

Big Data Analytics with Apache **Spark**



Matei Zaharia



Stanford



databricks

Outline

The big data problem

MapReduce

Apache Spark

How people are using Spark

The Big Data Problem

Data is growing faster than processor speeds

Growing data sources

- » Mostly machine generated

Cheap storage

Stalling CPU speeds



Examples

Facebook's daily event data: 4 PB

Google's web index: 100 PB

1000 genomes project: 200 TB

Cost of 1 TB of disk: \$16

Time to read 1 TB from disk: 6 hours (50 MB/s)

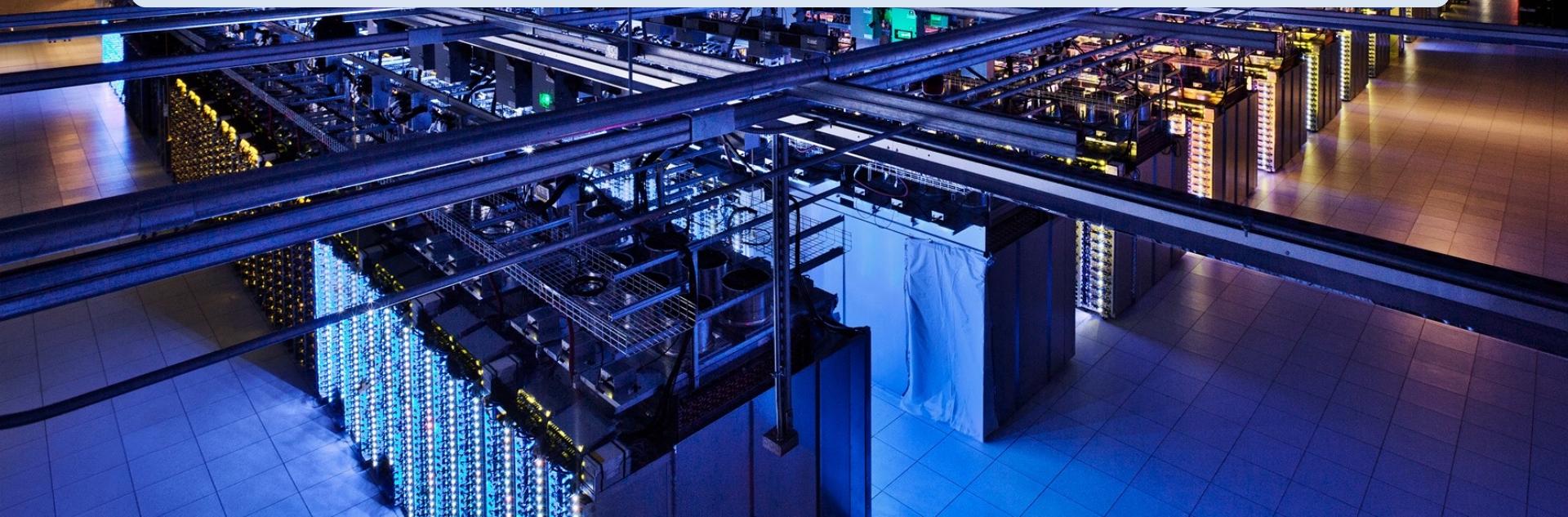
The Big Data Problem

Single machine can no longer process or even store all the data!

Only solution is to **distribute** over large clusters

Google Datacenter

How do we program these?



Traditional Network Programming

Message-passing between nodes

Really hard to do at scale:

- » How to split problem across nodes?
- » How to deal with failures?
- » Even worse: stragglers (node is not failed, but slow)

Data-Parallel Models

Restrict the programming interface so that the system can do more automatically

“Here’s an operation, run it on all of the data”

- » I don’t care *where* it runs (you schedule that)
- » In fact, feel free to run it *twice* on different nodes

Early example: MapReduce

MapReduce

First widely popular programming model for
data-intensive apps on clusters

Published by Google in 2004

Popularized by open-source Hadoop project

MapReduce Programming Model

Data type: key-value records

Map function:

$$(K_{in}, V_{in}) \rightarrow \text{list}(K_{inter}, V_{inter})$$

Reduce function:

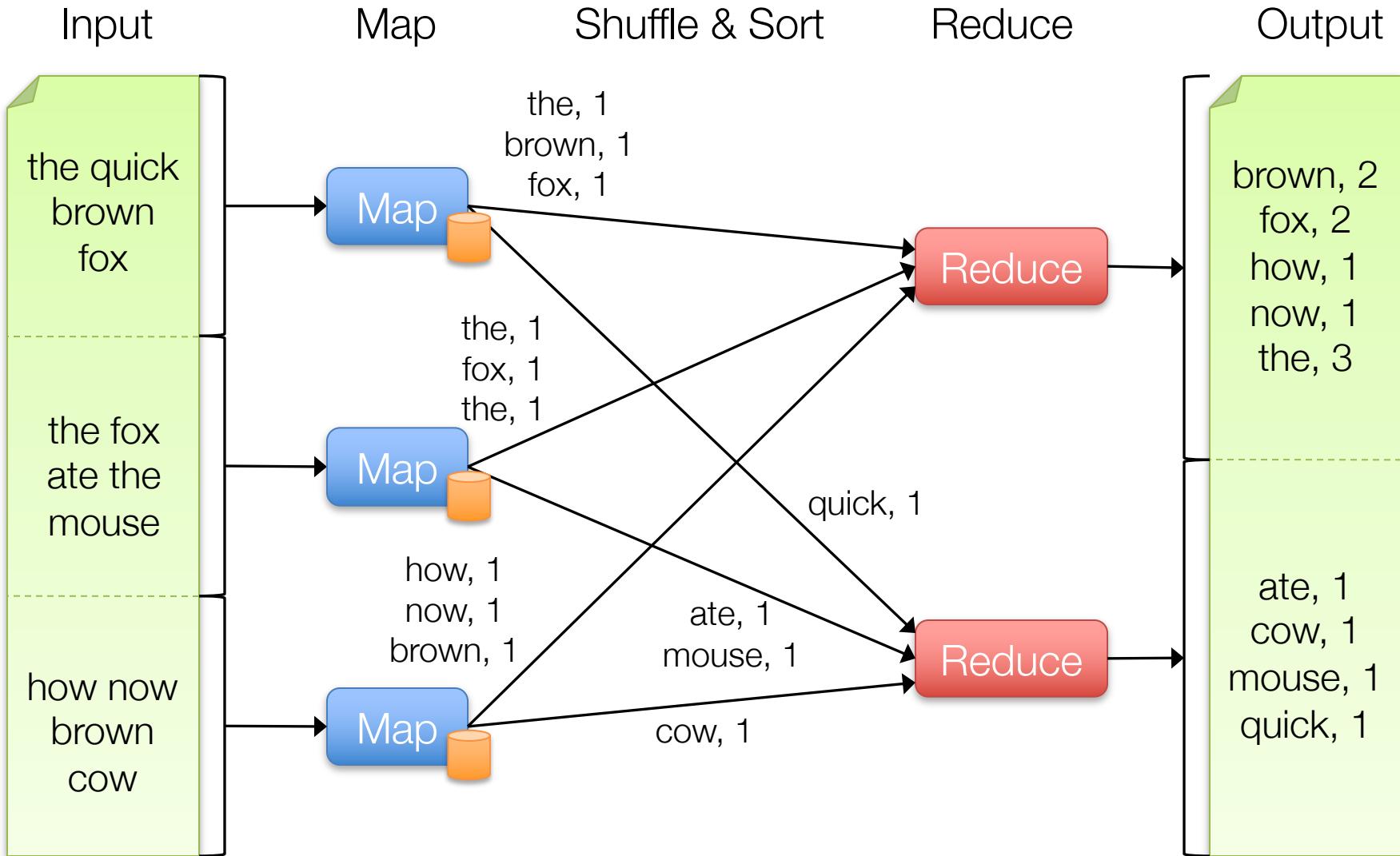
$$(K_{inter}, \text{list}(V_{inter})) \rightarrow \text{list}(K_{out}, V_{out})$$

Example: Word Count

```
def map(line):  
    foreach word in line.split():  
        output(word, 1)
```

```
def reduce(key, values):  
    output(key, sum(values))
```

Word Count Execution



Other MapReduce Apps

How would you write these with map & reduce?

Search: find all the records in a file that contain the word “Berkeley”

Top words: Find the top 10 most common words

MapReduce Execution

Mappers are scheduled on same node as their input block if possible

- » Minimize network use to improve performance

Mappers save outputs to local disk before serving to reducers

- » Allows recovery if a reducer crashes
- » Allows running more reducers than # of nodes

Fault Tolerance in MapReduce

1. If a task crashes:

» Retry on another node

- OK for a map because it had no dependencies
- OK for reduce because map outputs are on disk

» If the same task repeatedly fails, fail the job or ignore that input block

➤ Note: For fault tolerance to work, user tasks must be deterministic and side-effect-free

Fault Tolerance in MapReduce

2. If a node crashes:

- » Relaunch its current tasks on other nodes
- » Relaunch any maps the node previously ran
 - Necessary because their output files were lost along with the crashed node

Fault Tolerance in MapReduce

3. If a task is going slowly (straggler):
 - Launch second copy of task on another node
 - Take the output of whichever copy finishes first, and kill the other one

Critical for performance in large clusters (many possible causes of stragglers)

Summary

Data-parallel programming models let systems automatically manage much of execution:

- » Assigning work, load balancing, fault recovery

But... the story doesn't end here!

Outline

The big data problem

MapReduce

Apache Spark

How people are using Spark

Limitations of MapReduce

Programmability: most applications require higher level functions than map / reduce

- » E.g. statistics, matrix multiply, graph search
- » Google ads pipeline had 20 MR steps

Performance: inefficient to combine multiple MapReduce steps into complex programs

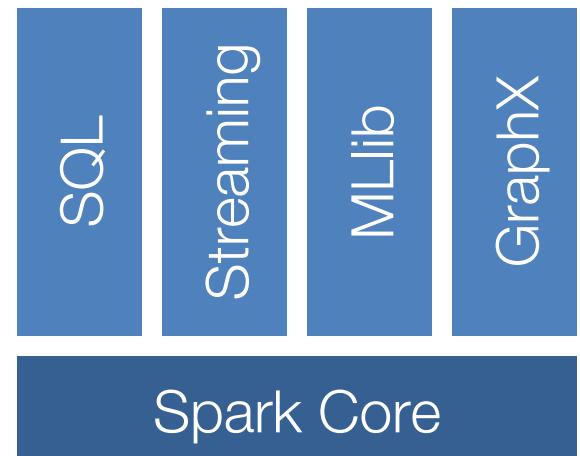
Apache Spark

Programming model that generalizes MapReduce
to support more applications

- » Adds efficient, in-memory data sharing

Large library of built-in functions

APIs in Python, Java, Scala, R



Spark Programmability

WordCount in MapReduce:

```
#include "mapreduce/mapreduce.h"

// User's map function
class Splitwords: public Mapper {
public:
virtual void Map(const MapInput& input)
{
    const string& text = input.value();
    const int n = text.size();
    for (int i = 0; i < n; ) {
        // Skip past leading whitespace
        while (i < n && isspace(text[i]))
            i++;
        // Find word end
        int start = i;
        while (i < n && !isspace(text[i]))
            i++;
        if (start < i)
            Emit(text.substr(
                start,i-start),"1");
    }
};

REGISTER_MAPPER(Splitwords);

// User's reduce function
class Sum: public Reducer {
public:
virtual void Reduce(ReduceInput* input)
{
    // Iterate over all entries with the
    // same key and add the values
    int64 value = 0;
    while (!input->done()) {
        value += StringToInt(
            input->value());
        input->NextValue();
    }
    // Emit sum for input->key()
    Emit(IntToString(value));
};

REGISTER_REDUCER(Sum);

int main(int argc, char** argv) {
ParseCommandLineFlags(argc, argv);
MapReduceSpecification spec;
for (int i = 1; i < argc; i++) {
    MapReduceInput* in= spec.add_input();
    in->set_format("text");
    in->set_filepattern(argv[i]);
    in->set_mapper_class("Splitwords");
}

// Specify the output files
MapReduceOutput* out = spec.output();
out->set_filebase("/gfs/test/freq");
out->set_num_tasks(100);
out->set_format("text");
out->set_reducer_class("Sum");

// Do partial sums within map
out->set_combiner_class("Sum");

// Tuning parameters
spec.set_machines(2000);
spec.set_map_megabytes(100);
spec.set_reduce_megabytes(100);

// Now run it
MapReduceResult result;
if (!MapReduce(spec, &result)) abort();
return 0;
}
```

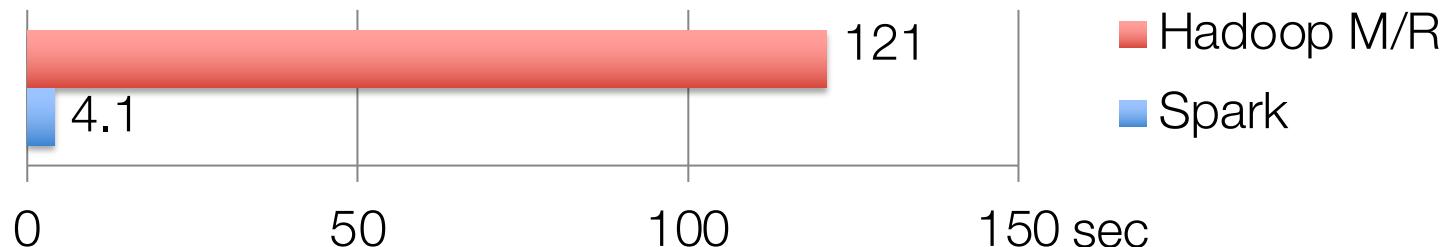
Spark Programmability

WordCount in Spark:

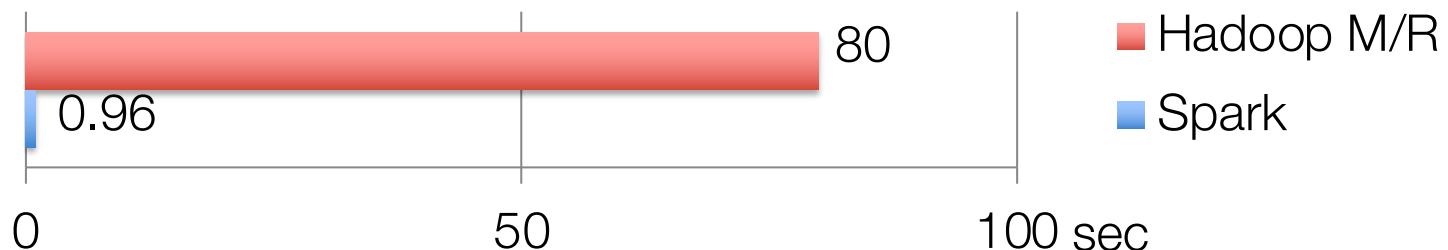
```
file = spark.textFile("hdfs://...")  
  
counts = file.flatMap(lambda line: line.split(" ")).  
         .map(lambda word: (word, 1)).  
         .reduceByKey(lambda a, b: a+b)  
  
counts.save("out.txt")
```

Spark Performance

K-means Clustering



Logistic Regression



Programming Model

Write programs in terms of transformations on distributed datasets

Low-level API: Resilient Distributed Dataset (RDD)

- » Collection of objects that can be stored in memory or disk across a cluster
- » Built via parallel transformations (map, filter, ...)
- » Automatically rebuilt on failure

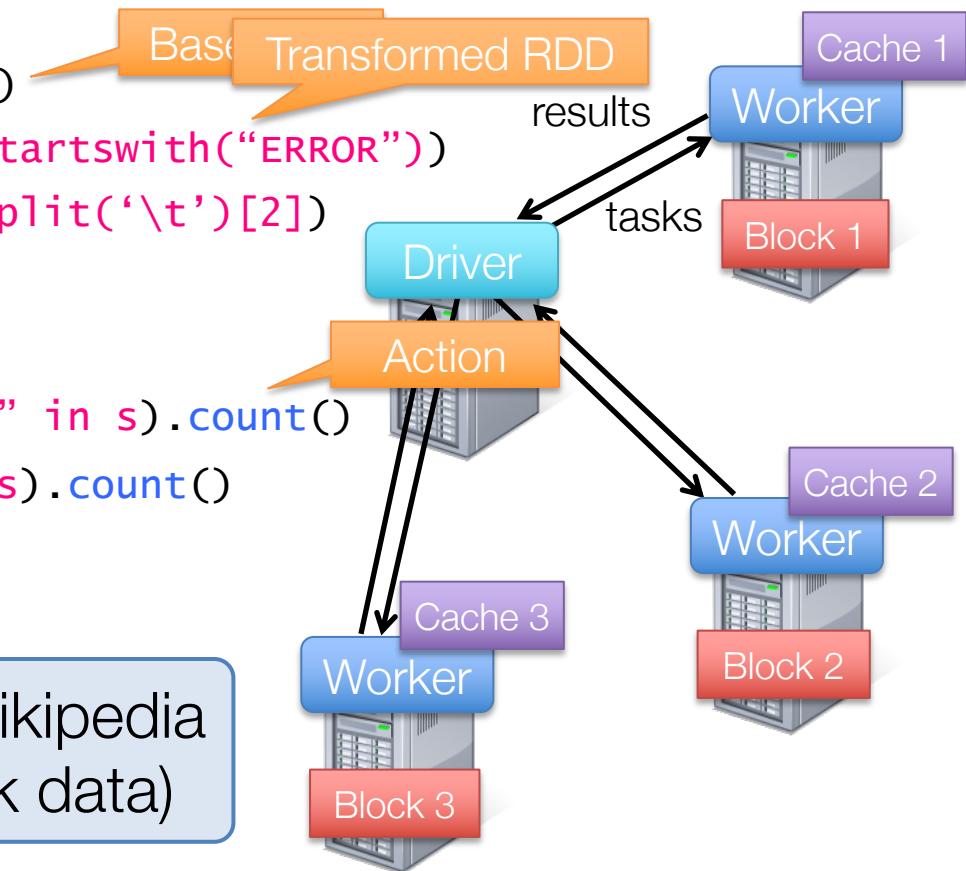
Example: Text Search

Load a large log file 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: "Berkeley" in s).count()  
messages.filter(lambda s: "MIT" in s).count()  
...
```

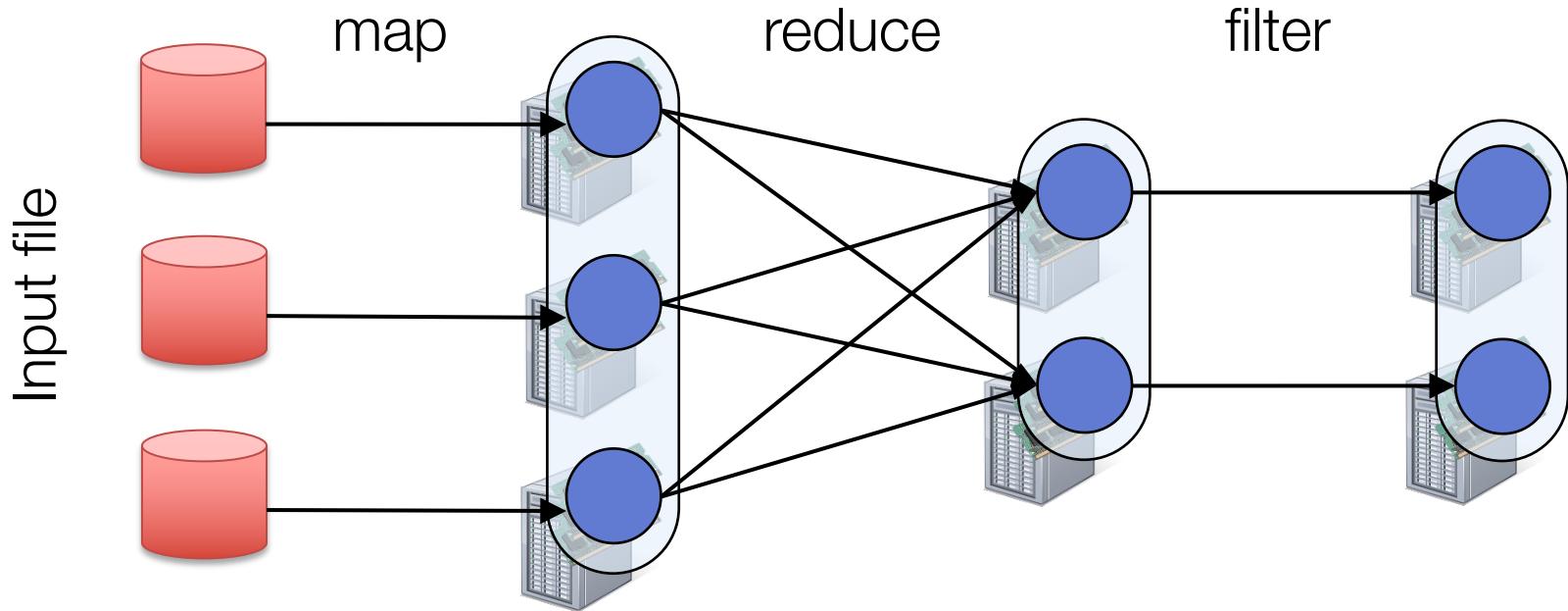
Example: full-text search of Wikipedia
in 0.5 sec (vs 20 s for on-disk data)



Fault Tolerance

RDDs track *lineage* info to rebuild lost data

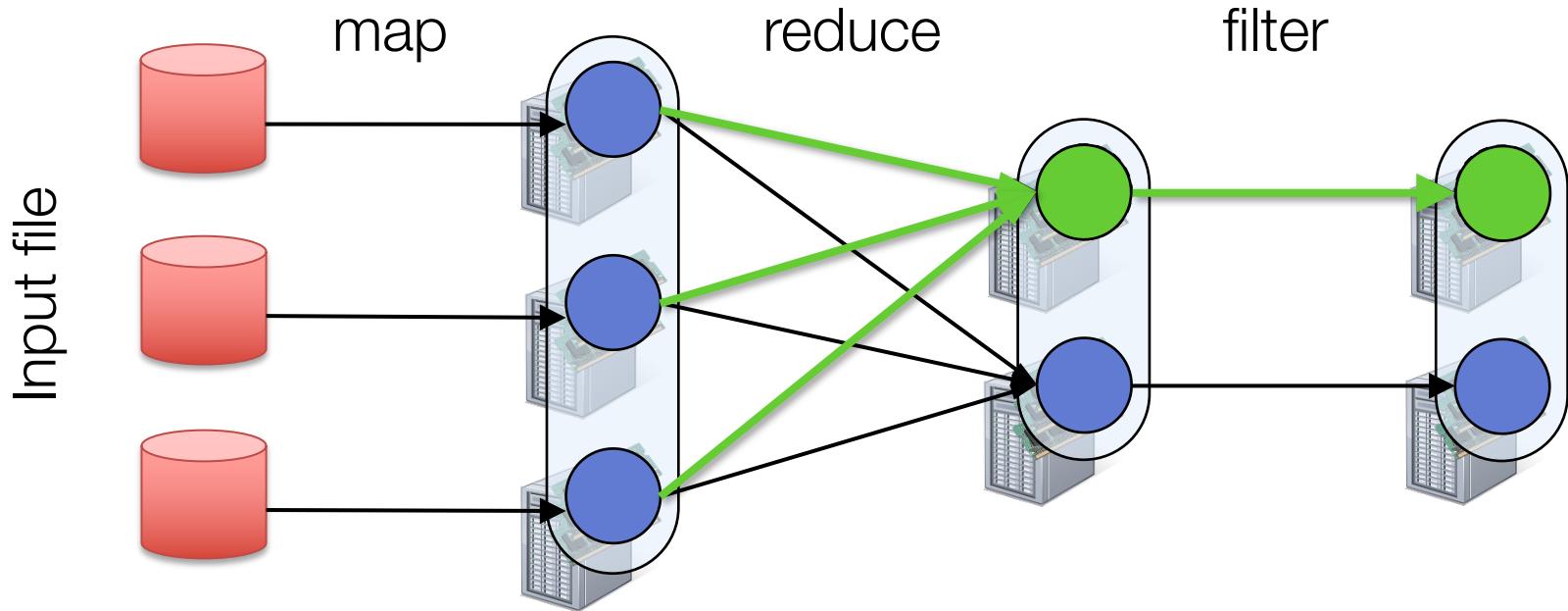
```
file.map(lambda record: (record.type, 1))  
    .reduceByKey(lambda x, y: x + y)  
    .filter(lambda type, count: count > 10)
```



Fault Tolerance

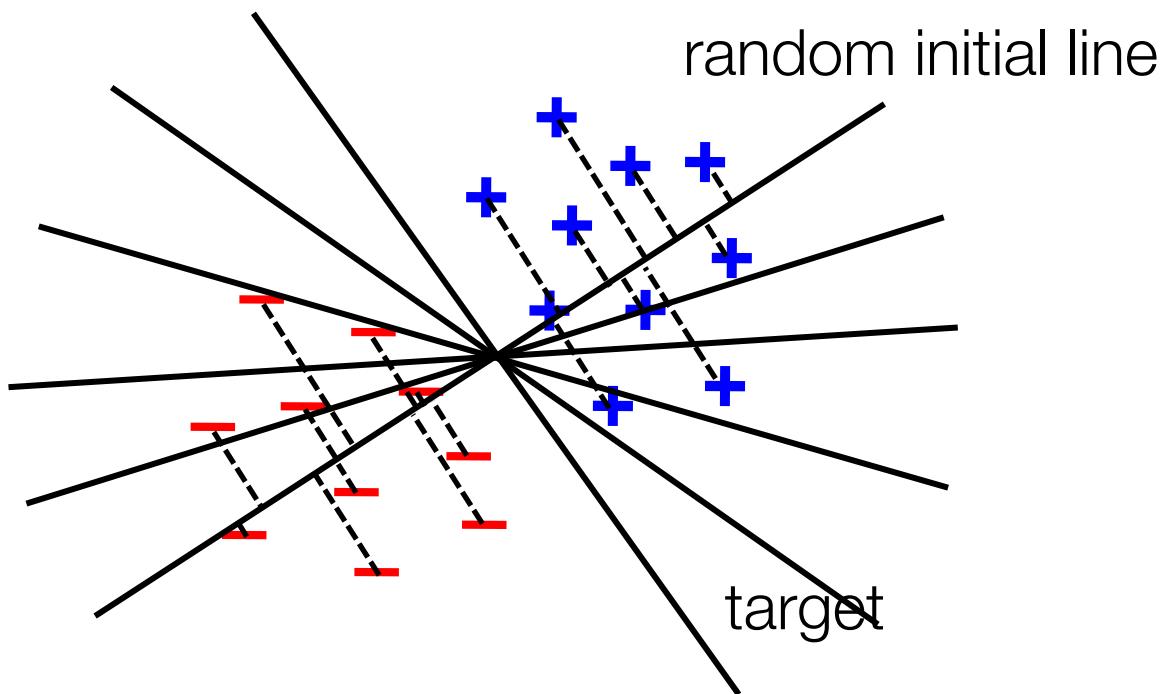
RDDs track *lineage* info to rebuild lost data

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

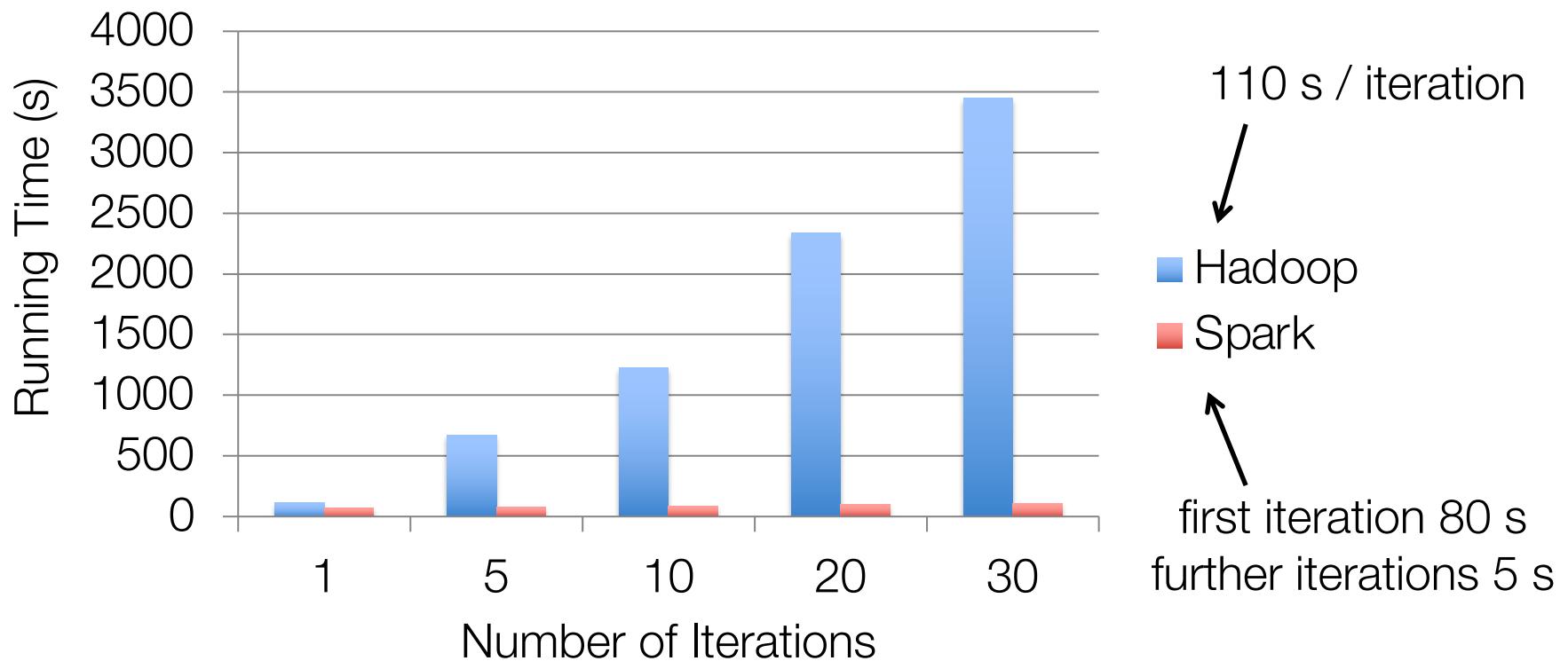


Example: Logistic Regression

Goal: find line separating two sets of points



Logistic Regression Performance



Built-in Libraries

DataFrames
and SQL

Structured
Streaming

MLlib

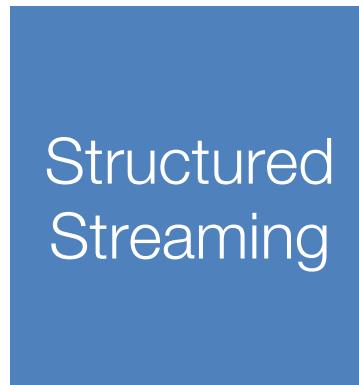
GraphX

Spark Core (RDDs)

Largest integrated standard library for big data

Built-in Libraries

Now the most common way to use Spark



Spark Core (RDDs)

Largest integrated standard library for big data

Challenges with RDD API

Looks high-level, but hides many semantics of computation from engine

- » Functions passed in are arbitrary blocks of code
- » Data stored is arbitrary Java/Python objects

Users can mix APIs in suboptimal ways

Example Problem

```
pairs = data.map(lambda word: (word, 1))
```

```
groups = pairs.groupByKey()
```

```
groups.map(lambda k, vs: (k, vs.sum))
```



Writes all groups as lists
of integers in memory

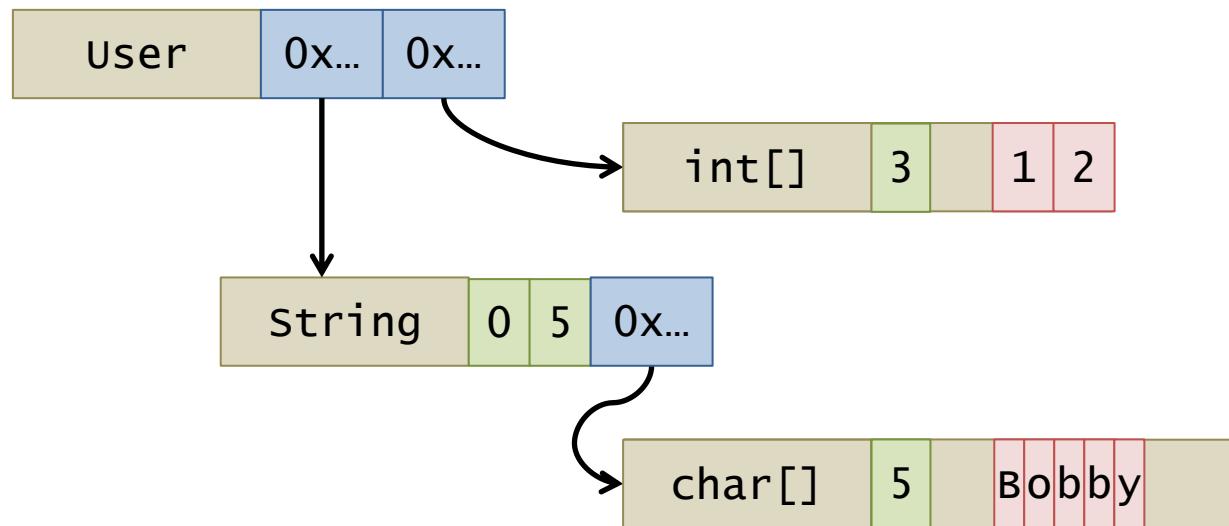


Then promptly
aggregates them

Challenge: Data Representation

Java & Python objects often have high overhead

```
class User(name: String, friends: Array<Int>)
User("Bobby", Array(1, 2))
```

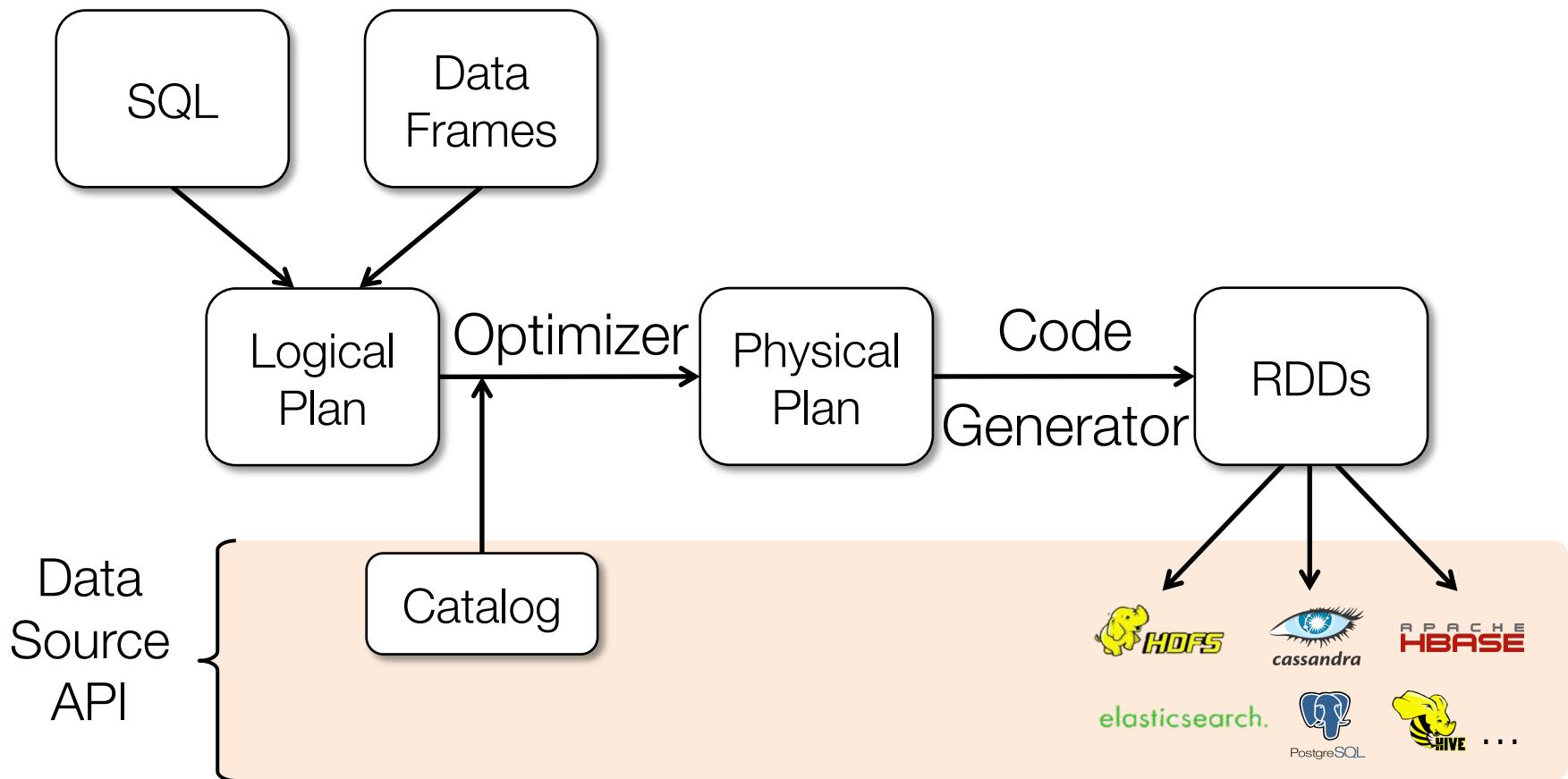


Spark SQL & DataFrames

Efficient library for working with structured data

- » Two interfaces: SQL for data analysts and external apps, DataFrames for complex programs
- » Automatically optimizes computation and storage underneath, as in a relational database

Spark SQL Architecture



DataFrame API

DataFrames hold rows with a known schema
and offer relational operations through a DSL

```
users = spark.sql("select * from users")
```

```
ma_users = users[users.state == "MA"]  
ma_users.count()
```

Expression tree (AST)

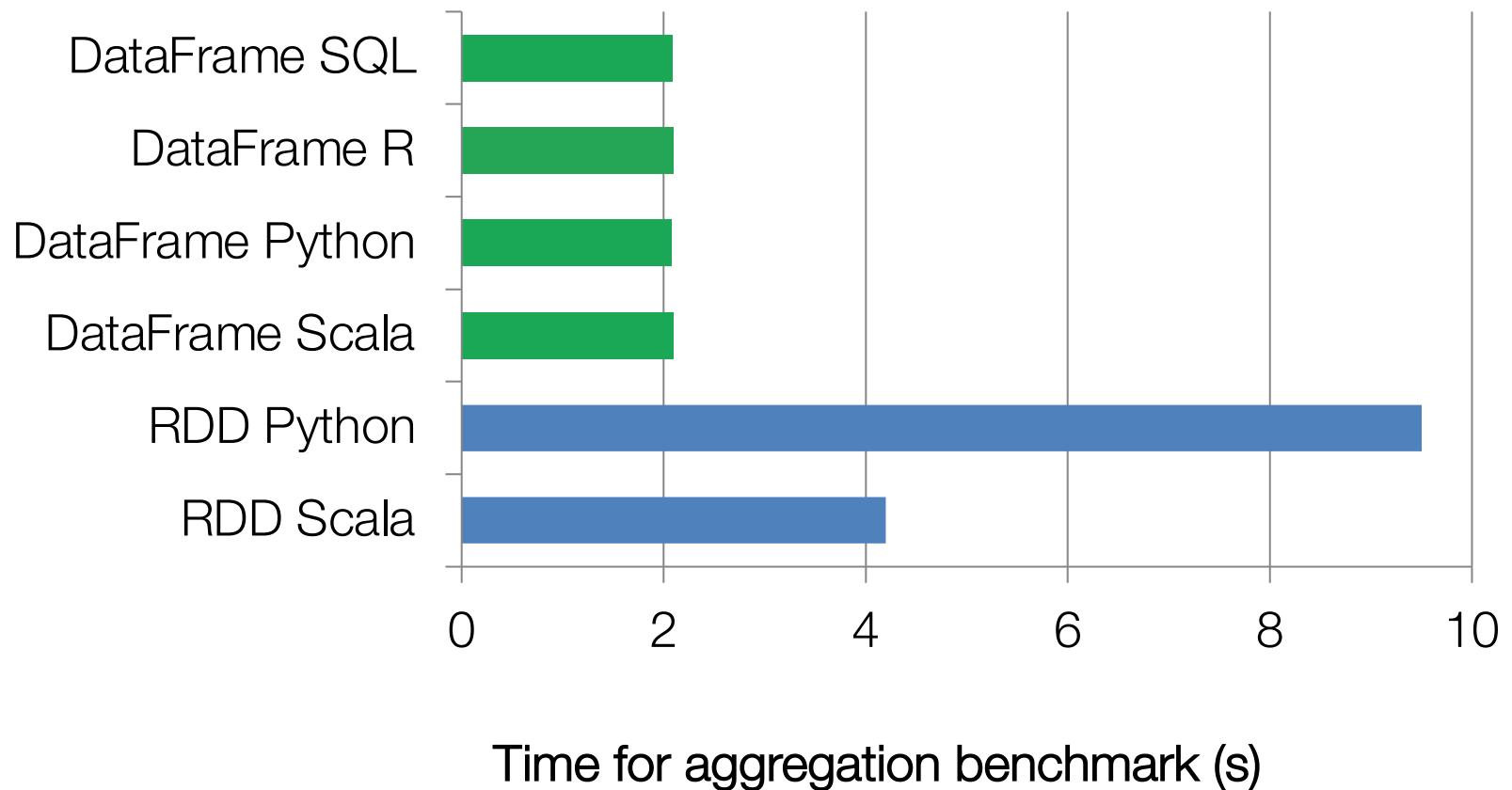
```
ma_users.groupBy("name").avg("age")
```

```
ma_users.map(lambda row: row.user.to_upper())
```

What DataFrames Enable

1. Compact binary representation
 - Columnar, compressed cache; rows for processing
2. Optimization across operators (join reordering, predicate pushdown, etc)
3. Runtime code generation

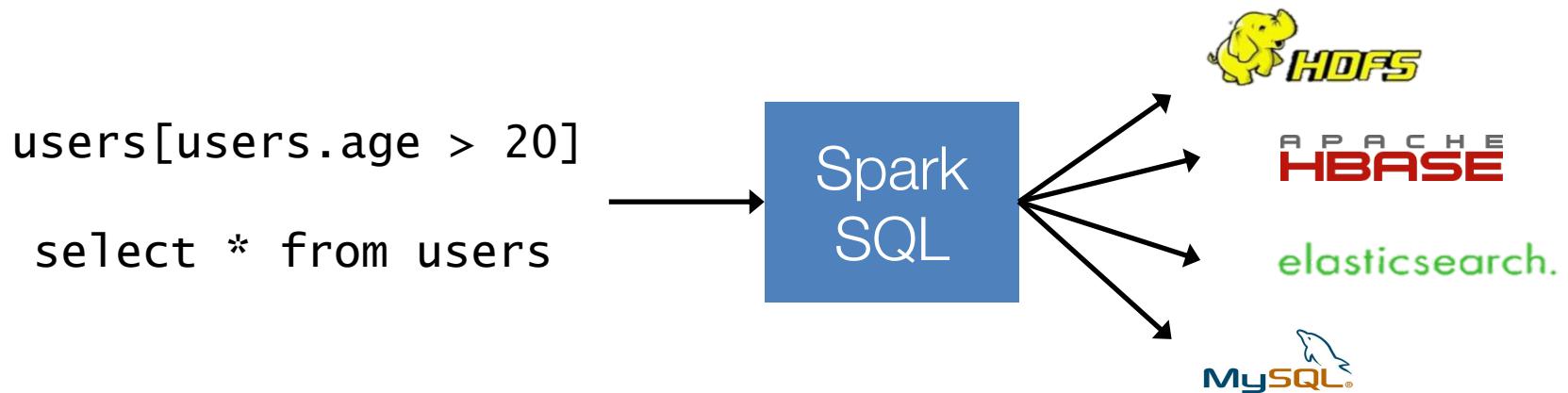
Performance



Data Sources

Uniform way to access structured data

- » Apps can migrate across Hive, Cassandra, JSON, ...
- » Rich semantics allows *query pushdown* into sources



Examples

JSON:

```
select user.id, text from tweets
```

```
{  
  "text": "hi",  
  "user": {  
    "name": "bob",  
    "id": 15  
  }  
}
```

JDBC:

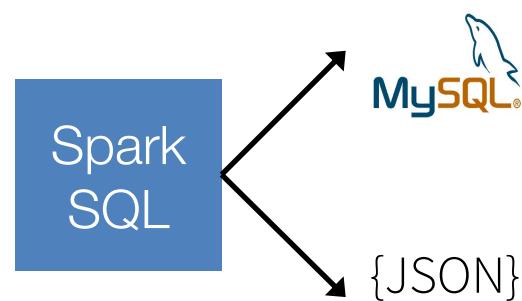
```
select age from users where lang = "en"
```

tweets.json

Together:

```
select t.text, u.age  
from tweets t, users u  
where t.user.id = u.id  
and u.lang = "en"
```

```
select id, age from  
users where lang="en"
```



Built-in Libraries

DataFrames
and SQL

Structured
Streaming

MLlib

GraphX

Spark Core (RDDs)

Largest integrated standard library for big data

Spark MLlib

Machine learning library with 50+ algorithms:

- » Classification/regression: linear models, decision trees, random forests, isotonic regression, ...
- » Clustering: k-means, Gaussian mixture, LDA, ...
- » Collaborative filtering: ALS, NMF
- » Dimensionality reduction: SVD, PCA

“Pipeline” API inspired by scikit-learn

MLlib Pipeline Example

```
# Configure an ML pipeline with three stages
tokenizer = Tokenizer(inputCol="text", outputCol="words")
hashingTF = HashingTF(inputCol=tokenizer.getOutputCol(),
                      outputCol="features")
lr = LogisticRegression(maxIter=10, regParam=0.001)

pipeline = Pipeline(stages=[tokenizer, hashingTF, lr])

# Fit the pipeline to training data
model = pipeline.fit(train_data)

# Make predictions on DataFrame of test data
prediction = model.transform(test_data)
```

Combining Libraries

```
# Prepare a dataset using SQL
ctx.jsonFile("tweets.json").registerTempTable("tweets")
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))
    .reduceByKey("5s", lambda a, b: a+b)
```

Performance of Composition

Separate computing frameworks:



Spark:



Summary

Libraries + function-based interface let users write parallel programs similar to sequential code

Can use Spark interactively in Python, R, etc

Outline

The big data problem

MapReduce

Apache Spark

How people are using Spark

Spark Community

20M downloads/month, clusters up to 8000 nodes

1800 contributors to open source project



Uber



airbnb

ebay

NETFLIX



SAMSUNG



REGENERON



TOYOTA



Telefonica



NVIDIA

Example Use Cases

REGENERON Correlate 500,000 patient records with DNA to design therapies



Query 170 PB of data that used to be in 14 databases



Optimize production using ML and SQL on petabyte-scale data

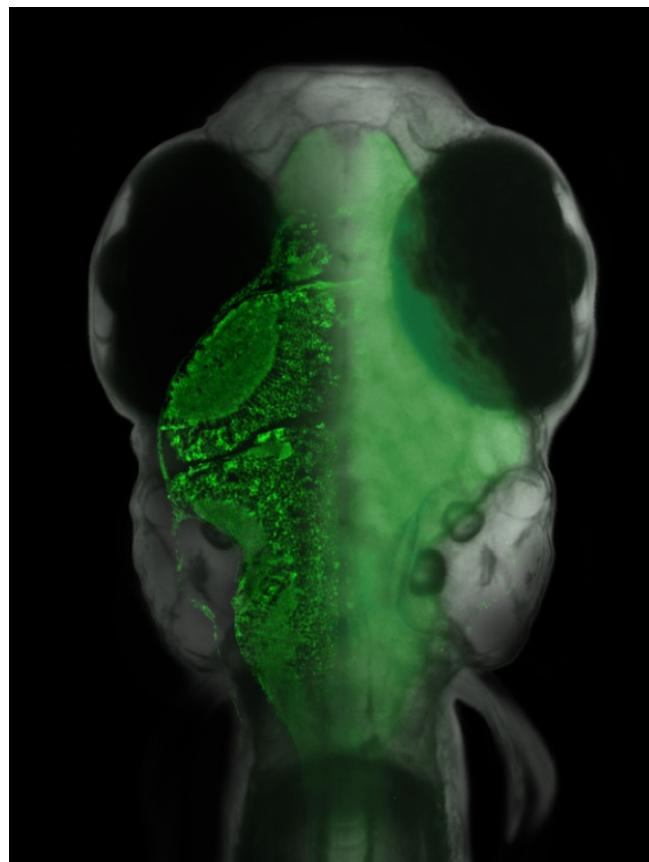


Identify securities fraud via ML on 30 PB of data

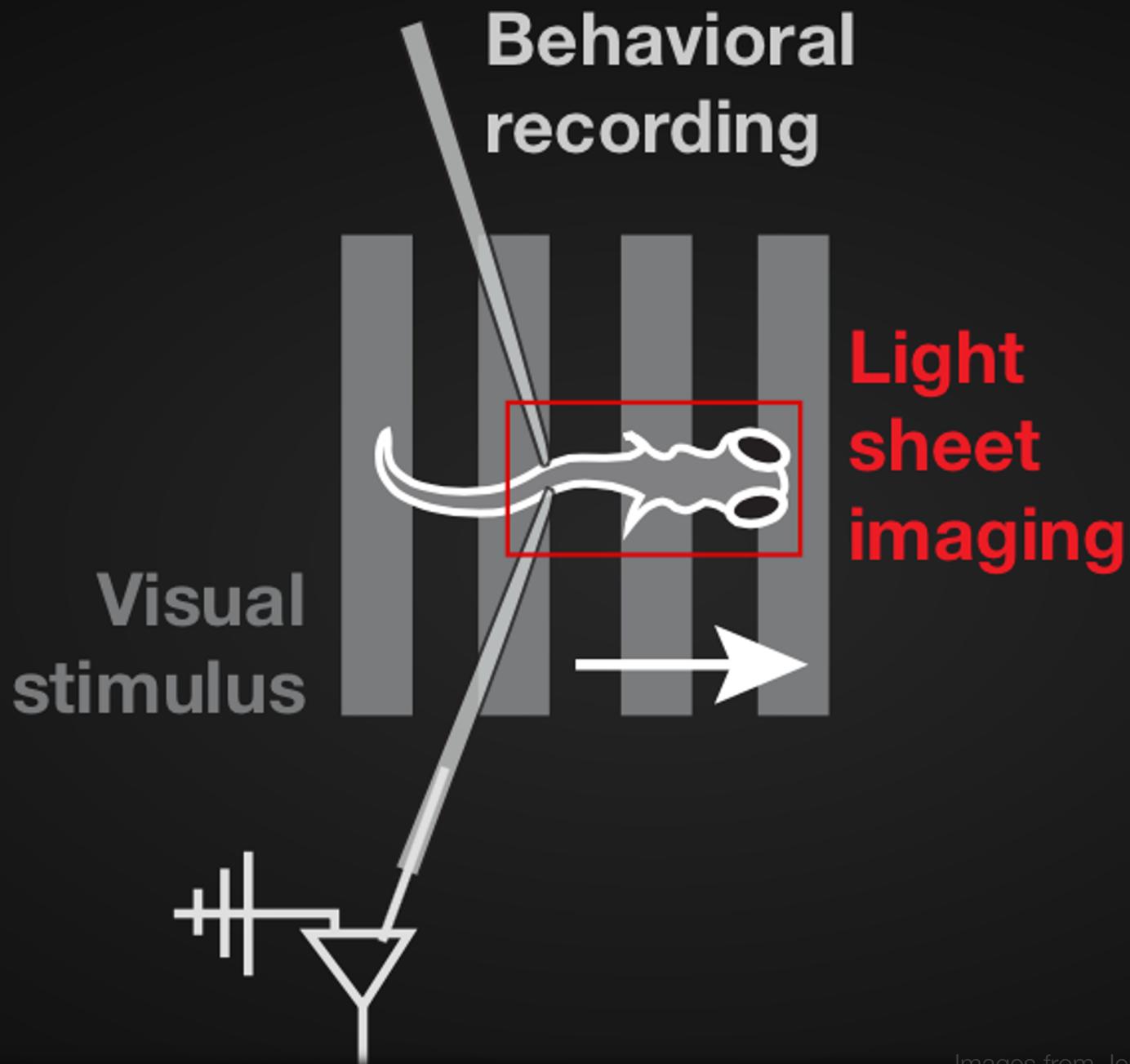
5000+ customers running 10 million VMs/day on Databricks alone

Early Use Case: Neuroscience

HHMI Janelia Farm analyzes data from full-brain imaging of neural activity

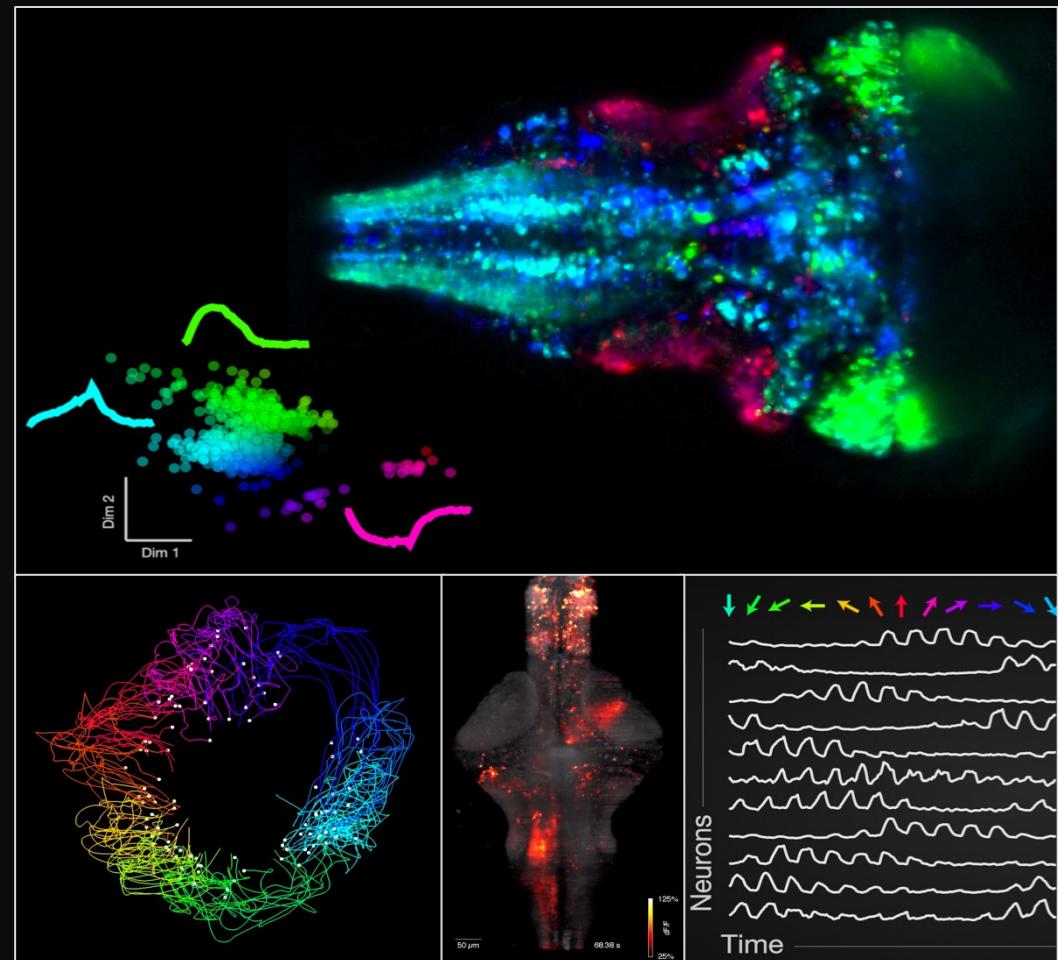


Larval zebrafish
+
Light-sheet imaging
=
2 TB / hour of data



Data Analysis

Streaming code
does clustering,
dimens. reduction
on 80-node cluster



Early Use Case: Genomics

Berkeley ADAM stores and processes reads with standard big data tools & formats

bdgenomics.org

Big Data Genomics Blog Archives Projects Mailing List CLAs Search 

 <https://github.com/bigdatagenomics/>  <https://twitter.com/bigdatagenomics/>

FEB 25TH, 2016 **ADAM 0.19.0 Released**

ADAM version 0.19.0 has been [released](#), built for both [Scala 2.10](#) and [Scala 2.11](#). The 0.19.0 release contains various concordance fixes and performance improvements for accessing read metadata. Schema changes, including a bump to version 0.7.0 of the Big Data Genomics [Avro data formats](#), were made to support the read metadata performance improvements. Additionally, the performance of exporting a single BAM file was improved, and this was made to be guaranteed correct for sorted data.

ADAM now targets Apache Spark 1.5.2 and Apache Hadoop 2.6.0 as the default build environment. ADAM and applications built on ADAM should run on a wide range of Apache Spark (1.3.1 up to and including the most recent, 1.6.0) and Apache Hadoop (currently 2.3.0 and 2.6.0) versions. A compatibility matrix of Spark, Hadoop, and Scala version builds in our [continuous integration system](#) verifies this. Please note, as of this release, support for Apache Spark 1.2.x and Apache Hadoop 1.0.x [has been dropped](#).

The full list of changes since version 0.18.2 is below.

[Read on →](#)

NOV 11TH, 2015 **ADAM 0.18.2 Released**

A few ADAM releases have been made since the last announcement; we'll attempt to catch up here.

The most recent is a version [0.18.2 bugfix release](#), built for both [Scala 2.10](#) and [Scala 2.11](#). It fixes a [minor issue](#) with the binary distribution artifact from version 0.18.1.

Prior to version 0.18.2, we made significant changes to support version 0.6.0 of the Big Data Genomics [Avro data formats](#). We also improved performance on core transforms (markdups,

About us...
This project is supported in part by NIH BD2K Award 1-U54HG007990-01 and NIH Cancer Cloud Pilot Award HHSN261201400006C with collaborators from the AMPLab at UC Berkeley, Genome Informatics Lab at UC Santa Cruz, Icahn School of Medicine at Mount Sinai, Microsoft Research, Cloudera, and the Broad Institute.

Chat with us... [gitter](#) [join chat](#)
If you're interested in contributing, take a look at the open "pick me up!" issues.

Recent Posts

ADAM 0.19.0 Released

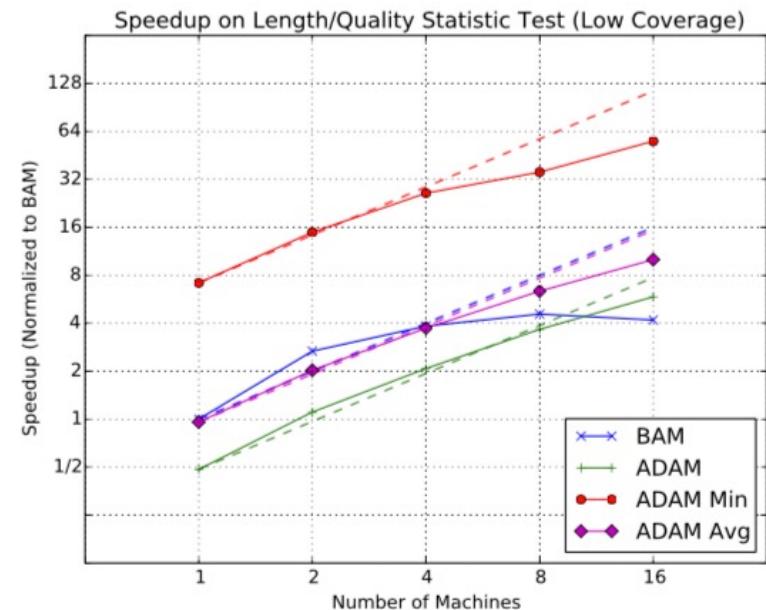
ADAM 0.18.2 Released

ADAM Results

25% reduction in file size over BAM, by using standard Apache Parquet format

Sorting, filtering, and base quality score recalibration (BQSR) scale up linearly

Interactive use in shell



Conclusion

Apache Spark offers a high-level interface to work with big data based on data-parallel model

Large set of existing libraries

Easy to try on just your laptop!

spark.apache.org

