# State of Spark, and where it is going

Reynold Xin 辛湜
@hashjoin
2015-12-10, Beijing BDTC

databricks™

# About Databricks
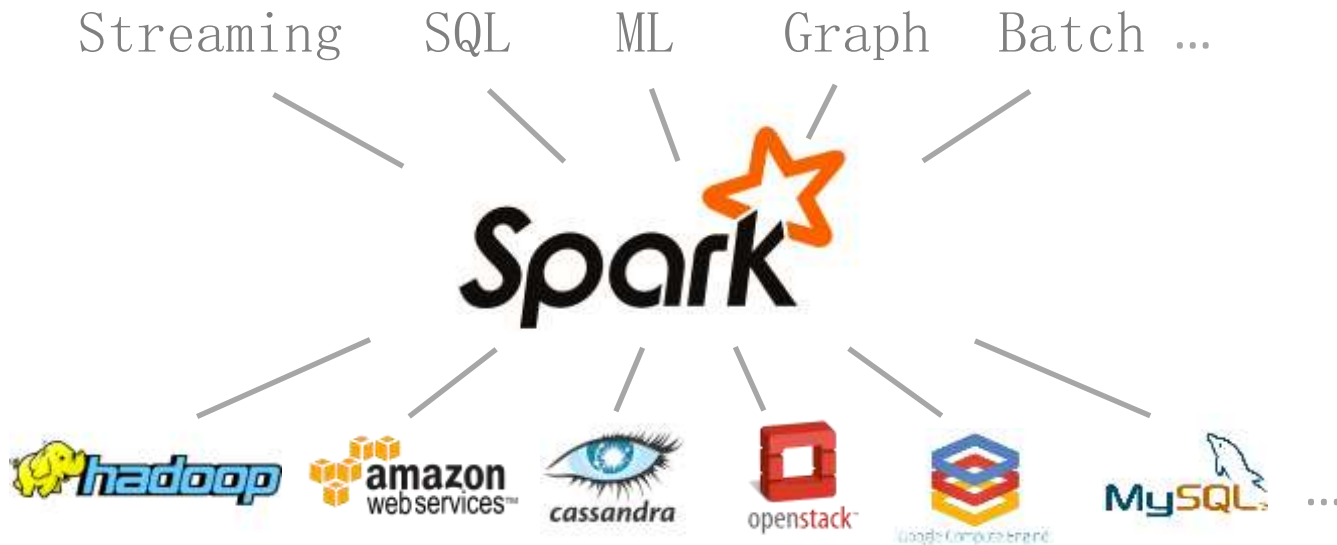
Founded by creators of Spark in 2013

Cloud service for end-to-end data processing

- Interactive notebooks, dashboards,
  production jobs, security, …

# Our Goal for Spark

Unified engine across data workloads and platforms

Streaming    SQL    ML    Graph    Batch …



databricks™

# Spark "Hall of Fame"

★

**LARGEST CLUSTER**

Tencent
(8000+ nodes)

★

**LARGEST SINGLE-DAY INTAKE**

Tencent
(1PB+ /day)

★

**LONGEST-RUNNING JOB**

Alibaba
(1 week on 1PB+ data)

★

**LARGEST SHUFFLE**

Databricks PB Sort
(1PB)

★

**MOST INTERESTING APP**

Jeremy Freeman
Mapping the Brain at Scale
(with lasers!)

databricks™

Based on Reynold's personal knowledge

# A Great Year for Spark

Most active open source project in big data

New language: R

Widespread industry support & adoption

databricks™

# IBM calls Apache Spark "most important new open source project in a decade"

June 15, 2015 Written by Business Cloud News

🖨 Print    ✉ Email

IBM said it will throw its weight behind Apache Spark, an open source community developing a processing engine for large-scale datasets, putting thousands of internal developers to work on Spark-related projects and contributing its machine learning technology to the code ecosystem.

Spark, an Apache open source project born in 2009, is essentially an engine that can process vast amounts of data very quickly. It runs in Hadoop clusters through YARN or as a standalone deployment and can process data in HDFS, HBase, Cassandra, Hive, and any Hadoop InputFormat; it currently supports Scala, Java and Python.

IBM is throwing its weight behind Apache Spark in a bid to bolster its IoT strategy

It is designed to perform general data

**databri**

"Spark is the Taylor
Swift
  of big data
software."

- Derrick Harris, Fortune

"Spark是大数据中的 Angelababy."

- Derrick Harris, Fortune

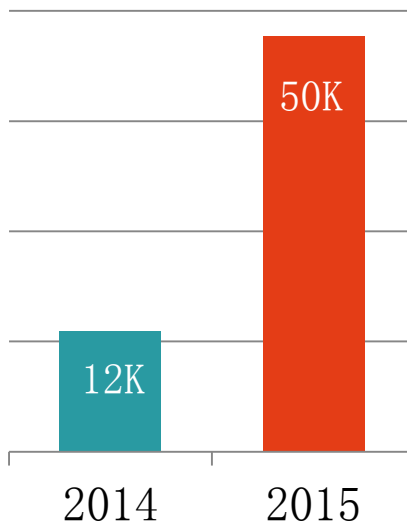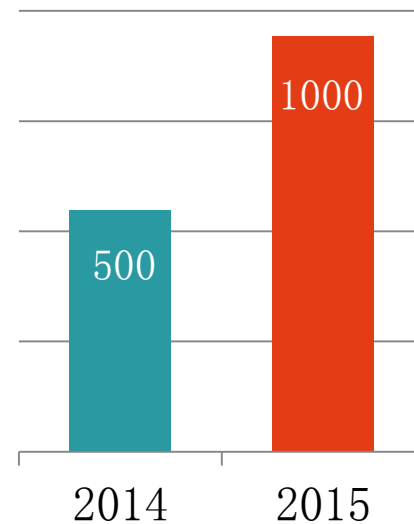# Meetup Groups: December 2014

databricks™

# Meetup Groups: December 2015

# Open Source Ecosystem

Applications

Sparkling

$H_2O$

IP[y]

Apache Ambari

thunder

mahout

sqoop

HIVE

Spark

hadoop

kubernetes

MESOS

docker

spring

openstack

HDFS

mongoDB

Parquet

cassandra

elasticsearch.

TACHYON

MySQL

PostgreSQL

HIVE

APACHE HBASE

SequoiaDB

kafka

Environments

Data Sources

# Users

## 1000+ companies



# Distributors + Apps

## 50+ companies

# Spark Survey 2015

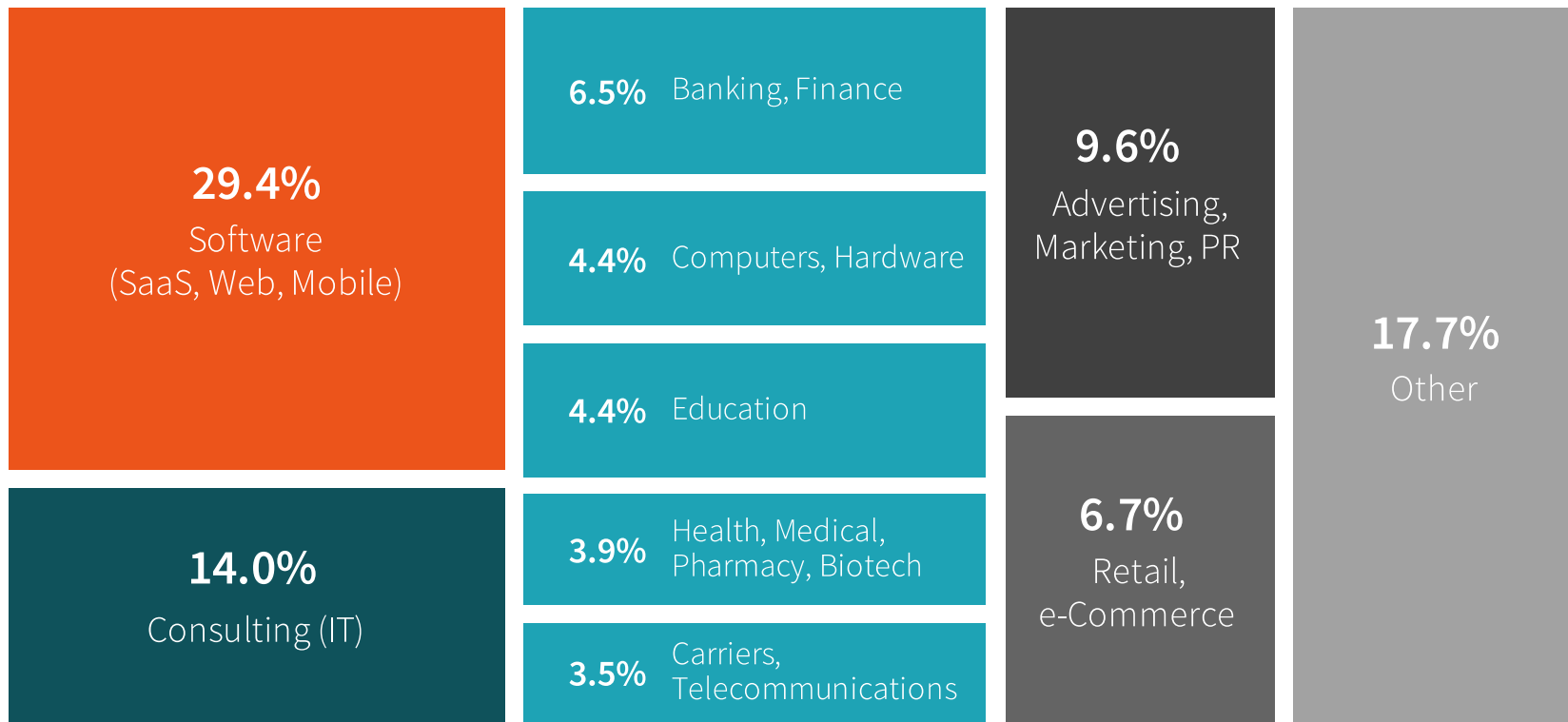# Databricks Survey

1400 respondents from 840 companies
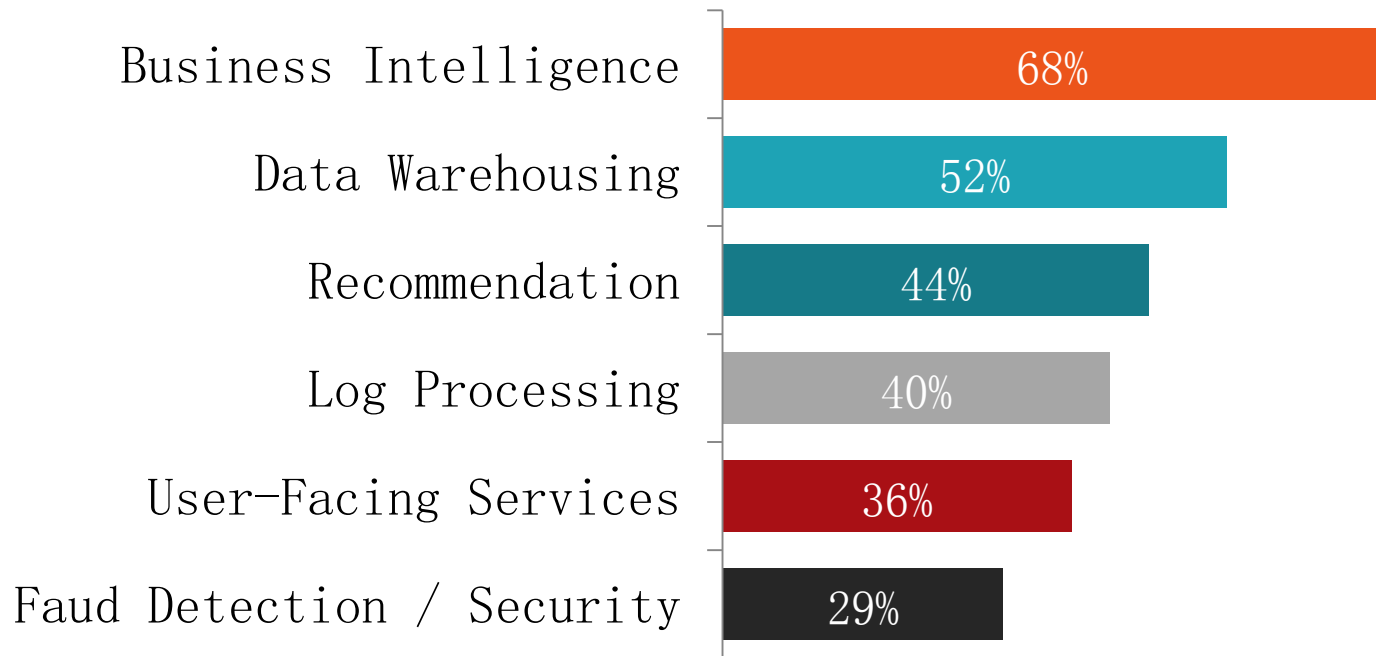
Three trends:
1) Diverse applications
2) More runtime environments
3) More types of users

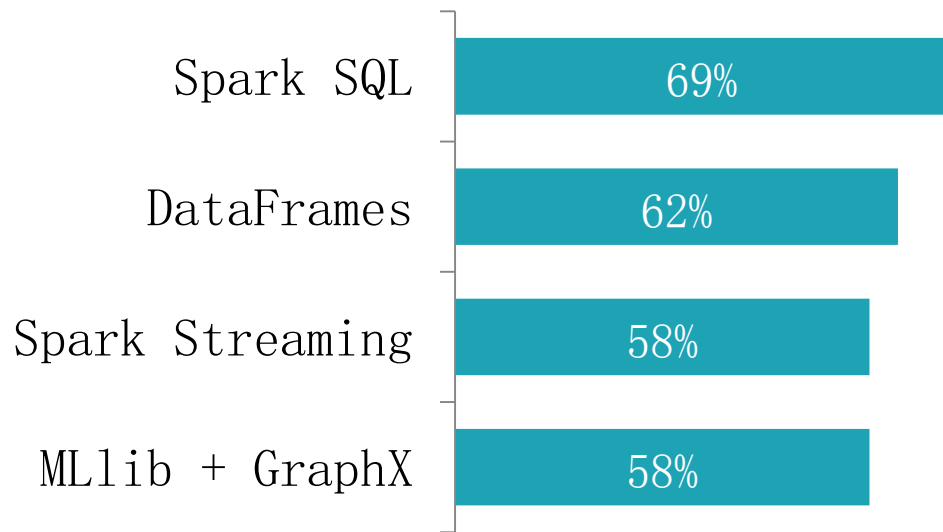databricks™

# Industries Using Spark

**29.4%**
Software
(SaaS, Web, Mobile)

**14.0%**
Consulting (IT)

**6.5%** Banking, Finance

**4.4%** Computers, Hardware

**4.4%** Education

**3.9%** Health, Medical, Pharmacy, Biotech

**3.5%** Carriers, Telecommunications

**9.6%**
Advertising, Marketing, PR

**6.7%**
Retail, e-Commerce

**17.7%**
Other

databricks™

# Top Applications



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

databricks™

# Spark Components Used

Spark SQL — 69%

DataFrames — 62%

Spark Streaming — 58%

MLlib + GraphX — 58%

75%

of users use more than one component

databricks™

# Diverse Runtime Environments

| MapReduce |
|-----------|
| HDFS |

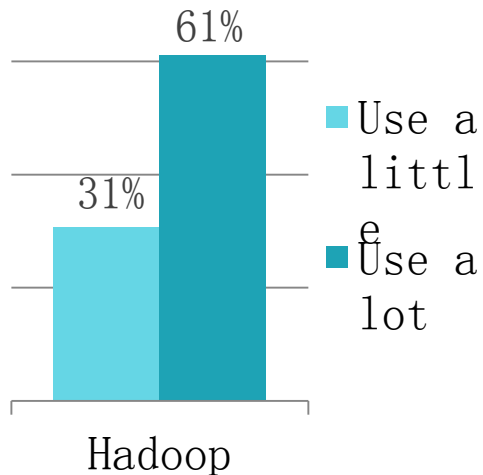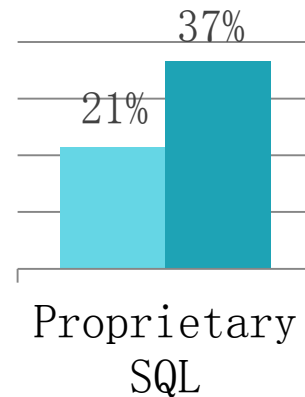| Spark | | |
|-------|---|---|
| HDFS | NoSQL e.g. Cassandra | SQL e.g. Oracle |

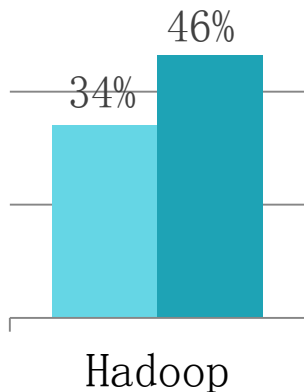Hadoop: combined
compute + storage

Spark: independent
of storage layer

databricks™

# Diverse Runtime Environments

# Diverse Runtime Environments
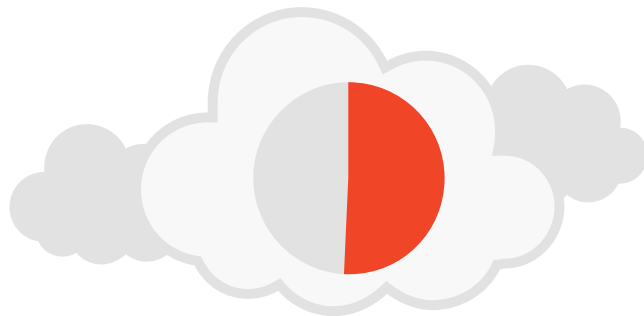
## Cluster Managers

**48%**
Standalone mode

**40%**
YARN

**11%**
Mesos

**51%**
on a public cloud

databricks™

# Diversity of Users



## Languages Used: 2014

84% Scala
38% Java
38% python

## Languages Used: 2015

71% Scala
31% Java
58% python
18% R

databricks™

# Fastest Growing Components

**+280%**

increase in
Windows users

**+56%**

production use
of Streaming

**+380%**

production
use of SQL

databricks™

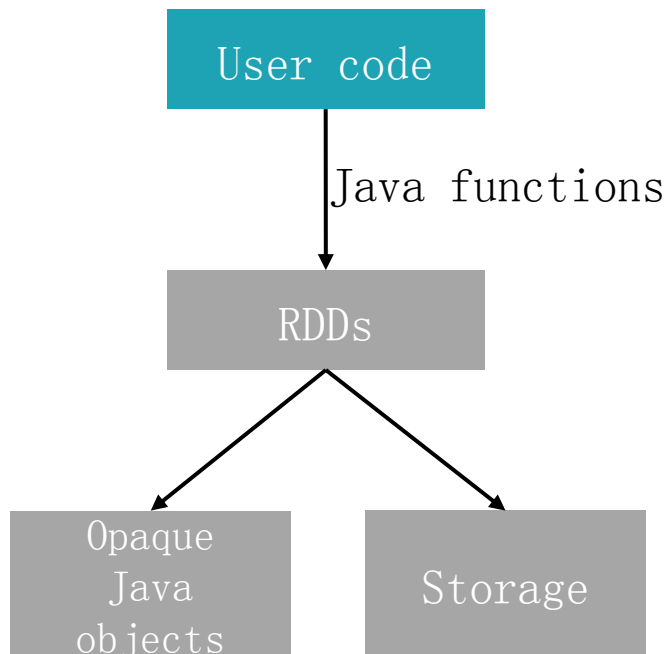# Are We Done?

No! Development is faster than ever.
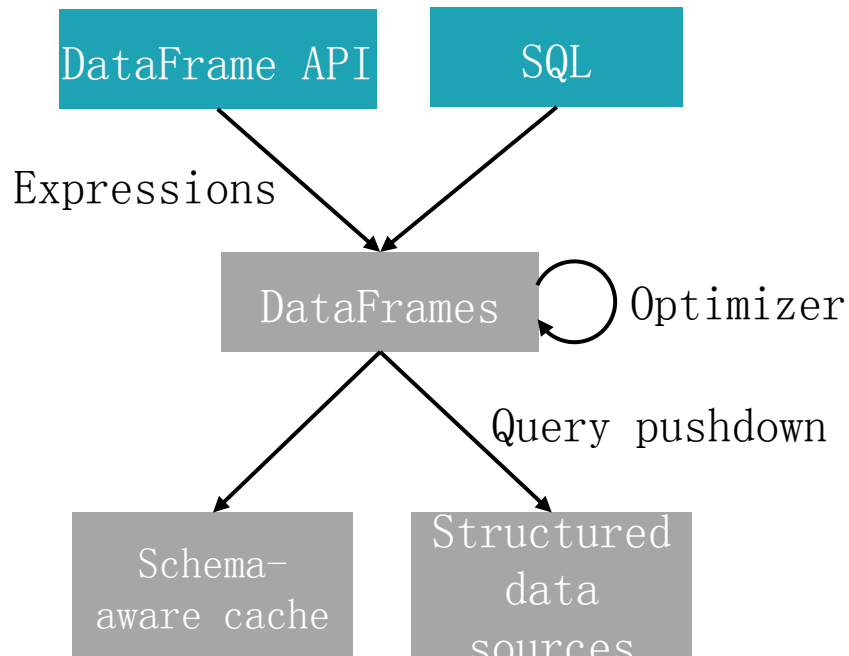
Biggest technical change in 2015 was DataFrames
- Moves many computations onto the relational Spark SQL optimizer

Enables both new APIs and more optimization, which is now happening through Project Tungsten

databricks™

# Traditional Spark

```
User code
```

Java functions

```
RDDs
```

```
Opaque
Java
objects
```

```
Storage
```

# DataFrames

```
DataFrame API
```

```
SQL
```

Expressions

```
DataFrames
```
Optimizer

Query pushdown

```
Schema-
aware cache
```

```
Structured
data
sources
```

# 3 Things to Look Forward To

# Dataset API in Spark 1.6 (SPARK-9999)

Typed interface over DataFrames / Tungsten

```scala
case class Person(name: String, age: Int)

val dataframe = read.json("people.json")
val ds: Dataset[Person] = dataframe.as[Person]

ds.filter(p => p.name.startsWith("M"))
    .groupBy("name")
    .avg("age")
```

databricks™

# Streaming DataFrames

Easier-to-use APIs (batch, streaming, and interactive)

```
val stream = read.kafka("...")
stream.window(5 mins, 10 secs)
    .agg(sum("sales"))
    .write.jdbc("mysql://...")
```

And optimizations:
- Tungsten backends
- native support for out-of-order data
- data sources and sinks

databricks™

# Microsoft Azure

BLOG > ANNOUNCEMENTS , VIRTUAL MACHINES

## Largest VM in the Cloud

THURSDAY, JANUARY 8, 2015

**DREW MCDANIEL**
Principal Program Manager, Azure

## G-Series Size Details

| VM Size | Cores | RAM |
|---------|-------|-----|
| Standard_G1 | 2 | 2 |
| Standard_G2 | 4 | 5 |
| Standard_G3 | 8 | 11 |
| Standard_G4 | 16 | 22 |
| Standard_G5 | 32 | 448 GiB | 6596 GB | 64 |

# AWS Announces X1 Instances For EC2 With 2TB Of Memory, Launching Next Year

Posted Oct 8, 2015 by **Frederic Lardinois** (@fredericl)

926
SHARES

**Introducing X1**

AVAILABLE IN THE FIRST ... 2016

Amazon today announced a massive new instance type for its AWS EC2 compute service. The

## CrunchBase

**Amazon**

FOUNDED
1994

OVERVIEW
Amazon is an e-commerce retailer former
to provide consumers with products in tw
It offers users with merchandise and con
purchased for resale from vendors and th
by third-party sellers. Operating in North
and international markets, Amazon prov
services through websites such as amazo
amazon.ca. It also enables authors, music
filmmakers, ...
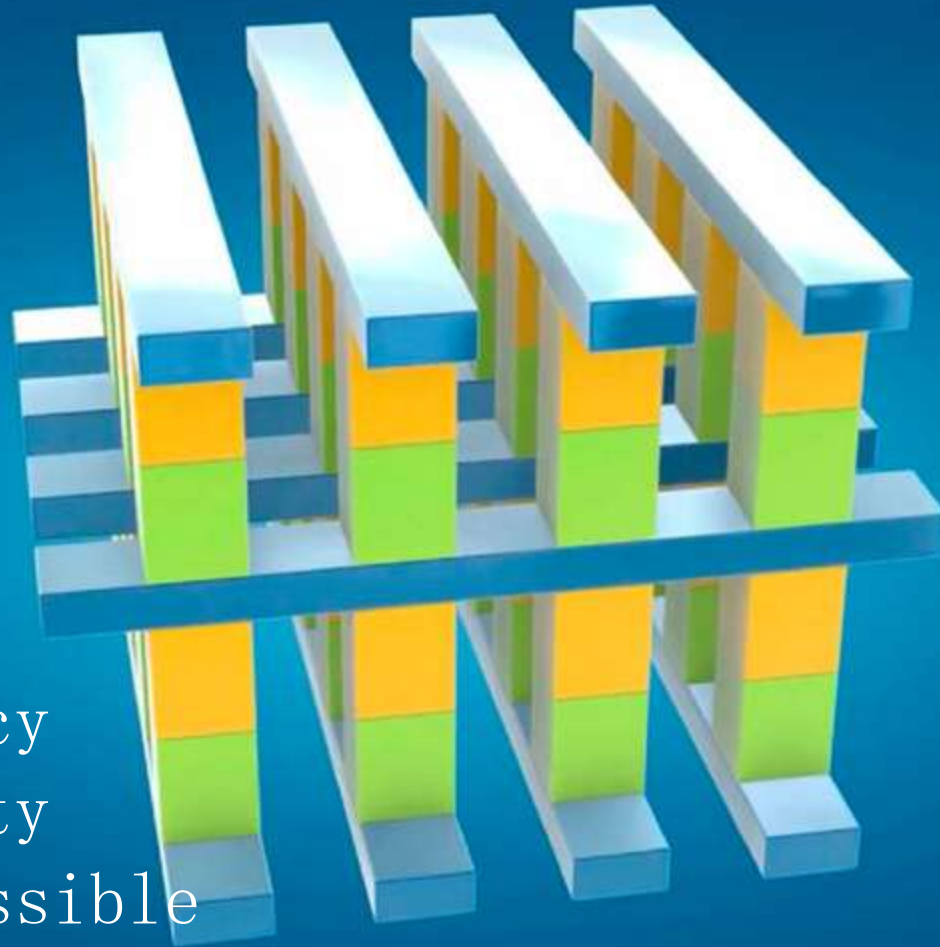
LOCATION
Seattle, WA

CATEGORIES
E-Commerce, Crowdsourcing, Groceries, Co
Goods, Delivery, Software, Retail, Internet

FOUNDERS
Jeff Bezos

3D XPoint

- DRAM latency
- SSD capacity
- Byte addressible

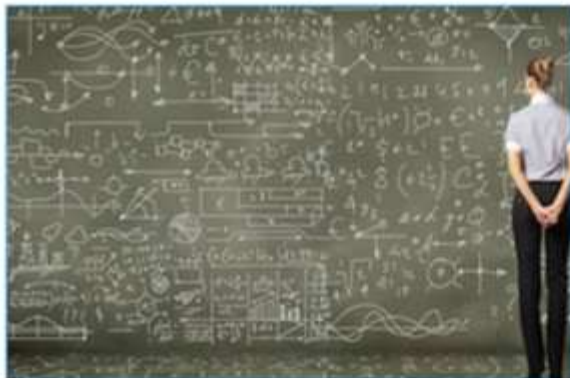# Unified API, One Engine, Automatically Optimized



language frontend

| SQL | Python | Java/Scala | R | ... |

DataFrame Logical Plan

Tungsten backend

| JVM | LLVM | SIMD | 3D XPoint | ... |

databricks™

# Introduction to Big Data with Apache Spark

Learn how to apply data science techniques using parallel programming in Apache Spark to explore big (and small) data.

## Berkeley
UNIVERSITY OF CALIFORNIA

## About this course

2 Reviews   4.5/5 ★★★★★

### This is an Archived Course

EdX keeps courses open for enrollment after they end to allow learners to explore content and continue learning. *All features and materials may not be all available*. Check back often to see when new course start dates are announced.

⊕ See more

## What you'll learn

| | | |
|---|---|---|
| ⏱ | Length: | 5 weeks |
| 🕐 | Effort: | 5 - 7 hours per week |
| 🏷 | Price: | FREE Verified Certificate option closed |
| 🏛 | Institution: | UC BerkeleyX |
| 🎓 | Subject: | Computer Science |

# 谢谢！

@rxin

databricks™