# Making Big Data Processing Simple with Spark

Matei Zaharia
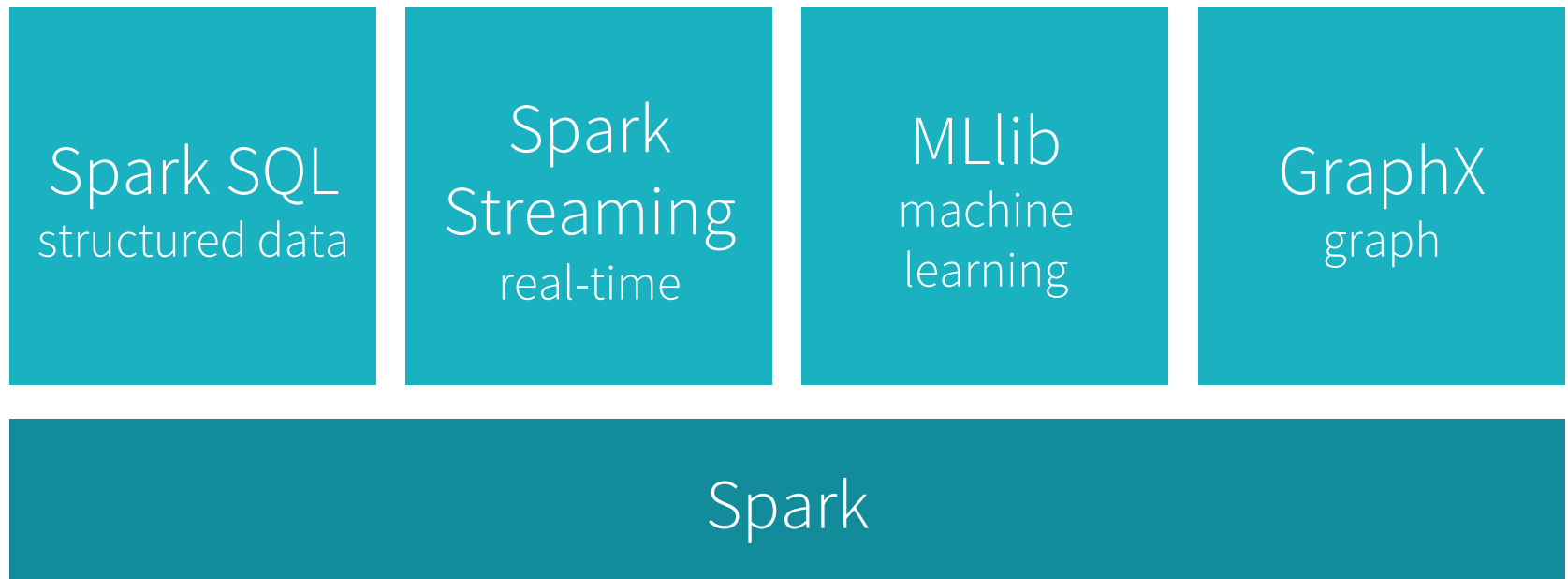
December 17, 2015

databricks™

MIT

# What is Apache Spark?

Fast and general cluster computing engine that generalizes the MapReduce model

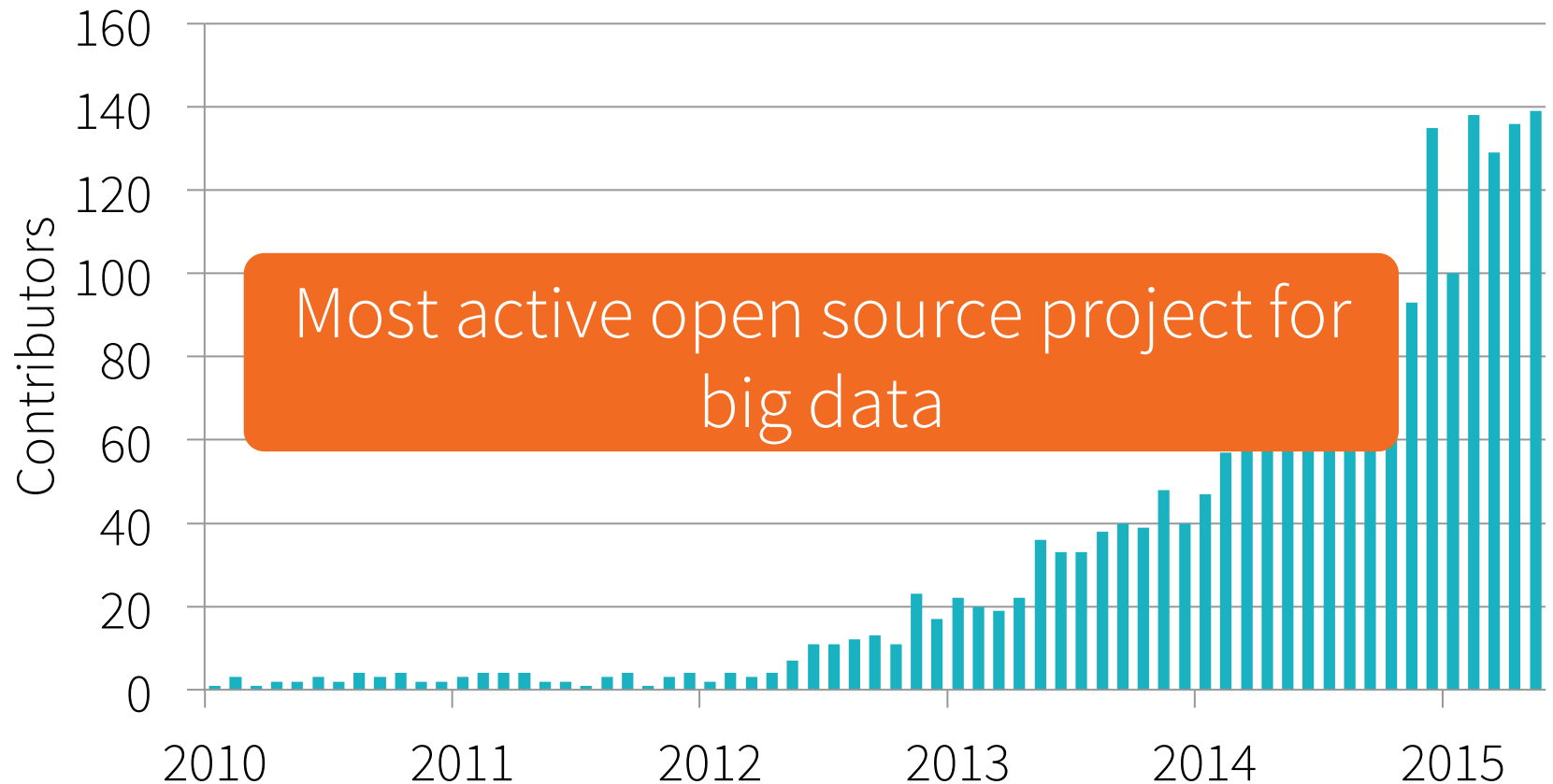Makes it easy and fast to process large datasets
- High-level APIs in Java, Scala, Python, R
- Unified engine that can capture many workloads

databricks™

# A Unified Engine

# A Large Community

## Contributors / Month to Spark



Most active open source project for big data

databricks™

# Overview

Why a unified engine?

Spark programming model

Built-in libraries

Applications

databricks™

# History: Cluster Computing

2004

## MapReduce: Simplified Data Processing on Large Clusters

Jeffrey Dean and Sanjay Ghemawat

jeff@google.com, sanjay@google.com

*Google, Inc.*

### Abstract

MapReduce is a programming model and an associated implementation for processing and generating large data sets. Users specify a *map* function that processes a key/value pair to generate a set of intermediate key/value pairs, and a *reduce* function that merges all intermediate values associated with the same intermediate key. Many

given day, etc. Most such computations are conceptually straightforward. However, the input data is usually large and the computations have to be distributed across hundreds or thousands of machines in order to finish in a reasonable amount of time. The issues of how to parallelize the computation, distribute the data, and handle failures conspire to obscure the original simple computation with large amounts of complex code to deal with

databricks™

# MapReduce

A *general* engine for batch processing

> We wrote the first version of the MapReduce library in February of 2003, and made significant enhancements to it in August of 2003, including the locality optimization, dynamic load balancing of task execution across worker machines, etc. Since that time, we have been pleasantly surprised at how broadly applicable the MapReduce library has been for the kinds of problems we work on. It has been used across a wide range of domains within Google, including:
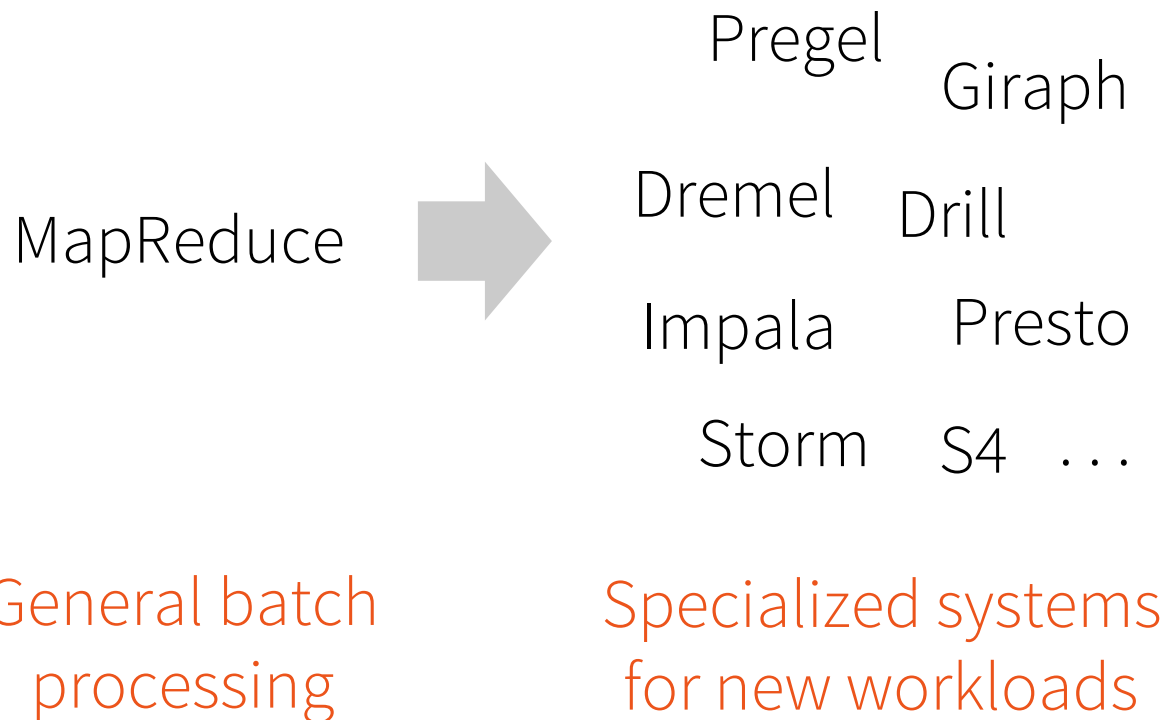
# Beyond MapReduce

MapReduce was great for batch processing, but users quickly needed to do more:

- More complex, multi-pass algorithms
- More interactive ad-hoc queries
- More real-time stream processing

Result: *specialized* systems for these workloads

databricks™

# Big Data Systems Today

MapReduce →

Pregel
Giraph
Dremel
Drill
Impala
Presto
Storm
S4   …

General batch
processing

Specialized systems
for new workloads

databricks™
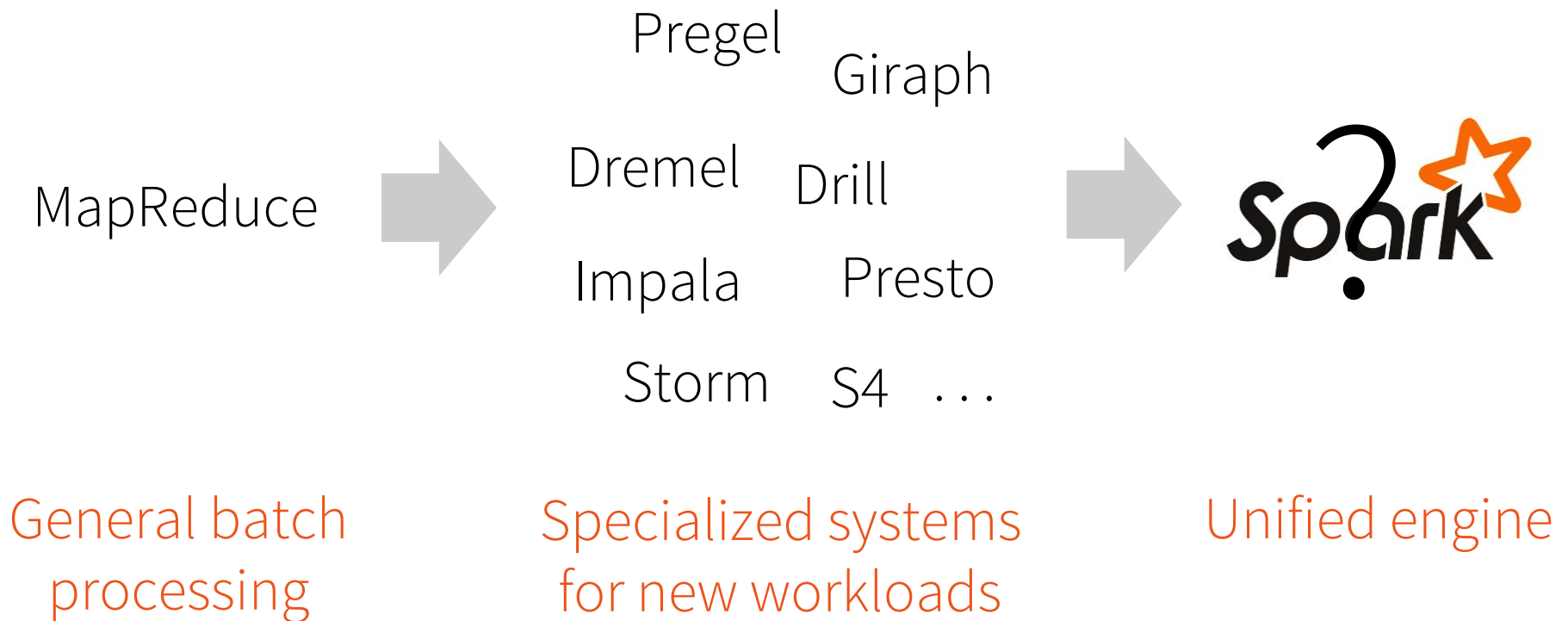
# Problems with Specialized Systems

More systems to manage, tune, deploy

Can't easily *combine* processing types
- Even though most applications need to do this!
- E.g. load data with SQL, then run machine learning

In many cases, data transfer between engines is a dominant cost!

databricks™

# Big Data Systems Today

MapReduce ➡ Pregel  Giraph
Dremel  Drill
Impala  Presto
Storm  S4  …
➡ **Spark**

General batch
processing

Specialized systems
for new workloads

Unified engine

databricks™

# Overview

Why a unified engine?

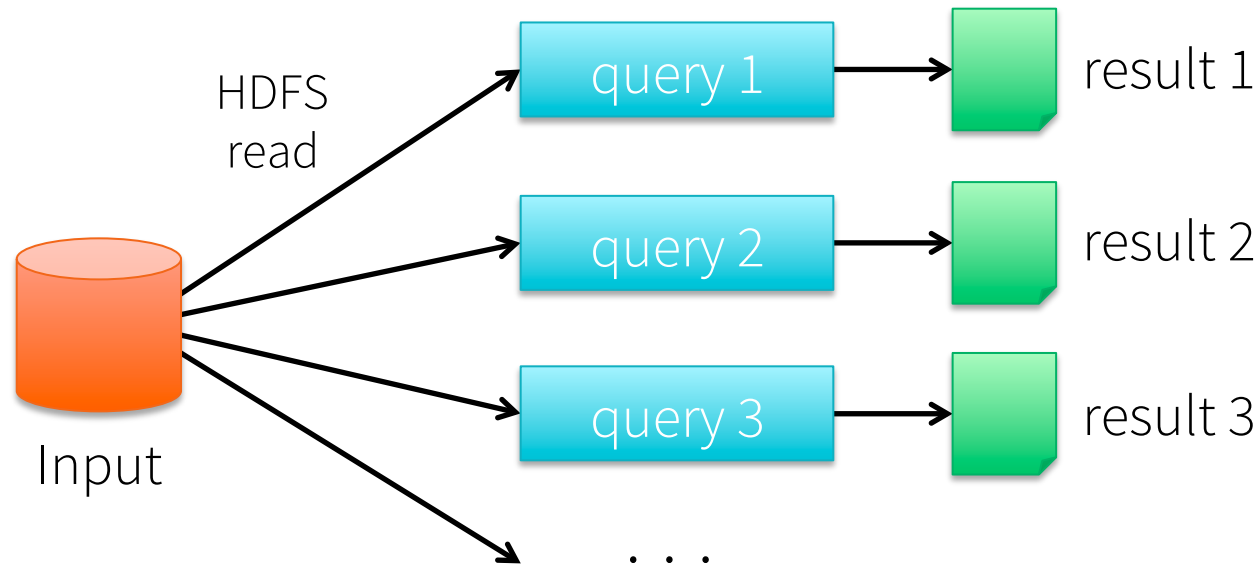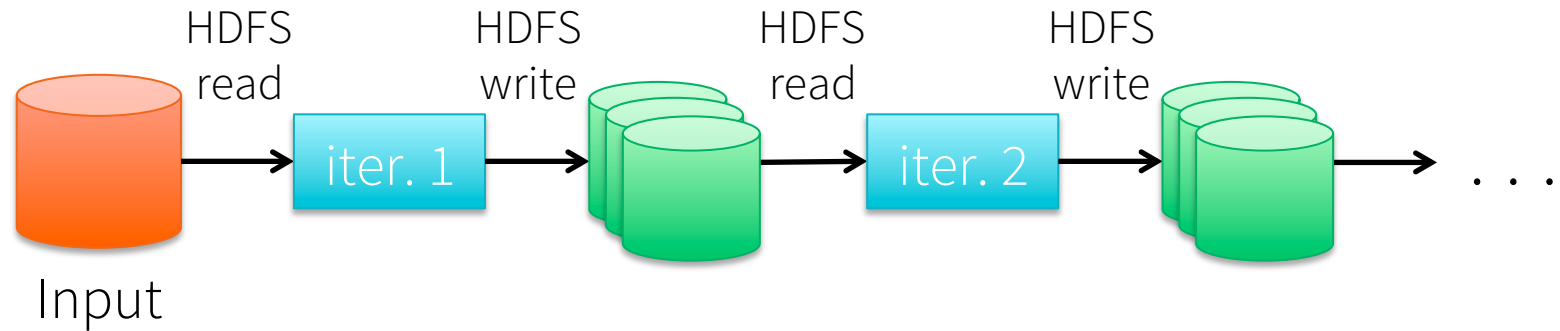Spark programming model

Built-in libraries

Applications

databricks™

# Background

Recall 3 workloads were issues for MapReduce:

- More complex, multi-pass algorithms
- More interactive ad-hoc queries
- More real-time stream processing

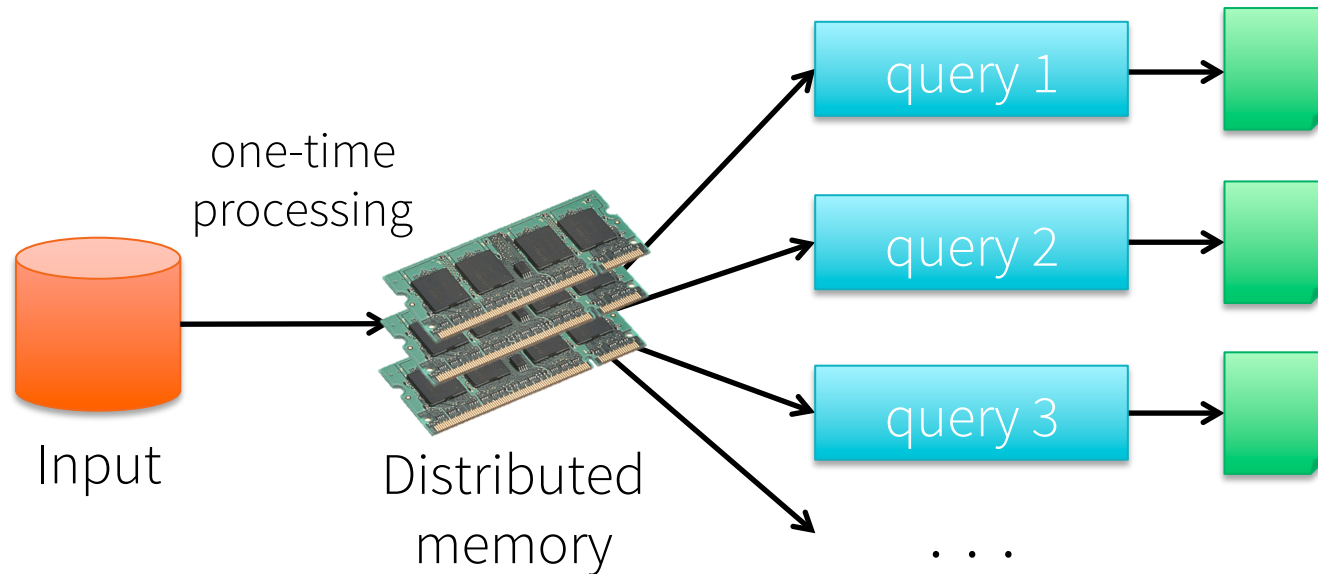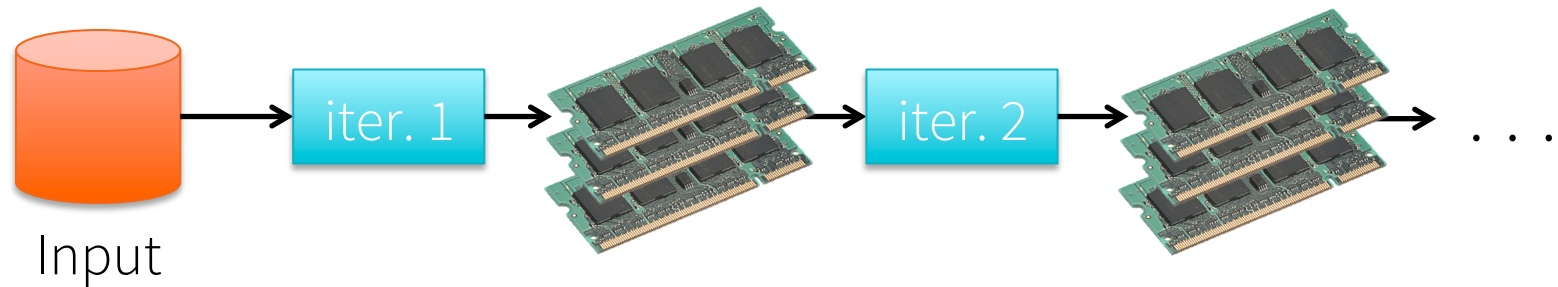While these look different, all 3 need one thing that MapReduce lacks: efficient data sharing

databricks™

# Data Sharing in MapReduce



Slow due to replication and disk I/O

# What We'd Like



Input → iter. 1 → iter. 2 → . . .

Input → one-time processing → Distributed memory → query 1 → 
Distributed memory → query 2 → 
Distributed memory → query 3 → 
. . .

**10-100x faster than network and disk**

databri

# Spark Programming Model

Resilient Distributed Datasets (RDDs)

- Collections of objects stored in RAM or disk across cluster
- Built via parallel transformations (map, filter, …)
- Automatically rebuilt on failure
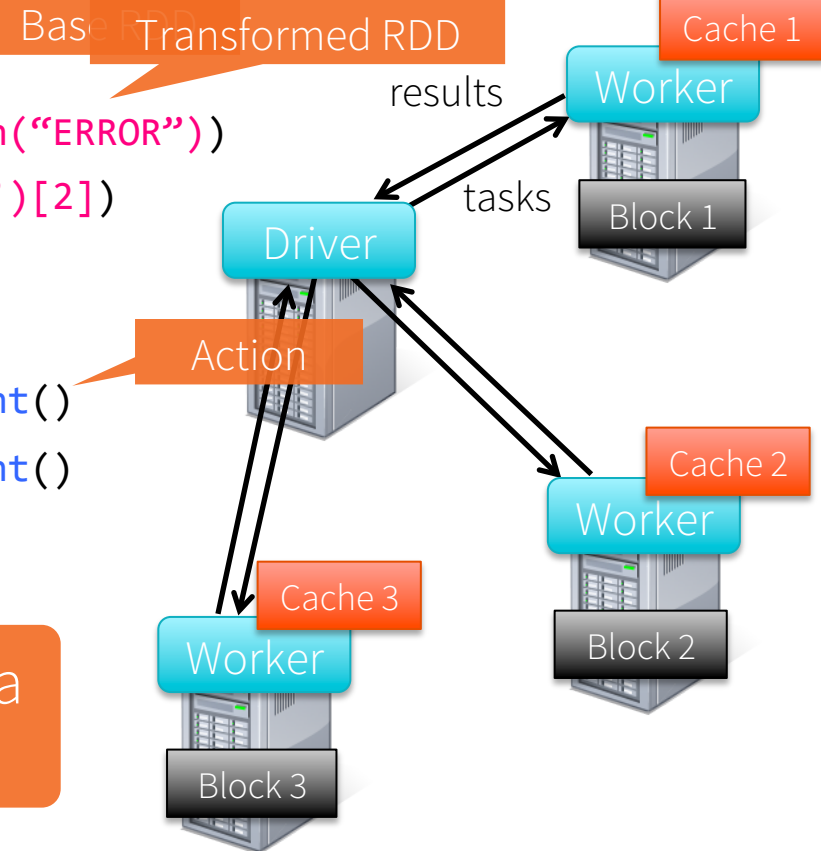
databricks™

# 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: "MySQL" in s).count()
messages.filter(lambda s: "Redis" in s).count()
. . .
```
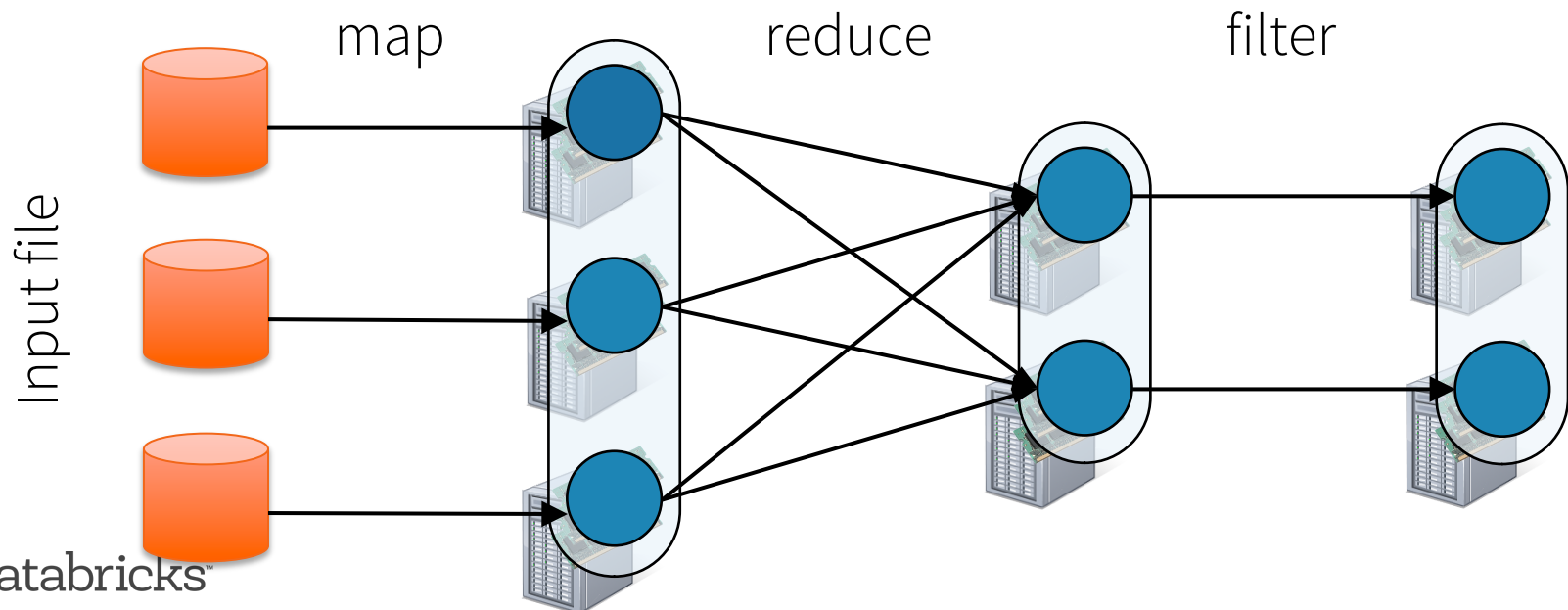
Base RDD

Transformed RDD

results

tasks

Action

Cache 1

Worker

Block 1

Driver

Cache 2

Worker

Block 2

Cache 3

Worker

Block 3

Example: full-text search of Wikipedia in 0.5 sec (vs 20s for on-disk data)

databricks™

# Fault Tolerance

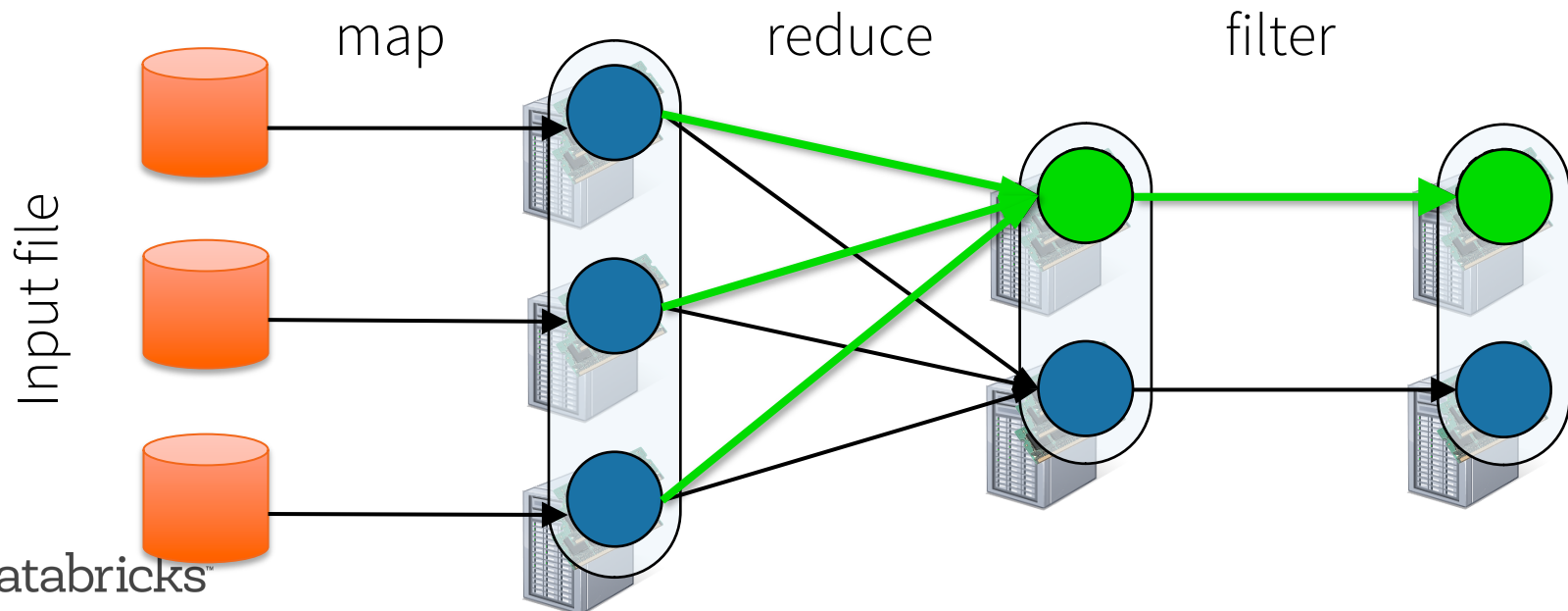RDDs track *lineage* info to rebuild lost data

```
file.map(lambda rec: (rec.type, 1))
    .reduceByKey(lambda x, y: x + y)
    .filter(lambda (type, count): count > 10)
```
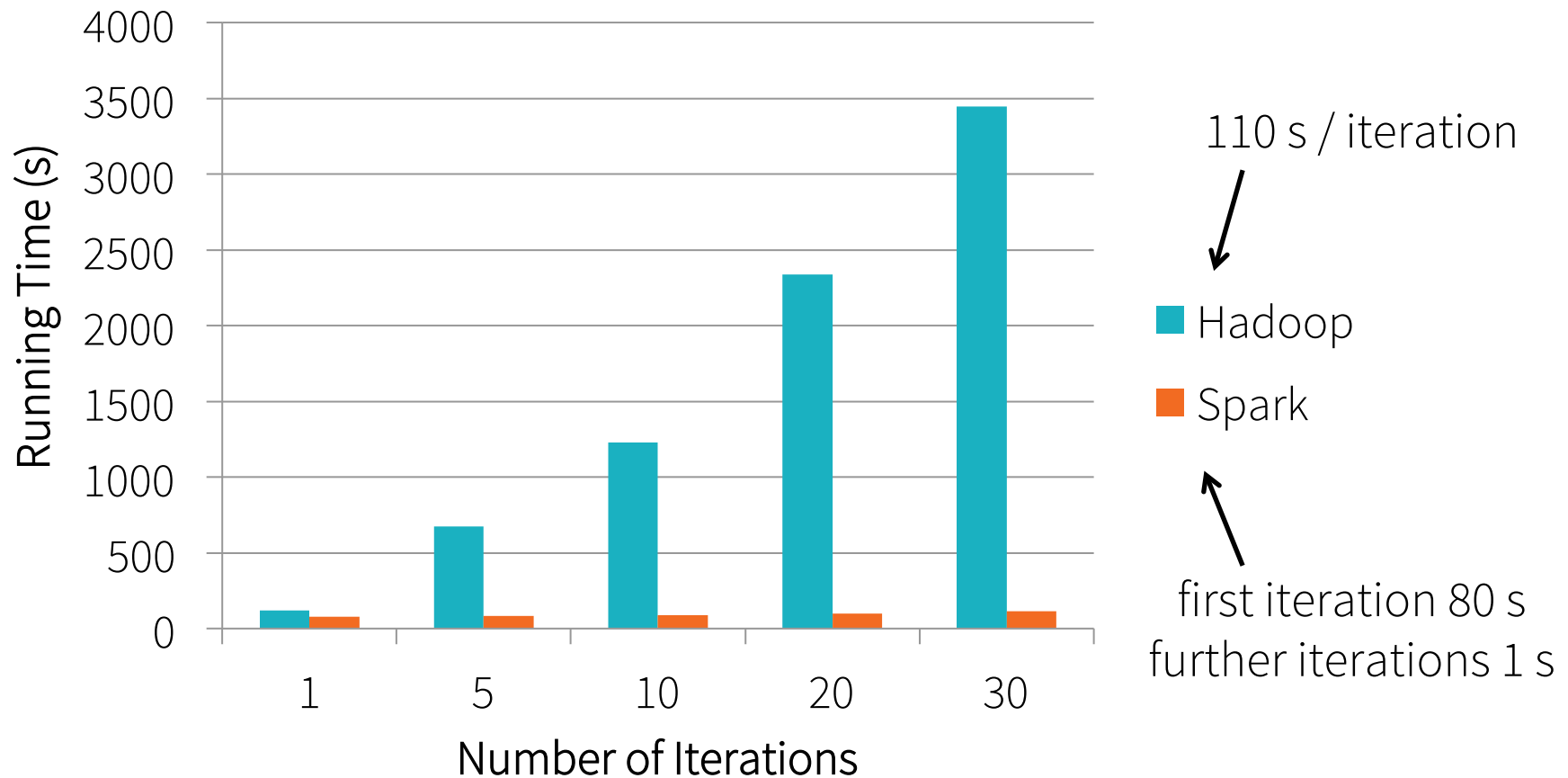
# Fault Tolerance

RDDs track *lineage* info to rebuild lost data

```
file.map(lambda rec: (rec.type, 1))
    .reduceByKey(lambda x, y: x + y)
    .filter(lambda (type, count): count > 10)
```

# Example: Logistic Regression



110 s / iteration

Hadoop

Spark

first iteration 80 s
further iterations 1 s

databricks™

# On-Disk Performance

Time to sort 100TB

2013 Record:
Hadoop

2100 machines

72 minutes

2014 Record:
Spark

207 machines

23 minutes

# Libraries Built on Spark



Spark SQL
structured data

Spark Streaming
real-time

MLlib
machine learning

GraphX
graph

Spark

databricks™

# Combining Processing Types

```python
// Load data using SQL
points = ctx.sql("select latitude, longitude from tweets")

// Train a machine learning model
model = KMeans.train(points, 10)

// Apply it to a stream
sc.twitterStream(...)
  .map(lambda t: (model.predict(t.location), 1))
  .reduceByWindow("5s", lambda a, b: a + b)
```
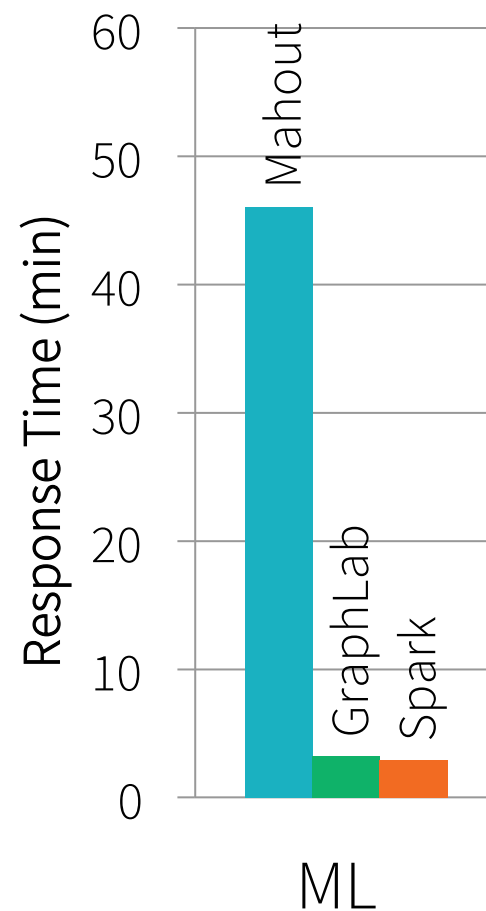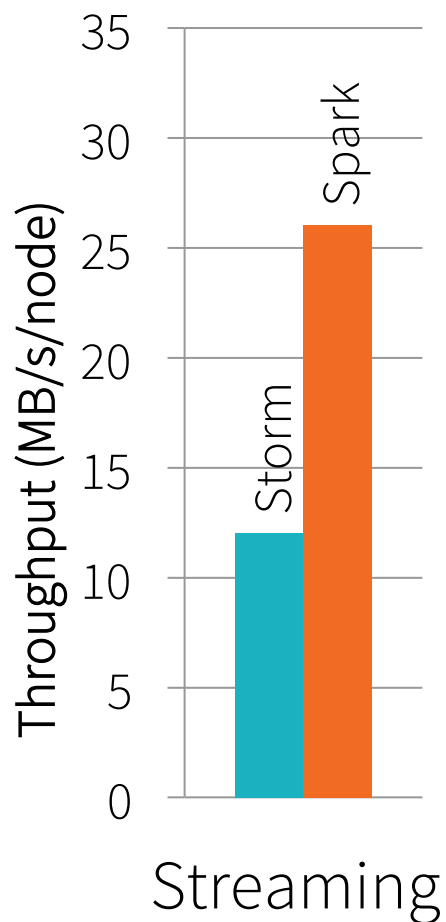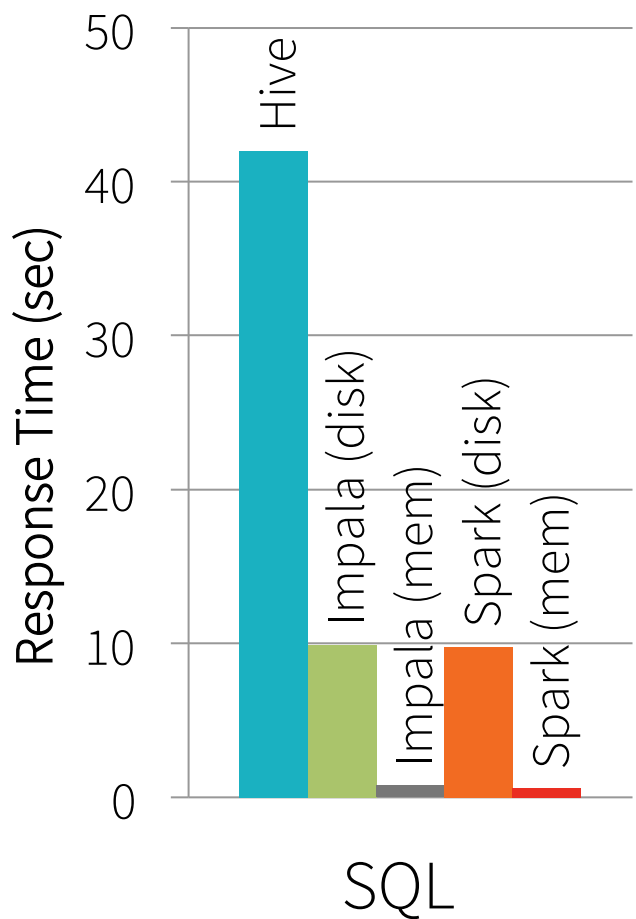
databricks™

# Combining Processing Types

Separate systems:

| HDFS read | ETL | HDFS write |   | HDFS read | train | HDFS write |   | HDFS read | query | HDFS write | ... |

Spark:

| HDFS read | ETL | train | query | HDFS write |

databricks™

# Performance vs Specialized Systems

**SQL** — Response Time (sec)
- Hive
- Impala (disk)
- Impala (mem)
- Spark (disk)
- Spark (mem)

**Streaming** — Throughput (MB/s/node)
- Storm
- Spark

**ML** — Response Time (min)
- Mahout
- GraphLab
- Spark

databricks™

# Some Recent Additions

DataFrame API (similar to R and Pandas)

- Easy programmatic way to work with structured data

R interface (SparkR)

Machine learning pipelines (like SciKit-learn)

databricks™

# Overview

Why a unified engine?

Spark programming model

Built-in libraries

Applications

databricks™

# Spark Community

Over 1000 deployments, clusters up to 8000 nodes

# Top Applications



| Application | Percentage |
|---|---|
| Business Intelligence | 68% |
| Data Warehousing | 52% |
| Recommendation | 44% |
| Log Processing | 40% |
| User-Facing Services | 36% |
| Faud Detection / Security | 29% |

# Spark Components Used

| Component | Percentage |
|---|---|
| Spark SQL | 69% |
| DataFrames | 62% |
| Spark Streaming | 58% |
| MLlib + GraphX | 58% |

**75%**

of users use more than one component

databricks™

# Learn More

Get started on your laptop: [spark.apache.org](spark.apache.org)

Resources and MOOCs: [sparkhub.databricks.com](sparkhub.databricks.com)

Spark Summit: [spark-summit.org](spark-summit.org)

Thank You