# Introduction to Apache Spark

Patrick Wendell - Databricks

# What is Spark?

Fast and **Expressive** Cluster Computing Engine Compatible with Apache Hadoop

*Up to* **10x** *faster on disk,* **100x** *in memory*

## Efficient

- General execution graphs
- In-memory storage

**2-5x** *less code*

## Usable

- Rich APIs in Java, Scala, Python
- Interactive shell

**Spark**

# Spark Programming Model

# Key Concept: RDD's

Write programs in terms of **operations** on distributed datasets

## Resilient Distributed Datasets

- Collections of objects spread across a cluster, stored in RAM or on Disk

- Built through parallel transformations

- Automatically rebuilt on failure

## Operations

- Transformations (e.g. map, filter, groupBy)

- Actions (e.g. count, collect, save)

Spark

# **Example**: Log Mining

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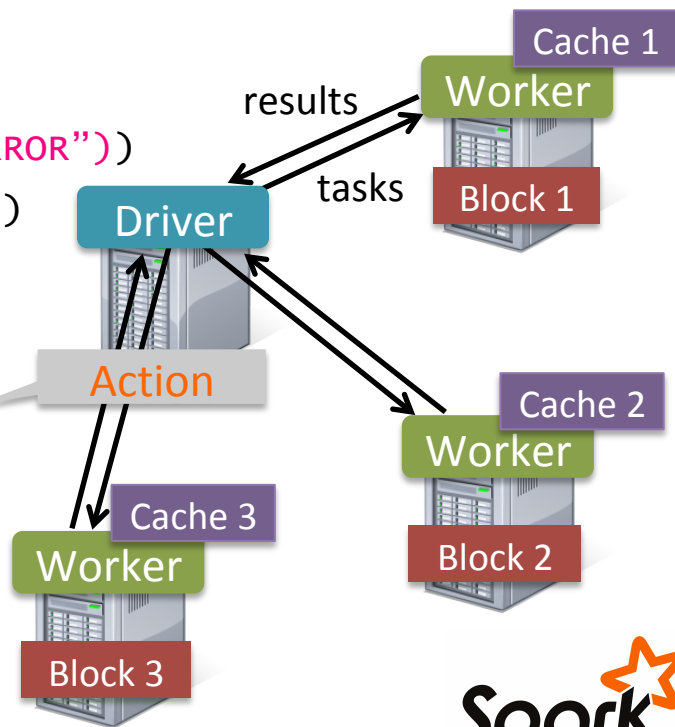


Transformed RDD

```
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: "mysql" in s).count()
messages.filter(lambda s: "php" in s).count()

. . .
```

Action

Driver

results

tasks

Worker
Cache 1
Block 1

Worker
Cache 2
Block 2

Worker
Cache 3
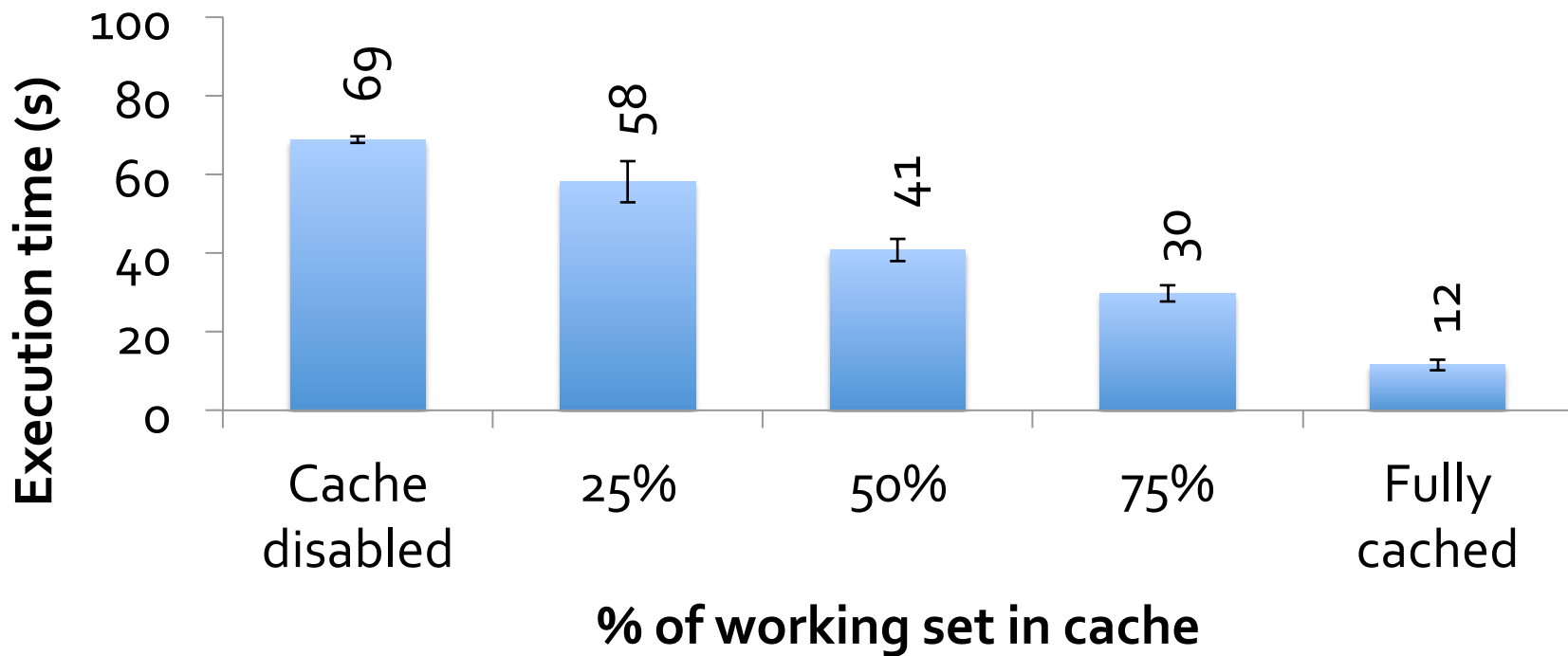Block 3

**Full-text search of Wikipedia**
- 60GB on 20 EC2 machine
- 0.5 sec vs. 20s for on-disk

*Spark*

# More RDD Operators

- map
- filter
- groupBy
- sort
- union
- join
- leftOuterJoin
- rightOuterJoin

- reduce
- count
- fold
- reduceByKey
- groupByKey
- cogroup
- cross
- zip

sample

take

first

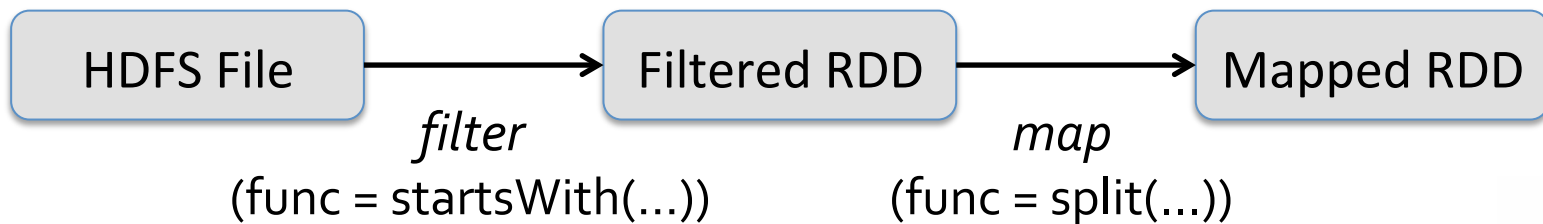partitionBy

mapWith

pipe

save     ...

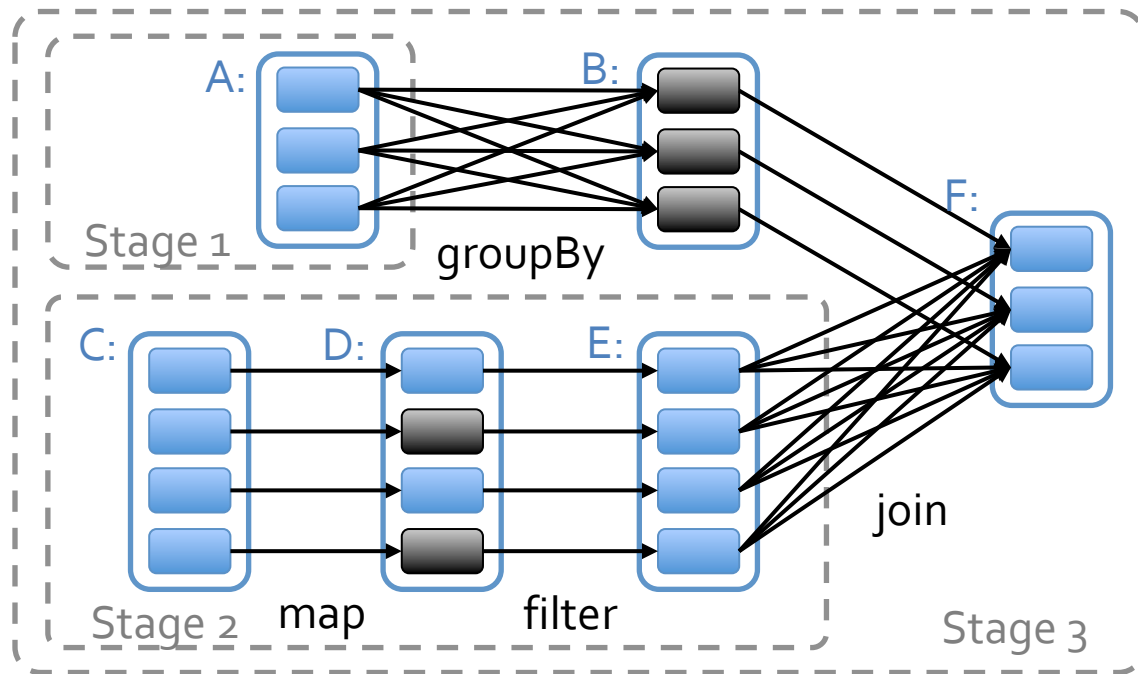# Scaling Down

# Fault Recovery

RDDs track *lineage* information that can be used to efficiently recompute lost data

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

HDFS File → *filter* (func = startsWith(…)) → Filtered RDD → *map* (func = split(…)) → Mapped RDD

# Under The Hood: DAG Scheduler

- General task graphs
- Automatically pipelines functions
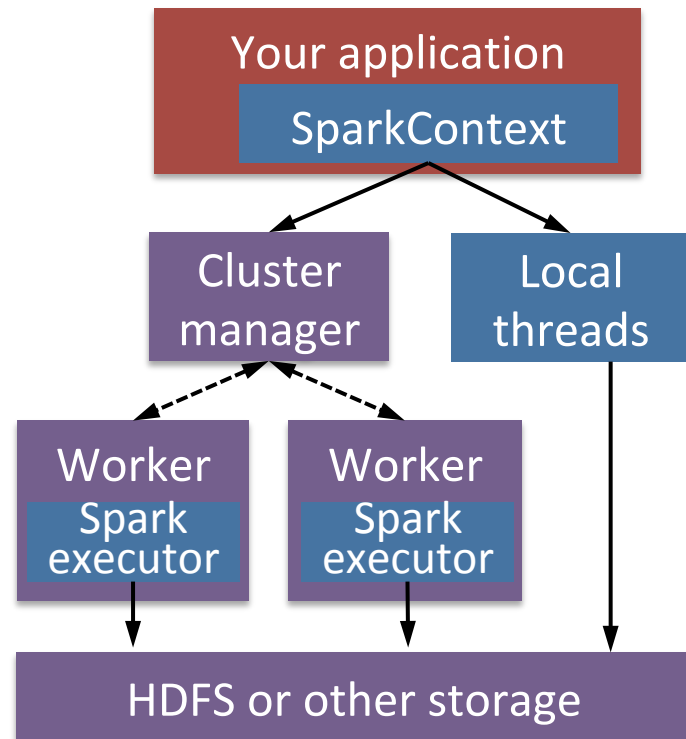- Data locality aware
- Partitioning aware to avoid shuffles

# Software Components

- Spark runs as a library in your program (1 instance per app)
- Runs tasks locally or on cluster
  – Mesos, YARN or standalone mode
- Accesses storage systems via Hadoop InputFormat API
  – Can use HBase, HDFS, S3, …



Your application
SparkContext

Cluster manager

Local threads

Worker
Spark executor

Worker
Spark executor

HDFS or other storage

Spark

# CONCLUSION

# Conclusion

- Spark offers a rich API to make data analytics *fast*: both fast to write and fast to run

- Achieves 100x speedups in real applications

- Growing community with 25+ companies contributing

**Spark**