples

Demo

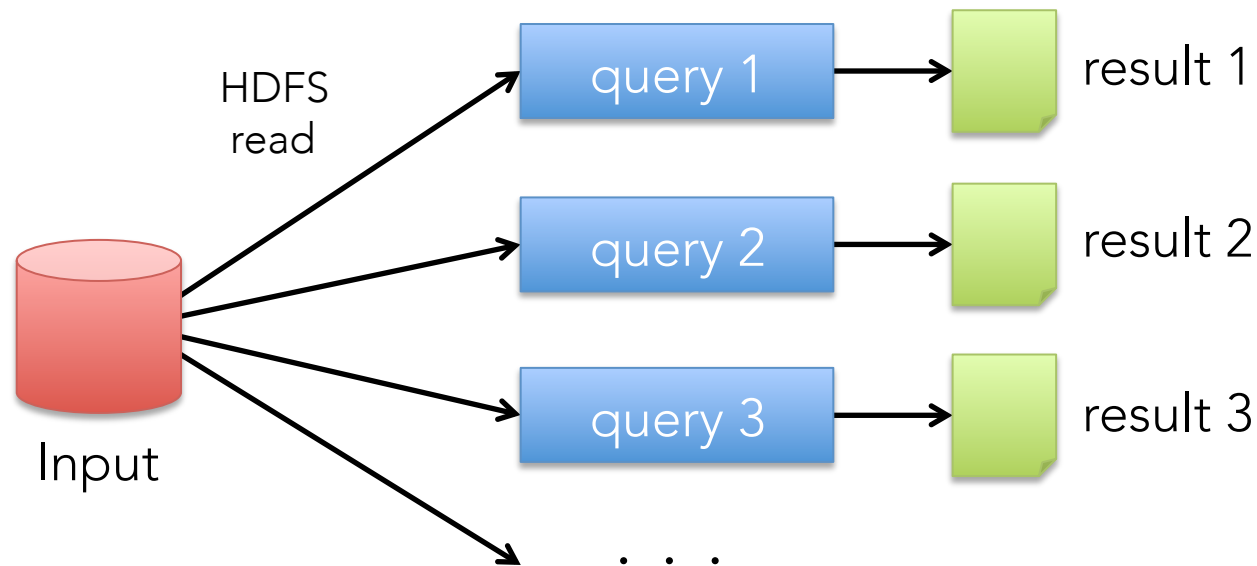Implementation

Trying it out

# Why a New Programming Model?
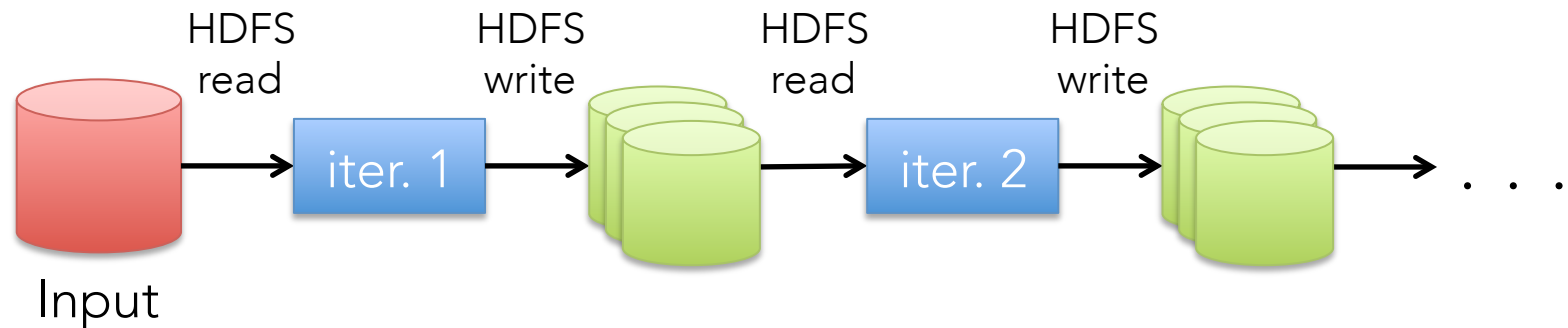
MapReduce simplified big data processing, but users quickly found two problems:

**Programmability:** tangle of map/red functions

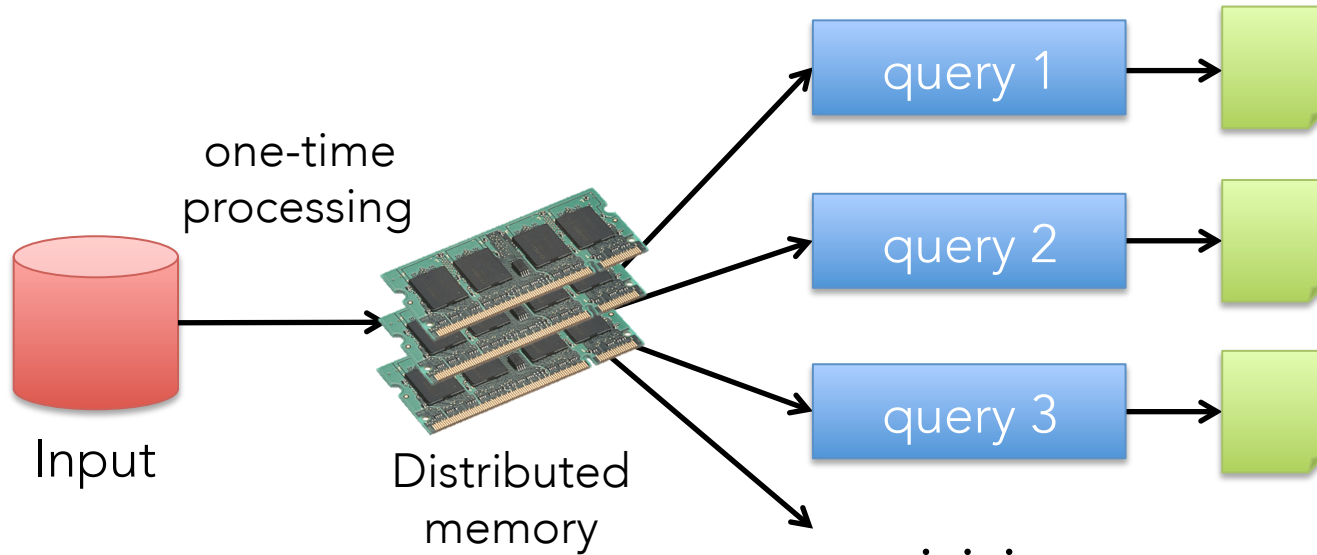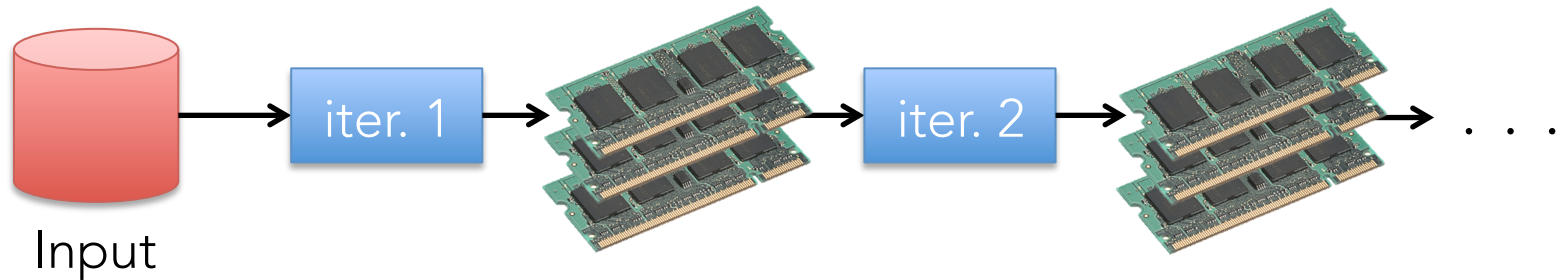**Speed:** MapReduce inefficient for apps that share data across multiple steps
  » Iterative algorithms, interactive queries

# Data Sharing in MapReduce



Slow due to data replication and disk I/O

# What We'd Like



Input

iter. 1 → iter. 2 → . . .

one-time processing

Input

Distributed memory

query 1 →

query 2 →

query 3 →

. . .

10-100× faster than network and disk

# Spark Model

*Write programs in terms of transformations on distributed datasets*

Resilient Distributed Datasets (RDDs)
- » Collections of objects that can be stored in memory or disk across a cluster
- » Built via parallel transformations (map, filter, …)
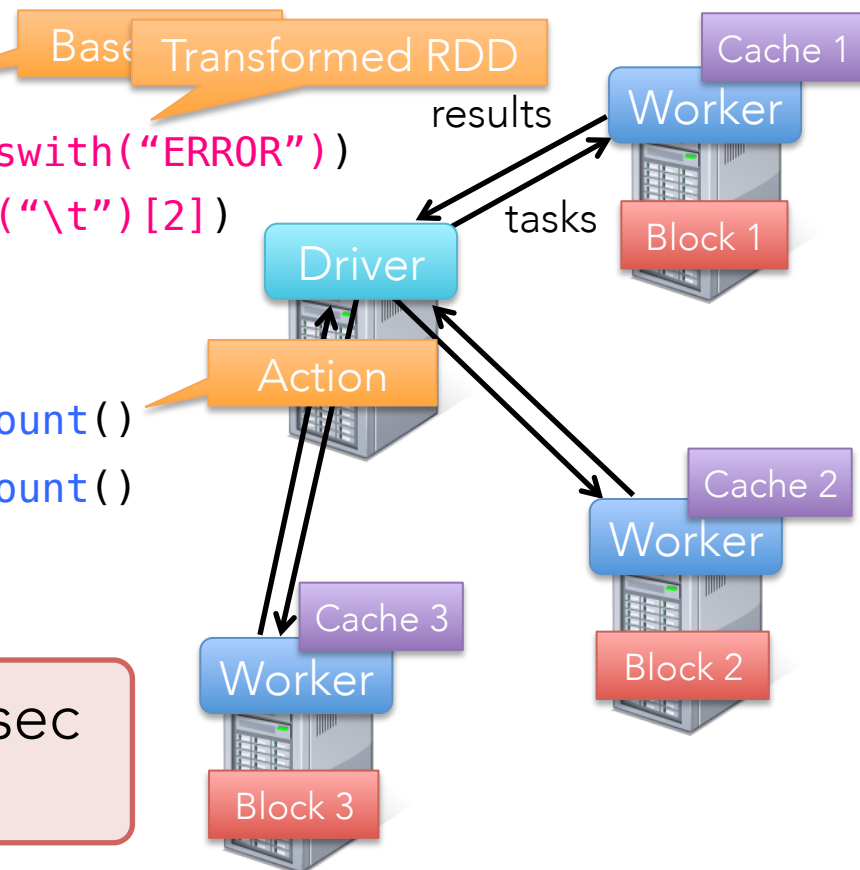- » Automatically rebuilt on failure

# Example: Log Mining

Load error messages from a log into memory, then interactively search for various patterns



```
lines = spark.textFile("hdfs://...")
errors = lines.filter(lambda s: s.startswith("ERROR"))
messages = errors.map(lambda s: s.split("\t")[2])
messages.cache()


messages.filter(lambda s: "foo" in s).count()
messages.filter(lambda s: "bar" in s).count()
. . .
```

Base

Transformed RDD

Action

results

tasks

Driver

Worker — Cache 1 — Block 1

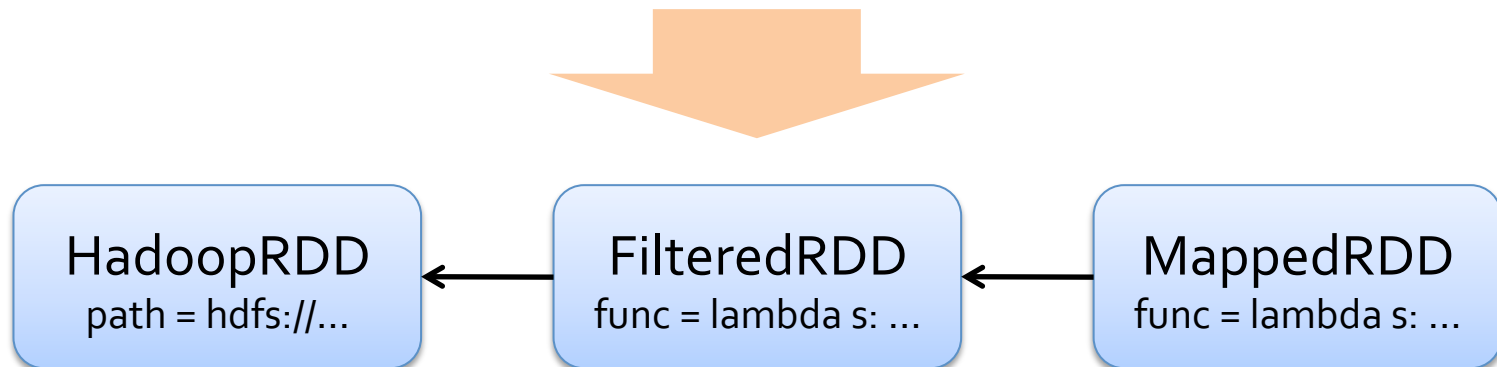Worker — Cache 2 — Block 2

Worker — Cache 3 — Block 3

**Result:** scaled to 1 TB data in 7 sec
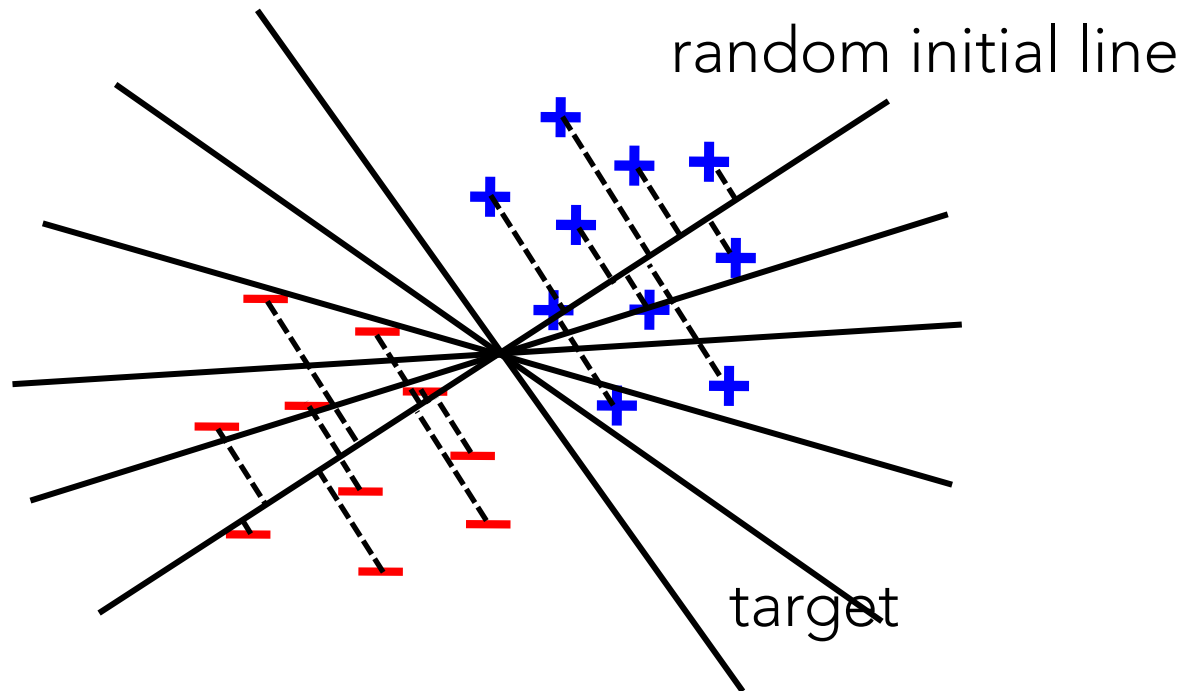(vs 180 sec for on-disk data)

# Fault Tolerance

RDDs track the transformations used to build them (their *lineage*) to recompute lost data

```
messages = textFile(...).filter(lambda s: "ERROR" in s)
                        .map(lambda s: s.split("\t")[2])
```

| HadoopRDD | ← | FilteredRDD | ← | MappedRDD |
|---|---|---|---|---|
| path = hdfs://... | | func = lambda s: ... | | func = lambda s: ... |

# Example: Logistic Regression

Goal: find line separating two sets of points



random initial line

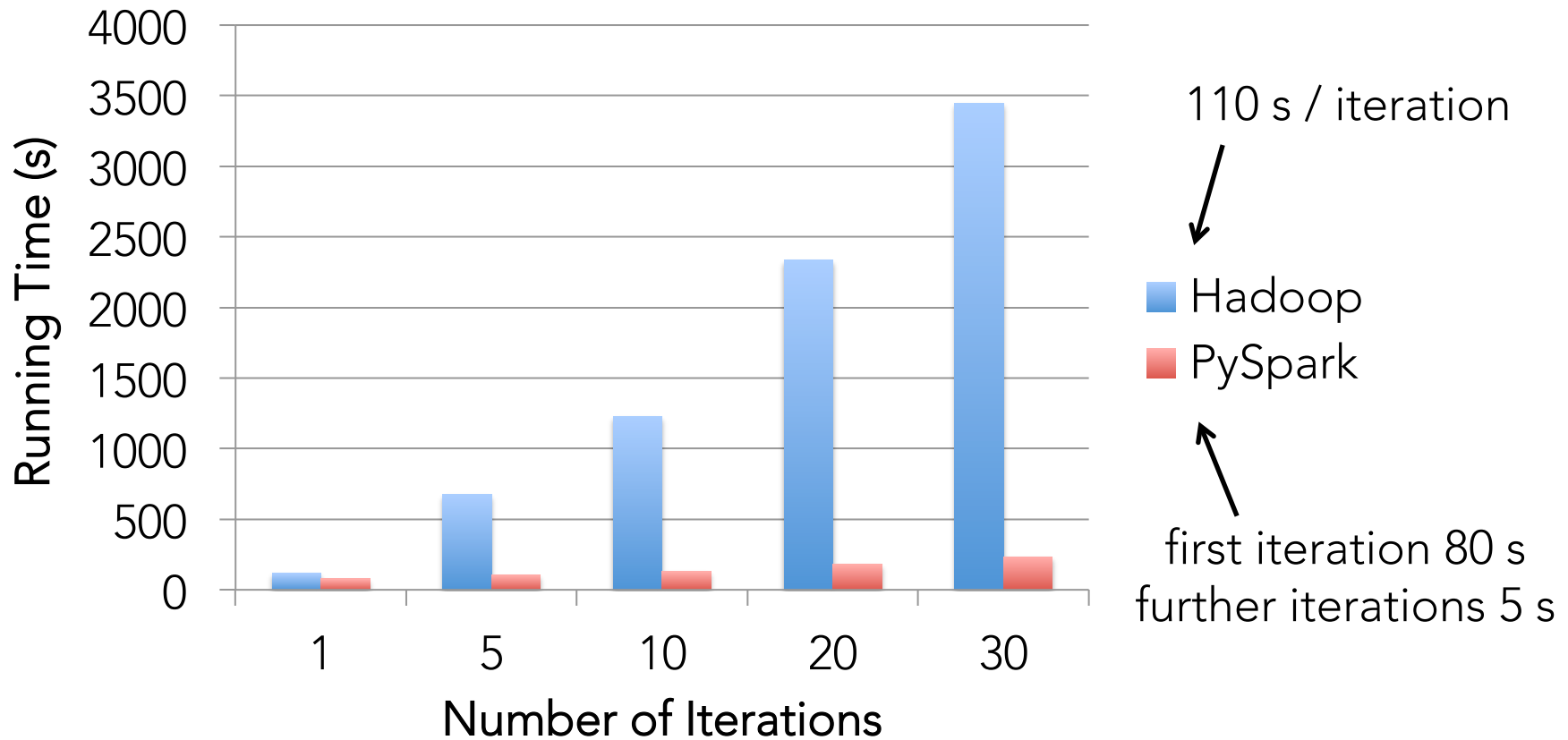target

# Example: Logistic Regression

```python
data = spark.textFile(...).map(readPoint).cache()

w = numpy.random.rand(D)

for i in range(iterations):
    gradient = data.map(lambda p:
        (1 / (1 + exp(-p.y * w.dot(p.x)))) * p.y * p.x
    ).reduce(lambda x, y: x + y)
    w -= gradient

print "Final w: %s" % w
```

# Logistic Regression Performance



**Running Time (s)** vs **Number of Iterations**

Legend: Hadoop, PySpark

110 s / iteration

first iteration 80 s
further iterations 5 s

# Demo

# Supported Operators

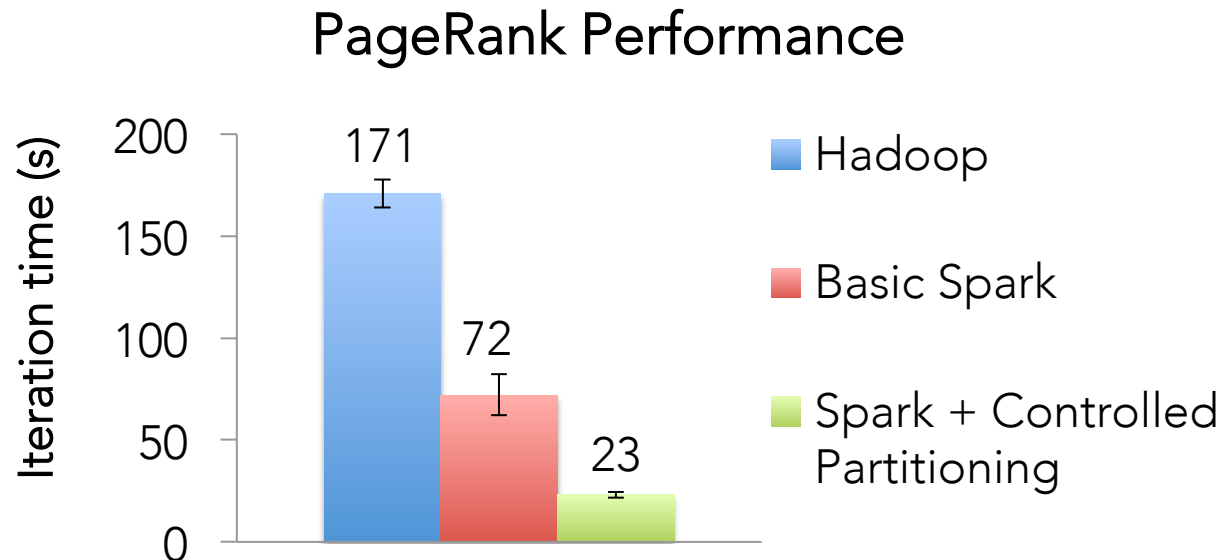| | | |
|---|---|---|
| map | reduce | take |
| filter | count | first |
| groupBy | fold | partitionBy |
| union | reduceByKey | pipe |
| join | groupByKey | distinct |
| leftOuterJoin | cogroup | save |
| rightOuterJoin | flatMap | ... |

# Other Engine Features

General operator graphs (not just map-reduce)

Hash-based reduces (faster than Hadoop's sort)

Controlled data partitioning to save communication

## PageRank Performance

# Spark Community



1000+ meetup members

60+ contributors

17 companies contributing

# This Talk

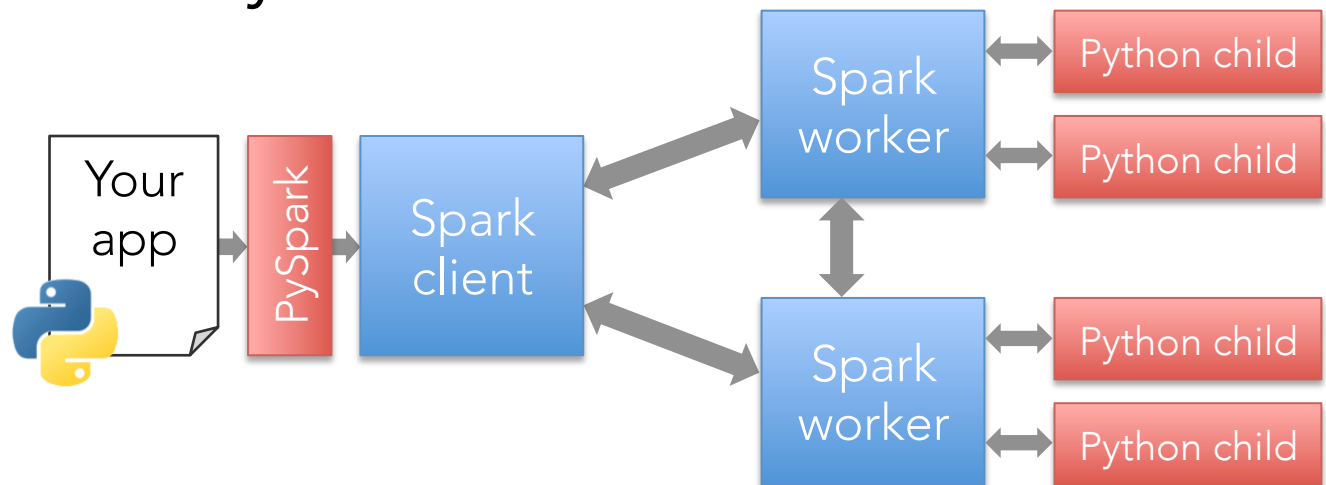Spark programming model

Examples

Demo

Implementation

Trying it out

# Overview

Spark core is written in Scala

PySpark calls existing scheduler, cache and networking layer (2K-line wrapper)

No changes to Python

# Overview

Spark core is written in Scala

PySpark calls existing scheduler, cache and networking layer (2K-line wrapper)

No changes to Python

Python child

Python child

Main PySpark author:

Josh Rosen

Python child

cs.berkeley.edu/~joshrosen

Python child

# Object Marshaling

Uses pickle library for both communication and cached data
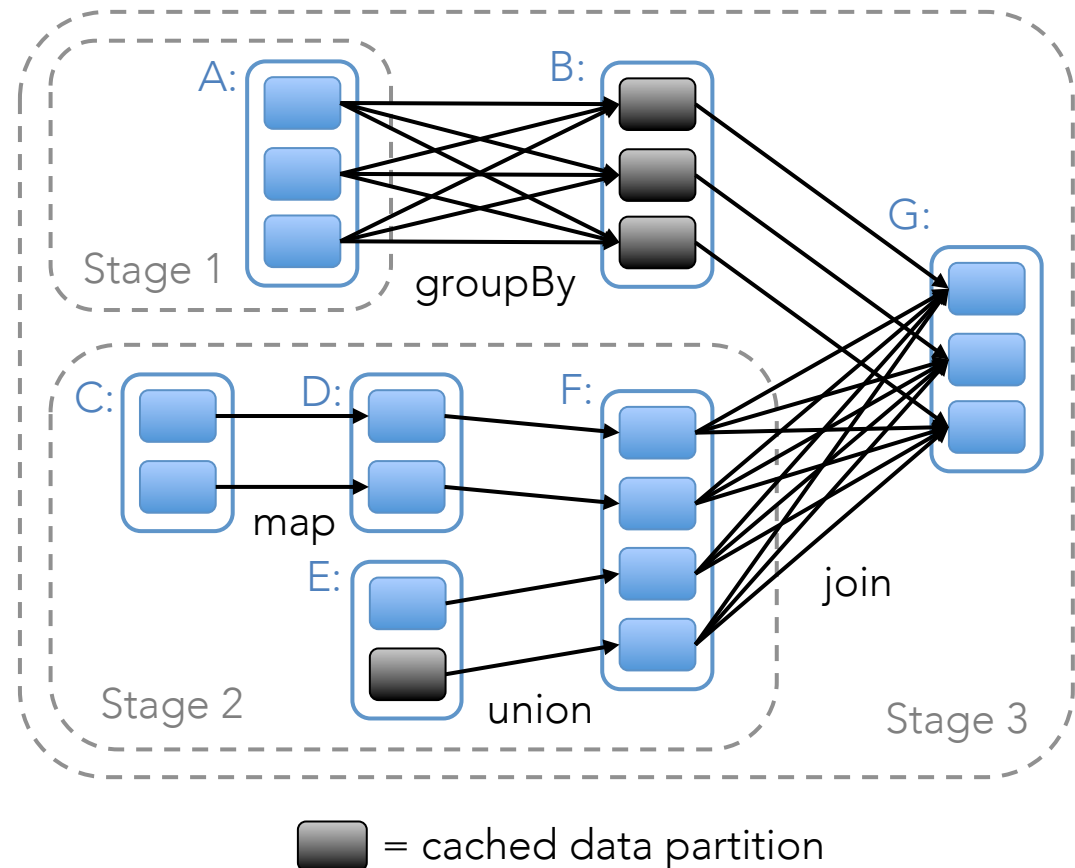  » Much cheaper than Python objects in RAM

Lambda marshaling library by PiCloud

# Job Scheduler

Supports general operator graphs

Automatically pipelines functions

Aware of data locality and partitioning



= cached data partition

# Interoperability

Runs in standard CPython, on Linux / Mac
  » Works fine with extensions, e.g. NumPy

Input from local file system, NFS, HDFS, S3
  » Only text files for now

Works in IPython, including notebook

Works in doctests – see our tests!

# Getting Started

Visit spark-project.org for video tutorials, online exercises, docs

Easy to run in local mode (multicore), standalone clusters, or EC2

Training camp at Berkeley in August (free video): ampcamp.berkeley.edu

# Getting Started

Easiest way to learn is the shell:

```
$ ./pyspark

>>> nums = sc.parallelize([1,2,3]) # make RDD from array

>>> nums.count()
3

>>> nums.map(lambda x: 2 * x).collect()
[2, 4, 6]
```

# Writing Standalone Jobs

```python
from pyspark import SparkContext

if __name__ == "__main__":
    sc = SparkContext("local", "WordCount")
    lines = sc.textFile("in.txt")

    counts = lines.flatMap(lambda s: s.split()) \
                  .map(lambda word: (word, 1)) \
                  .reduceByKey(lambda x, y: x + y)

    counts.saveAsTextFile("out.txt")
```

# Conclusion

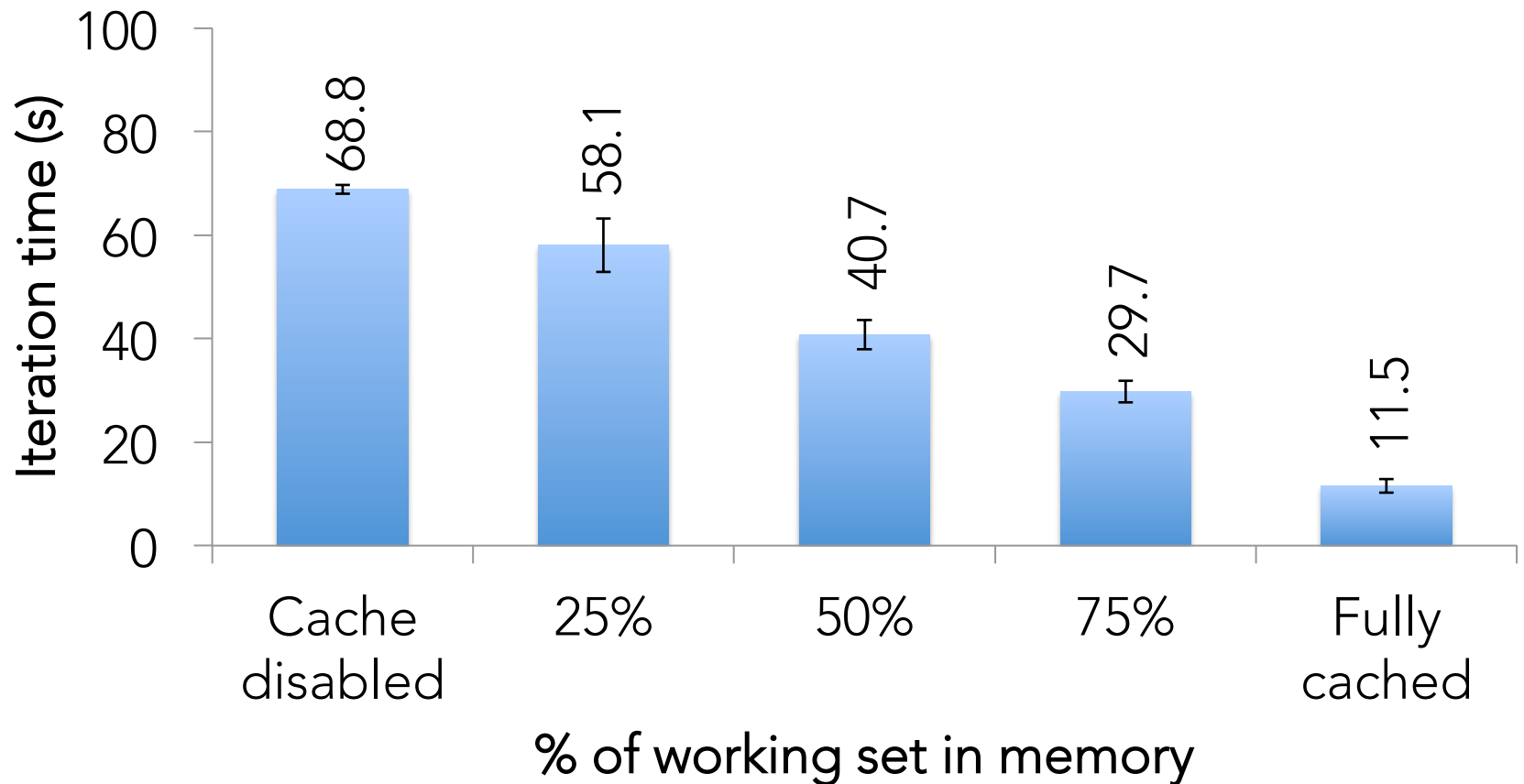PySpark provides a fast and simple way to analyze big datasets from Python

Learn more or contribute at [spark-project.org](spark-project.org)

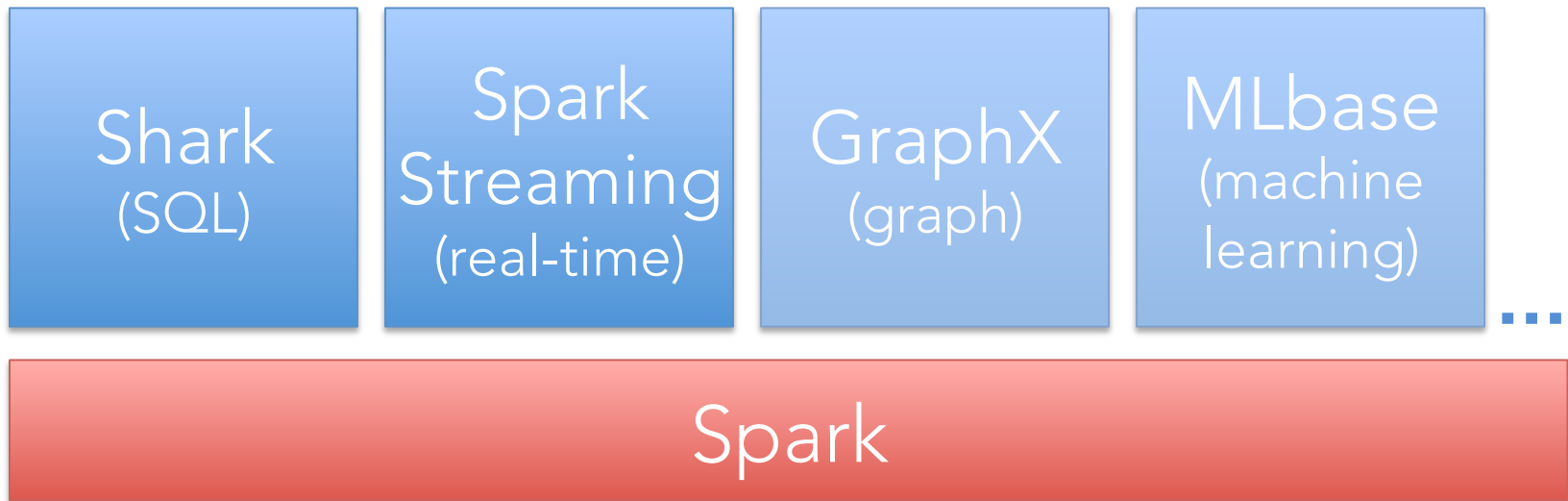> Look for our training camp
> on August 29-30!



My email: matei@berkeley.edu

# Behavior with Not Enough RAM

# The Rest of the Stack

Spark is the foundation for wide set of projects in the Berkeley Data Analytics Stack (BDAS)



More details: amplab.berkeley.edu

# Performance Comparison