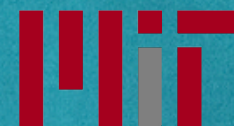# What's Changing in Big Data?

Matei Zaharia
June 21, 2016
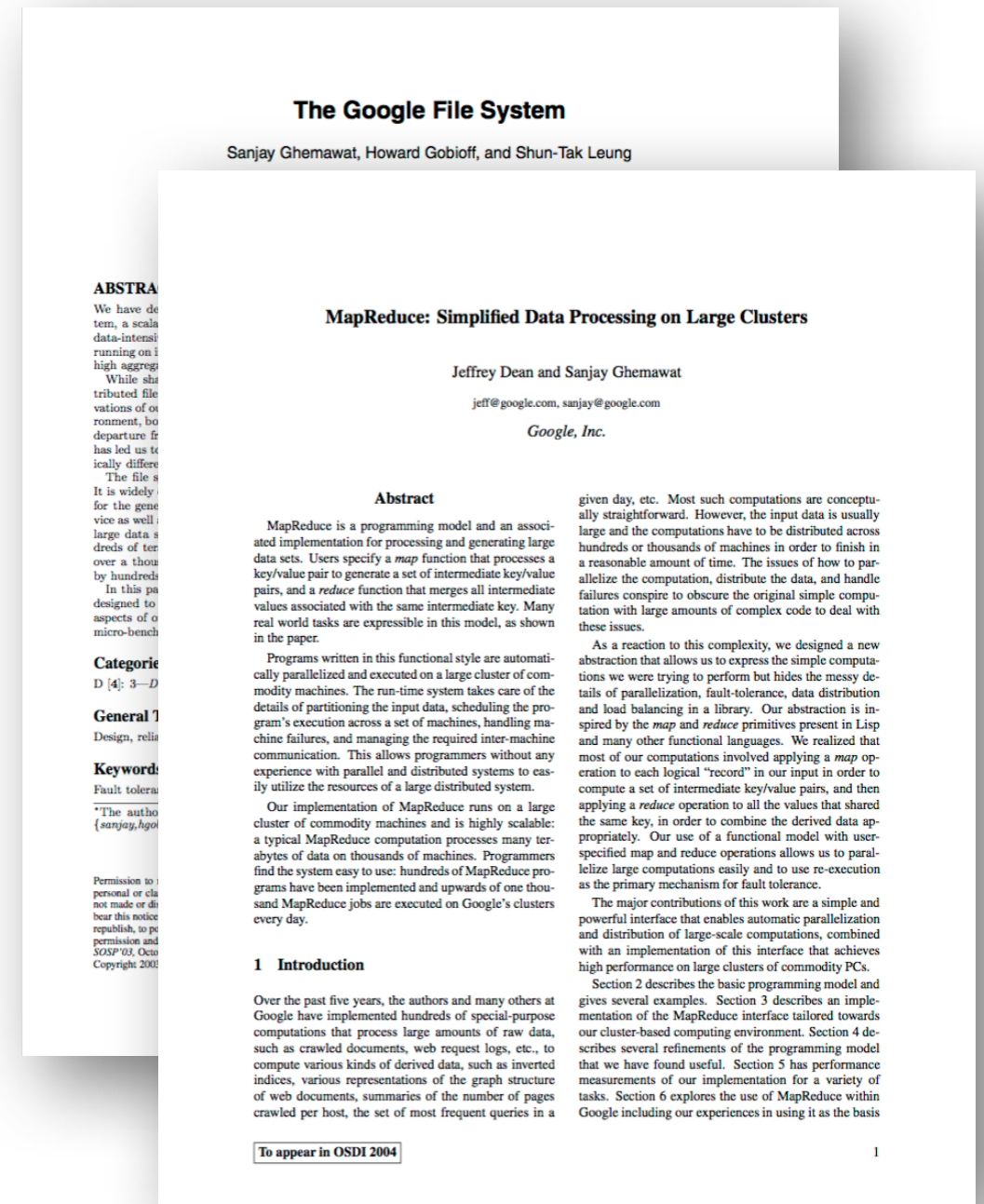
# Background

Big data systems became a popular research topic nearly 10 years ago

- Large-scale, commodity clusters

What has changed since then?



databricks

# My Perspective



Open source processing engine and set of libraries



Cloud data processing service based on Apache Spark

# Three Key Changes

① **Users:** engineers ➡ analysts

② **Hardware:** I/O bottleneck ➡ compute

③ **Delivery:** strong trend toward cloud

databricks™

# Changing Users

# Initial Big Data Users

lying runtime system automatically parallelizes the computation across large-scale clusters of machines, handles machine failures, and schedules inter-machine communication to make efficient use of the network and disks. Programmers find the system easy to use: more than ten thousand distinct MapReduce programs have been implemented internally at Google over the past four years, and an average of one hundred thousand MapReduce jobs are executed on

## Software engineers:

- Use Java, C++, etc to create large projects
- Build applications out of low-level operators

databricks™

# Expanding the User Base

Scripting / query languages inspired by SQL, awk, etc

Used by new roles:
- **Data scientists** (technical domain experts, e.g. ML)
- **Analysts** (business)



databricks™

# Challenges for Non-Engineers

API familiarity

Performance predictability & debugging

Can't hide that it's large-scale

Access from small data tools

E.g. Excel, Tableau
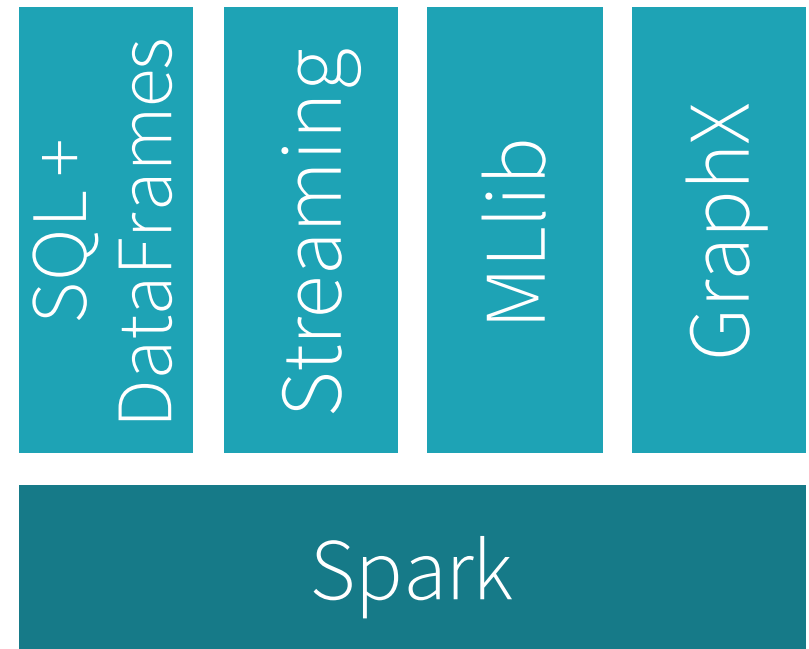
Worse with more familiar APIs!

databricks™

# Case Study: Apache Spark

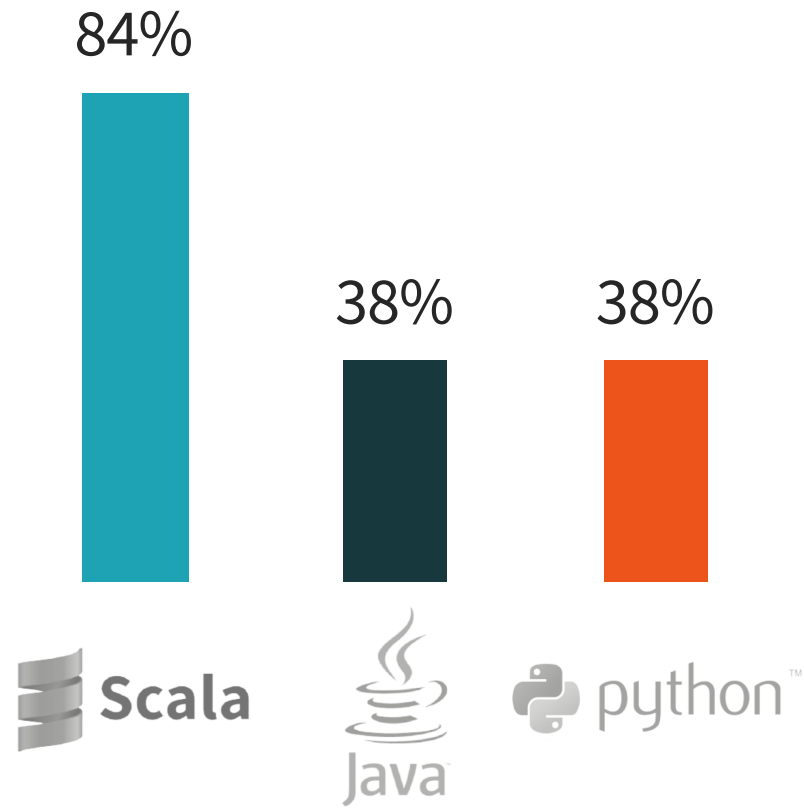Cluster computing engine that generalizes MapReduce

Collection of APIs and libraries

- APIs in Scala, Java, Python and R
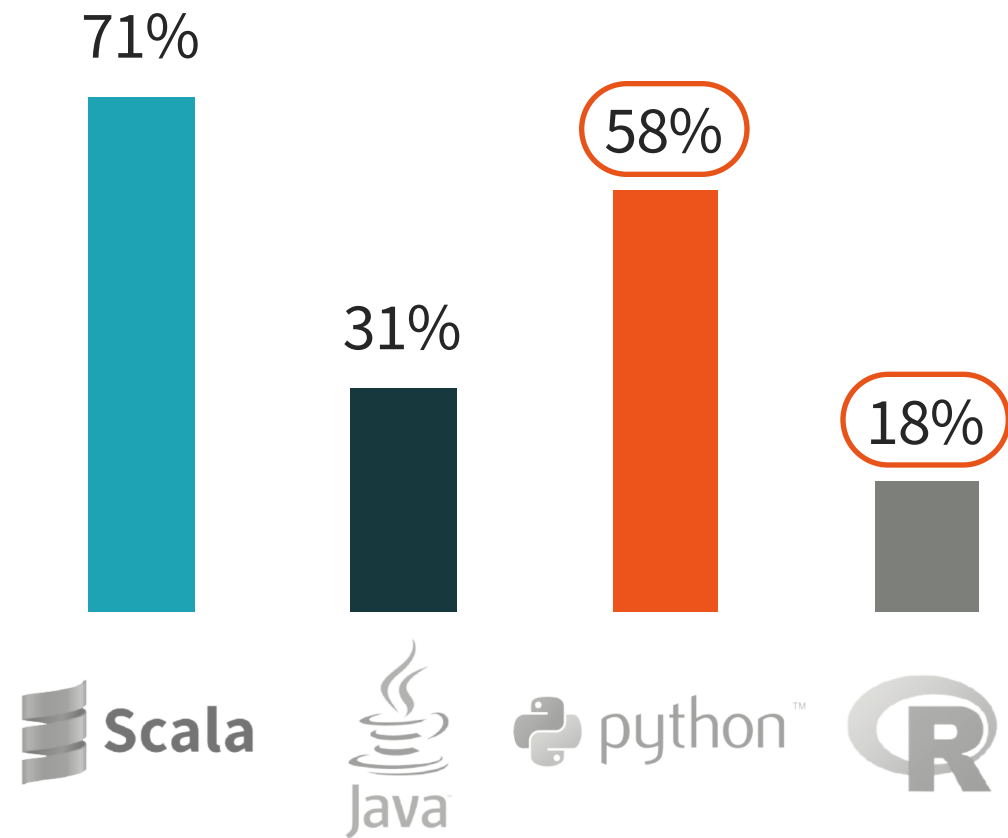- Streaming, SQL, ML, graph, …

1000+ deployments, max > 8000 nodes

| SQL + DataFrames | Streaming | MLlib | GraphX |
|---|---|---|---|
| Spark | | | |

# Languages Used for Spark

2014 Languages Used

84%
38%    38%

Scala    Java    python™

2015 Languages Used

71%
58%
31%
18%

Scala    Java    python™    R

# Original Spark API

Functional API aimed at Java / Scala developers

Resilient Distributed Datasets (RDDs): distributed collections with functional transformations

```
lines = spark.textFile("hdfs://...")          // RDD[String]
points = lines.map(line => parsePoint(line))  // RDD[Point]
points.filter(p => p.x > 100).count()
```

databricks

# Challenge with Functional API

Looks high-level, but hides many semantics of computation

- Functions are arbitrary blocks of Java bytecode
- Data stored is arbitrary Java objects

Users can mix APIs in suboptimal ways

databricks™

# Which Operator Causes Most Tickets?

map                reduce              sample

filter             count               take

groupBy            fold                first

sort               reduceByKey         partitionBy

union              groupByKey          mapWith

join               cogroup             pipe

leftOuterJoin      cross               save

rightOuterJoin     zip                 ...

databricks

# Example Problem

```
pairs = data.map(word => (word, 1))

groups = pairs.groupByKey()

groups.map((k, vs) => (k, vs.sum))
```

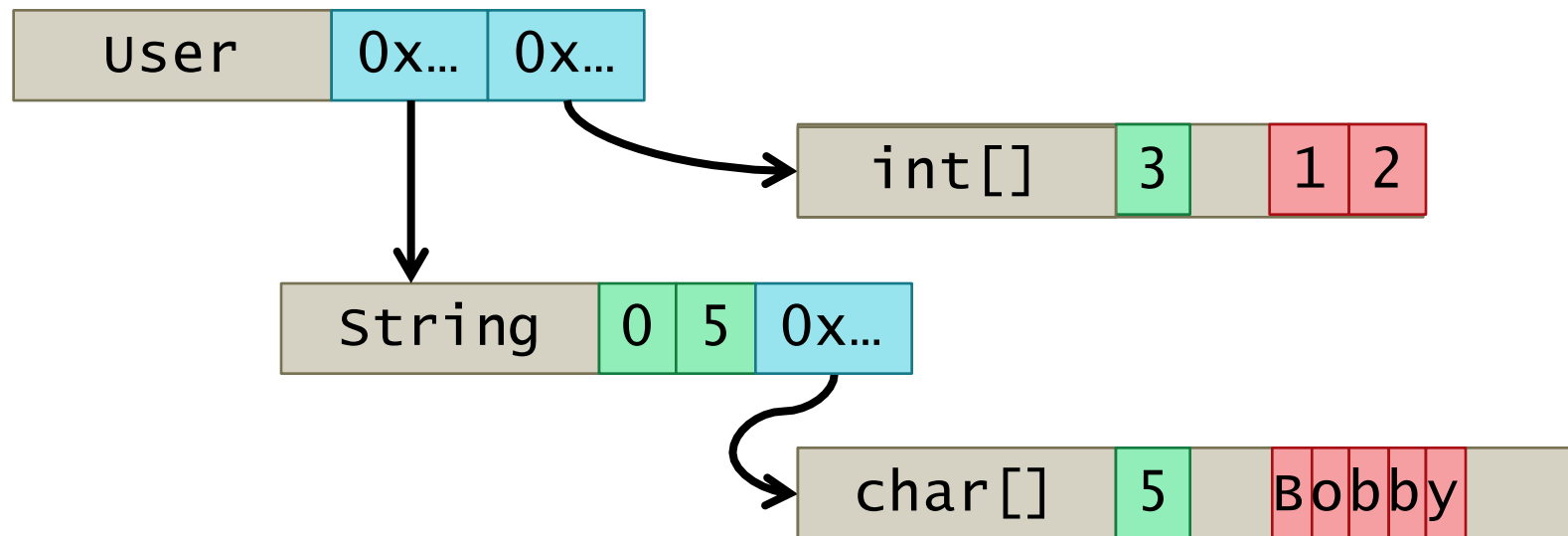Materializes all groups as Seq[Int] objects

Then promptly aggregates them

# Challenge: Data Representation

Java objects often many times larger than underlying fields

```
class User(name: String, friends: Array[Int])
new User("Bobby", Array(1, 2))
```

# Structured APIs: DataFrames + Spark SQL

databricks™

# DataFrames and Spark SQL

Efficient library for structured data (data with a known schema)

- Two interfaces: SQL for analysts + apps, DataFrames for programmers

Optimized computation and storage, similar to RDBMS

## Spark SQL: Relational Data Processing in Spark

Michael Armbrust[†], Reynold S. Xin[†], Cheng Lian[†], Yin Huai[†], Davies Liu[†], Joseph K. Bradley[†], Xiangrui Meng[†], Tomer Kaftan[‡], Michael J. Franklin[††], Ali Ghodsi[†], Matei Zaharia[†*]

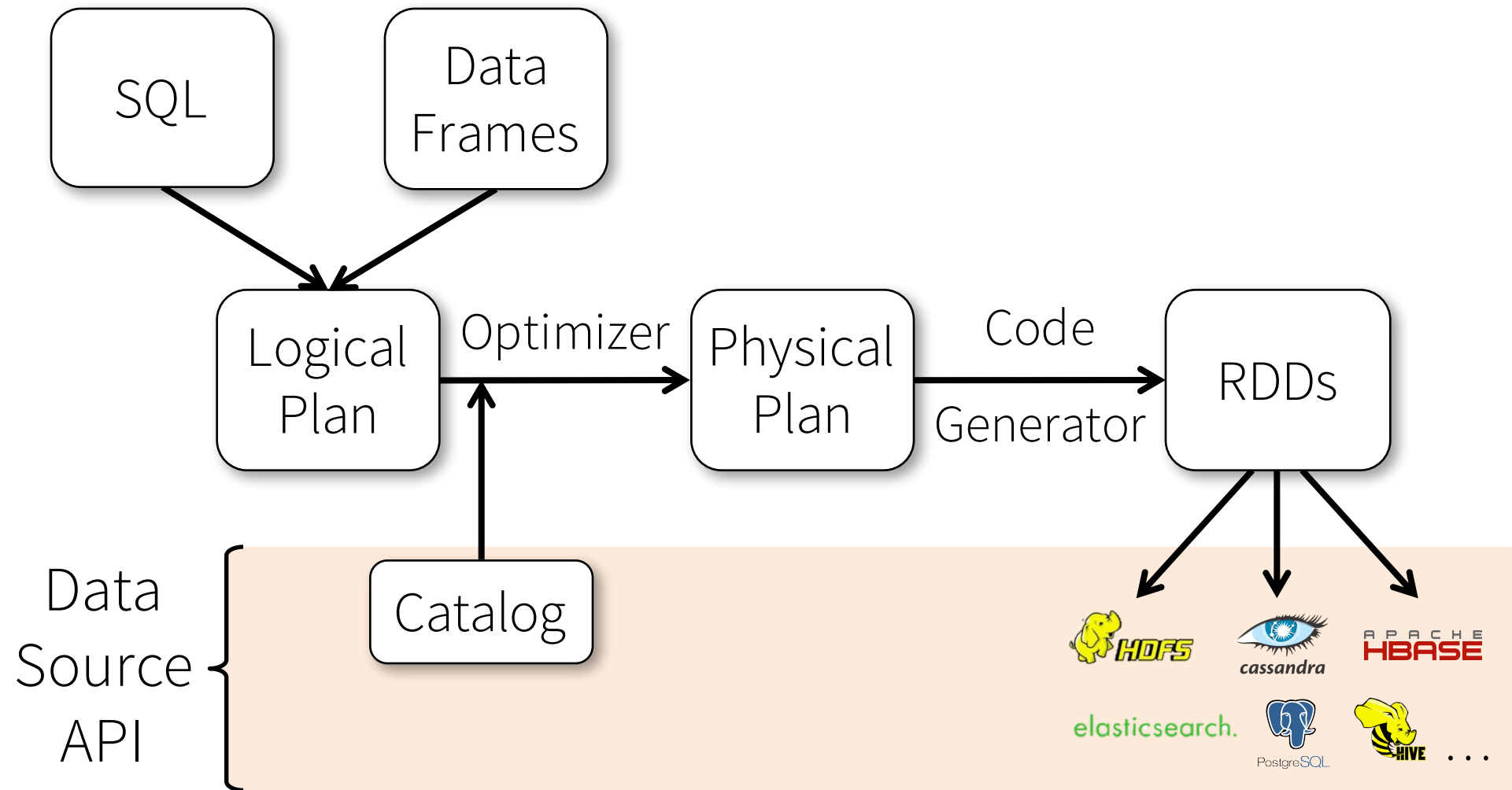[†]Databricks Inc.        [*]MIT CSAIL        [‡]AMPLab, UC Berkeley

**ABSTRACT**

Spark SQL is a new module in Apache Spark that integrates relational processing with Spark's functional programming API. Built on our experience with Shark, Spark SQL lets Spark program-

While the popularity of relational systems shows that users often prefer writing declarative queries, the relational approach is insufficient for many big data applications. First, users want to perform ETL to and from various data sources that might be semi- or un-

databricks™

# Execution Steps

# DataFrame API

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

```
val c = new HiveContext()
val users = c.sql("select * from users")

val massUsers = users(users("state") === "MA")
                     └─────── Expression AST ───────┘
massUsers.count()

massUsers.groupBy("name").avg("age")

massUsers.map(row => row.getString(0).toUpper())
```

# Why DataFrames?

Based on data frame concept in R and Python

- Spark is the first to make this a declarative API

Integrates with other data science libraries

- MLlib, GraphFrames, …



2005    2007    2009    2011    2013    2015
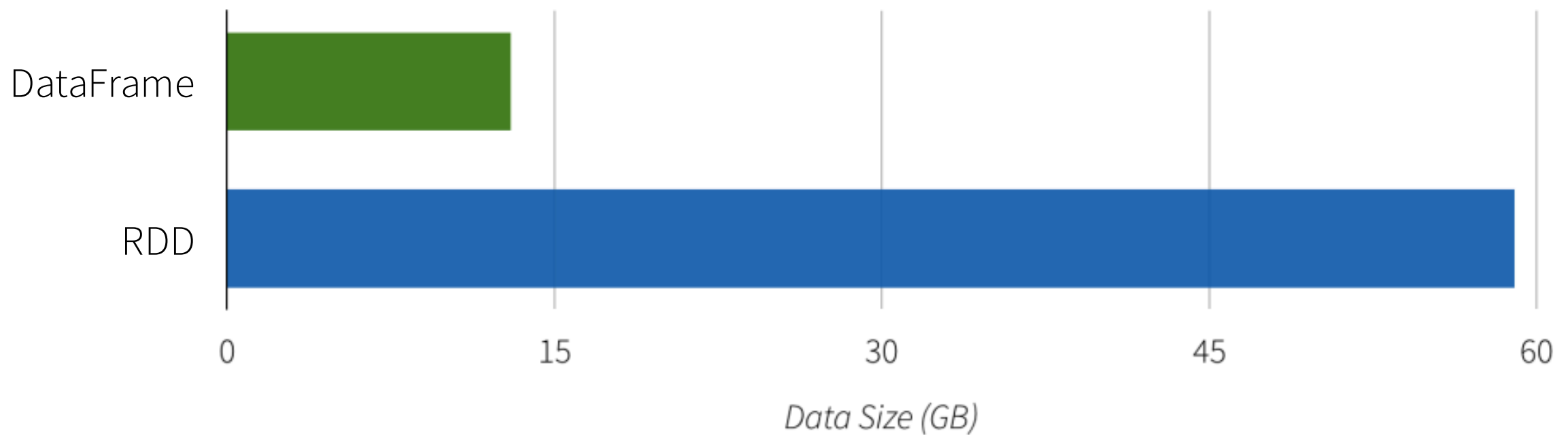
databricks™

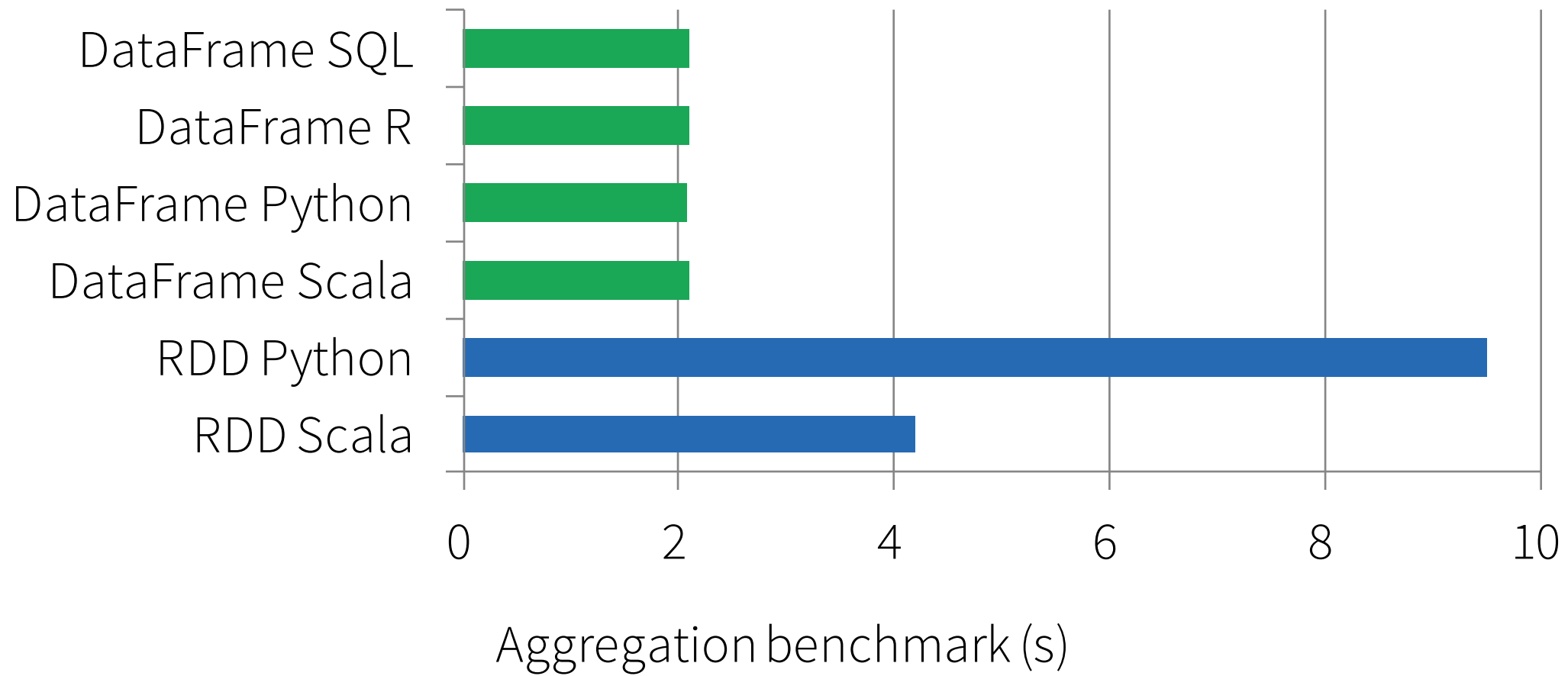Google trends for "data frame"

# What Structured APIs Enable

1. Compact binary representation
   - Columnar, compressed format for caching; rows for processing

2. Optimization across operators (join ordering, pushdown, etc)

3. Runtime code generation

databricks™

# Space Usage



Memory Usage when Caching

Data Size (GB)

# Performance



Aggregation benchmark (s)

# Uptake

DataFrames were released in March 2015, but already see high use:
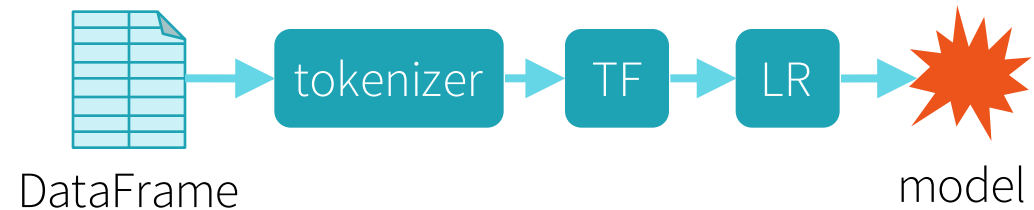
62% of users in 2015 survey use DataFrames

69% of users use Spark SQL

SQL & Python are the top languages on Databricks

databricks™

# Other High-Level APIs

## Machine Learning Pipelines
Modular API based on scikit-learn



DataFrame → tokenizer → TF → LR → model

## GraphFrames
Relational + graph operations

## Structured Streaming
Declarative streaming API in Spark 2.0

Many high-level data science APIs can be declarative

databricks™

# Changing Hardware

# Hardware Trends

Storage

Network

CPU

databricks™

# Hardware Trends

|  | 2010 |
|---|---|
| Storage | 50+MB/s (HDD) |
| Network | 1Gbps |
| CPU | ~3GHz |

databricks™

# Hardware Trends

|         | 2010              | 2016               |
|---------|-------------------|--------------------|
| Storage | 50+MB/s (HDD)     | 500+MB/s (SSD)     |
| Network | 1Gbps             | 10Gbps             |
| CPU     | ~3GHz             | ~3GHz              |

databricks™

# Hardware Trends

|  | 2010 | 2016 |  |
|---|---|---|---|
| Storage | 50+MB/s (HDD) | 500+MB/s (SSD) | 10x |
| Network | 1Gbps | 10Gbps | 10x |
| CPU | ~3GHz | ~3GHz | ☹ |

databricks™

# Summary

In 2005-2010, I/O was the name of the game

- Network locality, compression, in-memory caching

Now, compute efficiency matters even for data-intensive apps

- Getting harder with more diverse hardware, e.g. GPUs, FPGAs
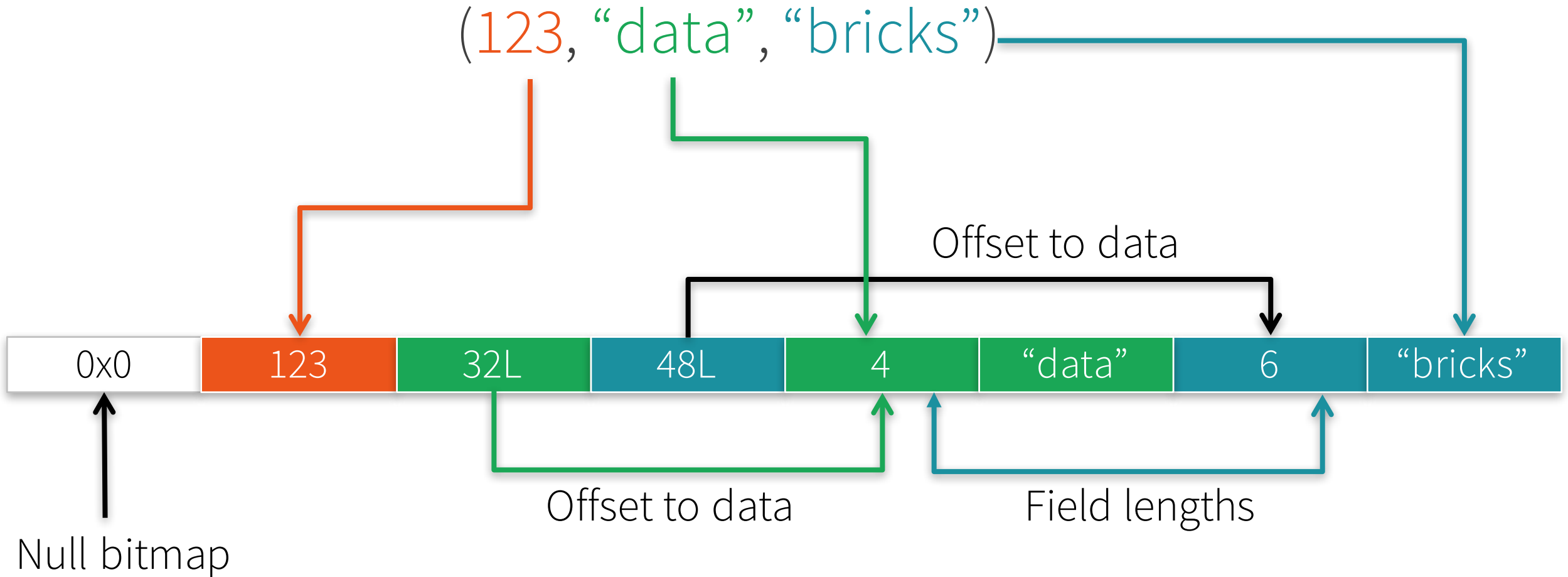
Future: network cards $\cong$ DRAM bandwidth

databricks™

# Spark Effort: Project Tungsten

Optimize Apache Spark's CPU and memory usage, via:

(1) Runtime code generation

(2) Exploiting cache locality

(3) Off-heap memory management

databricks™

# Tungsten's Binary Encoding

(123, "data", "bricks")

| 0x0 | 123 | 32L | 48L | 4 | "data" | 6 | "bricks" |
|-----|-----|-----|-----|---|--------|---|----------|

Offset to data

Offset to data

Field lengths

Null bitmap

databricks™

# Runtime Code Generation

DataFrame Code / SQL

```
df.where(df("year") > 2015)
```

Logical Expressions

```
GreaterThan(year#234, Literal(2015))
```

Low-level Bytecode

```
bool filter(Object baseObject) {
    int offset = baseOffset + bitSetWidthInBytes + 3*8L;
    int value = Platform.getInt(baseObject, offset);
    return value34 > 2015;
}
```

JVM intrinsic JIT-ed to pointer arithmetic

databricks™

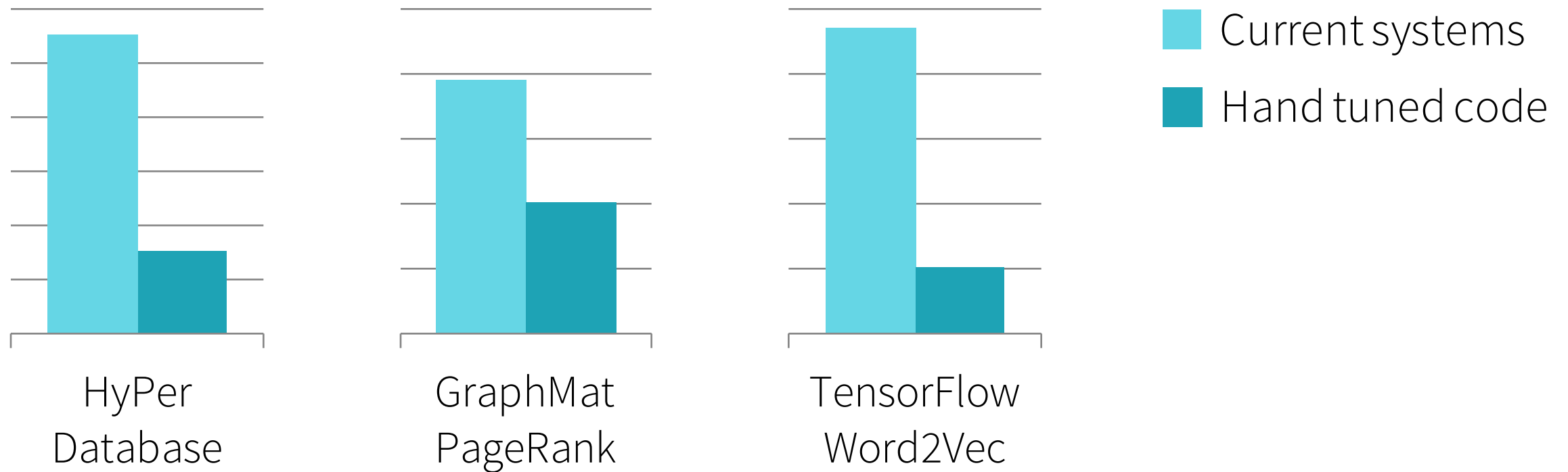# Recent Additions

Whole-stage code generation

- Fuse across multiple operators

| | |
|---|---|
| Spark 1.6 | 14M rows/s |
| Spark 2.0 | 125M rows/s |

Optimized input / output

- Apache Parquet + built-in cache

| | |
|---|---|
| Parquet in 1.6 | 11M rows/s |
| Parquet in 2.0 | 90M rows/s |

# Not Limited to Spark

Results from Nested Vector Language (NVL) project at MIT



HyPer
Database

GraphMat
PageRank

TensorFlow
Word2Vec

Current systems

Hand tuned code

databricks™

# Challenges

How to get this high performance while keeping the ease of use for non-programmers?

Can optimizations compose across libraries / systems?

databricks™

# Cloud Delivery

# The Public Cloud is Here

Many Fortune 100 companies have multiple PB of data in S3

Amazon Web Services up to $10B revenue

Especially attractive for big data
  • 51% of respondents in 2015 Spark survey run on public cloud

# Benefits

## For cloud users:

- Purchase an end-to-end experience, not just bits
- Rapidly experiment with new solutions (same data & infrastructure)

## For software vendors:

- Better products: end-to-end service, high visibility
- Fast iteration and uniform adoption

databricks™

# Challenges

Requires new way to build software that is not well understood by researchers (or traditional software companies)

- Multi-tenant: with untrusted tenants
- Highly available, yet with continuous updates
- Highly monitored for billing and security

databricks™

# Example Challenges

Deploying updates while keeping the service up
- And rolling back if needed!

Knowing whether the service is up

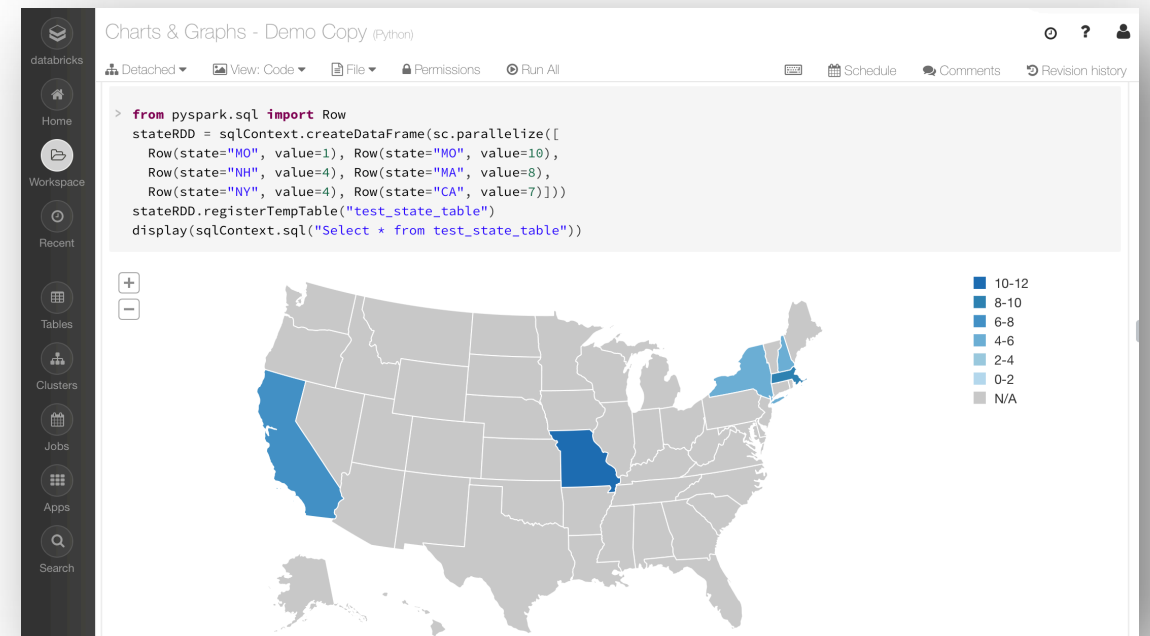Unexpected use, especially by code calling APIs

Performance isolation of tenants at all levels

Little academic research these

databricks™

# Example: Databricks

End-to-end data processing platform based around Apache Spark

Access control, collaboration, auditing, production workflows

200+ customers and thousands of individual users

# Lessons

## Cloud development model is superior

- Two week releases, immediate feedback, visibility

## State management is very hard at scale

- Per-tenant configuration, local data, VM images, etc

## Careful testing strategy is crucial

- Feature flags, stress tests, 70/20/10 testing pyramid

## Design to maximize dogfooding

databricks™

# Research Perspective

Computer systems is largely a social field: about interactions between users ↔ machines, users ↔ users, and machines ↔ machines

Cloud greatly changes the way users develop and consume software

Not much research beyond using it to parallelize stuff

databricks™

# Example Research Problems

Composing security interfaces of different cloud providers

- E.g. Databricks access controls + Amazon IAM

Deterministic updates and rollback for complex systems

"Elastic-first" systems for price and demand variability

databricks™

# Conclusion

Big data systems made great strides since they first came out

    **+** They're used well beyond tech companies

    **−** Not fully keeping up with new users & hardware

The cloud offers fantastic opportunities for research

    **+** People can try your new thing in production right away!

    **−** Not much research fully embraces it

databricks™

# Thanks!

## Databricks is Hiring

Full-timers and interns

matei@databricks.com

databricks™