



SPARK DEVELOPMENT BOOTCAMP

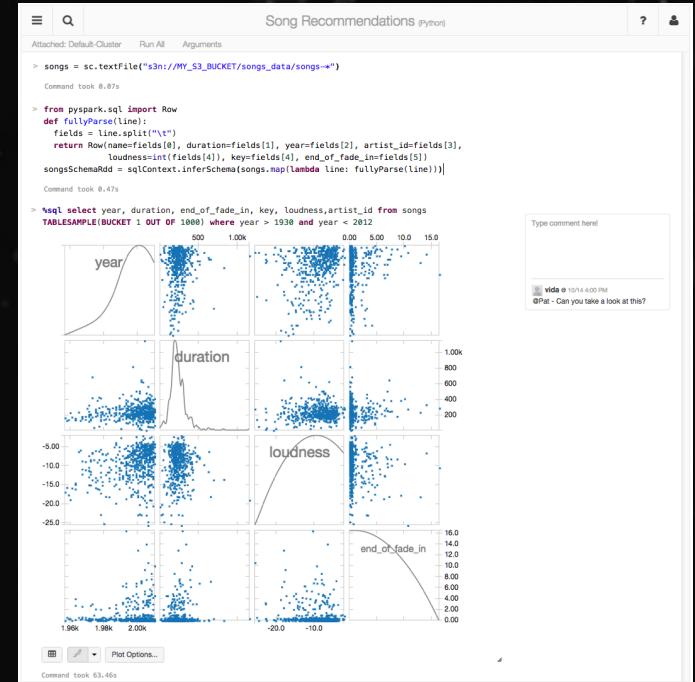
September 2015 Strata + Hadoop World NYC





making big data simple

- Founded in late 2013
- by the creators of Apache Spark
- Original team from UC Berkeley AMPLab
- Raised \$47 Million in 2 rounds
- ~60 employees
- We're hiring! (<http://databricks.workable.com>)
- Level 2/3 support partnerships with
 - Hortonworks
 - MapR
 - DataStax



Databricks Cloud:
“A unified platform for building Big Data pipelines
– from ETL to Exploration and Dashboards, to
Advanced Analytics and Data Products.”

Databricks has contributed more than **75%** of the code added to Spark in the past year



AGENDA

Day 1

- Welcome
- Login + Introductions
- Just enough Scala for Spark 
- DataFrames 
- RDD Fundamentals
- RDD lab 

Day 2

- Review
- Stages
- Accumulators & Broadcast Variables 
- Spark Runtime Architecture (w/ YARN)
- Memory & Persistence
- Transformations & Actions 
- Spark UI

Day 3

- Review
- Spark Streaming 
- MLlib
- GraphX 
- Config, tunings, settings
- Putting it all together 
- Certification Prep.

INSTRUCTOR: LAURENT WEICHLBERGER



Homepage: www.ompoint.com

LinkedIn: [@meherfalcon](https://www.linkedin.com/profile/view?id=2218412)

- 25+ years experience building & maintaining software systems
- Hadoop (at Hortonworks, and Cloudera), Cassandra (at DataStax), Couchbase (at Couchbase), Java & C# (at JP Morgan Chase, Pacific Controls, etc.) now Scala and much more ...
- Spark instructor for Databricks

100TB Daytona Sort Competition 2014



	Hadoop MR Record	Spark Record	Spark 1 PB
Data Size	102.5 TB	100 TB	1000 TB
Elapsed Time	72 mins	23 mins	234 mins
# Nodes	2100	206	190
# Cores	50400 physical	6592 virtualized	6080 virtualized
Cluster disk throughput	3150 GB/s (est.)	618 GB/s	570 GB/s
Sort Benchmark Daytona Rules	Yes	Yes	No
Network	dedicated data center, 10Gbps	virtualized (EC2) 10Gbps network	virtualized (EC2) 10Gbps network
Sort rate	1.42 TB/min	4.27 TB/min	4.27 TB/min
Sort rate/node	0.67 GB/min	20.7 GB/min	22.5 GB/min

Spark sorted the same data **3X faster**
using 10X fewer machines
than Hadoop MapReduce in 2013.

All the sorting took place on disk (HDFS) without
using Spark's in-memory cache!

More info:

<http://sortbenchmark.org>

<http://databricks.com/blog/2014/11/05/spark-officially-sets-a-new-record-in-large-scale-sorting.html>

Work by Databricks engineers: Reynold Xin, Parviz Deyhim, Xiangrui Meng, Ali Ghodsi, Matei Zaharia



GEAR SCIENCE ENTERTAINMENT BUSINESS SECURITY DESIGN OPINION MAGAZINE

ENTERPRISE

big data databricks google Hadoop

Startup Crunches 100 Terabytes of Data in a Record 23 Minutes

BY KLINT FINLEY 10.13.14 | 2:36 PM | PERMALINK



1.1k



789



75



565



565

GIGAOM

EVENTS RESEARCH

SIGN IN SUBSCRIBE

Cloud Data Media Mobile Science & Energy Social & Web Podcasts

Gigaom Research. Get unlimited market intelligence from over 200 in-depth reports.

MUST READS



Google launches Contributor, a crowdfunding tool for publishers



Net neutrality looks doomed in Europe before it even gets started

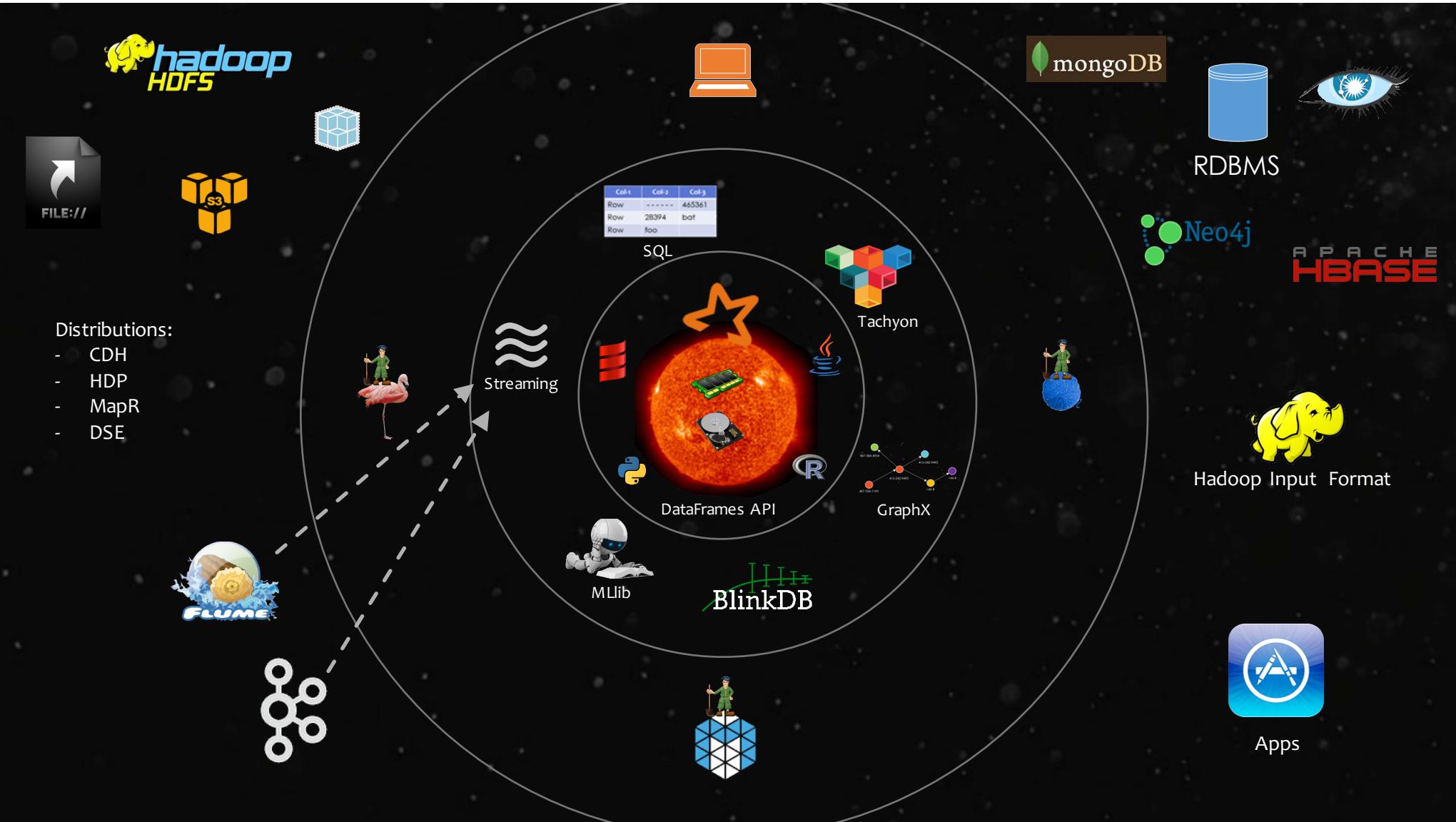


Five tech products that designers have fallen in love with

Databricks demolishes big data benchmark to prove Spark is fast on disk, too

by Derrick Harris Oct. 10, 2014 - 1:49 PM PST

1 Comment



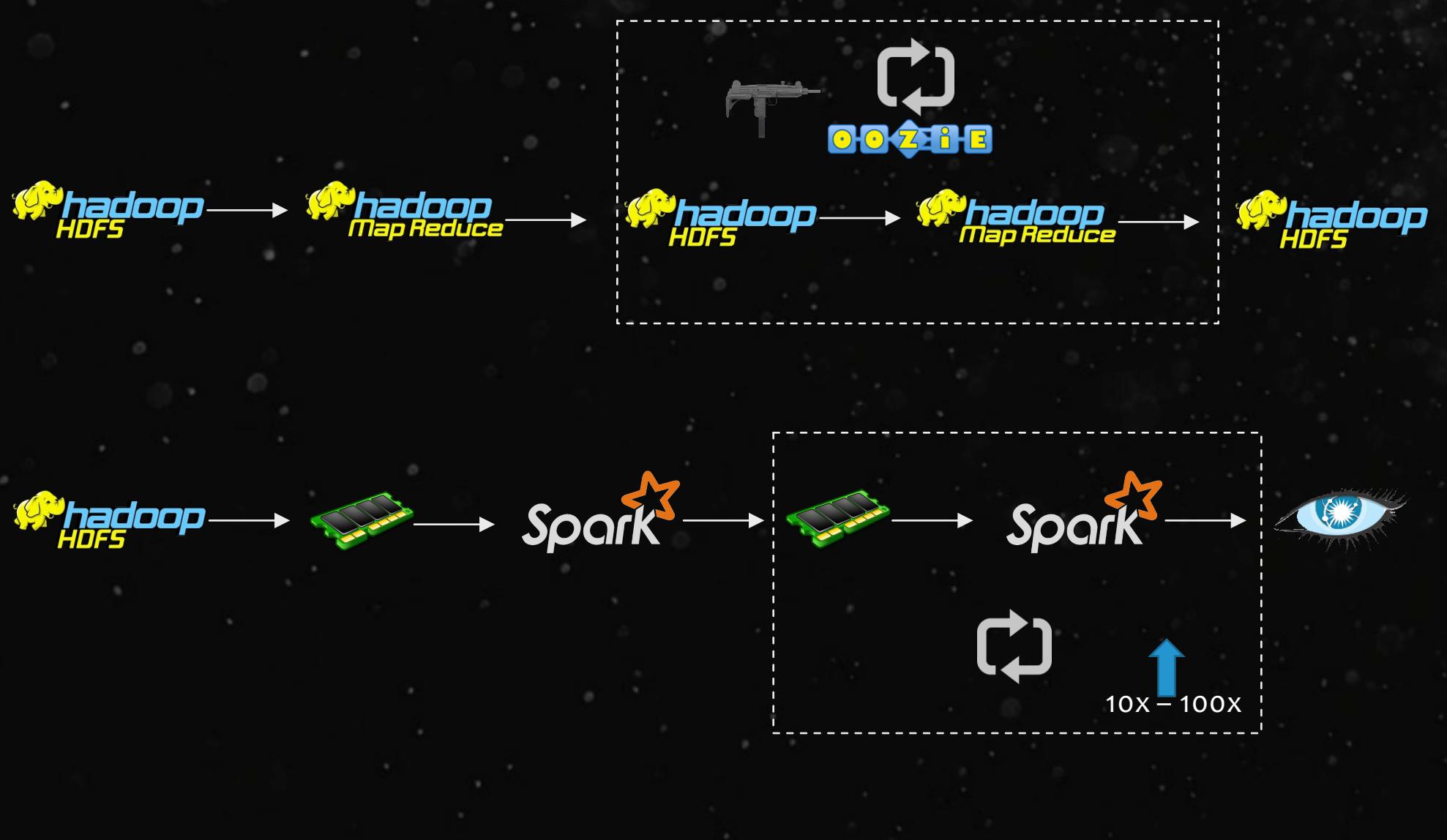


vs

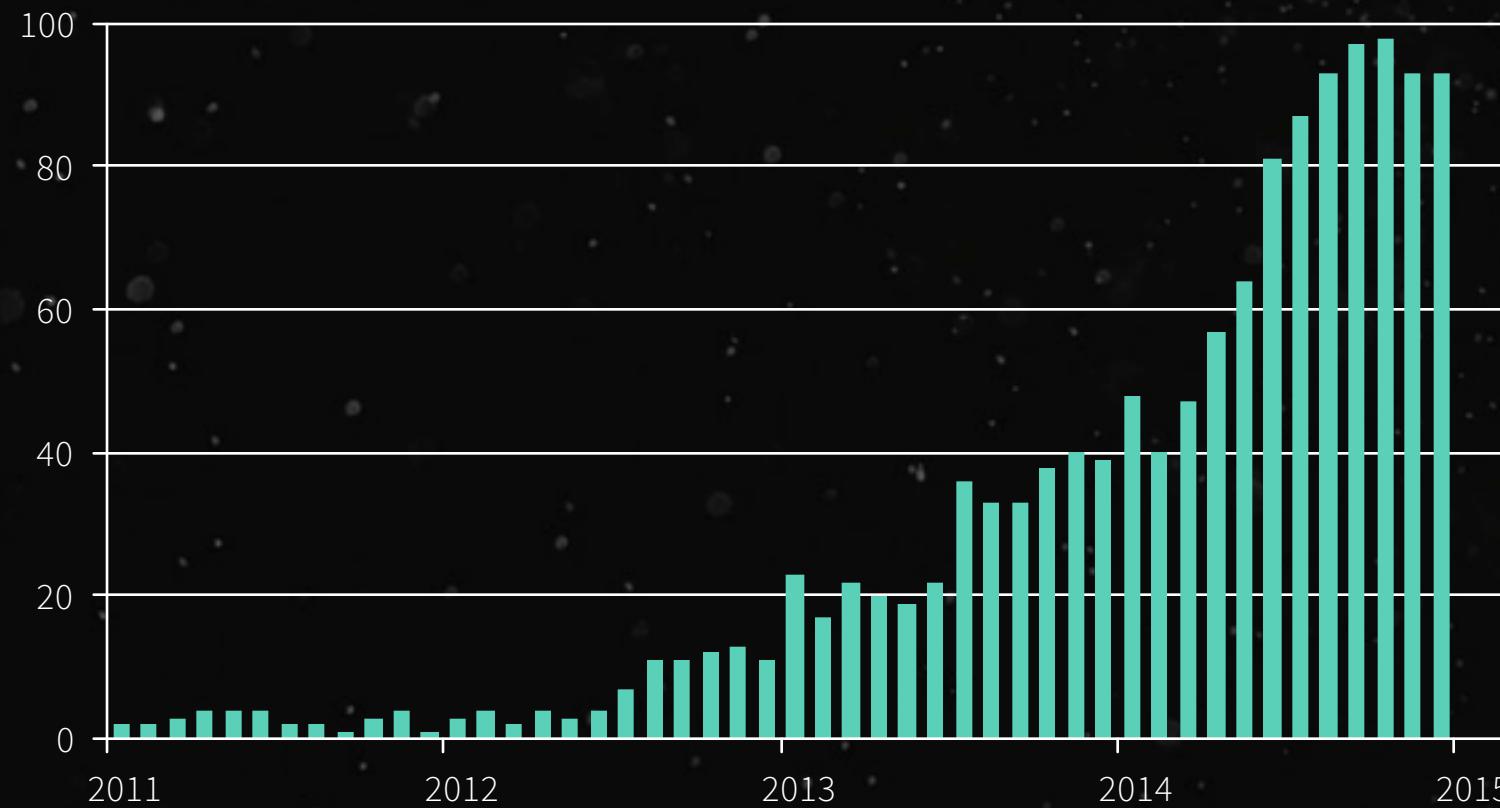


YARN





Contributors per month to Spark



COMMUNITY STATS

- Most active Apache project in contributors per month
- Most active open source project in a functional language

In a Nutshell, Apache Spark...

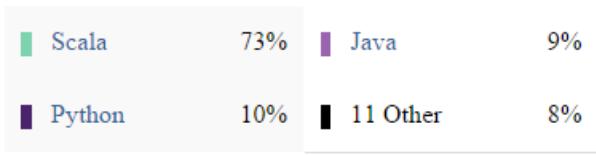
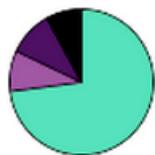
... has had 23,706 commits made by 797 contributors representing 563,770 lines of code

... is mostly written in Scala with a well-commented source code

... has a well established, mature codebase maintained by a very large development team with stable Y-O-Y commits

... took an estimated 154 years of effort (COCOMO model) starting with its first commit in ~~March, 2010~~ Aug 2009 ending with its most recent commit 1 day ago

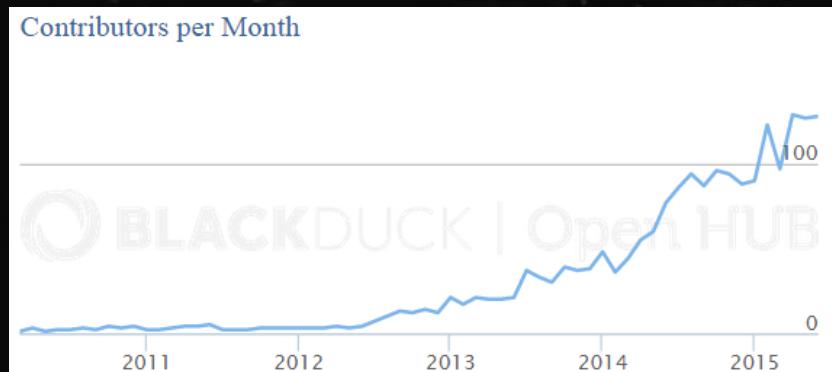
Languages



Lines of Code



Contributors per Month



...in June 2013

Source: openhub.net

INTRODUCTIONS

- 1. Your Name?**
- 2. Last programming language you used?**
- 3. When you wrote that code?**

4. Go ahead and login:

<https://classeast04.cloud.databricks.com/>

The screenshot shows the Databricks Home page. On the left, there is a sidebar with icons for Home, Workspace, Cluster, Jobs, Accounts, Recent, and Introductions. A large green arrow labeled '1' points down to the 'Introductions' link in the sidebar. Another green arrow labeled '2' points from the 'Introductions' link in the sidebar to the 'Introductions' section in the main content area. A third green arrow labeled '3' points from the 'Introductions' section in the main content area to the text 'Attached: community-1.4'. In the main content area, there is a heading 'Introductions (Scala)' and a paragraph of text. At the bottom of the main content area, there is a note about running cells and a 'shortcuts' link.

Home

Create New Features Getting Started

Introductions (Scala)

Attached: community-1.4 View: Notebook Run All Actions Permissions

This is just an Introductions sand box (play area). If you start a cell by typing "%md" you can write text in the cell. Scala code in it, to run the cell in the REPL against the attached Community Cluster. When we do lab work you will on it. Try writing something about yourself in a new cell. Click on the (+) sign to make a new cell here below.

Hi, I am Laurent Weichberger, your instructor. The last language I wrote in was Scala. That was on Friday, September 14, 2018. I wrote a StrataTest notebook, you can see it for yourself in Workspace >> Labs >> streaming >> lab04-Twtiter-ML

Shift+Enter to run [shortcuts](#)