

Install Apache Hadoop for Development/Production using Open source Distributions



Install Pure Apache Hadoop

- 1. Ubuntu or Centos (Linux)**
- 2. Installing JDK**
- 3. Download/Extract Hadoop**
- 4. Installing Hadoop**
- 5. Configure xml files**
- 6. Formatting HDFS**
- 7. Start Hadoop and then see Hadoop Web Console.**
- 8. Install Services (Hive , Sqoop ,HBASE,.....)**
- 9. Stop Hadoop (if needed.)**

```
<configuration>
<property>
  <name>yarn.resourcemanager.hostname</name>
  <value>localhost</value>
</property>
<property>
  <name>yarn.resourcemanager.scheduler.address</name>
  <value>localhost:8030</value>
</property>
<property>
  <name>yarn.resourcemanager.resource-tracker.address</name>
  <value>localhost:8031</value>
</property>
<property>
  <name>yarn.resourcemanager.address</name>
  <value>localhost:8032</value>
</property>
```

```
<configuration>
<property>
  <name>fs.defaultFS</name>
  <value>hdfs://localhost:9000</value>
</property>
</configuration>
```

Hadoop User Web Interfaces

Namenode Information - Mozilla Firefox

Namenode information

http://localhost:50070/dfshealth.html#tab-overview

Overview 'localhost:54310' (active)

Started:	5st Apr 18 15:53:55 PDT 2015
Versions:	2.6.0, re3496499ec08d220fb099dc5ed4c99cbf9e33cb1
Compiled:	2014-11-13T21:10Z by jenkins from (detached from e349649)
Cluster ID:	CID-c2f515ac-3300-45bc-8466-50110002bf7f
Block Pool ID:	BP-130729900-192.168.1.1-1429393391595

Summary

Security is off.

SafeMode is off.

1 files and directories, 0 blocks = 1 total filesystem object(s).

Heap Memory used 58.43 MB of 167.5 MB Heap Memory. Max Heap Memory is 889 MB.

Non-Heap Memory used 28.34 MB of 23.94 MB Committed Non-Heap Memory. Max Non-Heap Memory is 214 MB.

<http://localhost:50070/dfshealth.html#tab-startup-progress>

Namenode Information - Mozilla Firefox

Namenode information

http://localhost:50070/dfshealth.html#tab-datanode

Datanode Information

In operation

Node	Last contact	Admin State	Capacity	Used	New DFS Used	Remaining	Blocks	Block pool used	Failed Volumes	Version
ip-10.0.0.1:50010	1	In Service	414.29 GB	28 EB	125.83 GB	328.47 GB	0	28 EB (0%)	0	2.6.0

Decommissioning

Node	Last contact	Under replicated blocks	Blocks with no replicas	Under Replicated Blocks In Block under construction

Hadoop, 2014.

Hadoop SecondaryNameNode - Mozilla Firefox

Hadoop SecondaryNa... +

<http://localhost:50090/status.jsp> C Search

SecondaryNameNode

Version: 2.6.0. e3496499ecb8d220fba99dc5ed4c99c8f9e33bb1

Compiled: 2014-11-13T21:10Z by Jenkins from (detached from e349649)

SecondaryNameNode Status

```
Name Node Address      : localhost/127.0.0.1:54310
Start Time              : Sat Apr 18 16:43:38 PDT 2015
Last Checkpoint         : 79 seconds ago
Checkpoint Period       : 3600 seconds
Checkpoint Transactions: 1000000
Checkpoint Dirs         : [file:///app/hadoop/tmp/dfs/namesecondary]
Checkpoint Edits Dirs   : [file:///app/hadoop/tmp/dts/namesecondary]
```

Logs



Cluster	
About	Nodes
Node Labels	
Applications	
NEW	
NEW SAVING	
SUBMITTED	
ACCEPTED	
RUNNING	
FINISHED	
FAILED	
KILLED	
Scheduler	
Tools	

Cluster Metrics

Aps Submitted	Apps Pending	Apps Running	Apps Completed	Containers Running	Memory Used	Memory Total	Memory Reserved	VCores Used	VCores Total	VCores Reserved	Active Nodes	Decommissioned Nodes	Lost Nodes	Unhealthy Nodes	Rebooted Nodes
0	0	0	0	0	0 B	3 GB	0 B	0	8	0	1	0	0	0	0

Scheduler Metrics

Scheduler Type		Scheduling Resource Type		Minimum Allocation		Maximum Allocation	
Capacity Scheduler		[MEMORY]		<memory:1024, vCores:1>		<memory:8192, vCores:8>	

Cluster overview

Cluster ID: 1465353454396
ResourceManager state: STARTED
ResourceManager HA state: active
ResourceManager HA zookeeper connection state: ResourceManager HA is not enabled.
ResourceManager RMStateStore: org.apache.hadoop.yarn.server.resourcemanager.recovery.NullRMStateStore
ResourceManager started on: Wed Jun 08 09:37:34 +0700 2016
ResourceManager version: 2.7.2 from b165c4fe3a74265c792ce23f546c64004acf0e41 by Jenkins source checksum c63f7cc71b8f03249e35120fd7492d on 2016-01-26T00:08Z
Hadoop version: 2.7.2 from b165c4fe3a74265c792ce23f546c64004acf0e41 by Jenkins source checksum d0fd1a26633fa762bf87ec759ebe669c on 2016-01-26T00:08Z

Directory: /logs/ - Mozilla Firefox

Directory: /logs/ +

<http://localhost:50090/logs/> C Search

Directory: /logs/

[SecurityAuth-hduser.audit](#) 0 bytes Apr 18, 2015 3:40:58 PM
[hadoop-hduser-datanode-laptop.log](#) 72879 bytes Apr 18, 2015 4:44:13 PM
[hadoop-hduser-datanode-laptop.out](#) 718 bytes Apr 18, 2015 4:43:21 PM
[hadoop-hduser-datanode-laptop.out.1](#) 718 bytes Apr 18, 2015 3:53:49 PM
[hadoop-hduser-datanode-laptop.out.2](#) 718 bytes Apr 18, 2015 3:41:03 PM
[hadoop-hduser-namenode-laptop.log](#) 121216 bytes Apr 18, 2015 4:52:23 PM
[hadoop-hduser-namenode-laptop.out](#) 718 bytes Apr 18, 2015 4:43:16 PM
[hadoop-hduser-namenode-laptop.out.1](#) 718 bytes Apr 18, 2015 3:53:44 PM
[hadoop-hduser-namenode-laptop.out.2](#) 718 bytes Apr 18, 2015 3:40:58 PM
[hadoop-hduser-secondarynamenode-laptop.log](#) 51913 bytes Apr 18, 2015 4:52:38 PM
[hadoop-hduser-secondarynamenode-laptop.out](#) 718 bytes Apr 18, 2015 4:43:37 PM
[hadoop-hduser-secondarynamenode-laptop.out.1](#) 718 bytes Apr 18, 2015 3:54:06 PM
[hadoop-hduser-secondarynamenode-laptop.out.2](#) 718 bytes Apr 18, 2015 3:42:52 PM
[userlogs/](#) 4096 bytes Apr 18, 2015 4:52:22 PM
[yarn-hduser-nodemanager-laptop.log](#) 81625 bytes Apr 18, 2015 4:44:32 PM

Logged in as: dr.who :44:02 PM
:54:32 PM
:43:10 PM
:44:32 PM
:44:00 PM
:54:29 PM
:43:08 PM

About the Cluster

Hadoop Distribution

On-Premise

- Pure Apache Hadoop
- Cloudera
- Hortonworks
- MapR
- Pivotal
- IBM InfoSphere BigInsight

On-Cloud (Hadoop as a Service)

- Amazon EMR
- Microsoft Azure HDInsight
- Google Cloud Platform

The Forrester Wave: Big Data Hadoop

Distributions Q1 2016



The Forrester Wave: Big Data Hadoop

Cloud Q1 2016



Hadoop Development: Sandbox

Recommended

Cloudera Quickstart

Hortonworks Sandbox

On PC/Mac

Require to install VMWare Player or VirtualBox

On Linux Virtual Server

Need to install Docker Engine

<http://hortonworks.com/downloads/#sandbox>

Hortonworks Sandbox

Hortonworks Sandbox on a VM

HDP® 2.5 on Hortonworks Sandbox

[Tutorials](#) | [Release Notes](#) | [Import on Virtual Box](#) | MD5 : b08f1dce17ab3ae2431532be74bdbbbb

[DOWNLOAD FOR VIRTUALBOX](#)

[Tutorials](#) | [Release Notes](#) | [Import on VMware](#) | MD5 : fe1e91bc26e6879fdc6dde6e3778c262

[DOWNLOAD FOR VMWARE](#)

[Tutorials](#) | [Release Notes](#) | [Import on Docker](#) | MD5 : 2a710f236135e620ec8488a1229af07e

[DOWNLOAD FOR DOCKER](#)

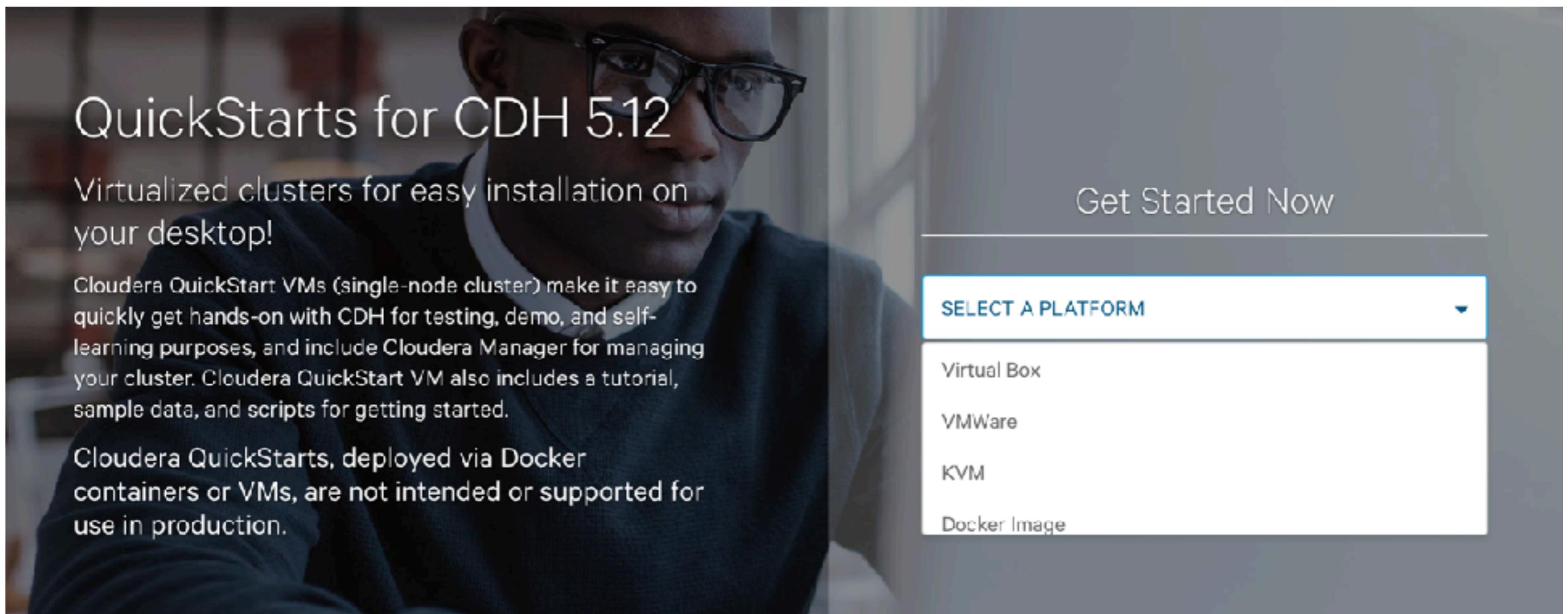
Hortonworks Sandbox in the Cloud

HDP 2.4 on Azure with Hortonworks Sandbox

[Tutorial: Sandbox on Azure](#) | Try it one month for free

[ONE MONTH TRIAL](#)

<http://www.cloudera.com/downloads.html>



QuickStarts for CDH 5.12

Virtualized clusters for easy installation on your desktop!

Cloudera QuickStart VMs (single-node cluster) make it easy to quickly get hands-on with CDH for testing, demo, and self-learning purposes, and include Cloudera Manager for managing your cluster. Cloudera QuickStart VM also includes a tutorial, sample data, and scripts for getting started.

Cloudera QuickStarts, deployed via Docker containers or VMs, are not intended or supported for use in production.

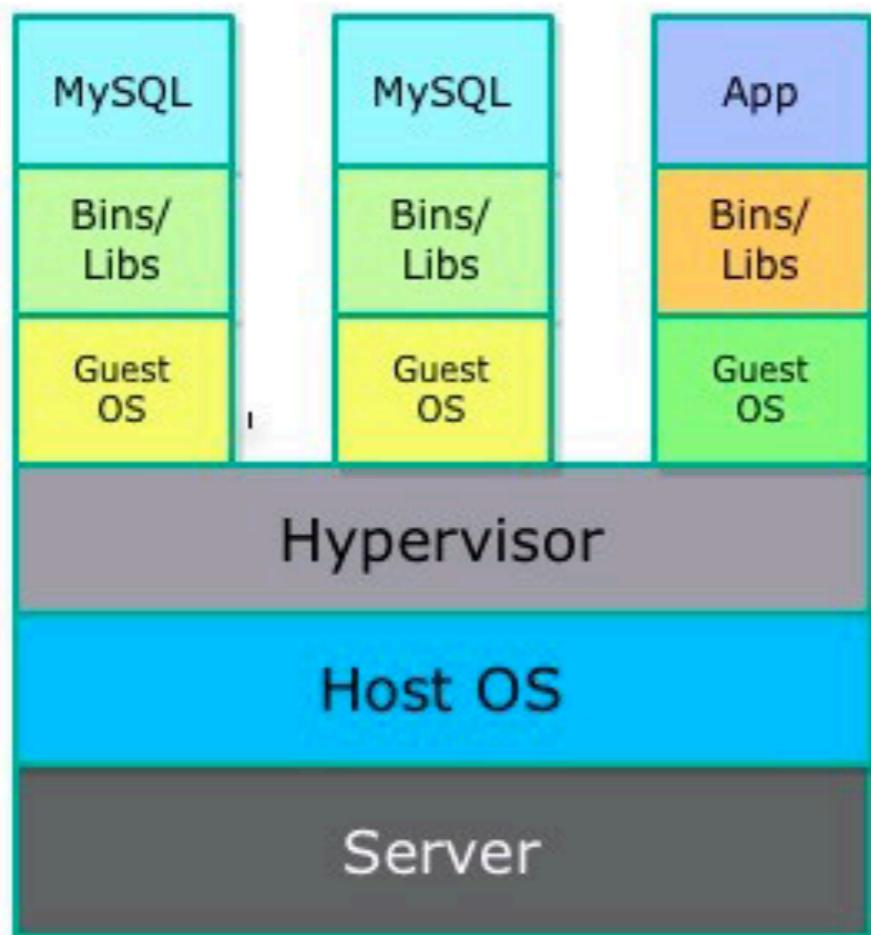
Get Started Now

SELECT A PLATFORM

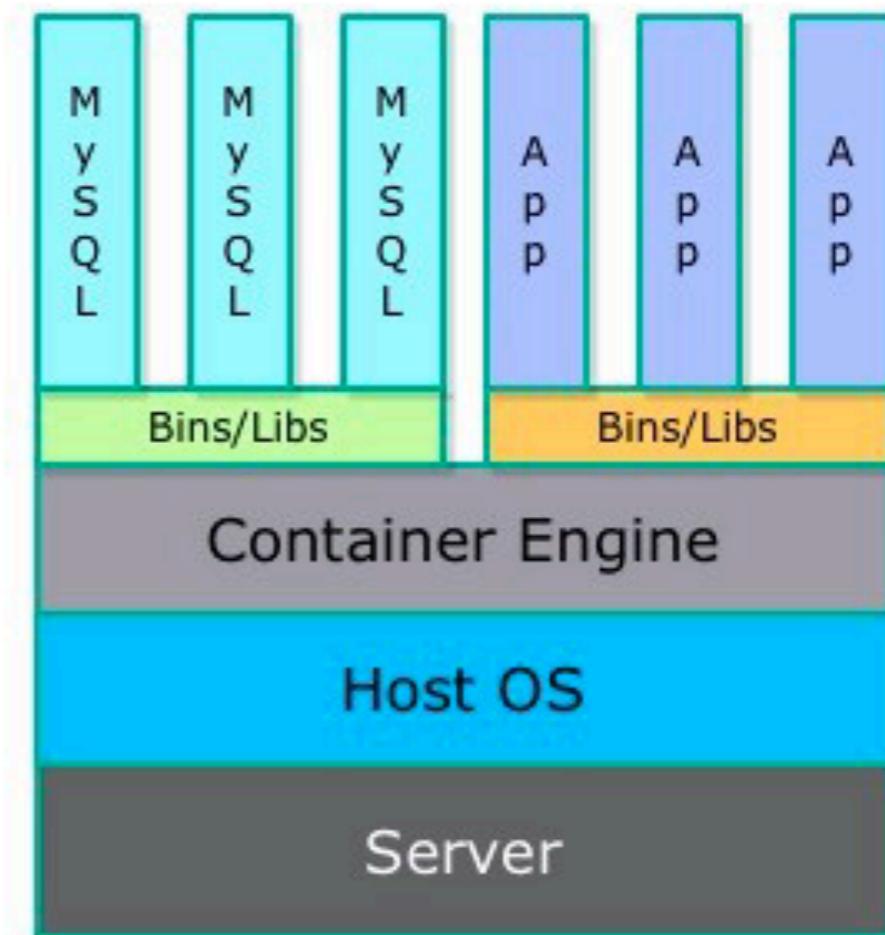
- Virtual Box
- VMWare
- KVM
- Docker Image

Install a Docker Engine

Virtual Machines



Containers



Pull Cloudera Quickstart

```
$ sudo docker pull cloudera/quickstart:latest
```

```
ubuntu@ip-172-31-30-238:~$ sudo docker pull cloudera/quickstart:latest
latest: Pulling from cloudera/quickstart
2cda82941cb7: Already exists
Digest: sha256:f91bee4cdfa2c92ea3652929a22f729d4d13fc838b00f120e630f91c941acb63
Status: Downloaded newer image for cloudera/quickstart:latest
ubuntu@ip-172-31-30-238:~$ █
```

```
$ sudo docker images
```

REPOSITORY	TAG	IMAGE ID	CREATED
VIRTUAL SIZE			
cloudera/quickstart	latest	2cda82941cb7	9 weeks ago
6.336 GB		—	

Run Cloudera quickstart

```
$ sudo docker run --hostname=quickstart.cloudera  
--privileged=true -t -i [OPTIONS] [IMAGE]  
/usr/bin/docker-quickstart
```

Example: \$ sudo docker run
--hostname=quickstart.cloudera --privileged=true -t -i -p
8888:8888 cloudera/quickstart /usr/bin/docker-quickstart

```
ubuntu@ip-172-31-30-238:~$ sudo docker run --hostname=quickstart.cloudera --privileged=true -t -i -p 8888:8888 -p 7180:7180 cloudera/quickstart /usr/bin/docker-quickstart  
Starting mysqld: [ OK ]  
  
if [ "$1" == "start" ] ; then  
  if [ "${EC2}" == 'true' ]; then  
    FIRST_BOOT_FLAG=/var/lib/cloudera-quickstart/.ec2-key-installed  
    if [ ! -f "${FIRST_BOOT_FLAG}" ]; then  
      METADATA_API=http://169.254.169.254/latest/meta-data  
      KEY_URL=${METADATA_API}/public-keys/0/openssh-key
```

Install Cloudera Cluster on AWS

Launch a virtual server on EC2 Amazon Web Services

Amazon Web Services

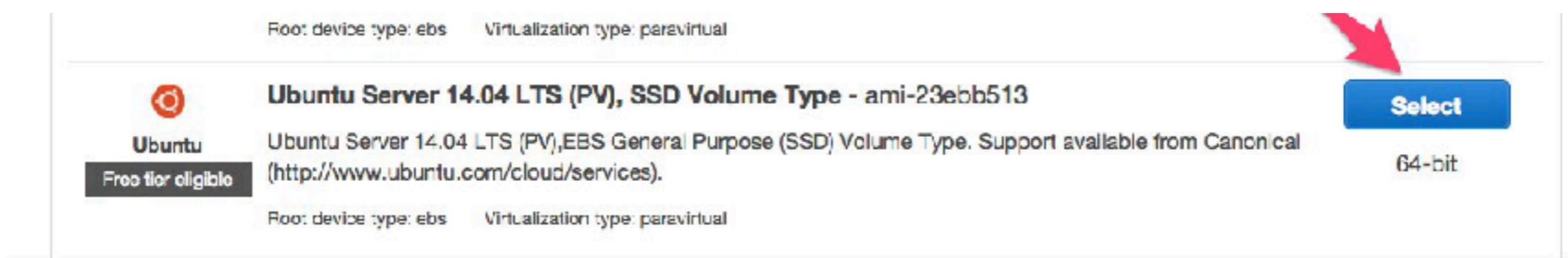
Compute	Administration & Security	Application Services
 EC2 Virtual Servers in the Cloud	 Directory Service Managed Directories in the Cloud	 SQS Message Queue Service
 Lambda PREVIEW Run Code in Response to Events	 Identity & Access Management Access Control and Key Management	 SWF Workflow Service for Coordinating Application Components
Storage & Content Delivery	 Trusted Advisor AWS Cloud Optimization Expert	 AppStream Low Latency Application Streaming
 S3	 CloudTrail	 Elastic Transcoder

Resource Groups

A resource group is a collection of resources that share one or more tags. Create a group for each project, application, or environment in your account.

[Create a Group](#) [Tag Editor](#)

Select an Amazon Machine Image (AMI) and Ubuntu Server 14.04 LTS (PV)



The screenshot shows the AWS Lambda console interface. A red arrow points to the 'Select' button next to the 'Ubuntu Server 14.04 LTS (PV), SSD Volume Type - ami-23ebb513' option. The details for this AMI are displayed:

- Root device type: ebs
- Virtualization type: paravirtual
- Provider: Ubuntu
- Image ID: ami-23ebb513
- Description: Ubuntu Server 14.04 LTS (PV), EBS General Purpose (SSD) Volume Type. Support available from Canonical (<http://www.ubuntu.com/cloud/services>).
- Architecture: 64-bit
- Status: Free tier eligible

Install Cloudera Cluster on AWS

Choose m3.xlarge Type virtual server

1. Choose AMI 2. Choose Instance Type 3. Configure Instance 4. Add Storage 5. Tag Instance 6. Configure Security Group 7. Review

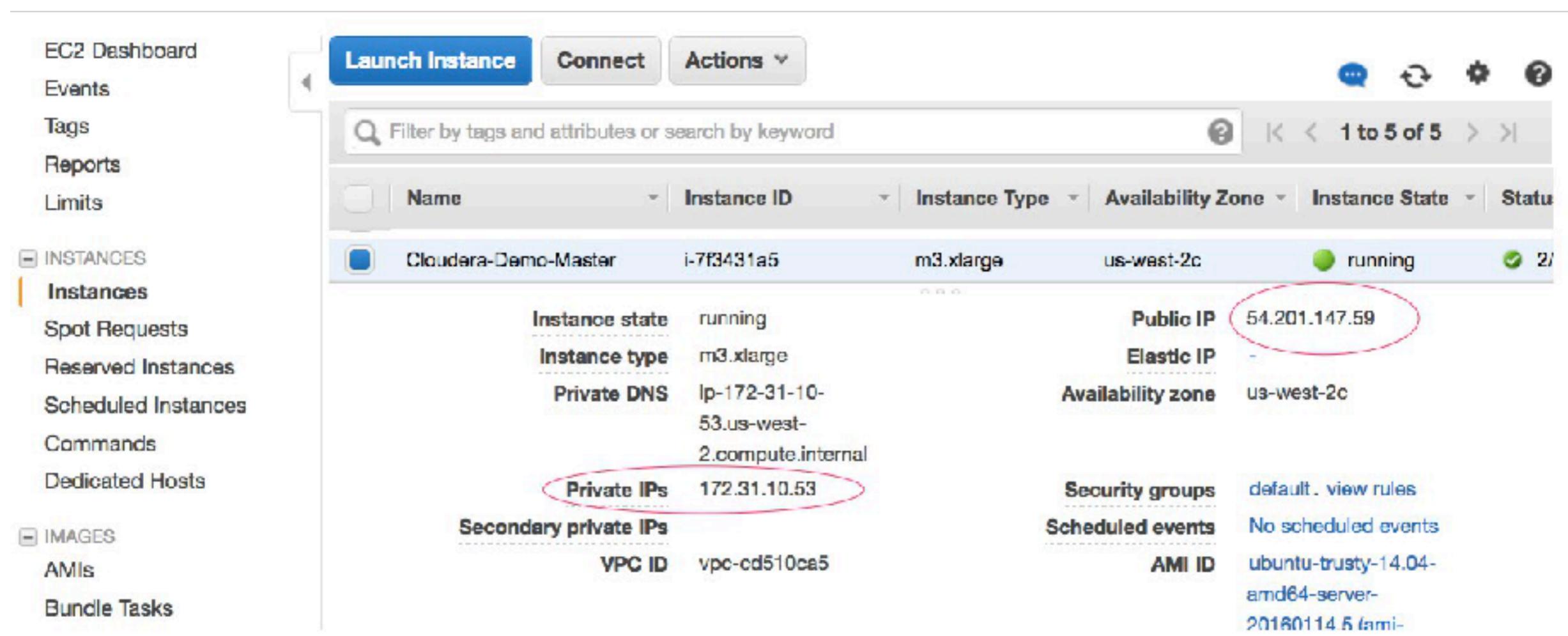
Step 2: Choose an Instance Type

					Avg. Price	Available	
<input type="checkbox"/>	Micro instances	t1.micro <small>Free tier eligible</small>	1	0.613	EBS only	-	Very Low
<input type="checkbox"/>	General purpose	t2.micro <small>Free tier eligible</small>	1	1	EBS only	-	Low to Moderate
<input type="checkbox"/>	General purpose	t2.small	1	2	EBS only	-	Low to Moderate
<input type="checkbox"/>	General purpose	t2.medium	2	4	EBS only	-	Low to Moderate
<input type="checkbox"/>	General purpose	m3.medium	1	3.75	1 x 4 (SSD)	-	Moderate
<input type="checkbox"/>	General purpose	m3.large	2	7.5	1 x 32 (SSD)	-	Moderate
<input checked="" type="checkbox"/>	General purpose	m3.xlarge	4	15	2 x 40 (SSD)	Yes	High
<input type="checkbox"/>	General purpose	m3.2xlarge	8	30	2 x 80 (SSD)	Yes	High

[Cancel](#) [Previous](#) [Review and Launch](#) [Next: Configure Instance Details](#)

Install Cloudera Cluster on AWS

Can also view details of the instance such as Public IP and Private IP



The screenshot shows the AWS EC2 Instances page. On the left, there's a sidebar with links like EC2 Dashboard, Events, Tags, Reports, Limits, INSTANCES (with Instances selected), Spot Requests, Reserved Instances, Scheduled Instances, Commands, Dedicated Hosts, and IMAGES (with AMIs selected). The main area has tabs for Launch Instance, Connect, and Actions. Below that is a search bar and a table with columns: Name, Instance ID, Instance Type, Availability Zone, Instance State, and Status. One row is selected, showing "Cloudera-Demo-Master" with Instance ID i-7f3431a5, m3.xlarge instance type, us-west-2c availability zone, running instance state, and a status of 2/2. To the right of the table are detailed instance settings: Instance state (running), Instance type (m3.xlarge), Private DNS (ip-172-31-10-53.us-west-2.compute.internal), Public IP (54.201.147.59), Elastic IP (-), Availability zone (us-west-2c), Security groups (default, view rules), Scheduled events (No scheduled events), and AMI ID (ubuntu-trusty-14.04-amd64-server-20160114.5 ami-). Two specific fields, "Private IPs" (172.31.10.53) and "Secondary private IPs", are circled in red.

Name	Instance ID	Instance Type	Availability Zone	Instance State	Status
Cloudera-Demo-Master	i-7f3431a5	m3.xlarge	us-west-2c	running	2/2

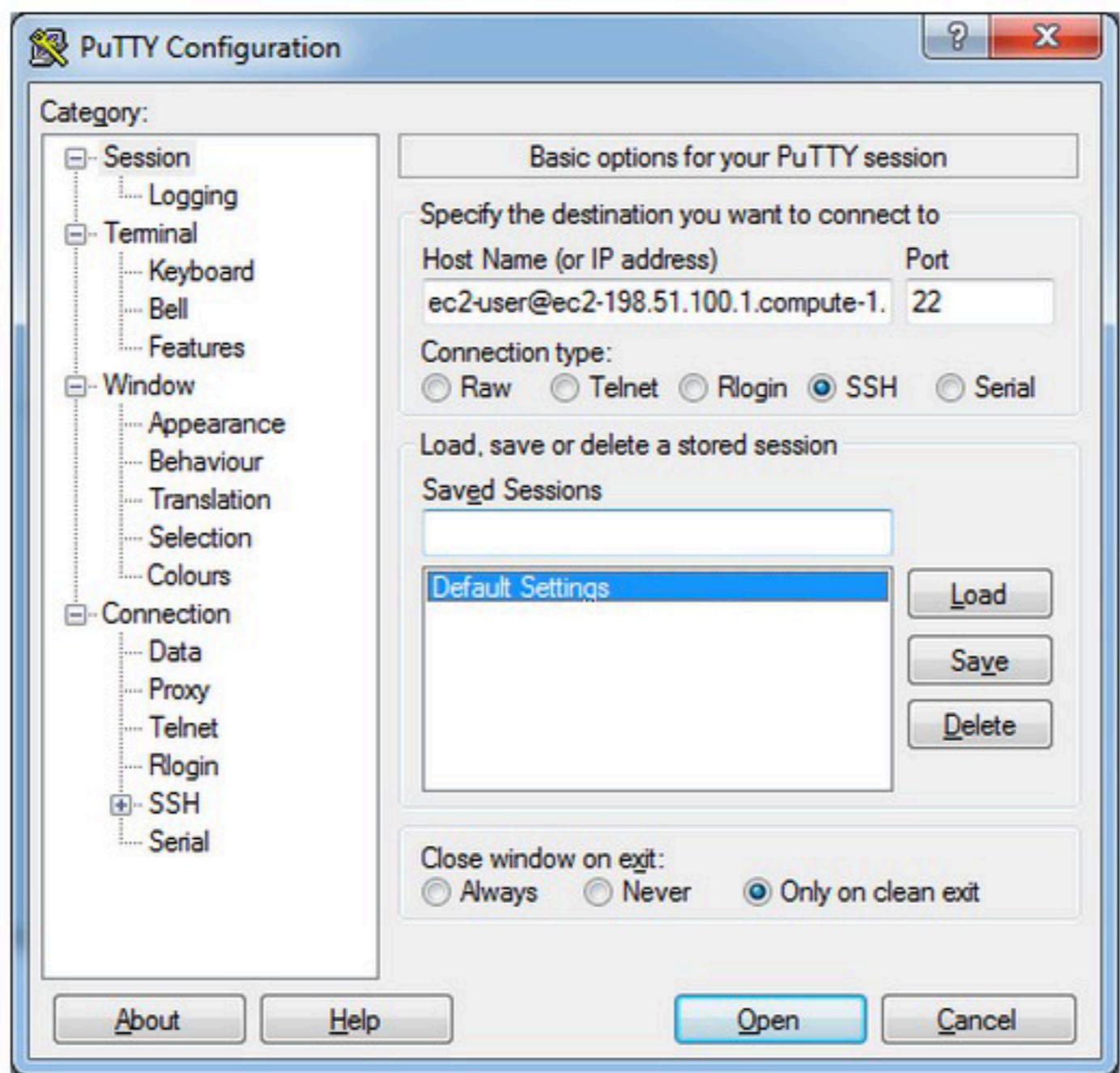
Private IPs: 172.31.10.53

Secondary private IPs

VPC ID: vpc-cd510ca5

Install Cloudera Cluster on AWS

Connect to an instance from Windows using Putty



Install Cloudera Cluster on AWS

Connect to the instance

Ubuntu comes with ABSOLUTELY NO WARRANTY, to the extent permitted by applicable law.

WARNING! Your environment specifies an invalid locale.

This can affect your user experience significantly, including the ability to manage packages. You may install the locales by running:

```
sudo apt-get install language-pack-UTF-8
```

or

```
sudo locale-gen UTF-8
```

To see all available language packs, run:

```
apt-cache search "^language-pack-[a-z] [a-z]$"
```

To disable this message for all users, run:

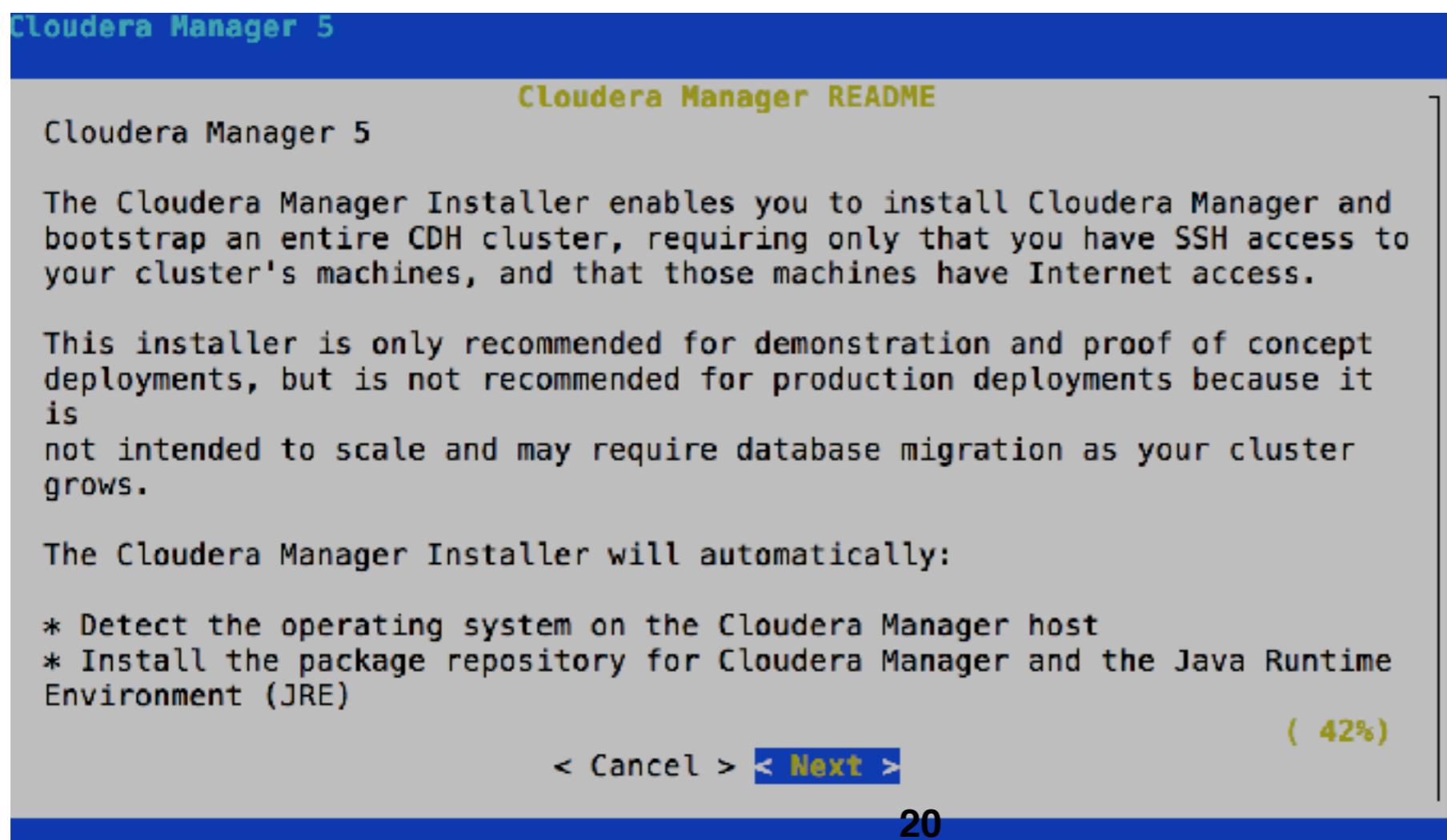
```
sudo touch /var/lib/cloud/instance/locale-check.skip
```

```
ubuntu@ip-172-31-1-242:~$
```

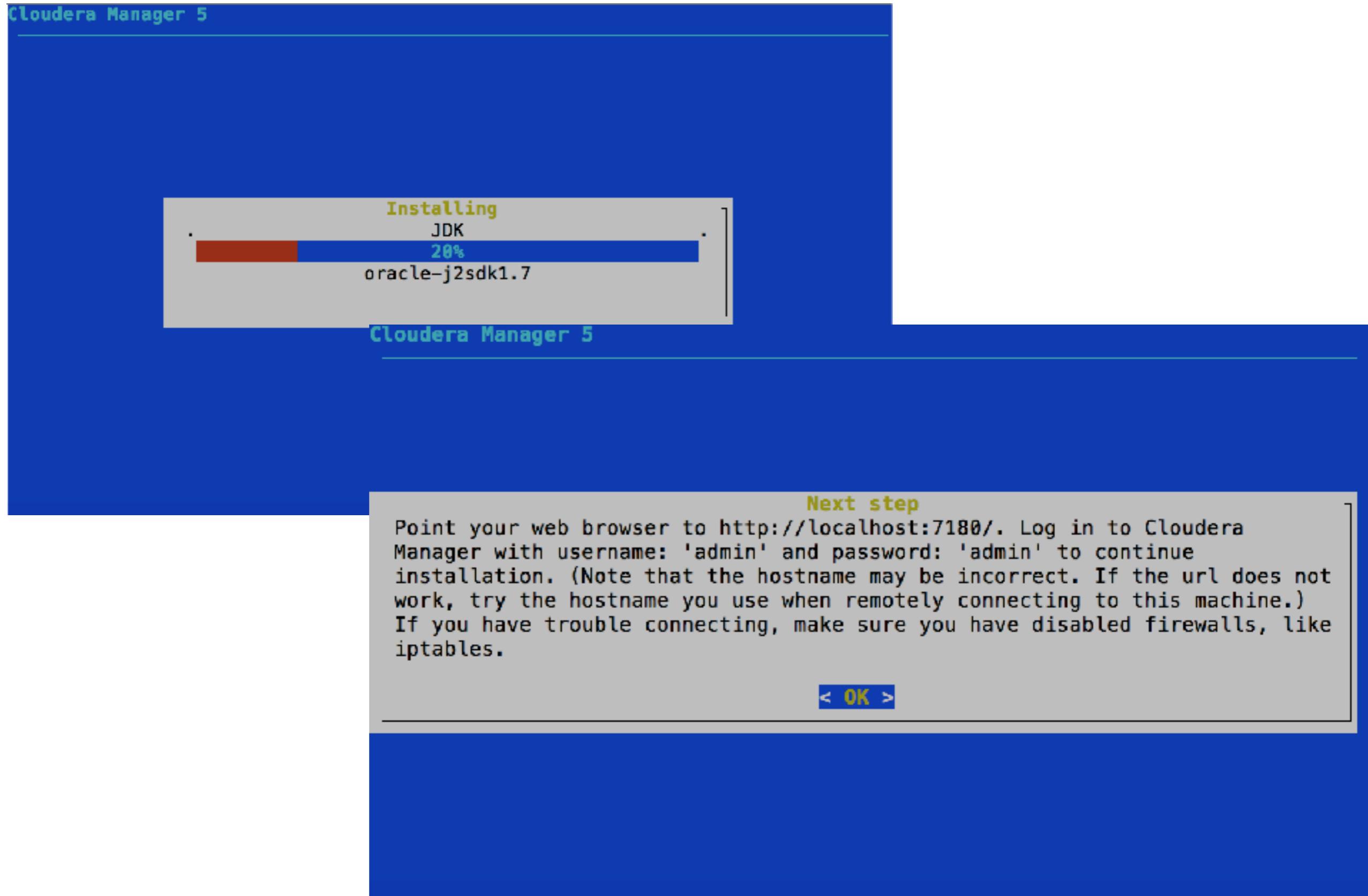
Installing Cloudera on EC2

Download Cloudera Manager

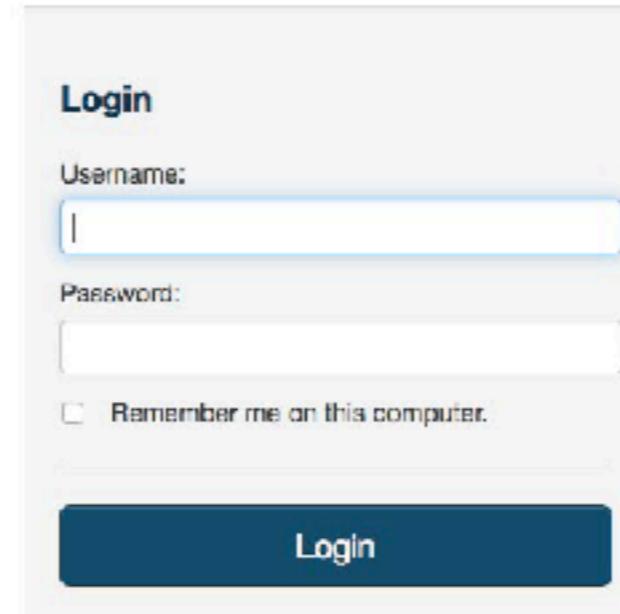
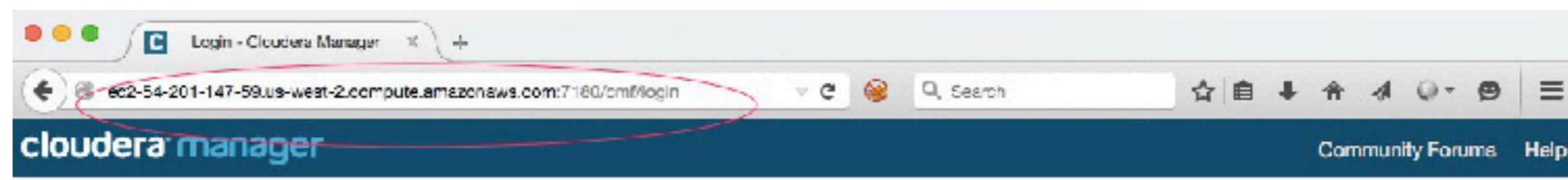
- 1) Type command > **wget http://archive.cloudera.com/cm5/installer/latest/cloudera-manager-installer.bin**
- 2) Type command > **chmod u+x cloudera-manager-installer.bin**
- 3) Type command > **sudo ./cloudera-manager-installer.bin**



Installing Cloudera on EC2



Wait several minutes for the Cloudera Manager Server to complete its startup. Then running web browser: [http:// public-dns: 7180](http://public-dns:7180)



The login form is titled "Login". It contains fields for "Username" and "Password", both of which are currently empty. Below these fields is a checkbox labeled "Remember me on this computer." At the bottom of the form is a large blue "Login" button.

Welcome to Cloudera Manager

End User License Terms and Conditions

Cloudera Standard License

Version 2015-08-06

END USER LICENSE TERMS AND CONDITIONS

THESE TERMS AND CONDITIONS (THESE "TERMS") APPLY TO YOUR USE OF THE PRODUCTS (AS DEFINED BELOW) PROVIDED BY CLOUDERA, INC. ("CLOUDERA").

PLEASE READ THESE TERMS CAREFULLY.

IF YOU ("YOU" OR "CUSTOMER") PLAN TO USE ANY OF THE PRODUCTS ON BEHALF OF A COMPANY OR OTHER ENTITY, YOU REPRESENT THAT YOU ARE THE EMPLOYEE OR AGENT OF SUCH COMPANY (OR OTHER ENTITY) AND YOU HAVE THE AUTHORITY TO ACCEPT ALL OF THE TERMS AND CONDITIONS SET FORTH IN AN ACCEPTED REQUEST FOR QUOTE ("RFQ"), ORDER AGREEMENT ("ORDER AGREEMENT") OR SIMILAR DOCUMENT ("AGREEMENT") ON BEHALF OF SUCH COMPANY (OR OTHER ENTITY).

BY USING ANY OF THE PRODUCTS, YOU ACKNOWLEDGE AND AGREE THAT:

- (A) YOU HAVE READ ALL OF THE TERMS AND CONDITIONS OF THIS AGREEMENT;
- (B) YOU UNDERSTAND ALL OF THE TERMS AND CONDITIONS OF THIS AGREEMENT;
- (C) YOU AGREE TO BE LEGALLY BOUND BY ALL OF THE TERMS AND CONDITIONS OF THIS AGREEMENT.

Yes, I accept the End User License Terms and Conditions.

cloudera manager

Support  admin

Welcome to Cloudera Manager. Which edition do you want to deploy?

Upgrading to Cloudera Enterprise Data Hub Edition provides important features that help you manage and monitor your Hadoop clusters in mission-critical environments.

	Cloudera Express	Cloudera Enterprise Data Hub Edition Trial	Cloudera Enterprise Data Hub Edition
License	Free	60 Days <small>After the trial period, the product will continue to function as Cloudera Express. Your cluster and your data will remain unaffected.</small>	Annual Subscription Upload License
Node Limit	Unlimited	Unlimited	Unlimited
CDH	✓	✓	✓
Care Cloudera Manager Features	✓	✓	✓
Advanced Cloudera Manager Features		✓	✓
Cloudera Navigator		✓	✓
			Continue

Thank you for choosing Cloudera Manager and CDH.

This installer will install **Cloudera Express 5.4.0** and enable you to later choose packages for the services below (there may be some license implications).

- Apache Hadoop (Common, HDFS, MapReduce, YARN)
- Apache HBase
- Apache ZooKeeper
- Apache Oozie
- Apache Hive
- Hue (Apache licensed)
- Apache Flume
- Cloudera Impala (Apache licensed)
- Apache Sentry
- Apache Sqoop
- Cloudera Search (Apache licensed)
- Apache Spark

You are using Cloudera Manager to install and configure your system. You can learn more about Cloudera Manager by clicking on the Support menu above.

 Continue

Provide your instances <private ip> addresses in the cluster



Specify hosts for your CDH cluster installation.

Hosts should be specified using the same hostname (FQDN) that they will identify themselves with.

Cloudera recommends including Cloudera Manager Server's host. This also enables health monitoring for that host.

Hint: Search for hostnames and/or IP addresses using [patterns](#).

172.31.10.50, 172.31.10.51, 172.31.10.52, 172.31.10.53

SSH Port:

Specify hosts for your CDH cluster installation.

Hosts should be specified using the same hostname (FQDN) that they will identify themselves with.

Cloudera recommends including Cloudera Manager Server's host. This also enables health monitoring for that host.

Hint: Search for hostnames and/or IP addresses using [patterns](#).

4 hosts scanned, 4 running SSH.

[New Search](#)

Expanded Query	Hostname (FQDN)
<input checked="" type="checkbox"/>	172.31.10.50 ip-172-31-10-50.us-west-2.compute.internal
<input checked="" type="checkbox"/>	172.31.10.51 ip-172-31-10-51.us-west-2.compute.internal
<input checked="" type="checkbox"/>	172.31.10.52 ip-172-31-10-52.us-west-2.compute.internal
<input checked="" type="checkbox"/>	172.31.10.53 ip-172-31-10-53.us-west-2.compute.internal

Cluster Installation

Select Repository

Cloudera recommends the use of parcels for installation over packages, because parcels enable Cloudera Manager to easily manage the software on your cluster, automating the deployment and upgrade of service binaries. Electing not to use parcels will require you to manually upgrade packages on all hosts in your cluster when software updates are available, and will prevent you from using Cloudera Manager's rolling upgrade capabilities.

Choose Method Use Packages [?](#)

Use Parcels (Recommended) [?](#) [More Options](#)

Select the version of CDH

CDH-5.6.0-1.cdh5.6.0.p0.45

CDH-4.7.1-1.cdh4.7.1.p0.47

Versions of CDH that are too new for this version of Cloudera Manager (5.6.0) will not be shown.

Additional Parcels ACCUMULO-1.6.0-1.cdh5.1.4.p0.116

1 2 3 4 5 6 7 8

[Back](#)

[Continue](#)



Cluster Installation

Enable Single User Mode

Only supported for CDH 5.8 and above.

By default, service processes run as distinct users on the system. For example, HDFS DataNodes run as user "hdfs" and HBase RegionServers run as user "hbase." Enabling "single user mode" configures Cloudera Manager to run service processes as a single user, by default "cloudera-scm", thereby prioritizing isolation between managed services and the rest of the system over isolation between the managed services.

The major benefit of this option is that the Agent does not run as root. However, directories which in the regular mode are created automatically by the Agent, must be set up for the configured user.

Switching back and forth between single user mode and regular mode is not recommended.

Single User Mode



Support - admin -

Cluster Installation

Installation completed successfully.

4 of 4 host(s) completed successfully.

Hostname	IP Address	Progress	Status	Details
ip-172-31-10-50.us-west-2.compute.internal	172.31.10.50	<div style="width: 100%; background-color: #2e6b2e;"></div>	Installation completed successfully.	Details
ip-172-31-10-51.us-west-2.compute.internal	172.31.10.51	<div style="width: 100%; background-color: #2e6b2e;"></div>	Installation completed successfully.	Details
ip-172-31-10-52.us-west-2.compute.internal	172.31.10.52	<div style="width: 100%; background-color: #2e6b2e;"></div>	Installation completed successfully.	Details
ip-172-31-10-53.us-west-2.compute.internal	172.31.10.53	<div style="width: 100%; background-color: #2e6b2e;"></div>	Installation completed successfully.	Details

[Back](#)

1 2 3 4 5 6 7 8

[Continue](#)



Cluster Setup

Choose the CDH 5 services that you want to install on your cluster.

Choose a combination of services to install.

- Core Hadoop**
HDFS, YARN (MapReduce 2 Included), ZooKeeper, Oozie, Hive, Hue, and Sqoop
- Core with HBase**
HDFS, YARN (MapReduce 2 Included), ZooKeeper, HBase, and HDFS
- Core with Impala**
HDFS, YARN (MapReduce 2 Included), ZooKeeper, Impala, and HDFS
- Core with Search**
HDFS, YARN (MapReduce 2 Included), ZooKeeper, Solr, and HDFS
- Core with Spark**
HDFS, YARN (MapReduce 2 Included), ZooKeeper, Spark, and HDFS
- All Services**
HDFS, YARN (MapReduce 2 Included), ZooKeeper, HBase, Impala, Solr, Spark, Oozie, Hive, Hue, and Sqoop
- Custom Services**
Choose your own services. Services required by other services will be installed automatically.

This wizard will also install the Cloudera Manager Agent.

[Back](#)

Cluster Setup

Customize Role Assignments

You can customize the role assignments for your new cluster here, but if assignments are made incorrectly, such as assigning too many roles to a single host, this can impact the performance of your services. Cloudera does not recommend altering assignments unless you have specific requirements, such as having pre-selected a specific host for a specific role.

You can also view the role assignments by host. [View By Host](#)

H Base

<input type="checkbox"/> Master x 1 New ip-172-31-26-220.us-west-2.compute...	<input type="checkbox"/> HBase REST Server Select hosts	<input type="checkbox"/> HBase Thrift Server Select hosts	<input type="checkbox"/> RegionServer x 3 New Same As DataNode ▾
--	--	--	---

HDFS

<input type="checkbox"/> NameNode x 1 New ip-172-31-26-220.us-west-2.compute...	<input type="checkbox"/> SecondaryNameNode x 1 New ip-172-31-26-220.us-west-2.compute...	<input type="checkbox"/> Balancer x 1 New ip-172-31-26-220.us-west-2.compute...	<input type="checkbox"/> HttpFS Select hosts
--	---	--	---

NFS Gateway

<input type="checkbox"/> Select hosts	ip-172-31-26-[221-223].us-west-2.compute...
---------------------------------------	---

Hive

[Back](#)

1 2 3 4 5 6

[Continue](#)

Cluster Setup

First Run Command

Status: Running Start Time: Jan 20, 4:41:25 PM

[Abort](#)

[Details](#) [Completed in 8 steps](#)

Step	Context	Start Time	Duration	Actions
1 ✓ Deploy Client Configuration Successfully deployed all client configurations.	Cluster 1.2	Jan 20, 4:41:25 PM	15.96s	
2 ✓ Start Cloudera Management Service, ZooKeeper Successfully completed 2 steps.		Jan 20, 4:41:41 PM	25.72s	
3 ✓ Start HDFS Successfully completed 1 steps.				
4 ✓ Start HiveServer, Sqoop Successfully completed 2 steps.				
5 ✓ Start YARN (MRP Included), Key-Value Store Indexer Successfully completed 2 steps.				
6 ✓ Start Spark 8/11 steps completed.				
Back				
✓ Deploy Client Configuration Successfully deployed all client configurations.				
✓ Creating Hive Metastore Database Created Hive Metastore Database. Details				
✓ Creating Hive user directory Successfully created HDFS directory. Details				
✓ Creating Hive warehouse directory Successfully created HDFS directory. Details				
✓ Starting Hive Service Service started successfully. Details				
✓ Creating Oozie database Oozie database created successfully. Details				
✓ Installing Oozie ShareLib in HDFS Successfully installed Oozie ShareLib. Details				
✓ Starting Oozie Service Service started successfully. Details				
✓ Starting Hue Service Service started successfully. Details				
✓ Deploying Client Configuration Successfully deployed all client configurations. Details				

Cloudera Manager

cloudera manager Clusters Hosts Diagnostics Audits Charts Administration Search (Hotkey: /) Support admin

Home

30 minutes preceding January 20, 2016, 4:49 PM UTC

Status All Health Issues Configuration X 5 All Recent Commands Add Cluster

Try Cloudera Enterprise Data Hub Edition for 60 Days

Cluster 1 (CDH 5.5.1, Parcels)

- Hosts
- HBase
- HDFS
- Hive
- Hue
- Impala
- Key-Value Store...
- Oozie
- Solr
- Spark
- YARN (MR2 Incl...)
- ZooKeeper

Charts

Cluster CPU



Percent

04:30 04:45

Cluster 1, Host CPU Usage Across Hosts 15%

Cluster Disk IO



Bytes / second

04:30 04:45

Total Disk Bytes Read 0 Total Disk Bytes Written 1.7M/s

Cluster Network IO



Bytes / second

04:30 04:45

Total Bytes Read 151K/s Total Bytes Written 7.8M/s

HDFS IO



Bytes / second

04:30 04:45

Total Bytes Read 1b/s Total Bytes Written 166b/s

Completed Impala Queries

Query ID	Start Time	End Time	Duration	State
Q1	2016-01-20T04:30:00Z	2016-01-20T04:45:00Z	15m	Completed
Q2	2016-01-20T04:30:00Z	2016-01-20T04:45:00Z	15m	Completed

Cloudera Management Service

- Cloudera...

X 2

Hadoop as a Service

Recommended

Amazon EMR

Azure HDInsight

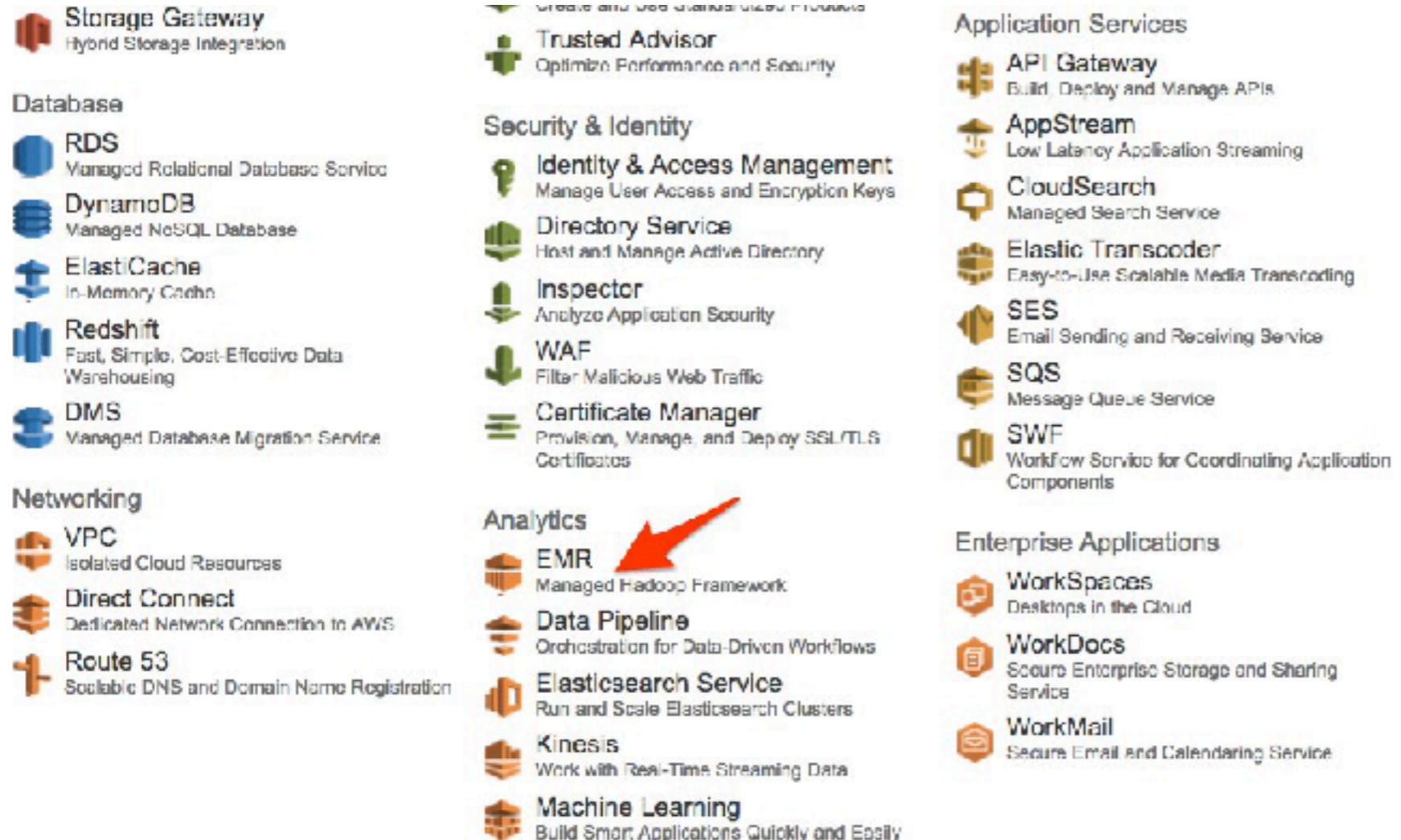
Amazon EMR

Install a cluster on 4 EC2 AWS

Azure HD Insight

Install a cluster on 4 Virtual Servers Microsoft Azure

Launch a EMR Cluster



The screenshot shows the AWS Services Catalog interface. On the left, there's a sidebar with categories like Storage, Database, and Networking. The main area lists various AWS services under different sections: CloudWatch Metrics, Trusted Advisor, Security & Identity, Application Services, AWS Marketplace, AWS re:Invent Announcements, Service Health, and a status message.

- Storage**
 - Storage Gateway
 - Hybrid Storage Integration
- Database**
 - RDS
 - Managed Relational Database Service
 - DynamoDB
 - Managed NoSQL Database
 - ElastiCache
 - In-Memory Cache
 - Redshift
 - Fast, Simple, Cost-Effective Data Warehousing
 - DMS
 - Managed Database Migration Service
- Networking**
 - VPC
 - Isolated Cloud Resources
 - Direct Connect
 - Dedicated Network Connection to AWS
 - Route 53
 - Social DNS and Domain Name Registration
- CloudWatch Metrics**
- Trusted Advisor**
- Security & Identity**
 - Identity & Access Management
 - Manage User Access and Encryption Keys
 - Directory Service
 - Host and Manage Active Directory
 - Inspector
 - Analyze Application Security
 - WAF
 - Filter Malicious Web Traffic
 - Certificate Manager
 - Provision, Manage, and Deploy SSL/TLS Certificates
- Analytics**
 - EMR
 - Managed Hadoop Framework
 - Data Pipeline
 - Orchestration for Data-Driven Workflows
 - Elasticsearch Service
 - Run and Scale Elasticsearch Clusters
 - Kinesis
 - Work with Real-Time Streaming Data
 - Machine Learning
 - Build Smart Applications Quickly and Easily
- Application Services**
 - API Gateway
 - Build, Deploy and Manage APIs
 - AppStream
 - Low Latency Application Streaming
 - CloudSearch
 - Managed Search Service
 - Elastic Transcoder
 - Easy-to-Use Scalable Media Transcoding
 - SES
 - Email Sending and Receiving Service
 - SQS
 - Message Queue Service
 - SWF
 - Workflow Service for Coordinating Application Components
- AWS Marketplace**
- AWS re:Invent Announcements**
- Service Health**
 - All services operating normally.
- Updated: Oct 02 2016 07:40:00 GMT+0700
- [Service Health Dashboard](#)

Create Cluster - Quick Options [Go to advanced options](#)

General Configuration

Cluster name

Logging [i](#)

S3 folder 

Launch mode Cluster [i](#) Step execution [i](#)

Software configuration

Vendor Amazon MapR

Release  

Applications Core Hadoop: Hadoop 2.7.2 with Ganglia 3.7.2, Hive 2.1.0, Hue 3.10.0, Mahout 0.12.2, Pig 0.16.0, and Taz 0.8.4
 HBase: HBase 1.2.2 with Ganglia 3.7.2, Hadoop 2.7.2, Hive 2.1.0, Hue 3.10.0, Phoenix 4.7.0, and ZooKeeper 3.4.8
 Presto: Presto 0.150 with Hadoop 2.7.2 HDFS and Hive 2.1.0 Metastore
 Spark: Spark 2.0.0 on Hadoop 2.7.2 YARN with Ganglia 3.7.2 and Zeppelin 0.6.1

Hardware configuration

Instance type 

Number of instances (1 master and 3 core nodes)

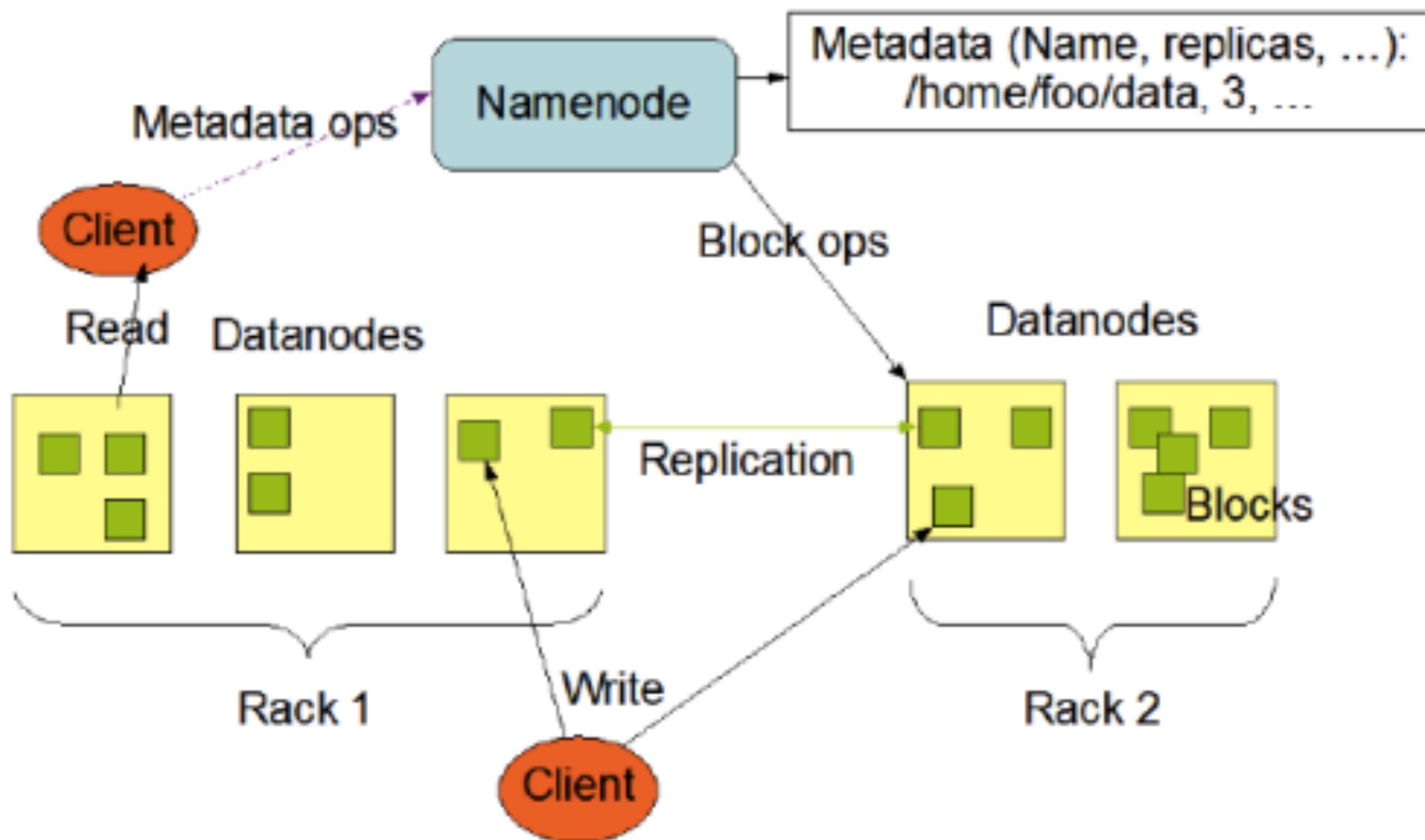
Connect to the instance

```
[https://aws.amazon.com/amazon-linux-ami/2016.03-release-notes/  
26 package(s) needed for security, out of 34 available  
Run "sudo yum update" to apply all updates.  
Amazon Linux version 2016.09 is available.  
-bash: warning: setlocale: LC_CTYPE: cannot change locale (UTF-8): No such file  
or directory  
  
EEEEEEEEEEEEEEEEEE MMMMMMMMM          MMMMMMMMM RRRRRRRRRRRRRRRR  
E:::::::::::E M:::::M          M:::::M R:::::::::::R  
EE:::::EEEEEEE:::E M:::::M          M:::::M R:::::RRRRRR:::R  
 E:::E     EEEEE M:::::M          M:::::M RR:::R      R:::R  
 E:::E          M:::::M:::M          M:::M:::::M R:::R      R:::R  
 E:::::EEEEEEEEE M:::::M M:::M M:::M M:::::M R:::RRRRRR:::R  
 E:::::::::::E M:::::M M:::M:::M M:::::M R:::::::::::RR  
 E:::::EEEEEEEEE M:::::M M:::::M M:::::M R:::RRRRRR:::R  
 E:::E          M:::::M M:::M          M:::::M R:::R      R:::R  
 E:::E     EEEEE M:::::M     MMM M:::::M R:::::R      R:::R  
EE:::::EEEEEEE:::E M:::::M          M:::::M R:::R      R:::R  
E:::::::::::E M:::::M          M:::::M RR:::R      R:::R  
EEEEEEEEEEEEEEEEEE MMMMMMMMM          MMMMMMMMM RRRRRRRR      RRRRRR  
  
[hadoop@ip-172-31-39-165 ~]$
```

Hadoop File System (HDFS)

- Default storage for the Hadoop cluster
- Data is distributed and replicated over multiple machines
- Designed to handle very large files with streaming data access patterns.
- NameNode/DataNode
- Master/slave architecture (1 master 'n' slaves)
- Designed for large files (64 MB default, but configurable) across all the nodes

HDFS Architecture

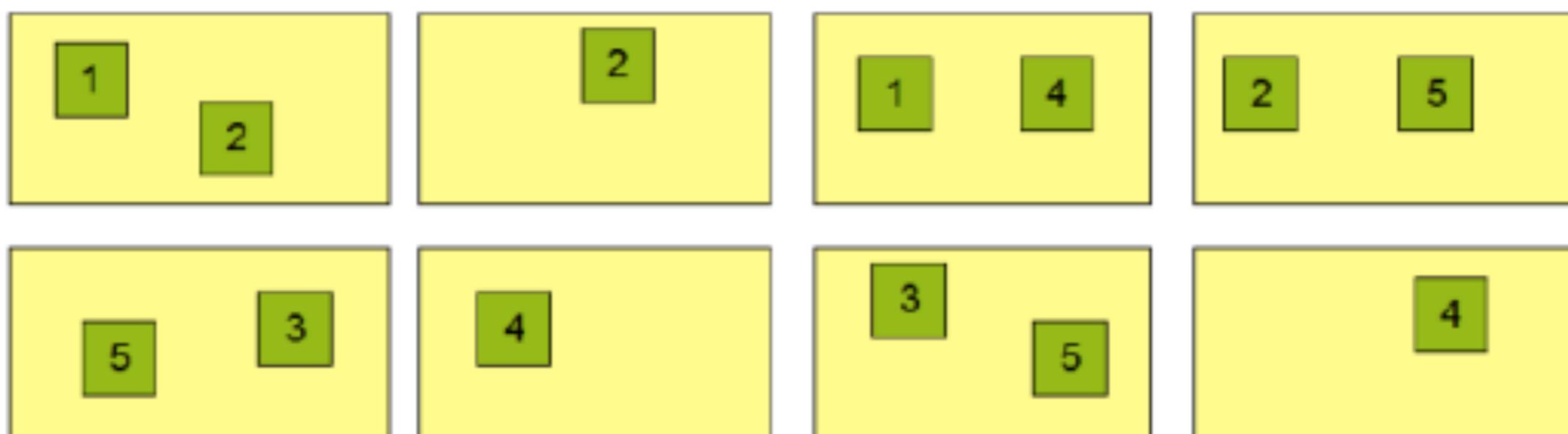


Data Replication in HDFS

Block Replication

```
Namenode (Filename, numReplicas, block-ids, ...)  
/users/sameerp/data/part-0, r:2, {1,3}, ...  
/users/sameerp/data/part-1, r:3, {2,4,5}, ...
```

Datanodes



How does HDFS work?

A file we want to store on HDFS ...

600 MB

We're raising the question because no one else wants to, because no one else wants to say what needs to be said.

And let's be real, it's the two-ton elephant in the room with nearly every other star's name on the trade rumor radar these days.

We've read over and over again about Nash refusing to ask for a trade, refusing to play the game that so many others have late in their careers.

How does HDFS work?

HDFS Splits file into blocks ...

256 MB

We're raising the question because no one else wants to, because no one else wants to say what needs to be said.

256 MB

And let's be real, it's the two-ton elephant in the room with nearly every other star's name on the trade rumor radar these days.

88 MB

We've read over and over again about Nash refusing to play the game that so many others have late in their careers.

How does HDFS work?

HDFS will create **3replicas** of each block ...

3 copies

We're raising the question because no one else wants to, because no one else wants to say what needs to be said.

3 copies

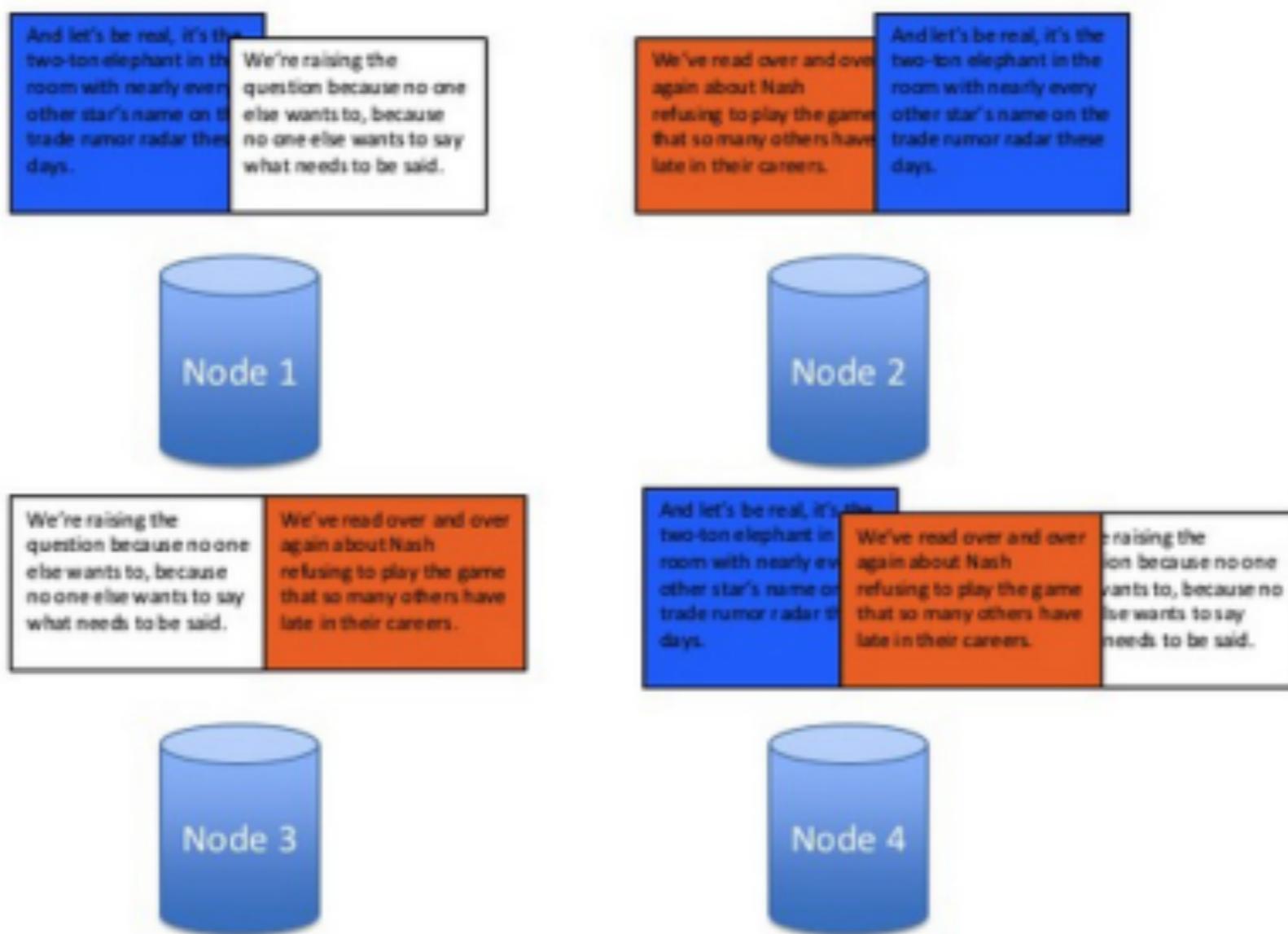
And let's be real, it's the two-ton elephant in the room with nearly every other star's name on the trade rumor radar these days.

3 copies

We've read over and over again about Nash refusing to play the game that so many others have late in their careers.

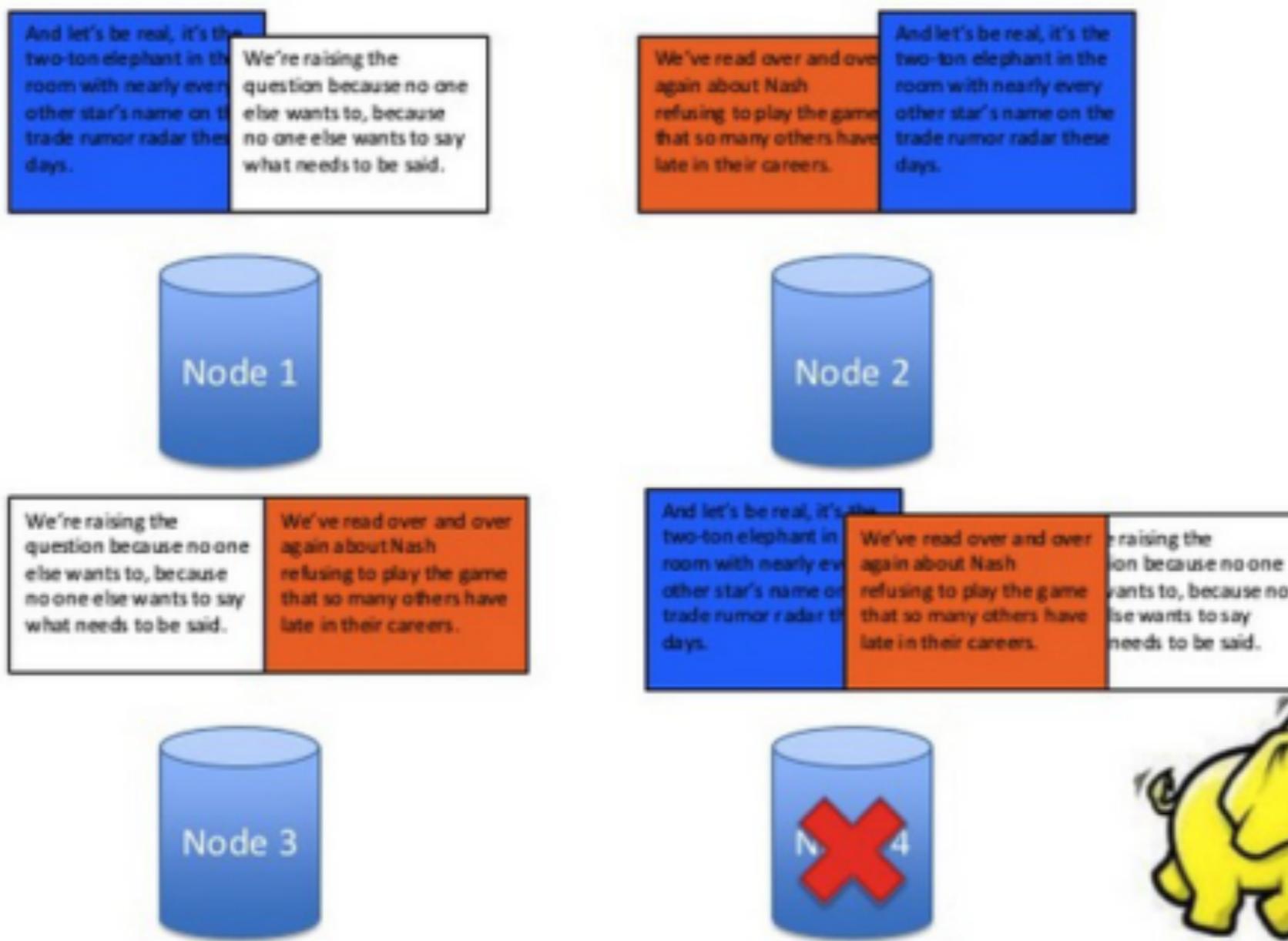
How does HDFS work?

HDFS **distributes** these replicas across the cluster ...



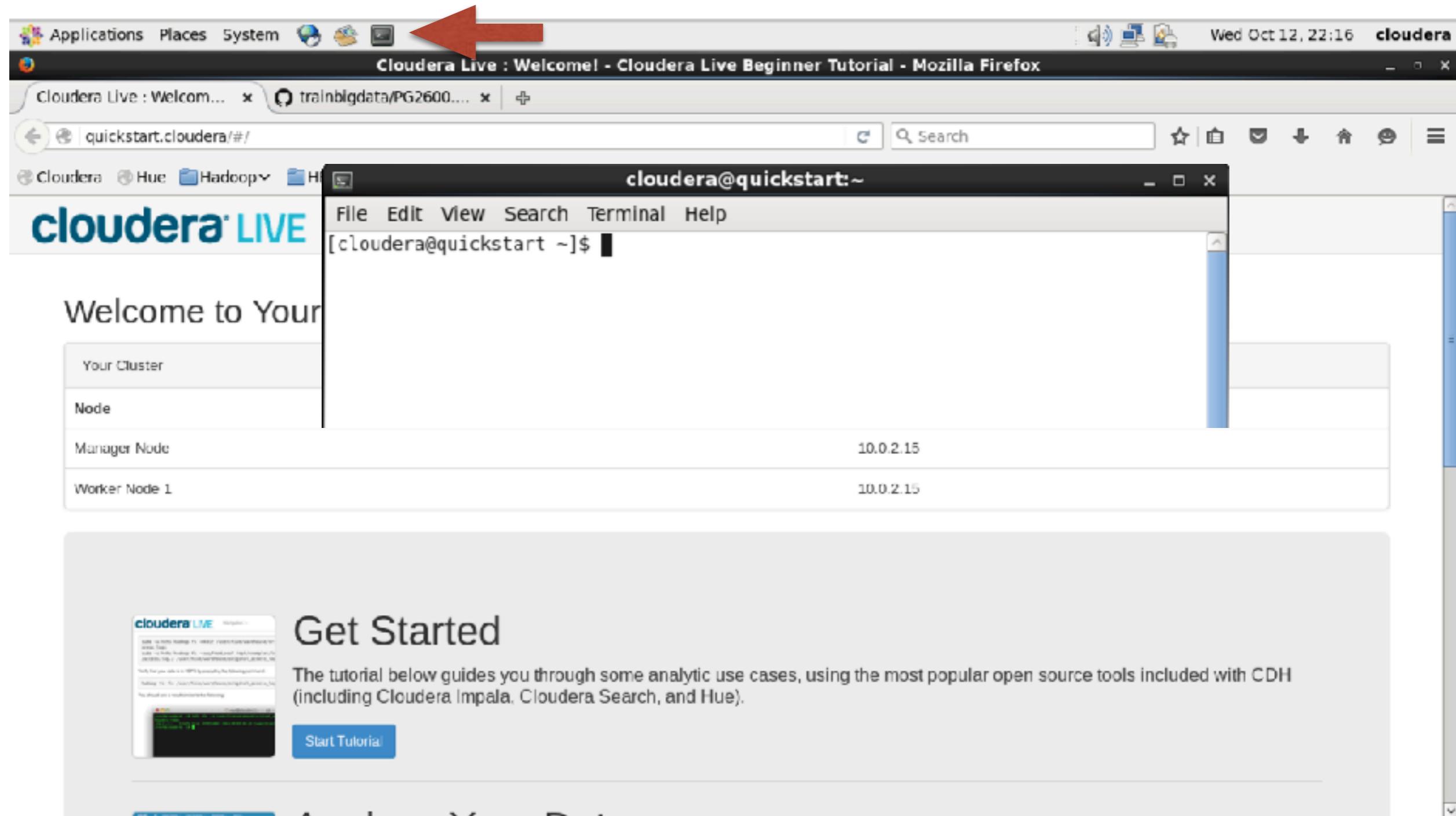
How does HDFS work?

If a node goes down, we have copies elsewhere



Importing/Exporting Data to HDFS

Download an example text file via SSH

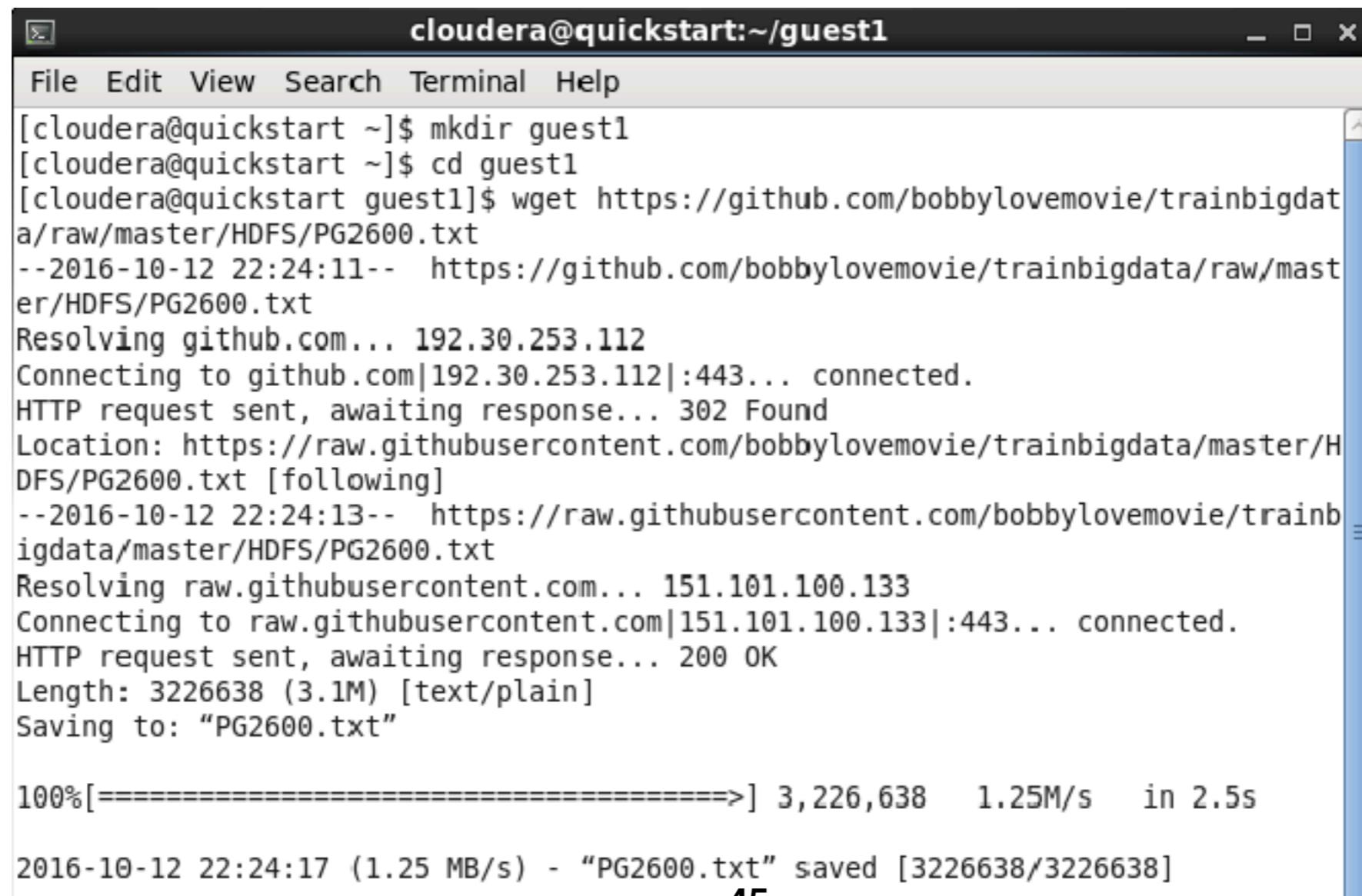


Importing/Exporting Data to HDFS

```
$ mkdir guest1
```

```
$ cd guest1
```

```
$ wget https://github.com/bobbylovemovie/trainbigdata/raw/master/HDFS/PG2600.txt
```



```
cloudera@quickstart:~/guest1
File Edit View Search Terminal Help
[cloudera@quickstart ~]$ mkdir guest1
[cloudera@quickstart ~]$ cd guest1
[cloudera@quickstart guest1]$ wget https://github.com/bobbylovemovie/trainbigdata/raw/master/HDFS/PG2600.txt
--2016-10-12 22:24:11-- https://github.com/bobbylovemovie/trainbigdata/raw/master/HDFS/PG2600.txt
Resolving github.com... 192.30.253.112
Connecting to github.com|192.30.253.112|:443... connected.
HTTP request sent, awaiting response... 302 Found
Location: https://raw.githubusercontent.com/bobbylovemovie/trainbigdata/master/HDFS/PG2600.txt [following]
--2016-10-12 22:24:13-- https://raw.githubusercontent.com/bobbylovemovie/trainbigdata/master/HDFS/PG2600.txt
Resolving raw.githubusercontent.com... 151.101.100.133
Connecting to raw.githubusercontent.com|151.101.100.133|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 3226638 (3.1M) [text/plain]
Saving to: "PG2600.txt"

100%[=====] 3,226,638 1.25M/s in 2.5s

2016-10-12 22:24:17 (1.25 MB/s) - "PG2600.txt" saved [3226638/3226638]
```

Upload Data to Hadoop

```
$hadoop fs -mkdir /user/cloudera/input
```

```
$hadoop fs -ls /user/cloudera/input
```

```
$hadoop fs -rm /user/cloudera/input/*
```

```
$hadoop fs -put PG2600.txt /user/cloudera/input/
```

```
$hadoop fs -ls /user/cloudera/input
```

```
[cloudera@quickstart guest1]$ hadoop fs -mkdir /user/cloudera/input
[cloudera@quickstart guest1]$ hadoop fs -rm /user/cloudera/input/*
rm: `/user/cloudera/input/*': No such file or directory
[cloudera@quickstart guest1]$ hadoop fs -put PG2600.txt /user/cloudera/input/
[cloudera@quickstart guest1]$ hadoop fs -ls /user/cloudera/input/*
-rw-r--r-- 1 cloudera cloudera 3226638 2016-10-12 23:06 /user/cloudera/inpu
t/PG2600.txt
```

Hadoop syntax for HDFS

Command	Syntax
Listing of files in a directory	<code>hadoop fs -ls /user</code>
Create a new directory	<code>hadoop fs -mkdir /user/guest/newdirectory</code>
Copy a file from a local machine to Hadoop	<code>hadoop fs -put C:\Users\Administrator\Downloads\localfile.csv /user/rajn/newdirectory/hadoopfile.txt</code>
Copy a file from Hadoop to a local machine	<code>hadoop fs -get /user/rajn/newdirectory/hadoopfile.txt C:\Users\Administrator\Desktop\</code>
Tail last few lines of a large file in Hadoop	<code>hadoop fs -tail /user/rajn/newdirectory/hadoopfile.txt</code>
View the complete contents of a file in Hadoop	<code>hadoop fs -cat /user/rajn/newdirectory/hadoopfile.txt</code>
Remove a complete directory from Hadoop	<code>hadoop fs -rm -r /user/rajn/newdirectory</code>
Check the Hadoop filesystem space utilization	<code>hadoop fs -du /</code>

Importing/Exporting Data to HDFS



Hue - Hadoop User Experience - The Apache Hadoop UI

Screenshot of the Hue web interface showing a query editor and a query history.

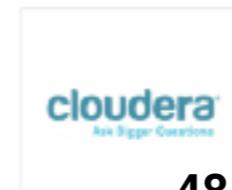
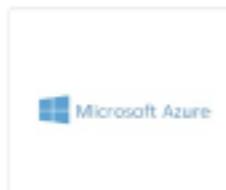
Query Editor:

```
18 -- Compute total amount per order for all customers
19 SELECT
20   c.id AS customer_id,
21   c.name AS customer_name,
22   ords.order_id AS order_id,
23   SUM(order_items.price * order_items.qty) AS total_amount
24 FROM
25   customers c
26   LATERAL VIEW EXPLODE(c.orders) o AS ords
27   LATERAL VIEW EXPLODE(ords.items) i AS order_items
28 GROUP BY c.id, c.name, ords.order_id;
29
```

Query History:

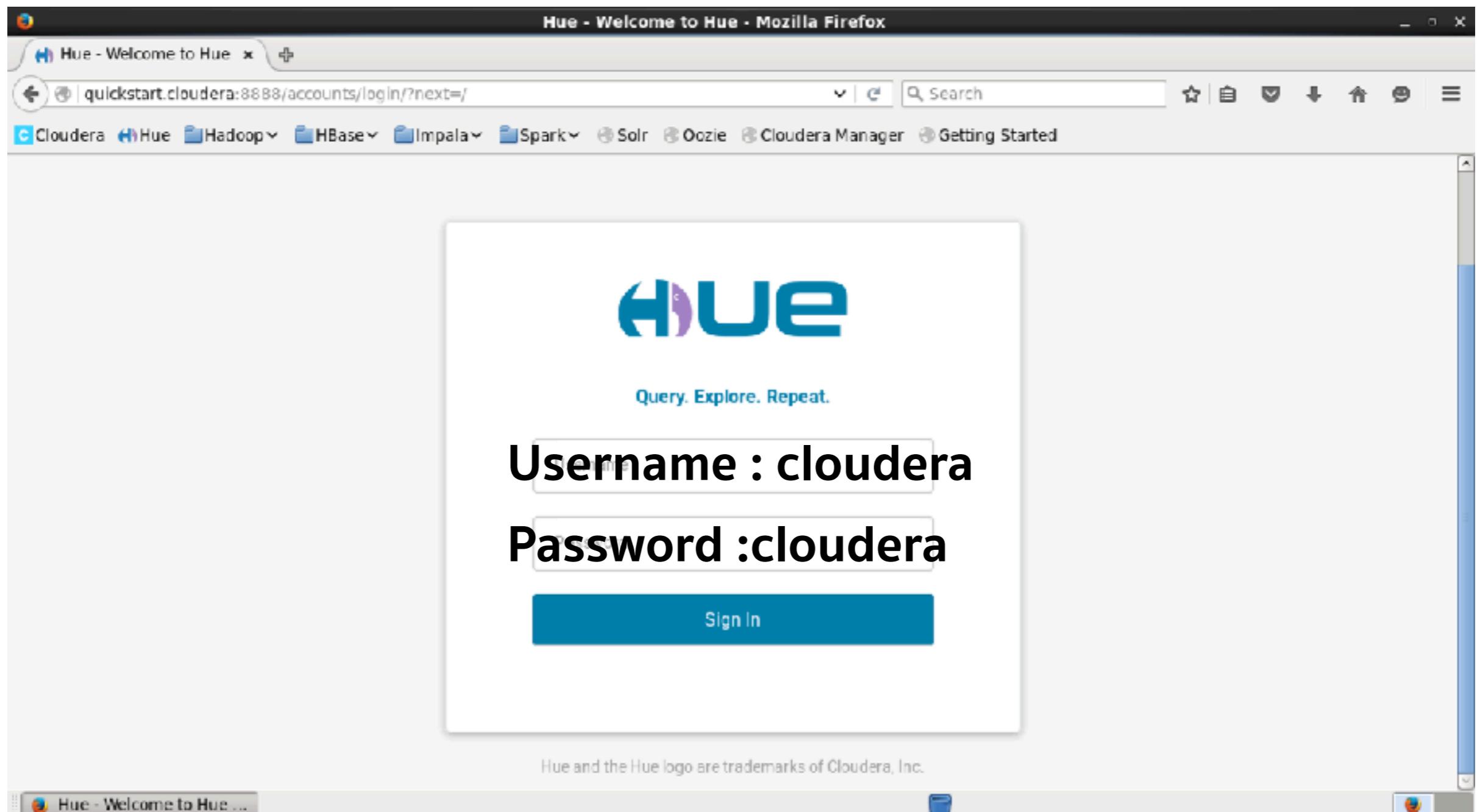
- a few seconds ago → Sample: Customers
- a few seconds ago → Sample: Customers
- a few seconds ago → Sample: Customers
- 15 hours ago → select * from web_logs;show tables;
- 15 hours ago → select * from web_logs;
- 15 hours ago → select * from web_logs;show tables;
- 15 hours ago → select * from web_logs;show tables;
- 15 hours ago → select * from web_logs;show tables;

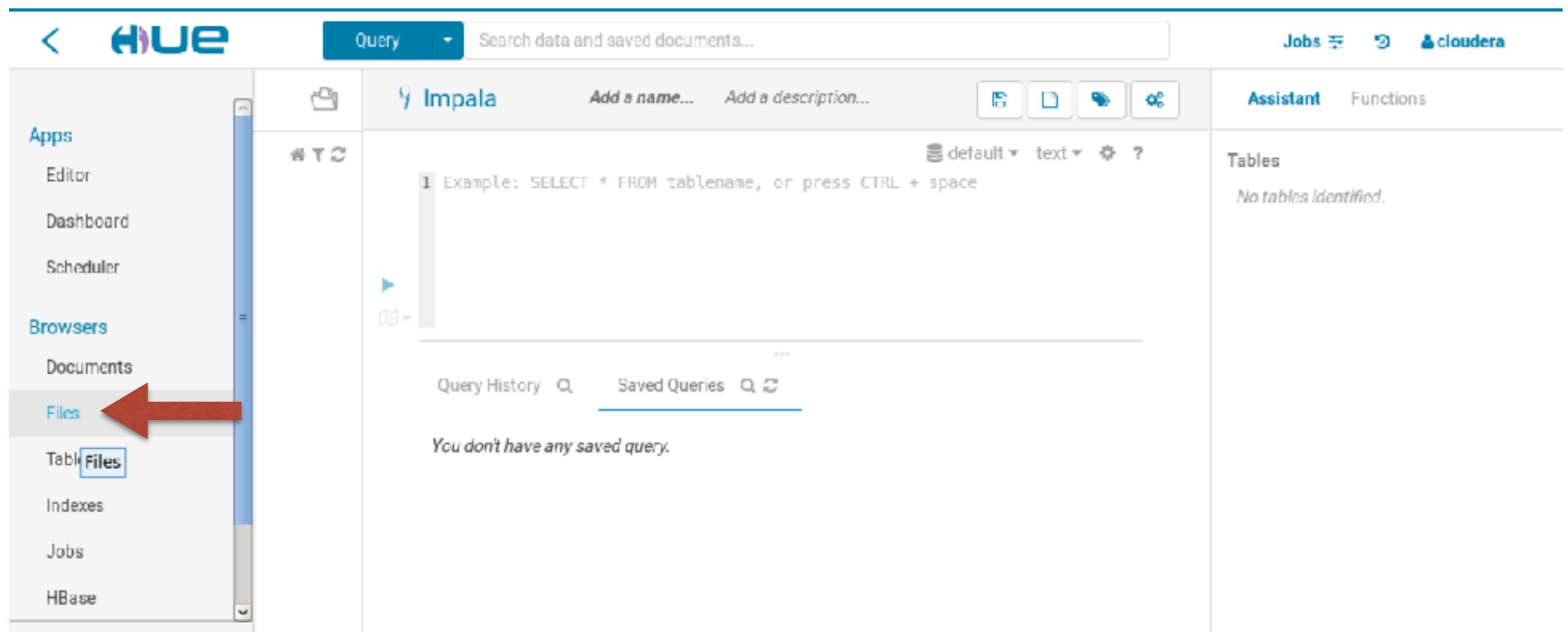
Available in



Review file in Hadoop HDFS using File Browse

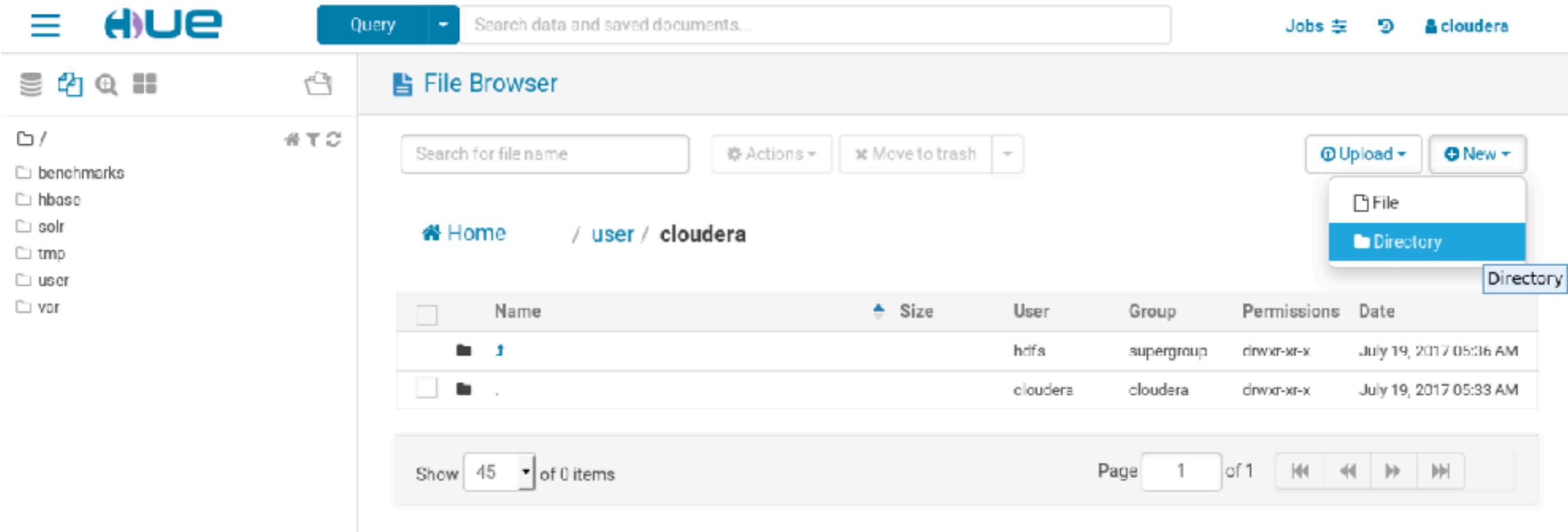
Open Web Browser : <http://quickstart.cloudera:8888>





The screenshot shows the Hue interface for managing Hadoop data. The left sidebar has a blue header "HUE" and a search bar "Search data and saved documents...". It lists several sections: Apps (Editor, Dashboard, Scheduler), Browsers (Documents, Files, Tables, Indexes, Jobs, HBase), and a dropdown menu. A red arrow points to the "Files" link under the Browsers section, which is highlighted with a blue border. The main content area is titled "Impala" and shows a query editor with placeholder text "Example: SELECT * FROM tablename, or press CTRL + space". Below the editor are tabs for "Query History" and "Saved Queries". A message at the bottom says "You don't have any saved query." To the right, there are sections for "Assistant" and "Functions", and a "Tables" section which displays "No tables identified".

Create a new directory name as: **inputByHue , output**



The screenshot shows the Hue File Browser interface. The top navigation bar includes 'Query' (selected), a search bar ('Search data and saved documents...'), 'Jobs' (with a dropdown arrow), and a user icon ('cloudera'). The left sidebar lists HDFS locations: benchmarks, hbase, solr, tmp, user, and var. The main area is titled 'File Browser' and shows the path '/user/cloudera'. It features a search bar ('Search for file name'), an 'Actions' dropdown, a 'Move to trash' button, and upload/new buttons ('Upload', 'New'). A context menu is open over the 'cloudera' directory, with 'File' and 'Directory' options; 'Directory' is highlighted in blue. Below is a table listing files and directories:

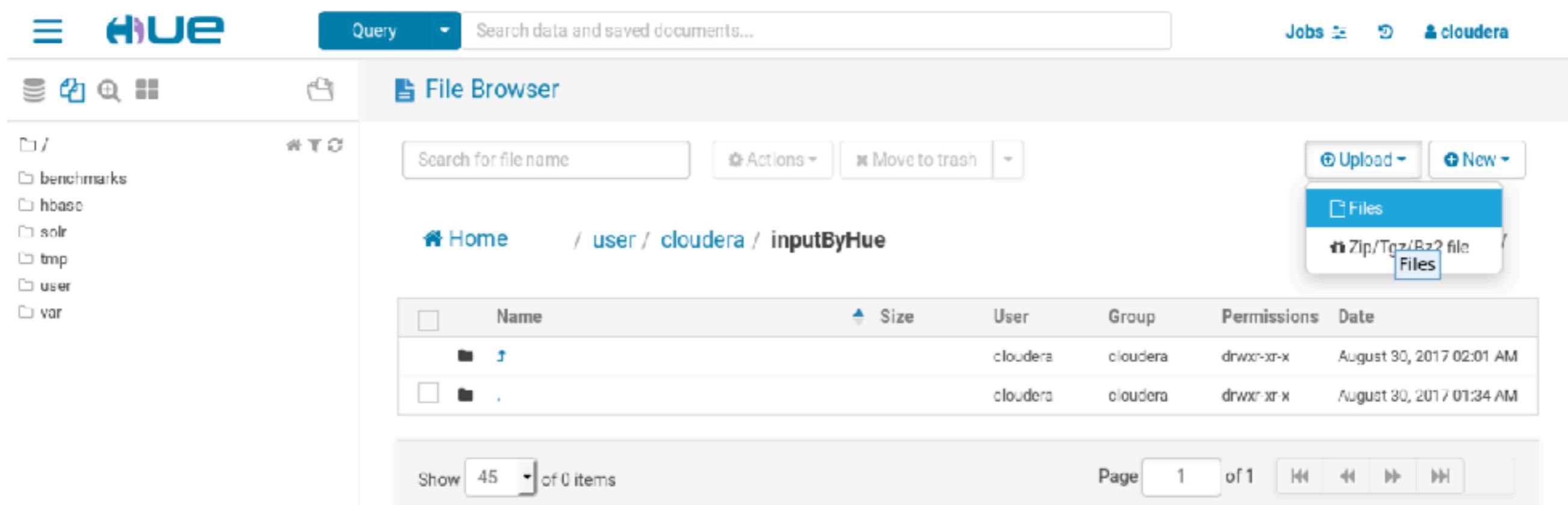
Name	Size	User	Group	Permissions	Date
inputByHue	0	hdfs	supergroup	drwxr-xr-x	July 19, 2017 05:36 AM
.	0	cloudera	cloudera	drwxr-xr-x	July 19, 2017 05:33 AM

At the bottom, there are pagination controls ('Show 45 of 0 items', 'Page 1 of 1', and navigation arrows) and a 'Create Directory' dialog box.

Create Directory

Directory Name:

Upload a local file to HDFS



The screenshot shows the Hue File Browser interface. The top navigation bar includes 'Query' dropdown, search bar, 'Jobs' link, and user 'cloudera'. The left sidebar shows a tree view of directory structure: /, benchmarks, hbase, solr, tmp, user, var. The main area is titled 'File Browser' with sub-titles 'Home / user / cloudera / inputByHue'. It features a search bar, actions dropdown, move to trash button, and upload/new buttons. A modal window is open for 'Upload' with tabs for 'Files' (selected) and 'Zip/Tgz/Bz2 file'. The table below lists files: 'input' (size 0, user cloudera, group cloudera, permissions drwxr-xr-x, date August 30, 2017 02:01 AM) and '.' (size 0, user cloudera, group cloudera, permissions drwxr-xr-x, date August 30, 2017 01:34 AM). Below the table are 'Show 45 of 0 items' and 'Page 1 of 1' navigation controls.

Name	Size	User	Group	Permissions	Date
input	0	cloudera	cloudera	drwxr-xr-x	August 30, 2017 02:01 AM
.	0	cloudera	cloudera	drwxr-xr-x	August 30, 2017 01:34 AM

Upload to /user/cloudera/input

Select files or drag and drop them here

03_Suitability test.pdf 99% from 0.3MB X

HUE

Query Search data and saved documents... Jobs cloudera

File Browser

Search for file name Actions Move to trash Upload New

Upload History

Home / user / cloudera / inputByHue

Name	Size	User	Group	Permissions	Date
..		cloudera	cloudera	drwxr-xr-x	August 30, 2017 02:01 AM
.		cloudera	cloudera	drwxr-xr-x	August 30, 2017 02:02 AM
PG2600.txt	3.1 MB	cloudera	cloudera	rwxr--r-	August 30, 2017 02:02 AM

Show 45 of 1 items Page 1 of 1

Big Data Processing

MapReduce

Before MapReduce...

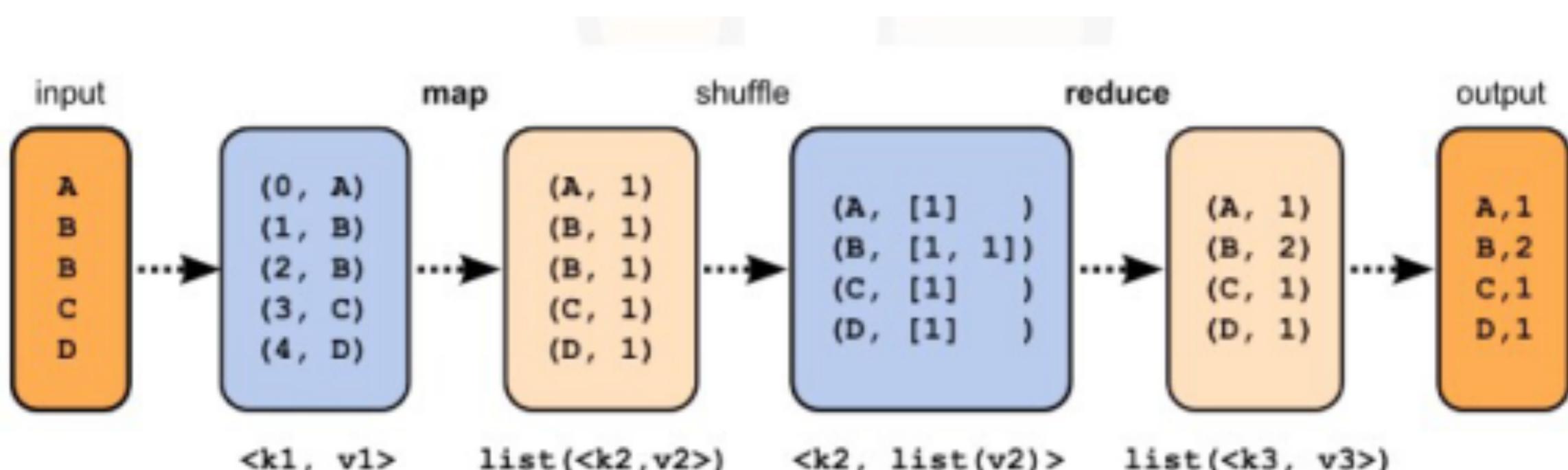
Large scale data processing was difficult!

- ➊ Managing hundreds or thousands of processors
- ➋ Managing parallelization and distribution
- ➌ I/O Scheduling
- ➍ Status and monitoring
- ➎ Fault/crash tolerance

MapReduce provides all of these, easily!

How Map and Reduce Work Together

- Map returns information
- Reduces accepts information
- Reduce applies a user defined function to reduce the amount of data



Example MapReduce: WordCount

```
$cd /guest1
```

```
$wget https://github.com/bobbylovemovie/trainbigdata/raw/master/HDFS/  
wordcount.jar
```

```
$hadoop jar wordcount.jar org.myorg.WordCount /user/cloudera/input/*  
/user/cloudera/output/wordcount
```

```
[cloudera@quickstart guest1]$ hadoop jar wordcount.jar org.myorg.WordCount /user/cloudera/input/* /user/cloudera/output/wordcount  
16/10/12 23:46:17 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.0:8032  
16/10/12 23:46:17 INFO client.RMProxy: Connecting to ResourceManager at /0.0.0.0:8032  
16/10/12 23:46:18 WARN mapreduce.JobResourceUploader: Hadoop command-line option parsing not performed. Implement the Tool interface and execute your application with ToolRunner to remedy this.  
16/10/12 23:46:18 INFO mapred.FileInputFormat: Total input paths to process : 1  
16/10/12 23:46:18 WARN hdfs.DFSClient: Caught exception  
java.lang.InterruptedException  
    at java.lang.Object.wait(Native Method)  
    at java.lang.Thread.join(Thread.java:1281)  
    at java.lang.Thread.join(Thread.java:1355)  
    at org.apache.hadoop.hdfs.DFSOutputStream$DataStreamer.closeResponder(DFSOutputStream.java:862)  
    at org.apache.hadoop.hdfs.DFSOutputStream$DataStreamer.endBlock(DFSOutputStream.java:606)  
    at org.apache.hadoop.hdfs.DFSOutputStream$DataStreamer.run(DFSOutputStream.java:789)  
16/10/12 23:46:18 WARN hdfs.DFSClient: Caught exception  
java.lang.InterruptedException  
    at java.lang.Object.wait(Native Method)  
    at java.lang.Thread.join(Thread.java:1281)  
    at java.lang.Thread.join(Thread.java:1355)  
    at org.apache.hadoop.hdfs.DFSOutputStream$DataStreamer.closeResponder(DFSOutputStream.java:862)  
    at org.apache.hadoop.hdfs.DFSOutputStream$DataStreamer.endBlock(DFSOutputStream.java:606)  
    at org.apache.hadoop.hdfs.DFSOutputStream$DataStreamer.run(DFSOutputStream.java:789)  
16/10/12 23:46:18 INFO mapreduce.JobSubmitter: number of splits:2  
16/10/12 23:46:18 INFO mapreduce.JobSubmitter: Submitting tokens for job: job_1476334425555_0001  
16/10/12 23:46:19 INFO impl.YarnClientImpl: Submitted application application_1476334425555_0001  
16/10/12 23:46:19 INFO mapreduce.Job: The url to track the job: http://quickstart.cloudera:8088/proxy/application_1476334425555_0001/  
16/10/12 23:46:19 INFO mapreduce.Job: Running job: job_1476334425555_0001  
16/10/12 23:46:32 INFO mapreduce.Job: Job job_1476334425555_0001 running in uber mode : false  
16/10/12 23:46:32 INFO mapreduce.Job: map 0% reduce 0%  
16/10/12 23:46:54 INFO mapreduce.Job: map 100% reduce 0%  
16/10/12 23:47:06 INFO mapreduce.Job: map 100% reduce 100%  
16/10/12 23:47:07 INFO mapreduce.Job: Job job_1476334425555_0001 completed successfully  
16/10/12 23:47:07 INFO mapreduce.Job: Counters: 49
```

Reviewing MapReduce Job in Hue



Job Browser

Jobs Workflows Schedules Bundles SLAs

user:cloudera Succeeded Running Failed in the last 7 days x Kill

<input type="checkbox"/>	Id	Name	User	Type	Status	Progress	Group	Started	Duration
<input type="checkbox"/>	application_1504064742599_0001	wordcount	cloudera	MAPREDUCE	SUCCEEDED	100	root.cloudera	August 30, 2017 2:05 AM	1m, 0s

Job Browser Jobs Workflows Schedules Bundles SLAs

job_1504064742599_0001

ID job_1504064742599_0001
TYPE MAPREDUCE
STATUS SUCCEEDED
USER cloudera
PROGRESS 100%
MAP 100% 2/2
REDUCE 100% 1/1

Logs Tasks Metadata Counters ✖ Kill

Filter by name Succeeded Running Failed Map Reduce

Type	Id	Elapsed Time	Progress	State	Start Time	Successful
MAP	task_1504064742599_0001_m_000000	28252	1	SUCCEEDED	1504083966966	attempt_1
MAP	task_1504064742599_0001_m_000001	27784	1	SUCCEEDED	1504083967448	attempt_1
REDUCE	task_1504064742599_0001_r_000000	11368	1	SUCCEEDED	1504083997798	attempt_1

◀ ▶

58

Reviewing MapReduce Output Result

File Browser

Search for file name Actions Move to trash Upload New

Home / user / cloudera / output History

<input type="checkbox"/>	Name	Size	User	Group	Permissions	Date
<input type="checkbox"/>	..		cloudera	cloudera	drwxr-xr-x	August 30, 2017 02:06 AM
<input type="checkbox"/>	.		cloudera	cloudera	drwxr-xr-x	August 30, 2017 02:06 AM
<input type="checkbox"/>	wordcount		cloudera	cloudera	drwxr-xr-x	August 30, 2017 02:06 AM

Show 45 of 1 items Page 1 of 1 << <> >>

Reviewing MapReduce Output Result

 File Browser

 View as binary  Home Page to of 1 

 Edit file / user / cloudera / output / wordcount / part-00000

 Download a 205807
 View file location e 315232
 Refresh i 174282
Last modified 08/30/2017 9:06 AM o 192879
User cloudera u 65433
Group cloudera
Size 44 B

Apache HBase

APACHE
HBASE

Understanding HBase

An open source, non-relational, distributed database

HBase is an open source, non-relational, distributed database modeled after Google's BigTable and is written in Java. It is developed as part of Apache Software Foundation's Apache Hadoop project and runs on top of HDFS providing BigTable-like capabilities for Hadoop. That is, it provides a fault-tolerant way of storing large quantities of sparse data.

HBase Features

- Hadoop database modelled after Google's Bigtable
- Column oriented data store, known as Hadoop Database
- Support random realtime CRUD operations
(unlike HDFS)
- No SQL Database
- Opensource, written in Java
- Run on a cluster of commodity hardware

When to use HBase ?

- When you need high volume data to be stored
- Un-structured data
- Sparse data
- Column-oriented data
- Versioned data (same data template, captured at various time, time-elapse data)
- When you need high scalability

Which one to use ?

HDFS

Only append dataset (no random write)

Read the whole dataset (no random read)

HBase

Need random write and/or read

Has thousands of operation per second on TB+ of data

RDBMS

Data fits on one big node

Need full transaction support

Need real-time query capabilities

HBase Components

Region

Row of table are stores

Region Server

Hosts the tables

Master

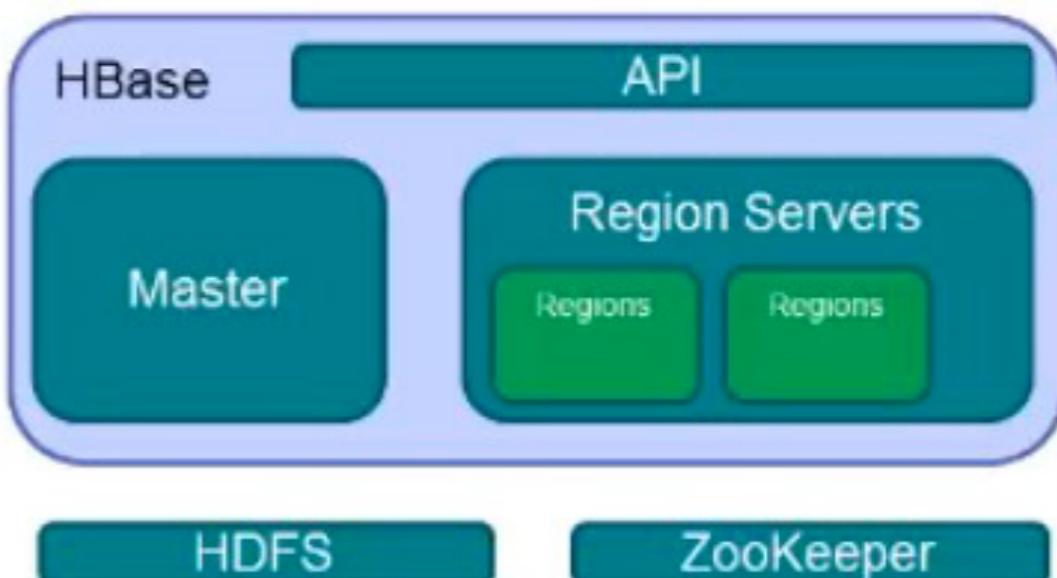
Coordinating the Region Servers

ZooKeeper

HDFS

API

The Java Client API



Hive.apache.org

HBase Shell Commands



List all tables in hbase.

```
hbase> list
```

Create table; pass table name

```
hbase> create 'ns1:t1', {NAME => 'f1', VERSIONS => 5}
```

```
hbase> create 't1', {NAME => 'f1'}, {NAME => 'f2'}, {NAME => 'f3'}
```

Put a cell ‘value’ at specified table/row/column

```
hbase> put 't1', 'r1', 'c1', 'value', ts1
```

Get row or cell contents

```
hbase> get 't1', 'r1'
```

Running HBase

\$ hbase shell

hbase(main):001:0> create 'employee', 'personal data', 'professional data'

hbase(main):002:0> list

```
[cloudera@quickstart guest1]$ hbase shell
2016-10-13 00:39:51,901 INFO  [main] Configuration.deprecation: hadoop.native.lib is deprecated. Instead, use io.native.lib.available
HBase Shell; enter 'help<RETURN>' for list of supported commands.
Type "exit<RETURN>" to leave the HBase Shell
Version 1.2.0-cdh5.8.0, rUnknown, Thu Jun 16 12:46:57 PDT 2016

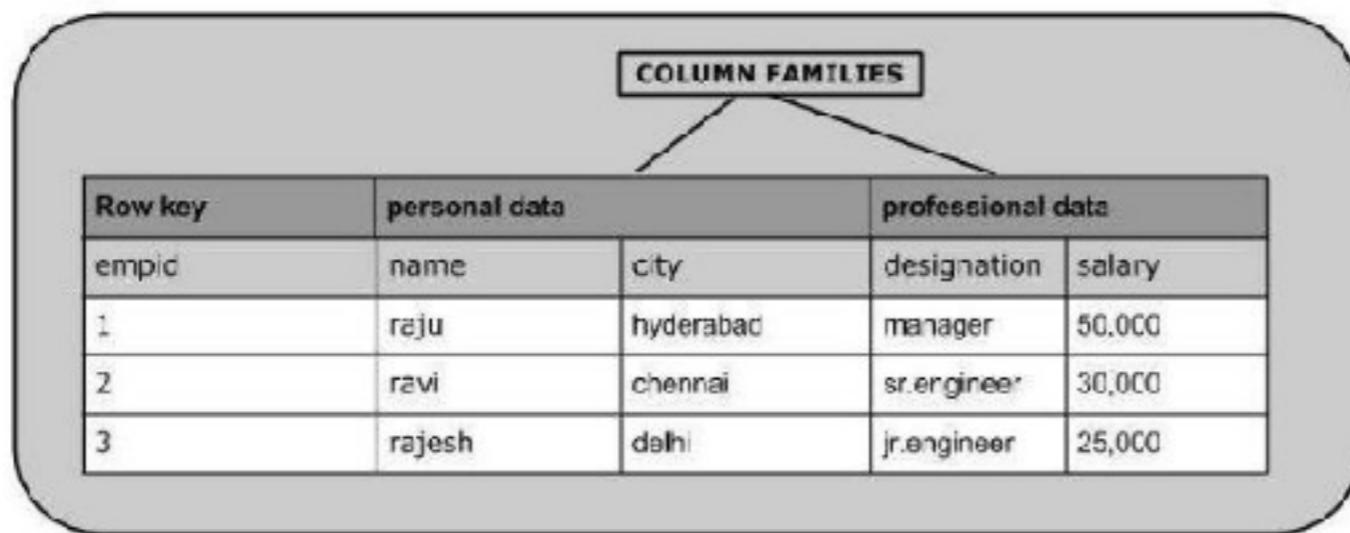
hbase(main):001:0> create 'employee', 'personal data', 'professional data'
0 row(s) in 1.6860 seconds

=> Hbase::Table - employee
hbase(main):002:0> list
TABLE
employee
1 row(s) in 0.0180 seconds

=> ["employee"]
hbase(main):003:0> ■
```

Row key	personal data	professional data

Create Data



```
hbase(main):004:0> put 'employee','1','personal data:name','raju'
```

```
hbase(main):005:0> put 'employee','1','personal data:city','hyderabad'
```

```
hbase(main):006:0> put 'employee','1','professional data:designation','manager'
```

```
hbase(main):007:0> put 'employee','1','professional data:salary','5000'
```

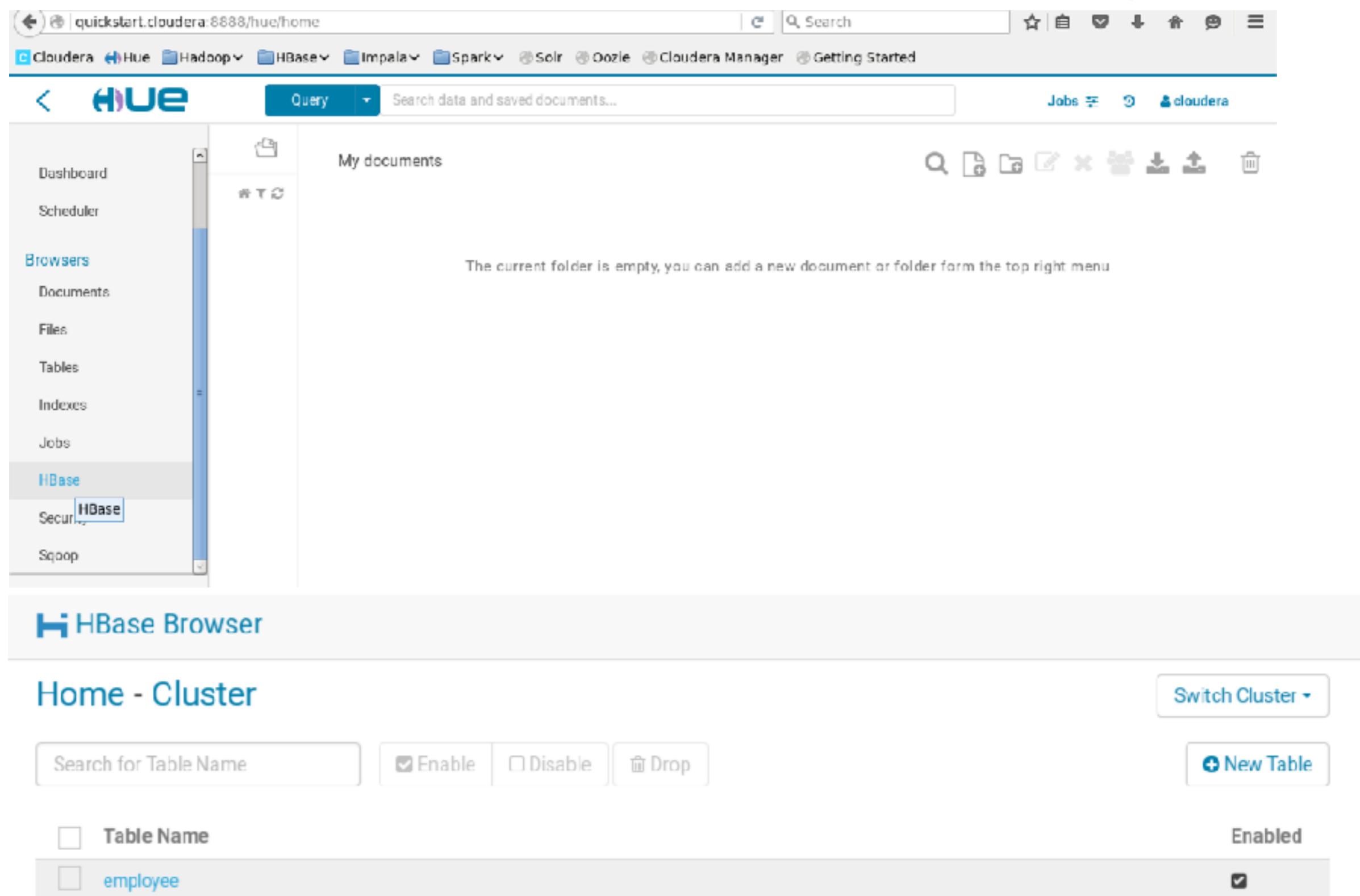
```
hbase(main):004:0> put 'employee','1','personal data:name','raju'  
0 row(s) in 0.0190 seconds
```

```
hbase(main):005:0> put 'employee','1','personal data:city','hyderabad'  
0 row(s) in 0.0050 seconds
```

```
hbase(main):006:0> put 'employee','1','professional data:designation','manager'  
0 row(s) in 0.0110 seconds
```

```
hbase(main):007:0> put 'employee','1','professional data:salary','5000'  
0 row(s) in 0.0050 seconds
```

Running HBase Browser



The screenshot shows the Cloudera Hue web interface. At the top, the URL is `quickstart.cloudera:8888/hue/home`. The navigation bar includes links for Cloudera, Hue, Hadoop, HBase, Impala, Spark, Solr, Oozie, Cloudera Manager, and Getting Started. The main header has a back arrow, the Hue logo, a search bar, and user icons. On the left, a sidebar menu lists Dashboard, Scheduler, Browsers (Documents, Files, Tables, Indexes, Jobs, HBase), HBase (selected), Security, and Sqoop. The main content area is titled "My documents" and displays a message: "The current folder is empty, you can add a new document or folder from the top right menu." Below this, the "HBase Browser" section is titled "Home - Cluster". It features a "Search for Table Name" input field, checkboxes for "Enable", "Disable", and "Drop", and a "New Table" button. A table lists tables: "Table Name" (checkbox) and "employee" (checkbox checked). The "Enabled" column shows "Enabled" for the first row and a checked checkbox for the "employee" row.

Table Name	Enabled
employee	<input checked="" type="checkbox"/>

Viewing Employee Table

HBase Browser

Home - Cluster / employee

Switch Cluster ▾

row_key, row_prefix* +scan_len [col1, family:col2, fam3:, col_prefix]

personal data: professional data: All

Filter Columns/Families

1	personal data: city	personal data: name	professional data: designation	professional data: salary	
	hyderabad	raju	manager	5000	

hbase(main):025:0> put 'employee','2','personal data:name','bobby'

1	personal data: city	personal data: name	professional data: designation	professional data: salary	
	hyderabad	raju	manager	5000	<input type="button" value=""/> <input type="button" value=""/>
2	personal data: name				
	bobby				

Create a table in HBase



HBase Browser

Home - HBase Switch Cluster ▾

Search for Table Name Enable Disable Drop

Table Name Enabled

+ New Table

Create New Table X

Table Name:

Column Families:

cf:name + Add a column property

+ Add an additional column family

Cancel Submit

HBase Browser

Home - HBase / Student

Switch Cluster 

row_key: row_prefix* +scan_len [col1, family:col2, fam3; col_prefix* +3, fam4]

No rows to display.

Fetched 123 in 0.074 seconds.



Add field into a new row

Insert New Row

Row Key	100
cf:firstname	Thaveewat
cf:lastname	khanan
+ Add Field	

[Cancel](#) [Submit](#)

Insert New Row

Row Key	101
cf:firstname	somchai
+ Add Field	

[Cancel](#) [Submit](#)

100	cf:firstname	cf:lastname	Filter Column Names/Family	Sort By ASC	Drop Columns	+
100	Thaveewat	khanan				
101	cf:firstname	somchai				
101						

APACHE HIVE



A Petabyte Scale Data Warehouse Using Hadoop

Hive is developed by Facebook, designed to enable easy data summarization, ad-hoc querying and analysis of large volumes of data. It provides a simple query language called Hive QL, which is based on SQL



What Hive is NOT

Hive is not designed for online transaction processing and does not offer real-time queries and row level updates. It is best used for batch jobs over large sets of immutable data (like web logs, etc.).

Sample HiveQL

The Query compiler uses the information stored in the metastore to convert SQL queries into a sequence of map/reduce jobs, e.g. the following query

`SELECT * FROM t where t.c = 'xyz'`

`SELECT t1.c2 FROM t1 JOIN t2 ON (t1.c1 = t2.c1)`

`SELECT t1.c1, count(1) from t1 group by t1.c1`

Hive Built in Functions

Return Type	Function Name (Signature)	Description
BIGINT	round(double a)	returns the rounded BIGINT value of the double
BIGINT	floor(double a)	returns the maximum BIGINT value that is equal or less than the double
BIGINT	ceil(double a)	returns the minimum BIGINT value that is equal or greater than the double
double	rand(), rand(int seed)	returns a random number (that changes from row to row). Specifying the seed will make sure the generated random number sequence is deterministic.
string	concat(string A, string B,...)	returns the string resulting from concatenating B after A. For example, concat('foo', 'bar') results in 'foobar'. This function accepts arbitrary number of arguments and return the concatenation of all of them.
string	substr(string A, int start)	returns the substring of A starting from start position till the end of string A. For example, substr('foobar', 4) results in 'bar'
string	substr(string A, int start, int length)	returns the substring of A starting from start position with the given length e.g. substr('foobar', 4, 2) results in 'ba'
string	upper(string A)	returns the string resulting from converting all characters of A to upper case e.g. upper('fOoBaR') results in 'FOOBAR'
string	ucase(string A)	Same as upper
string	lower(string A)	returns the string resulting from converting all characters of B to lower case e.g. lower('fOoBaR') results in 'foobar'
string	lcase(string A)	Same as lower
string	trim(string A)	returns the string resulting from trimming spaces from both ends of A e.g. trim(' foobar ') results in 'foobar'
string	ltrim(string A)	returns the string resulting from trimming spaces from the beginning(left hand side) of A. For example, ltrim(' foobar ') results in 'foobar'
string	rtrim(string A)	returns the string resulting from trimming spaces from the end(right hand side) of A. For example, rtrim(' foobar ') results in 'foobar'
string	regexp_replace(string A, string B, string C)	returns the string resulting from replacing all substrings in B that match the Java regular expression syntax(See Java regular expressions syntax) with C. For example, regexp_replace('foobar', 'oo ar',) returns 'fb'
string	from_unixtime(int unixtime)	convert the number of seconds from unix epoch (1970-01-01 00:00:00 UTC) to a string representing the timestamp of that moment in the current system time zone in the format of "1970-01-01 00:00:00"
string	to_date(string timestamp)	Return the date part of a timestamp string: to_date("1970-01-01 00:00:00") = "1970-01-01"
int	year(string date)	Return the year part of a date or a timestamp string: year("1970-01-01 00:00:00") = 1970, year("1970-01-01") = 1970
int	month(string date)	Return the month part of a date or a timestamp string: month("1970-11-01 00:00:00") = 11, month("1970-11-01") = 11
int	day(string date)	Return the day part of a date or a timestamp string: day("1970-11-01 00:00:00") = 1, day("1970-11-01") = 1
string	get_json_object(string json_string, string path)	Extract json object from a json string based on json path specified, and return json string of the extracted json object. It will return null if the input json string is invalid

Hive Commands

Command Line

Function	Hive
Run query	hive -e 'select a.col from tab1 a'
Run query silent mode	hive -S -e 'select a.col from tab1 a'
Set hive config variables	hive -e 'select a.col from tab1 a' -hiveconf hive.root.logger=DEBUG,console
Use initialization script	hive -i initialize.sql
Run non-interactive script	hive -f script.sql

Hive Shell

Function	Hive
Run script inside shell	source file_name
Run ls (dfs) commands	dfs -ls /user
Run ls (bash command) from shell	!ls
Set configuration variables	set mapred.reduce.tasks=32
TAB auto completion	set hive.<TAB>
Show all variables starting with hive	set
Revert all variables	reset
Add jar to distributed cache	add jar jar_path
Show all jars in distributed cache	list jars
Delete jar from distributed cache	delete jar jar_name

Hive Tables

Managed- CREATE TABLE

LOAD- File moved into Hive's data warehouse directory

DROP- Both data and metadata are deleted.

External- CREATE EXTERNAL TABLE

LOAD- No file moved

DROP- Only metadata deleted

**Use when sharing data between Hive and Hadoop applications
or you want to use multiple schema on the same data**

Hive External Table

- `CREATE EXTERNAL TABLE external_Table (dummy STRING)`
- `LOCATION '/user/notroot/external_table';`

Dropping External Table using Hive
Hive will delete metadata from metastore
Hive will NOT delete the HDFS file
You need to manually delete the HDFS file

HiveQL and MySQL Comparison

Metadata

Function	MySQL	HiveQL
Selecting a database	USE database;	USE database;
Listing databases	SHOW DATABASES;	SHOW DATABASES;
Listing tables in a database	SHOW TABLES;	SHOW TABLES;
Describing the format of a table	DESCRIBE table;	DESCRIBE (FORMATTED EXTENDED) table;
Creating a database	CREATE DATABASE db_name;	CREATE DATABASE db_name;
Dropping a database	DROP DATABASE db_name;	DROP DATABASE db_name (CASCADE);

HiveQL and MySQL Query Comparison

Query

Function	MySQL	HiveQL
Retrieving information	<code>SELECT from_columns FROM table WHERE conditions;</code>	<code>SELECT from_columns FROM table WHERE conditions;</code>
All values	<code>SELECT * FROM table;</code>	<code>SELECT * FROM table;</code>
Some values	<code>SELECT * FROM table WHERE rec_name = "value";</code>	<code>SELECT * FROM table WHERE rec_name = "value";</code>
Multiple criteria	<code>SELECT * FROM table WHERE rec1="value1" AND rec2="value2";</code>	<code>SELECT * FROM TABLE WHERE rec1 = "value1" AND rec2 = "value2";</code>
Selecting specific columns	<code>SELECT column_name FROM table;</code>	<code>SELECT column_name FROM table;</code>
Retrieving unique output records	<code>SELECT DISTINCT column_name FROM table;</code>	<code>SELECT DISTINCT column_name FROM table;</code>
Sorting	<code>SELECT col1, col2 FROM table ORDER BY col2;</code>	<code>SELECT col1, col2 FROM table ORDER BY col2;</code>
Sorting backward	<code>SELECT col1, col2 FROM table ORDER BY col2 DESC;</code>	<code>SELECT col1, col2 FROM table ORDER BY col2 DESC;</code>
Counting rows	<code>SELECT COUNT(*) FROM table;</code>	<code>SELECT COUNT(*) FROM table;</code>
Grouping with counting	<code>SELECT owner, COUNT(*) FROM table GROUP BY owner;</code>	<code>SELECT owner, COUNT(*) FROM table GROUP BY owner;</code>
Maximum value	<code>SELECT MAX(col_name) AS label FROM table;</code>	<code>SELECT MAX(col_name) AS label FROM table;</code>
Selecting from multiple tables (Join same table using alias w/"AS")	<code>SELECT pet.name, comment FROM pet, event WHERE pet.name = event.name;</code>	<code>SELECT pet.name, comment FROM pet JOIN event ON (pet.name = event.name);</code>

Loading Data using Hive

Start Hive

```
$ hive
```

```
[cloudera@quickstart guest1]$ hive
2016-10-13 02:08:13,500 WARN  [main] mapreduce.TableMapReduceUtil: The hbase-prefix-tree module jar containing PrefixTreeCodec is not present. Continuing without it.

Logging initialized using configuration in file:/etc/hive/conf.dist/hive-log4j.properties
WARNING: Hive CLI is deprecated and migration to Beeline is recommended.
hive> ■
```

Quit from Hive

```
hive> quit;
```

Create Hive Table

```
hive> CREATE TABLE TEST_TBL(ID INT,COUNTRY STRING) ROW FORMAT  
DELIMITED FIELDS TERMINATED BY ',' STORED AS TEXTFILE ;
```

```
hive> SHOW TABLES;
```

```
hive> describe test_tbl;
```

```
hive> CREATE TABLE TEST_TBL(ID INT,COUNTRY STRING) ROW FORMAT DELIMITED FIELDS TERMINATED BY ',' STORED AS TEXTFILE ;  
OK  
Time taken: 1.2 seconds  
hive> SHOW TABLES;  
OK  
test_tbl  
Time taken: 0.234 seconds, Fetched: 1 row(s)  
hive> describe test_tbl  
 > ;  
OK  
id          int  
country      string  
Time taken: 0.209 seconds, Fetched: 2 row(s)
```

Reviewing Hive Table in HDFS



File Browser

Search for file name Actions ▾ Move to trash ▾ Upload ▾ New ▾

Home / user / hive / warehouse History

<input type="checkbox"/>	Name	▲ Size	User	Group	Permissions	Date
<input type="checkbox"/>	↳		hive	supergroup	drwxrwxrwx	July 19, 2017 05:36 AM
<input type="checkbox"/>	.		hive	supergroup	drwxrwxrwx	August 30, 2017 09:53 PM
<input type="checkbox"/>	test_tbl		cloudera	supergroup	drwxrwxrwx	August 30, 2017 09:53 PM

Show 45 of 1 items Page 1 of 1 Previous page

Alter and Drop Hive Table

Hive > alter table test_tbl add columns (remarks STRING);

hive > describe test_tbl;

OK

id int

country string

remarks string

Time taken: 0.077 seconds

hive > drop table test_tbl;

OK

Time taken: 0.9 seconds

Preparing Large Dataset

<http://grouplens.org/datasets/movielens/>



The screenshot shows the top navigation bar of the GroupLens website. It features the "grouplens" logo on the left, followed by a horizontal menu with links: "about", "datasets", "publications", and "blog".

MovieLens

GroupLens Research has collected and made available rating data sets from the MovieLens web site (<http://movielens.org>). The data sets were collected over various periods of time, depending on the size of the set. Before using these data sets, please review their README files for the usage licenses and other details.

Help our research lab: Please take a short survey about the MovieLens datasets

MovieLens 100k

100,000 ratings from 1000 users on 1700 movies.

- [README.txt](#)
- [ml-100k.zip](#)
- [Index of unzipped files](#)

MovieLens 1M

1 million ratings from 6000 users on 4000 movies.

- [README.txt](#)

Datasets

[MovieLens](#)

[HetRec 2011](#)

[WikiLens](#)

[Book-Crossing](#)

[Jester](#)

[EachMovie](#)

MovieLens Dataset

1)Open Terminal

2)Type command > **mkdir movielens_dataset**

3)Type command >**cd movielens_dataset**

4)Type command > **wget http://files.grouplens.org/datasets/movielens/ml-100k.zip**

5)Type command > **unzip ml-100k.zip**

6)Type command > **more ml-100k/u.user**

```
[cloudera@quickstart ~]$ mkdir movielens_dataset
[cloudera@quickstart ~]$ cd movielens_dataset/
[cloudera@quickstart movielens_dataset]$ wget http://files.grouplens.org/datasets/movielens/ml-100k.zip
--2016-10-13 03:07:20--  http://files.grouplens.org/datasets/movielens/ml-100k.zip
Resolving files.grouplens.org... 128.101.34.146
Connecting to files.grouplens.org|128.101.34.146|:80... connected.
HTTP request sent, awaiting response... 200 OK
Length: 4924029 (4.7M) [application/zip]
Saving to: "ml-100k.zip"

100%[=====] 4,924,029    226K/s   in 36s

2016-10-13 03:08:02 (135 KB/s) - "ml-100k.zip" saved [4924029/4924029]
```

```
[cloudera@quickstart movielens_dataset]$ unzip ml-100k.zip
Archive:  ml-100k.zip
  creating: ml-100k/
  inflating: ml-100k/allbut.pl
  inflating: ml-100k/mku.sh
  inflating: ml-100k/README
  inflating: ml-100k/u.data
```

```
[cloudera@quickstart movielens_dataset]$ more ml-100k/u.user
1|24|M|technician|85711
2|53|F|other|94043
3|23|M|writer|32067
4|24|M|technician|43537
5|33|F|other|15213
6|42|M|executive|98101
7|57|M|administrator|91344
8|36|M|administrator|05201
9|29|M|student|01002
10|53|M|lawyer|99703
11|39|F|other|30329
12|28|F|other|06405
13|47|M|educator|29206
14|45|M|scientist|55106
15|49|F|educator|97301
16|21|M|entertainment|10309
17|30|M|programmer|06355
```

Moving dataset to HDFS

- 1) Type command > cd ml-100k
- 2) Type command > hadoop fs -mkdir /user/cloudera/movielens
- 3) Type command > hadoop fs -put u.user /user/cloudera/movielens
- 4) Type command > hadoop fs -ls /user/cloudera/movielens

```
[cloudera@quickstart movielens_dataset]$ cd ml-100k
[cloudera@quickstart ml-100k]$ hadoop fs -mkdir /user/cloudera/movielens
[cloudera@quickstart ml-100k]$ hadoop fs -put u.user /user/cloudera/movielens
[cloudera@quickstart ml-100k]$ hadoop fs -ls /user/cloudera/movielens
Found 1 items
-rw-r--r-- 1 cloudera cloudera 22628 2016-10-13 03:16 /user/cloudera/movi
elens/u.user
[cloudera@quickstart ml-100k]$ █
```

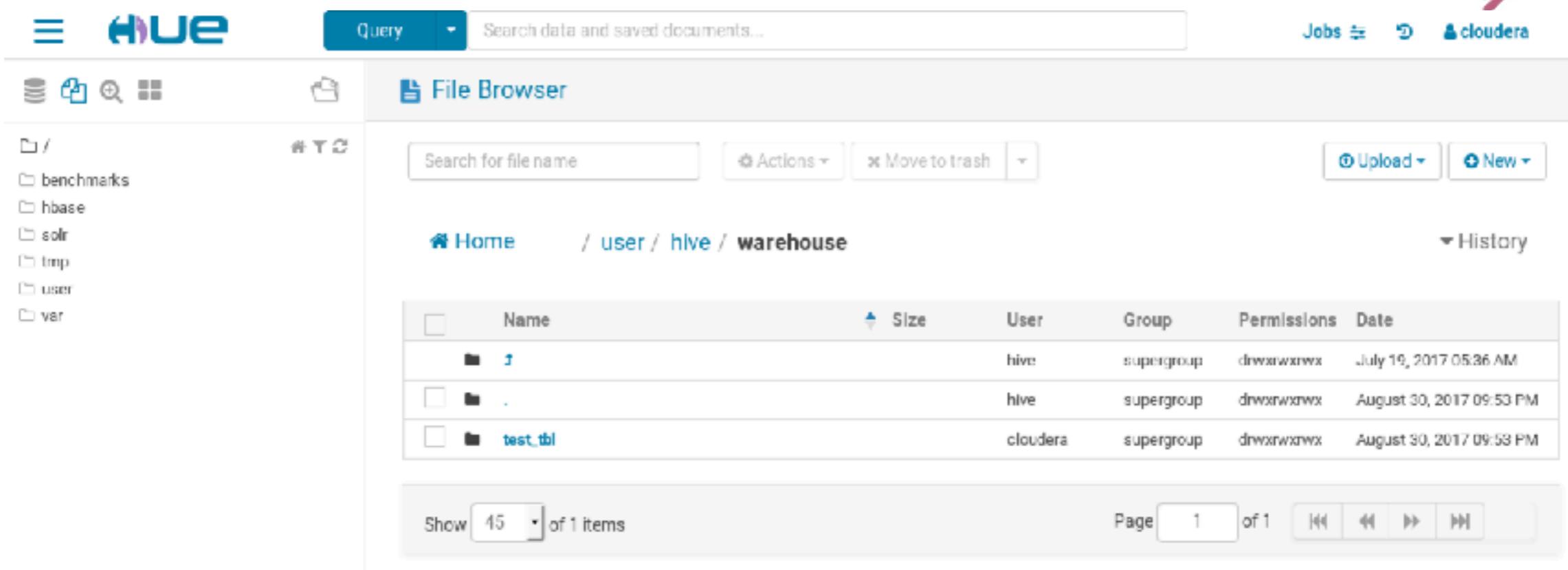
CREATE & SELECT TABLE

Type command > **hive**

```
hive> create external table users (userid int,age int,  
gender string,occupation string ,zipcode string)  
row format delimited fields terminated by '|'  
stored as textfile location '/user/cloudera/movielens';  
  
hive > select * from users;
```

```
hive> CREATE EXTERNAL TABLE users (userid INT, age INT,  
> gender STRING, occupation STRING, zipcode STRING) ROW FORMAT  
> DELIMITED FIELDS TERMINATED BY '|' STORED AS TEXTFILE  
> LOCATION '/user/cloudera/movielens';  
OK  
Time taken: 0.646 seconds  
hive> SELECT * FROM users;  
OK  
1 24 M technician 85711  
2 53 F other 94043  
3 23 M writer 32067  
4 24 M technician 43537  
5 33 F other 15213  
6 42 M executive 98101  
7 57 M administrator 91344  
8 36 M administrator 05201
```

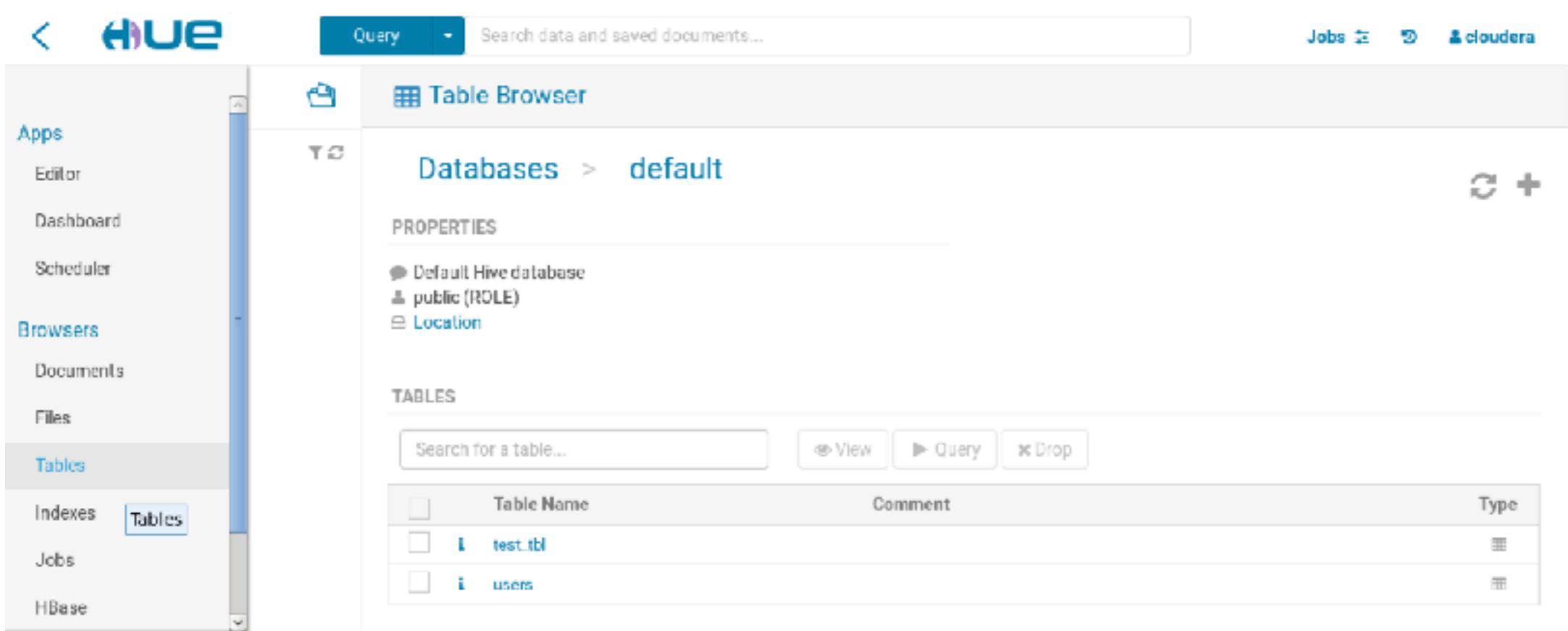
CREATE EXTERNAL TABLE Not File Moved



The screenshot shows the Hue File Browser interface. The left sidebar lists directories: /, benchmarks, hbase, solr, tmp, user, and var. The main area shows a list of files in the /user/hive/warehouse directory. The list includes:

Name	User	Group	Permissions	Date
test	hive	supergroup	drwxrwxrwx	July 16, 2017 05:36 AM
.	hive	supergroup	drwxrwxrwx	August 30, 2017 09:53 PM
test_tbl	cloudera	supergroup	drwxrwxrwx	August 30, 2017 09:53 PM

Below the table, it says "Show 45 of 1 items".



The screenshot shows the Hue Table Browser interface. The left sidebar has sections for Apps (Editor, Dashboard, Scheduler) and Browsers (Documents, Files, Tables, Indexes). The Tables section is selected. The main area shows the Databases > default view. Under PROPERTIES, it shows:

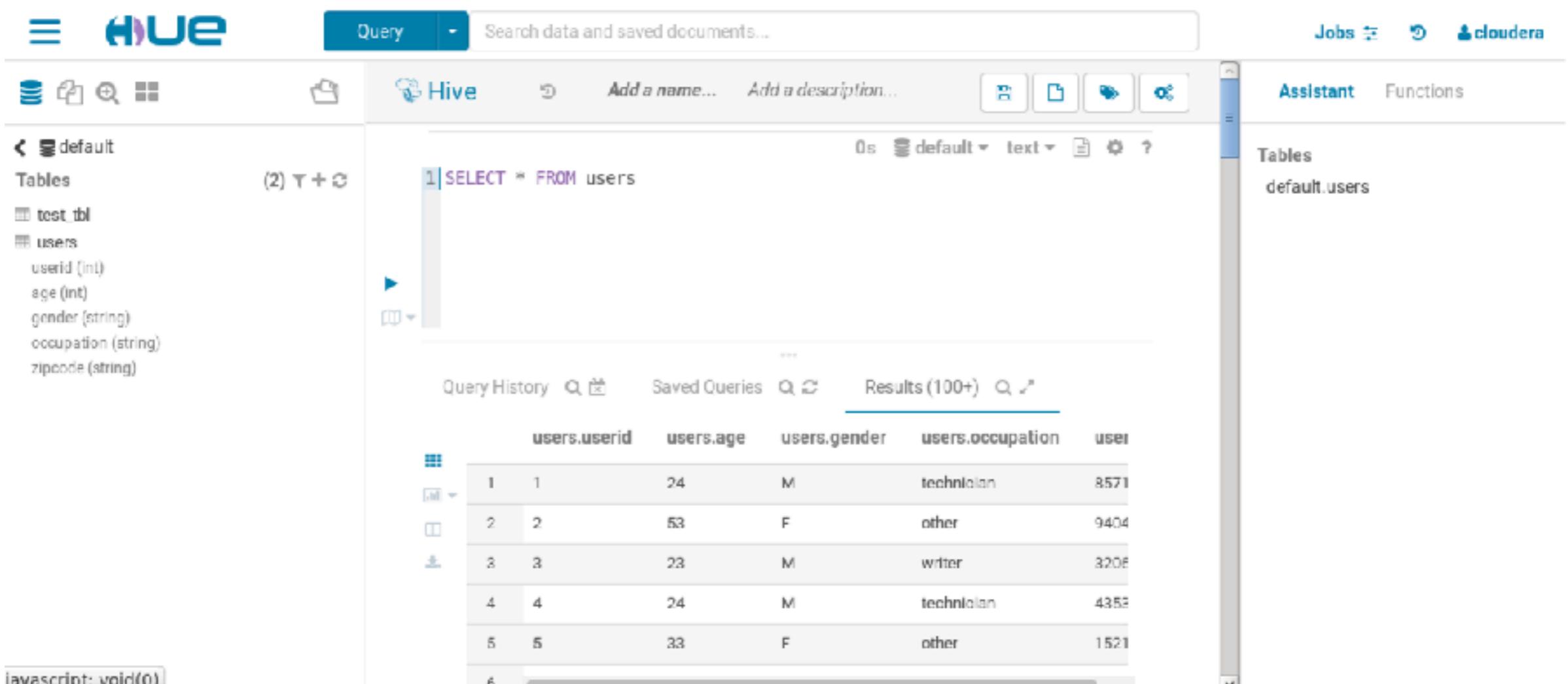
- Default Hive database
- public (ROLE)
- Location

Under TABLES, it shows:

Table Name	Comment	Type
test_tbl		
users		

At the bottom, there are buttons for View, Query, and Drop.

Starting Hive Editor on HUE



The screenshot shows the Hue web interface for managing Apache Hive. On the left, a sidebar lists tables in the 'default' database: 'test_tbl', 'users' (with columns 'userid', 'age', 'gender', 'occupation', 'zipcode'), and 'impressions' (partially visible). The main area is a query editor with the following details:

- Query:** SELECT * FROM users
- Database:** default
- Format:** text
- Results:** 100+ rows

The results table displays the following data:

	users.userid	users.age	users.gender	users.occupation	user
1	1	24	M	technician	8571
2	2	53	F	other	9404
3	3	23	M	writer	3208
4	4	24	M	technician	4353
5	5	33	F	other	1521
6					



IMPALA

open source massively parallel processing (MPP) SQL query engine



Cloudera Impala is a query engine that runs on Apache Hadoop. Impala brings scalable parallel database technology to Hadoop, enabling users to issue low-latency SQL queries to data stored in HDFS and Apache HBase without requiring data movement or transformation.

What is Impala ?

General--- purpose SQL engine

Real--time queries in Apache Hadoop

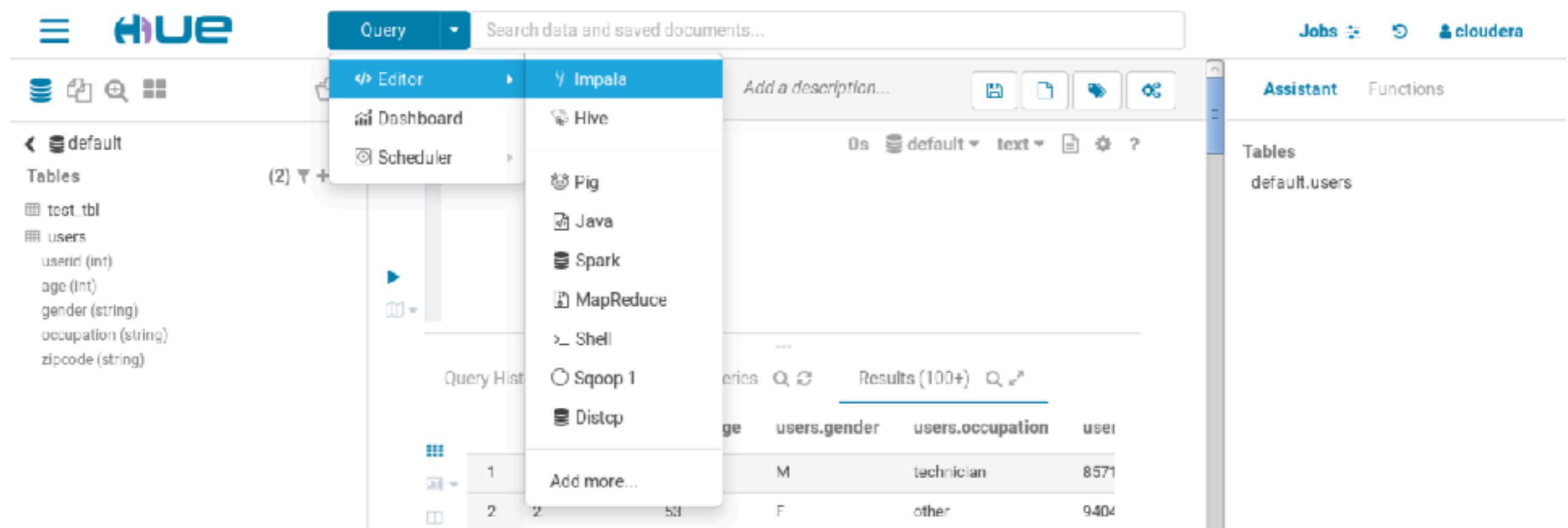
Opensource under Apache License

Runs directly within Hadoop

High performance

- C++ instead of Java**
- Runtime code generator**
- Roughly 4-100 x Hive**

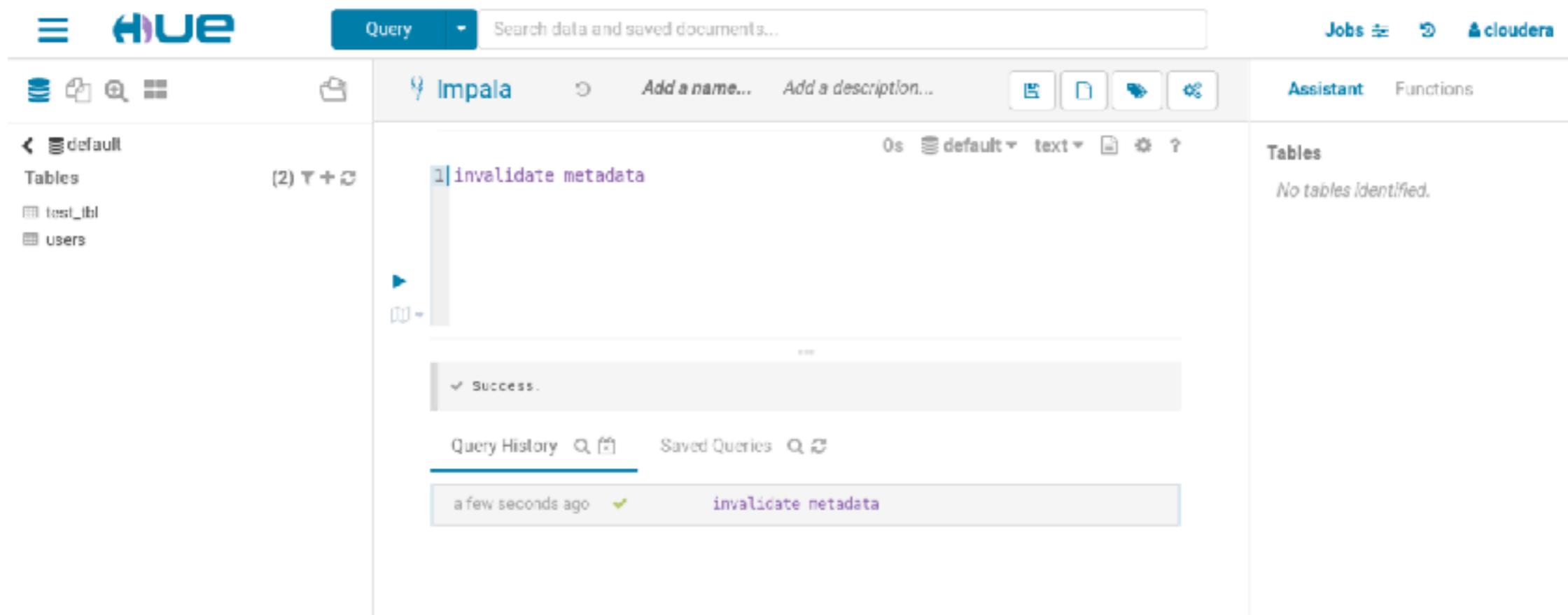
Start Impala Query Editor



The screenshot shows the Hue web interface. The top navigation bar has 'Query' selected. A dropdown menu is open under 'Query' with 'Editor' and 'Impala' options. The 'Impala' option is highlighted. Below this, there are links for Hive, Pig, Java, Spark, MapReduce, Shell, Sqoop 1, and Distcp. To the right, there's a search bar and a toolbar with icons for file operations. The main area displays a 'Query Hist' table with two rows. Row 1 shows 'eric' with gender 'M', occupation 'technician', and count '8571'. Row 2 shows 'other' with gender 'F', occupation 'other', and count '9404'. On the far right, there are tabs for 'Assistant' and 'Functions', and a sidebar showing 'Tables' with 'default.users' listed.

Index	Category	Gender	Occupation	Count
1	eric	M	technician	8571
2	other	F	other	9404

Update the list of tables/metadata by execute the command **invalidate metadata**



The screenshot shows the Hue interface for Apache Impala. In the top navigation bar, 'Query' is selected. The main area displays a single query entry:

```
1| invalidate metadata
```

The status bar at the bottom indicates 'Success.' and shows the query was run 'a few seconds ago'. On the left sidebar, under the 'Tables' section, there is a note: 'No tables identified.'

Run Impala On Terminal

Open Terminal

Type command > impala-shell

[quickstart.cloudera:21000] > select * from users;

```
[[quickstart.cloudera:21000] > select * from users;
Query: select * from users
+-----+-----+-----+-----+
| userid | age  | gender | occupation      | zipcode |
+-----+-----+-----+-----+
| 1      | 24   | M     | technician      | 85711   |
| 2      | 53   | F     | other           | 94043   |
| 3      | 23   | M     | writer          | 32067   |
| 4      | 24   | M     | technician      | 43537   |
| 5      | 33   | F     | other           | 15213   |
| 6      | 42   | M     | executive       | 98101   |
| 7      | 57   | M     | administrator   | 91344   |
| 8      | 36   | M     | administrator   | 05201   |
| 9      | 29   | M     | student         | 01002   |
| 10     | 53   | M     | lawyer          | 90703   |
| 11     | 39   | F     | other           | 30329   |
| 12     | 28   | F     | other           | 06405   |
| 13     | 47   | M     | educator        | 29206   |
| 14     | 45   | M     | scientist       | 55106   |
| 15     | 49   | F     | educator        | 97301   |
| 16     | 21   | M     | entertainment   | 10309   |
```

Apache Flume



Introduction



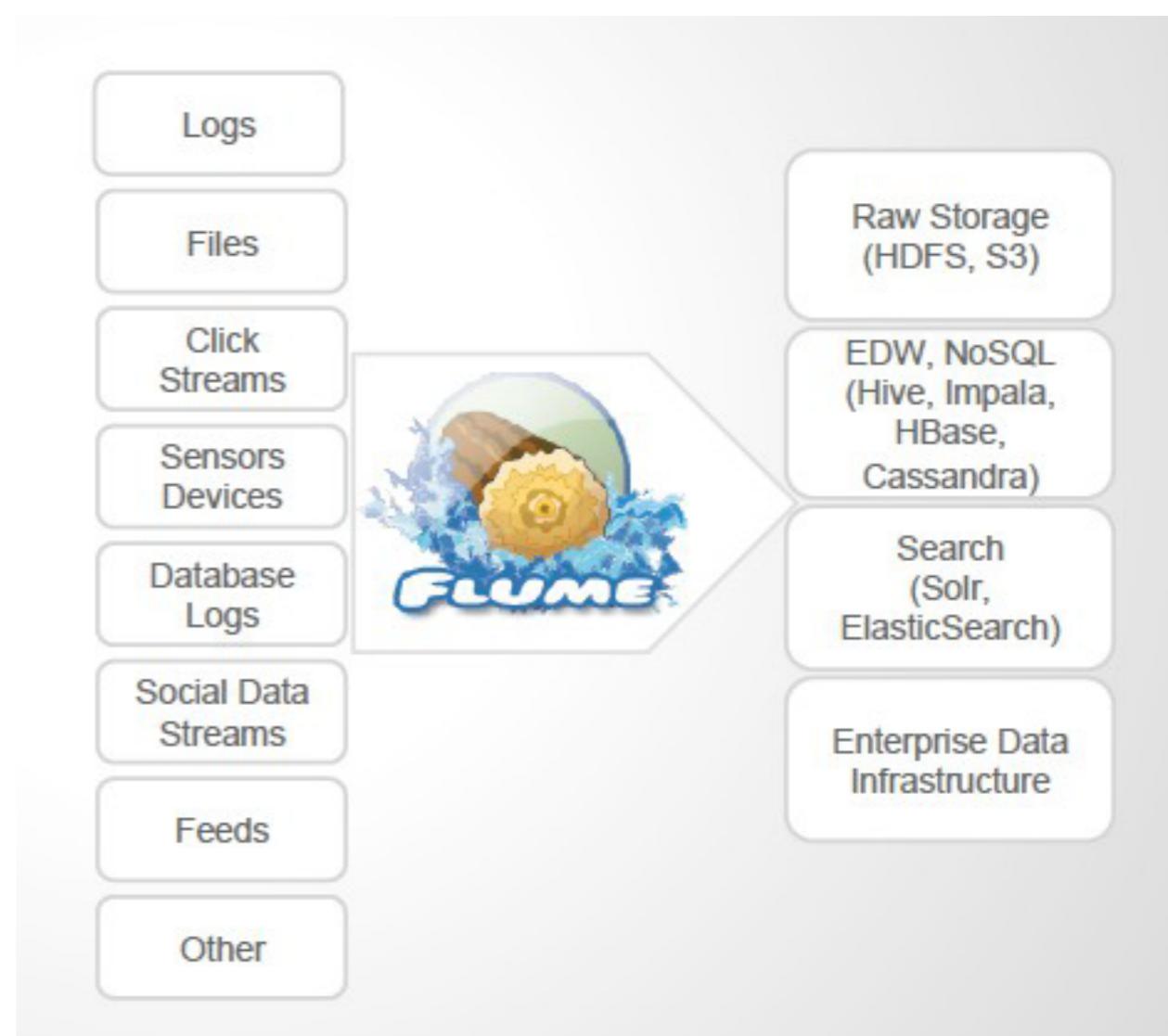
Apache Flume is:

- A distributed data transport and aggregation system for event- or log-structured data
- Principally designed for continuous data ingestion into Hadoop... But more flexible than that

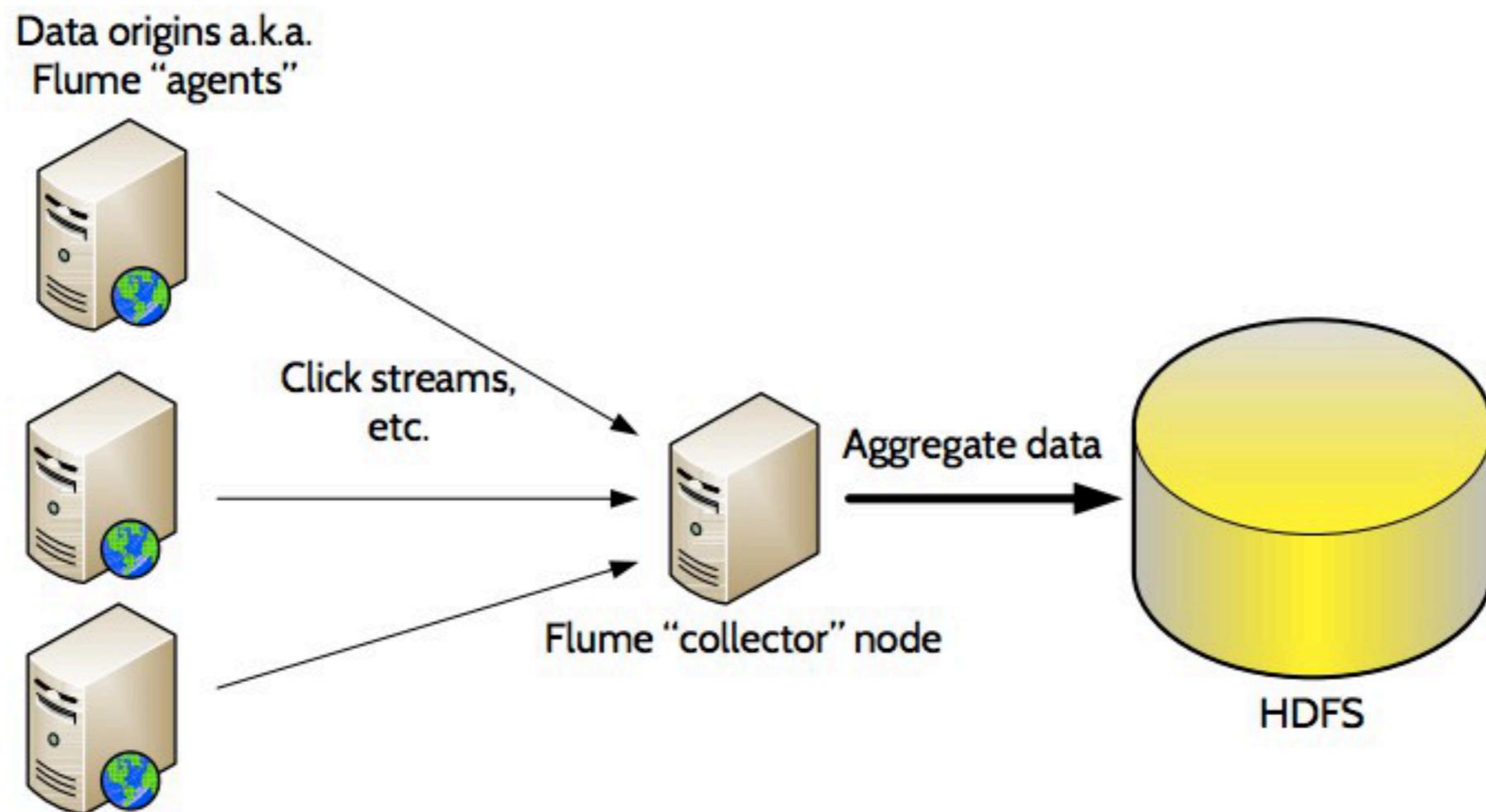
What is Flume?

Apache Flume is a continuous data ingestion system that is...

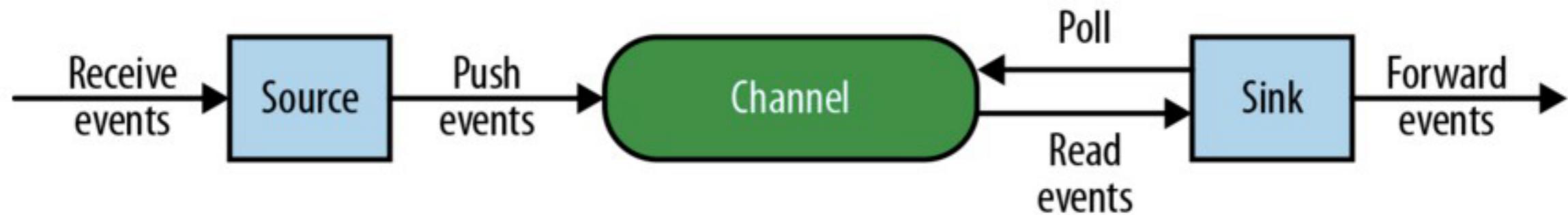
- **open-source,**
- **reliable,**
- **scalable,**
- **manageable,**
- **Customizable,**
- **and designed for**
- Big Data ecosystem**



Architecture Overview

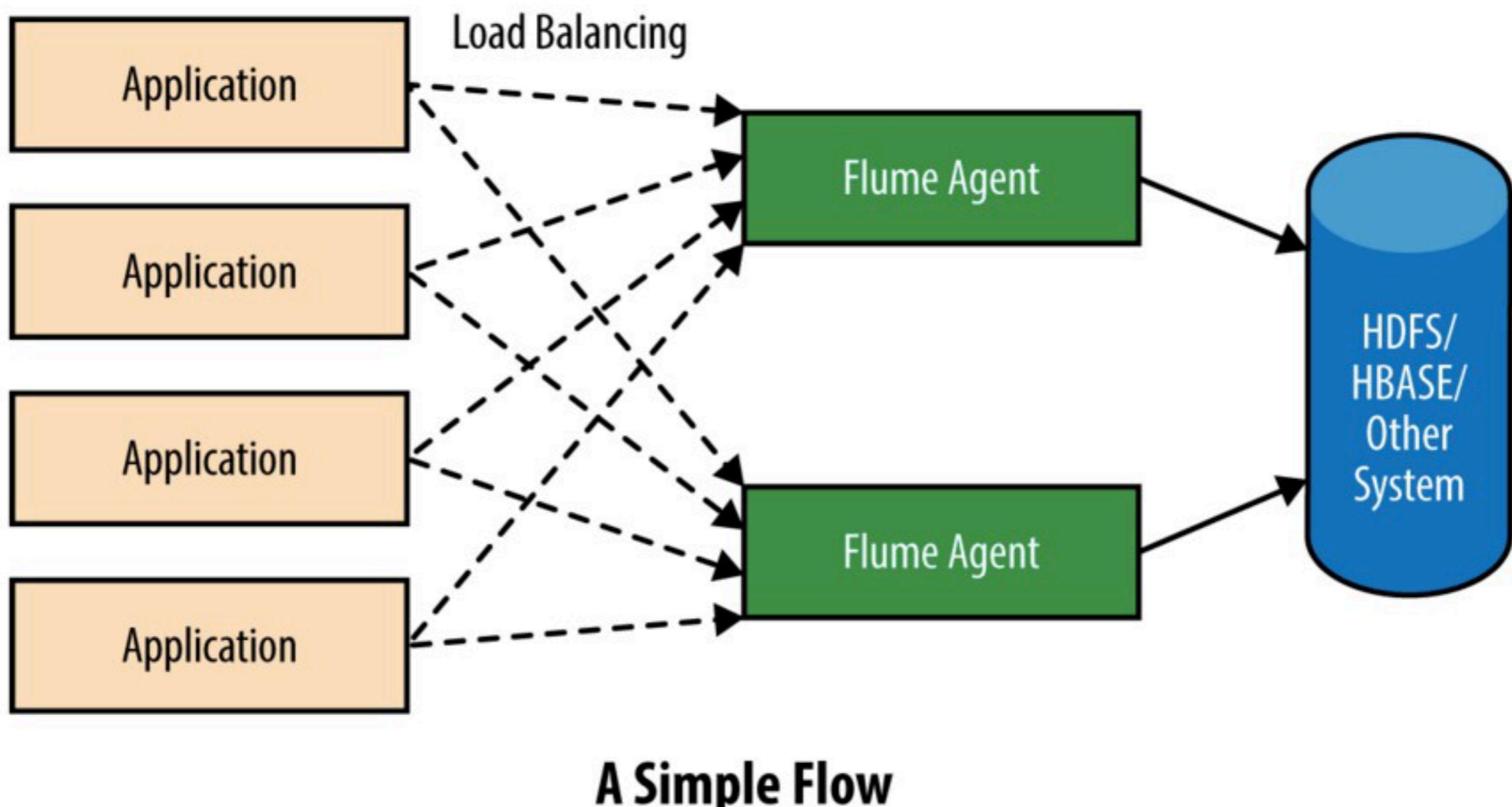


Flume Agent



- A source writes events to one or more channels.
- A channel is the holding area as events are passed from a source to a sink.
- A sink receives events from one channel only.
- An agent can have many channels.

Flow



Source: Using Flume, Hari Shreedharan, 2014

Flume Agent Configuration : Example

```
agent.sources = httpSrc
agent.channels = memory1 memory2
agent.sinks = hdfsSink hbaseSink

agent.sources.httpSrc.type = http
agent.sources.httpSrc.channels = memory1 memory2

# Bind to all interfaces
agent.sources.httpSrc.bind = 0.0.0.0
agent.sources.httpSrc.port = 4353

# Removing this line will disable SSL
agent.sources.httpSrc.ssl = true
agent.sources.httpSrc.keystore = /tmp/keystore
agent.sources.httpSrc.keystore-password = UsingFlume

agent.sources.httpSrc.handler = usingflume.ch03.HTTPSourceXMLHandler
agent.sources.httpSrc.handler.insertTimestamp = true

agent.sources.httpSrc.interceptors = hostInterceptor
agent.sources.httpSrc.interceptors.hostInterceptor.type = host
```

Flume Agent Configuration : Example

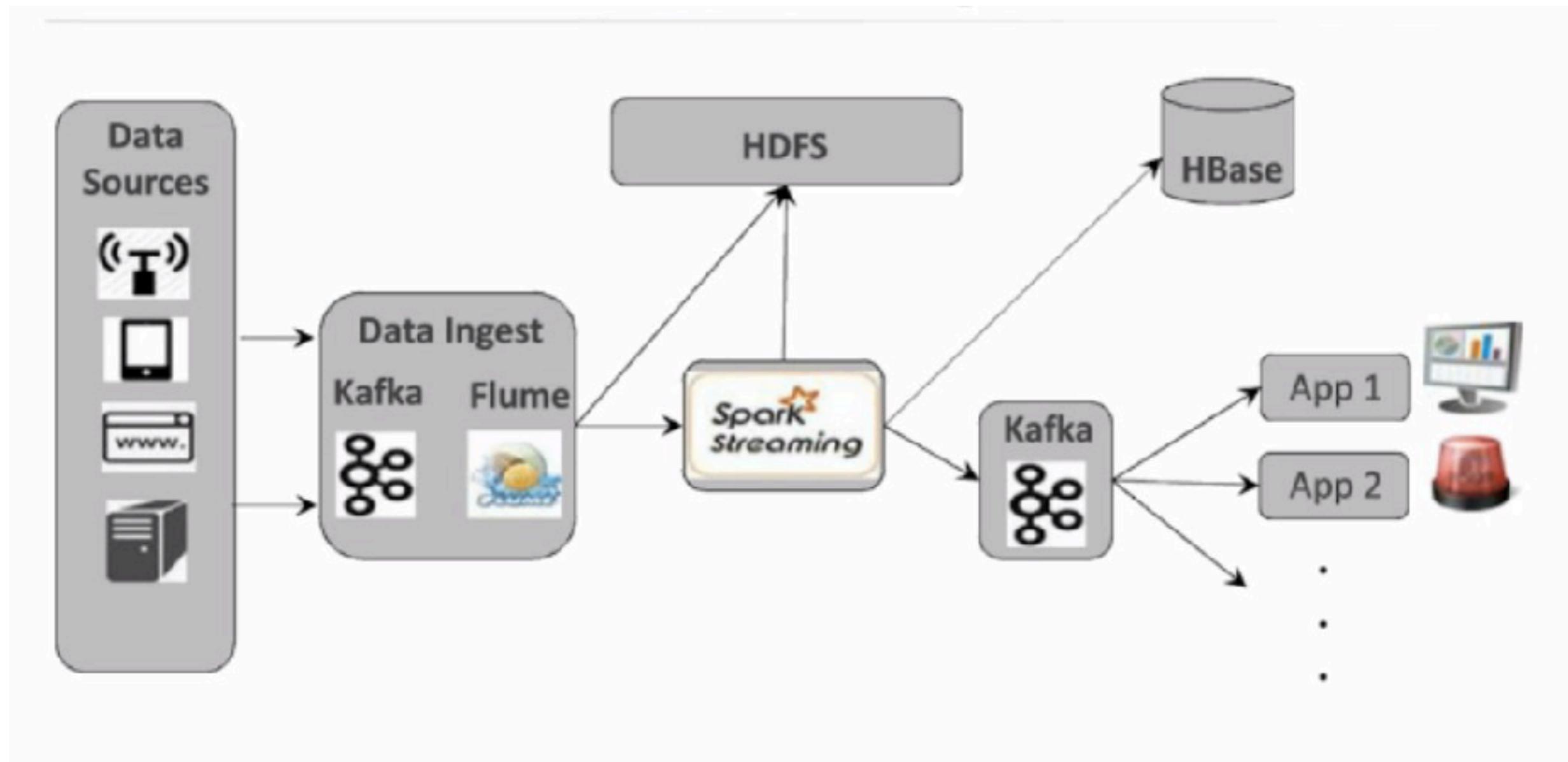
```
# Initializes a memory channel with default configuration
agent.channels.memory1.type = memory

# Initializes a memory channel with default configuration
agent.channels.memory2.type = memory

# HDFS Sink
agent.sinks.hdfsSink.type = hdfs
agent.sinks.hdfsSink.channel = memory1
agent.sinks.hdfsSink.hdfs.path = /Data/UsingFlume/{topic}/{Y}/{m}/{d}/{H}/{M}
agent.sinks.hdfsSink.hdfs.filePrefix = UsingFlumeData

agent.sinks.hbaseSink.type = asynchbase
agent.sinks.hbaseSink.channel = memory2
agent.sinks.hbaseSink.serializer = usingflume.ch05.AsyncHBaseDirectSerializer
agent.sinks.hbaseSink.table = usingFlumeTable
```

Stream Processing Architecture



Flume Loading Data to HDFS

```
$ cd /etc/flume-ng/conf/
```

```
$ sudo rm flume.conf
```

```
$sudo wget https://github.com/bobbylovemovie/trainbigdata/raw/master/flume/flume.conf
```

```
$cat flume.conf
```

```
agent.sources = netsource
agent.sinks = hdfssink
agent.channels = memorychannel
agent.sources.netsource.type = netcat
agent.sources.netsource.bind = localhost
agent.sources.netsource.port = 3030
agent.sources.netsource.interceptors = ts
agent.sources.netsource.interceptors.ts.type = org.apache.flume.interceptor.TimestampInterceptor$Builder
agent.sinks.hdfssink.type = hdfs
agent.sinks.hdfssink.hdfs.path = hdfs://localhost:8020/user/cloudera/flume/events
agent.sinks.hdfssink.hdfs.filePrefix = log
agent.sinks.hdfssink.hdfs.rollInterval = 0
agent.sinks.hdfssink.hdfs.rollCount = 5
agent.sinks.hdfssink.hdfs fileType = DataStream
agent.channels.memorychannel.type = memory
agent.channels.memorychannel.capacity = 100
agent.channels.memorychannel.transactionCapacity = 100
agent.sources.netsource.channels = memorychannel
agent.sinks.hdfssink.channel_ = memorychannel
```

Flume Loading Data to HDFS

start flume-service

```
$sudo service flume-ng-agent restart
```

```
[cloudera@quickstart ~]$ sudo service flume-ng-agent restart
Flume agent is not running                                     [ OK ]
Starting Flume NG agent daemon (flume-ng-agent):           [ OK ]
```

start flume Agent

```
$sudo flume-ng agent --conf /etc/flume-ng/conf/ --conf-file /etc/flume-ng/conf/flume.conf
--name agent -Dflume.root.logger=INFO,console
```

```
2016-10-31 03:20:39,476 (lifecycleSupervisor-1-3) [INFO - org.apache.flume.source.NetcatSource.start(NetcatSource.java:169)] Created serverSocket:sun.nio.ch.ServerSocketChannelImpl[/127.0.0.1:3030]
```

Flume Loading Data to HDFS

Datasource Connect By Telnet

Open New Terminal

```
$sudo yum install telnet
```

```
$telnet localhost 3030
```

```
[cloudera@quickstart ~]$ telnet localhost 3030
Trying 127.0.0.1...
Connected to localhost.
Escape character is '^]'.
Hello Bigdata from Bobby Thaveewat
OK
a,b,c
OK
d,e,f
OK
g,h,i
OK
j,k,l
OK
```

View Result



File Browser

View as
binary

Edit file

Download

View file
location

Refresh

Last modified
09/01/2017
4:01 AM

User
root

Group
cloudera

Size
28 B

Home

Page to of 1

/ user / cloudera / flume / events / log.1504238424110

a,b,c
d,e,f
g,h,i
j,k,l

APACHE SQOOP



Introduction



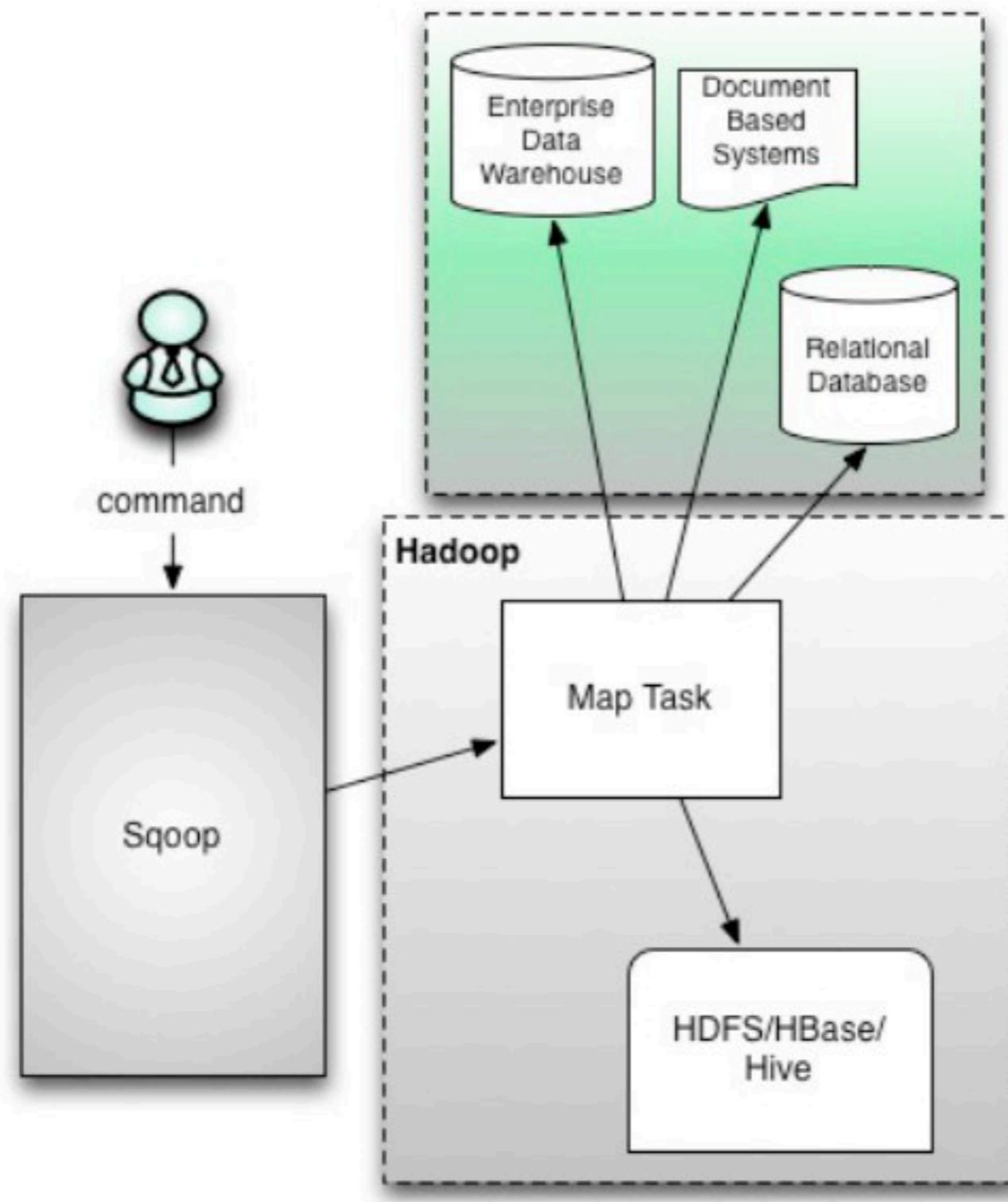
Sqoop (“SQL-to-Hadoop”) is a straightforward command-line tool with the following capabilities:

Imports individual tables or entire databases to files in HDFS

Generates Java classes to allow you to interact with your imported data

Provides the ability to import from SQL databases straight into your Hive data warehouse

Architecture Overview



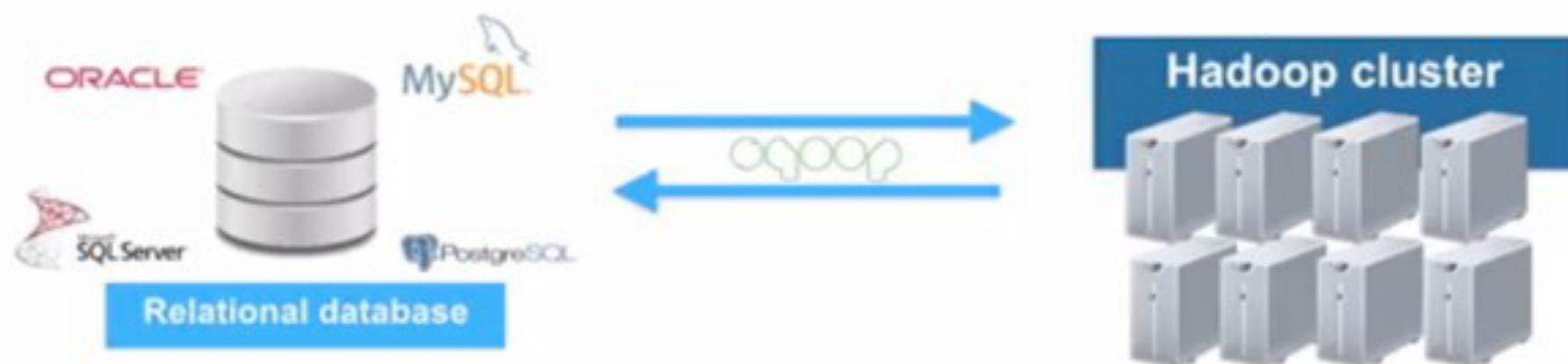
Sqoop Benefit

- Leverages RDBMS metadata to get the column data types
- It is simple to script and uses SQL
- It can be used to handle change data capture by importing daily transactional data to Hadoop
- It uses MapReduce for export and import that enables parallel and efficient data movement

Sqoop Mode

Sqoop import: Data moves from RDBMS to Hadoop

Sqoop export: Data moves from Hadoop to RDBMS



Use Case : Data Consolidation

- Integrate data from various organizational “data stores” to Hadoop for various data processing requirements



Import Commands

Parameters	Description
<code>--connect <jdbc-uri></code>	Specifies the server or database to connect to. It also specifies the port. For example: <code>--connect jdbc:mysql://host:port/databaseName</code>
<code>--connection-manager <class-name></code>	Specifies the connection manager class name.
<code>--driver <class-name></code>	Specifies the fully qualified name of the JDBC driver class.
<code>--hadoop-home <dir></code>	This parameter is used to override the <code>\$HADOOP_HOME</code> environment variable.
<code>-P</code>	If a user doesn't want to specify the database password along with the command, we can use the <code>-P</code> option to read the password from the console.
<code>--password <password></code>	Sets the authentication password required to connect to the input source.
<code>--username <username></code>	Sets the authentication username.
<code>--connection-param-file <properties-file></code>	Specifies the connection parameter's file.
<code>--help</code>	This option will provide the usage instructions.
<code>--verbose</code>	Prints more information during a query execution.

Export Commands

Parameters	Description
<code>--direct</code>	Use the direct mode to perform the export quickly. Note that it is only supported for MySQL.
<code>--export-dir<dir></code>	The location of input files in HDFS.
<code>--table <table-name></code>	Name of the output table (the RDBMS table).
<code>-m, --num-mappers <n></code>	Refers to the number of map tasks.
<code>--update-mode <mode></code>	Specifies how updates are performed when new rows are found with non-matching keys in the database. Legal values for the mode include <code>updateonly</code> (default) and <code>allowinsert</code> .
<code>--update-key <col-name></code>	The value of this column is used to identify the records that a user wants to update during the update mode. Use a comma-separated list of columns if there is more than one column.
<code>--staging-table <staging-table-name></code>	Specifies the name of the staging table. The staging table is used to stage the data before inserting it into the destination table.
<code>--clear-staging-table</code>	This argument is used to clean the data from the staging table.

Loading Data from RDBMS to Hadoop

Configuring MySQL On Cloudera.Quickstart

```
$ sudo /usr/bin/mysql_secure_installation
```

Enter current password for root (enter for none): **cloudera**

OK, successfully used password, moving on...

Set root password? [Y/n] **N**

Remove anonymous users? [Y/n] **Y**

Disallow root login remotely? [Y/n] **N**

Remove test database and access to it [Y/n] **Y**

Reload privilege tables now? [Y/n] **Y**

All done!

Running MySQL

\$ mysql -uroot -p"cloudera"

```
[cloudera@quickstart ~]$ mysql -uroot -p"cloudera"
Welcome to the MySQL monitor.  Commands end with ; or \g.
Your MySQL connection id is 389
Server version: 5.1.73 Source distribution

Copyright (c) 2000, 2013, Oracle and/or its affiliates. All rights reserved.

Oracle is a registered trademark of Oracle Corporation and/or its
affiliates. Other names may be trademarks of their respective
owners.

Type 'help;' or '\h' for help. Type '\c' to clear the current input statement.

mysql> █
```

mysql> show databases;

```
mysql> show databases;
+-----+
| Database |
+-----+
| information_schema |
| cm |
| firehose |
| hue |
| metastore |
| mysql |
| nav |
| navms |
| oozie |
| retail_db |
| rman |
| sentry |
+-----+
12 rows in set (0.01 sec)
```

Prepare a test database table

```
mysql> CREATE DATABASE test_mysql_db;  
mysql> USE test_mysql_db;  
mysql> CREATE TABLE country_tbl(id INT NOT NULL, country  
VARCHAR(50), PRIMARY KEY (id));  
mysql> INSERT INTO country_tbl VALUES(1, 'USA');  
mysql> INSERT INTO country_tbl VALUES(2, 'CANADA');  
mysql> INSERT INTO country_tbl VALUES(3, 'Mexico');  
mysql> INSERT INTO country_tbl VALUES(4, 'Brazil');  
mysql> INSERT INTO country_tbl VALUES(61, 'Japan');  
mysql> INSERT INTO country_tbl VALUES(65, 'Singapore');  
mysql> INSERT INTO country_tbl VALUES(66, 'Thailand');
```

View data in the table

```
mysql> SELECT * FROM country_tbl;
```

id	country
1	USA
2	CANADA
3	Mexico
4	Brazil
61	Japan
65	Singapore
66	Thailand

7 rows in set (0.00 sec)

```
mysql> exit;
```

Importing data from MySQL to HDFS

```
$ sqoop import --connect jdbc:mysql://localhost/test_mysql_db --username root --password cloudera --table country_tbl --target-dir /user/cloudera/test_table -m 1
```



The screenshot shows the Cloudera Manager File Browser interface. The top navigation bar includes 'File Browser' and a 'File' menu. Below the header, there's a 'ACTIONS' dropdown and a breadcrumb navigation path: Home / user / cloudera / test_table / part-m-00000. On the right, there are page navigation controls for 'Page 1 of 1' and icons for back, forward, and search. The main content area displays a list of data rows from the 'country_tbl' table:

Country ID	Country Name
1	USA
2	CANADA
3	Mexico
4	Brazil
61	Japan
65	Singapore
66	Thailand

Importing data from MySQL to Hive Table

```
$ sqoop import --connect jdbc:mysql://localhost/test_mysql_db --  
username root --password cloudera --table country_tbl --hive-import --  
hive-table country -m 1
```



The screenshot shows the Cloudera Manager File Browser interface. On the left, there's a sidebar with actions: View as binary, Edit file, Download, View file location, and Refresh. The main area shows a list of countries with IDs 1 through 6. The URL bar at the top shows the path: /user/hive/warehouse/country/part-m-00000, which is circled in red. The page number is 1 of 1, and there are navigation icons for previous and next pages.

ID	Country
1	USA
2	CANADA
3	Mexico
4	Brazil
5	Japan
6	Singapore
66	Thailand

Reviewing data from Hive Table

```
[cloudera@quickstart ~]$ hive
```

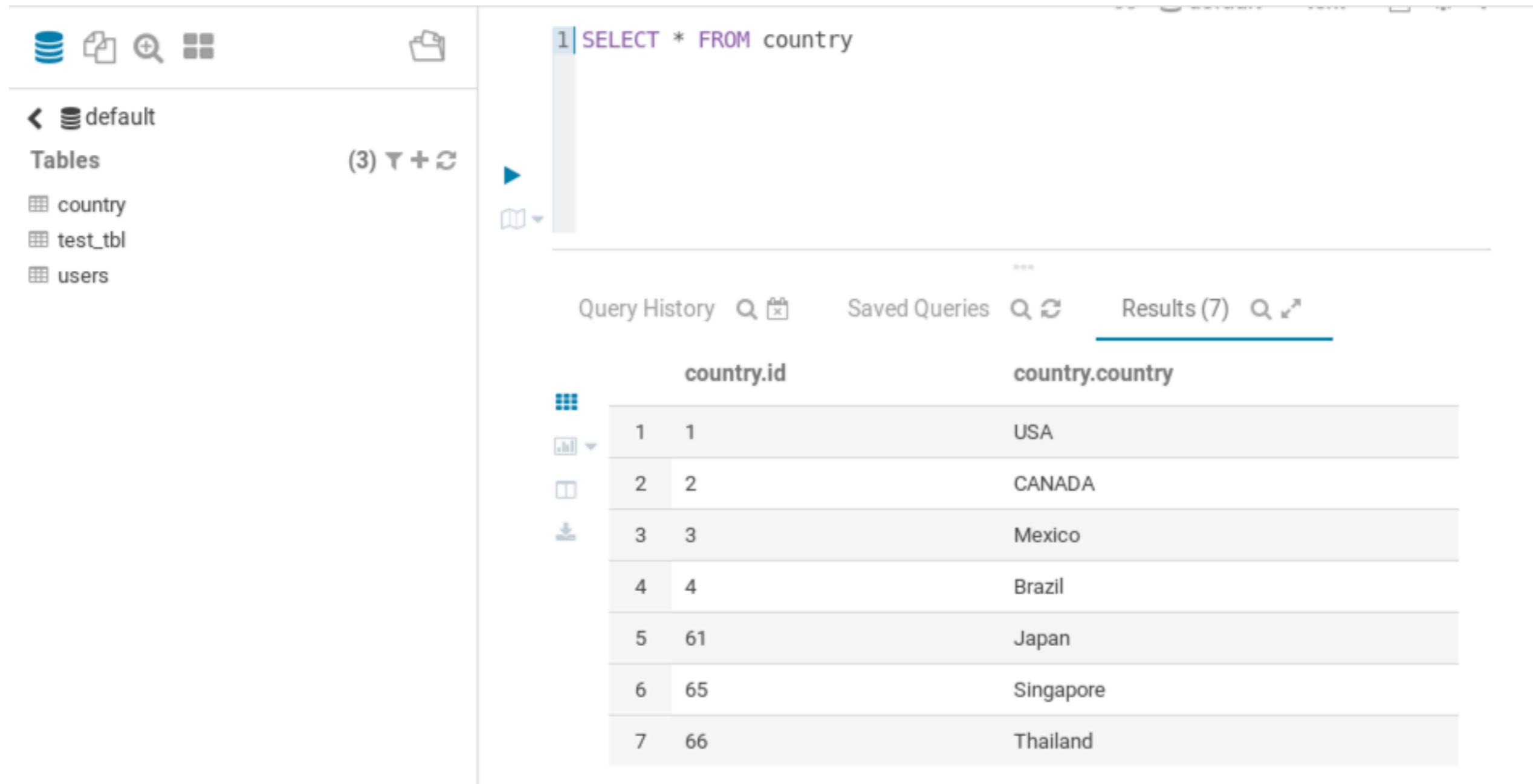
```
Logging initialized using configuration in file:/etc/hive/conf.dist/hive-log4j.properties
WARNING: Hive CLI is deprecated and migration to Beeline is recommended.
```

```
hive> show tables;
```

```
hive> select * from country;
```

```
...
1      USA
2      CANADA
3      Mexico
4      Brazil
61     Japan
65     Singapore
66     Thailand
Time taken: 0.587 seconds, Fetched: 7 row(s)
```

Running from Hue: Beewax



The screenshot shows a database interface with a sidebar on the left and a main query editor and results table on the right.

Left Sidebar:

- Icons for Tables, Views, and Queries.
- A folder icon.
- Text: < default
- Tables section:
 - country (selected)
 - test_tbl
 - users
- (3) T + C

Query Editor:

```
1| SELECT * FROM country
```

Results Table:

	country.id	country.country
1	1	USA
2	2	CANADA
3	3	Mexico
4	4	Brazil
5	61	Japan
6	65	Singapore
7	66	Thailand

Importing data from MySQL to HBase



```
$ sqoop import --connect jdbc:mysql://localhost/test_mysql_db --username root --password cloudera --table country_tbl --hbase-table country --column-family hbase_country_cf --hbase-row-key id --hbase-create-table -m 1
```

Start HBase

```
$ hbase shell
```

```
hbase(main):001:0> list
```

```
hbase(main):001:0> list
TABLE
country
employee
student
3 row(s) in 0.3720 seconds
```

```
=> ["country", "employee", "student"]
```

Viewing Hbase data

```
hbase(main):003:0> scan 'country'
ROW                                COLUMN+CELL
 1                                column=hbase_country_cf:country, timestamp=1468081466623, value=USA
 2                                column=hbase_country_cf:country, timestamp=1468081466623, value=CANADA
 3                                column=hbase_country_cf:country, timestamp=1468081466623, value=Mexico
 4                                column=hbase_country_cf:country, timestamp=1468081466623, value=Brazil
 61                               column=hbase_country_cf:country, timestamp=1468081466623, value=Japan
 65                               column=hbase_country_cf:country, timestamp=1468081466623, value=Singapore
 66                               column=hbase_country_cf:country, timestamp=1468081466623, value=Thailand
7 row(s) in 0.1670 seconds
```

Viewing data from Hbase browser

HBase Browser

Home - Cluster / country Switch Cluster ▾

row_key, row_prefix* +scan_len [col1, family:col2, fam3:, col_prefi] Search icon hbase_country_cf:

grid icon Filter Columns/Families All Sort By ASC ▾

1 x Filter Column Names/Family Sort By ASC ▾ Drop Columns +

hbase_country_cf: country refresh icon trash icon

USA

2 x Filter Column Names/Family Sort By ASC ▾ Drop Columns +

hbase_country_cf: country refresh icon trash icon

CANADA

3 x 75 seconds. Drop Rows Bulk Upload New Row

hbase_country_cf: country

Row Key	Column Family	Column Name	Value
1	hbase_country_cf	country	USA
2	hbase_country_cf	country	CANADA
3	hbase_country_cf	country	75 seconds.



Introduction

A fast and general engine for large scale data processing



An open source big data processing framework built around speed, ease of use, and sophisticated analytics. Spark enables applications in Hadoop clusters to run up to 100 times faster in memory and 10 times faster even when running on disk.

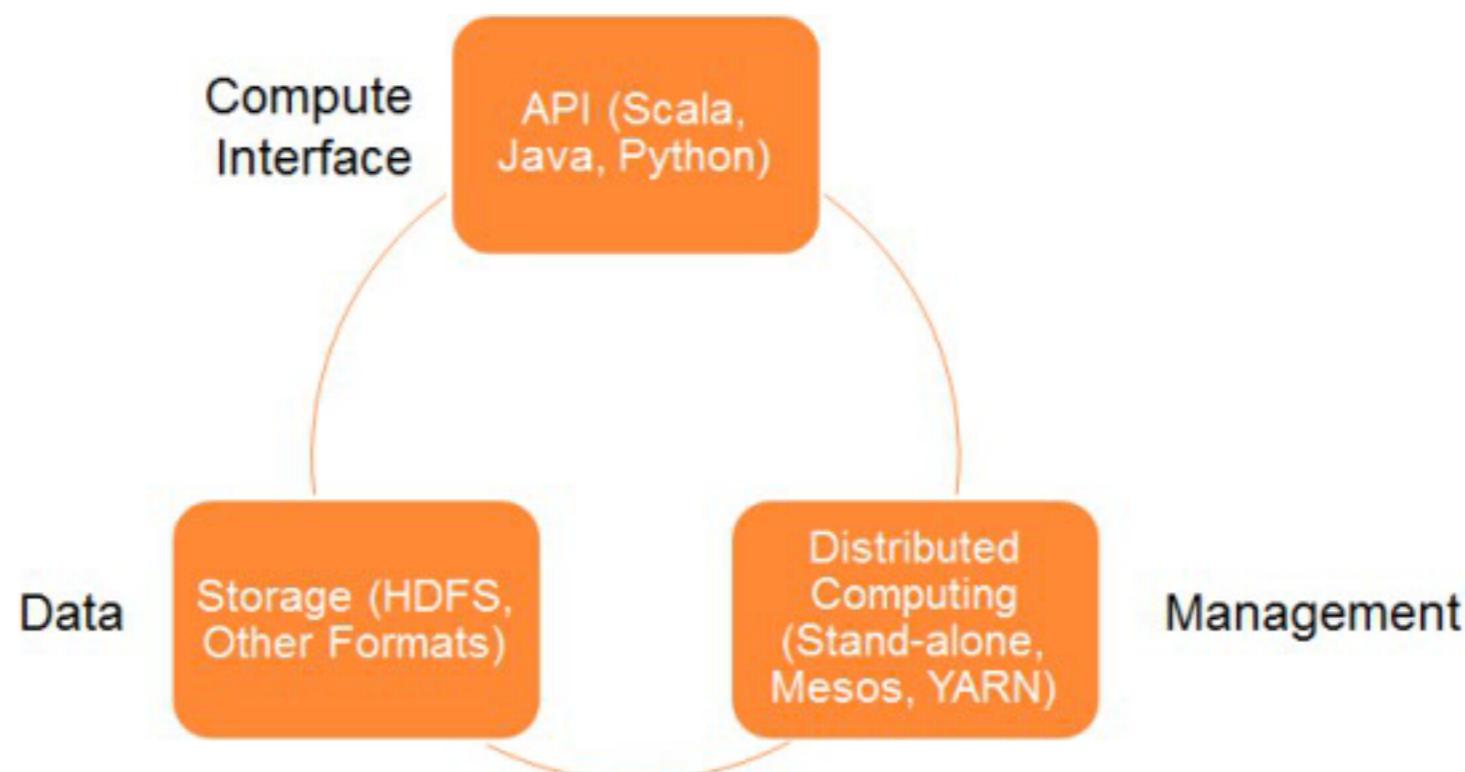
What is Spark ?

Framework for distributed processing.

In-memory, fault tolerant data structures

Flexible APIs in Scala, Java, Python, SQL, R

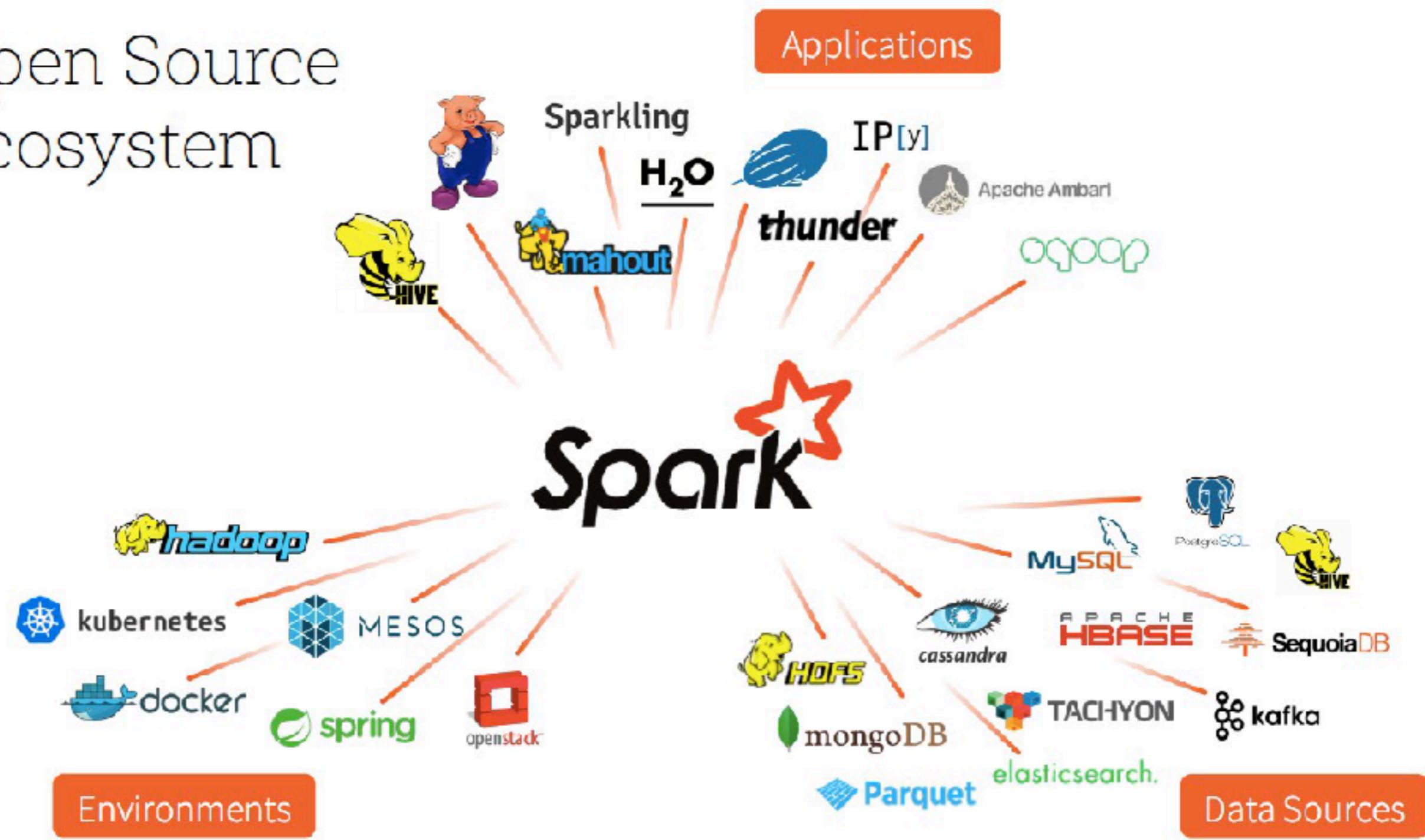
Open source



Why Spark ?

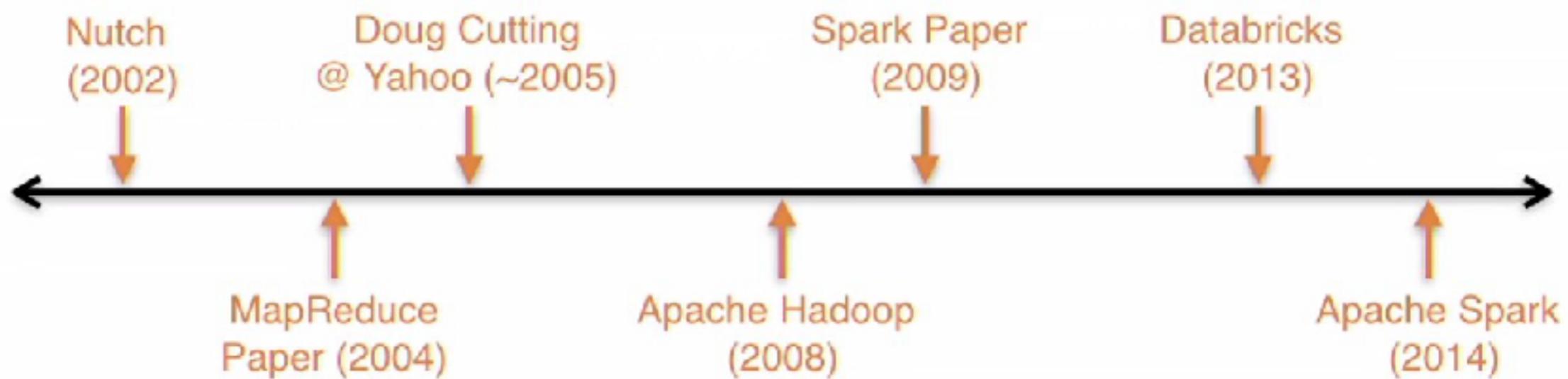
- Handle Petabytes of data
- Significant faster than MapReduce
- Simple and intuitive APIs
- General framework
 - Runs anywhere
 - Handles (most) any I/O
 - Interoperable libraries for specific use-cases

Open Source Ecosystem



Spark: History

- Founded by AMPLab, UC Berkeley
- Created by Matei Zaharia (PhD Thesis)
- Maintained by Apache Software Foundation
- Commercial support by Databricks





Data Science made easy, from ingest to production. Powered by Apache Spark™.

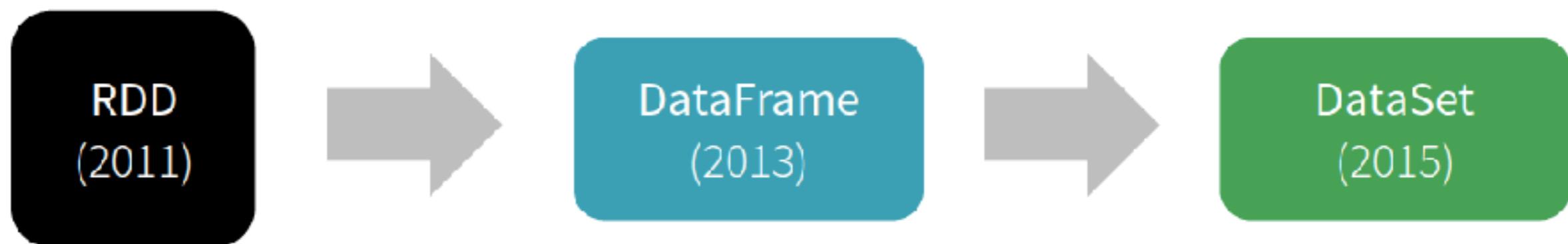
[SIGN UP FOR A 14-DAY FREE TRIAL](#)



LEARN SPARK

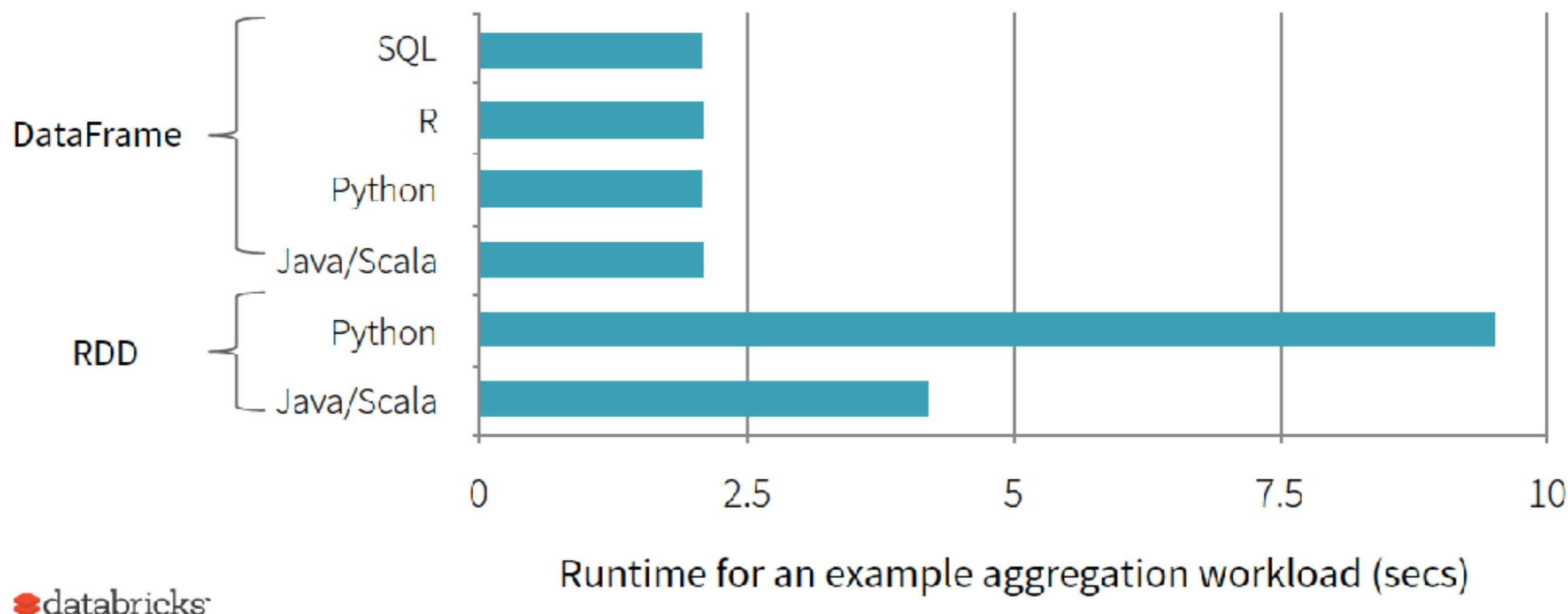
[Join the Community Edition Beta waitlist >](#)

History of Spark APIs



- Distribute collection of JVM objects
- Functional Operators (map, filter, etc.)
- Distribute collection of Row objects
 - Expression-based operations and UDFs
 - Logical plans and optimizer
 - Fast/efficient internal representations
- Internally rows, externally JVM objects
 - “Best of both worlds” **type safe + fast**

Benefit of Logical Plan: Performance Parity Across Languages



What is a RDD ?

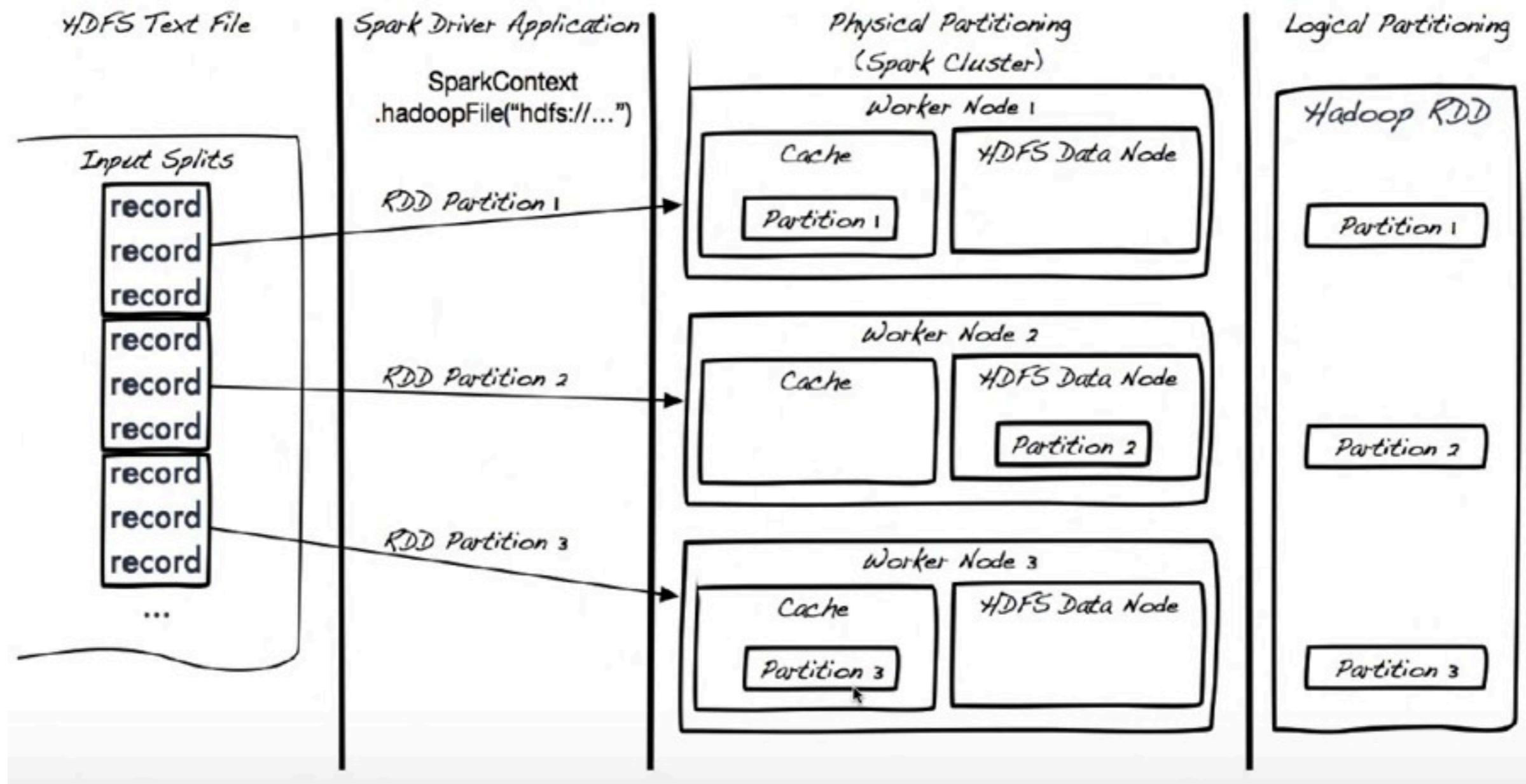
Resilient: if the data in memory (or on a node) is lost, it can be recreated.

Distributed: data is chunked into partitions and stored in memory across the cluster.

Dataset: initial data can come from a table or be created programmatically

RDD Creation

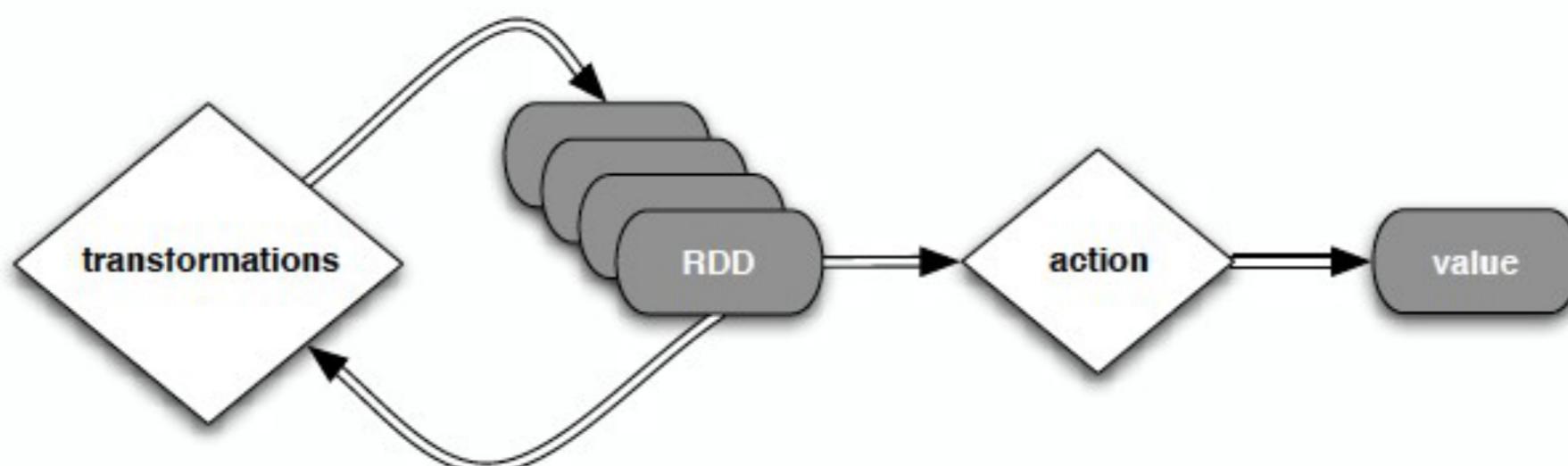
`hdfsData = sc.textFile("hdfs://data.txt")`



RDD: Operations

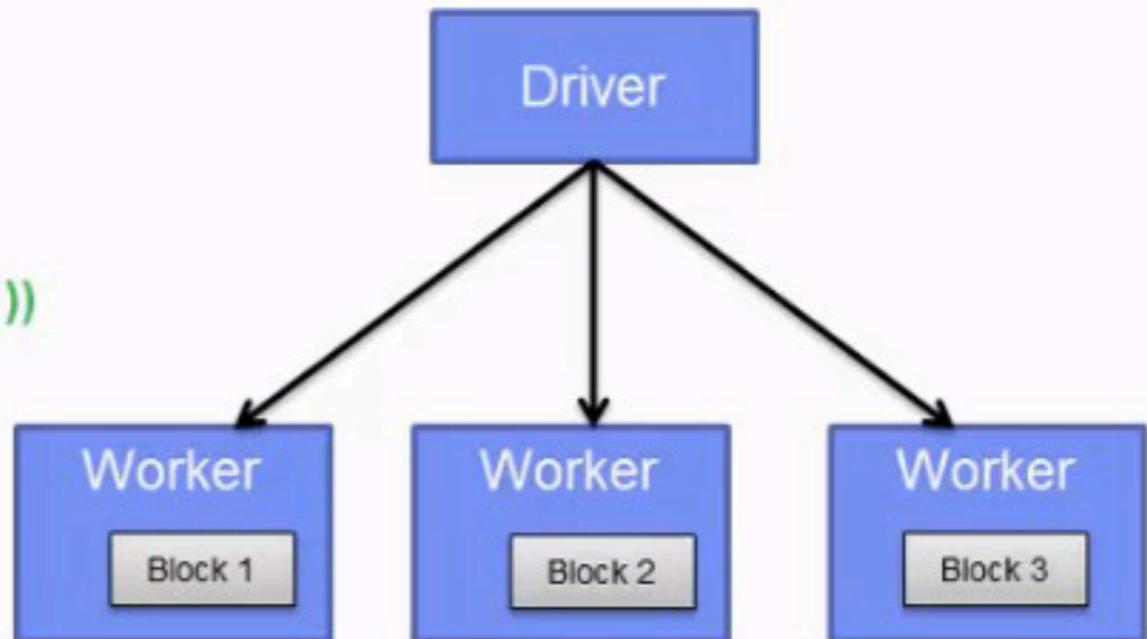
Transformations: transformations are lazy (not computed immediately)

Actions: the transformed RDD gets recomputed when an action is run on it (default)



What happens when an action is executed

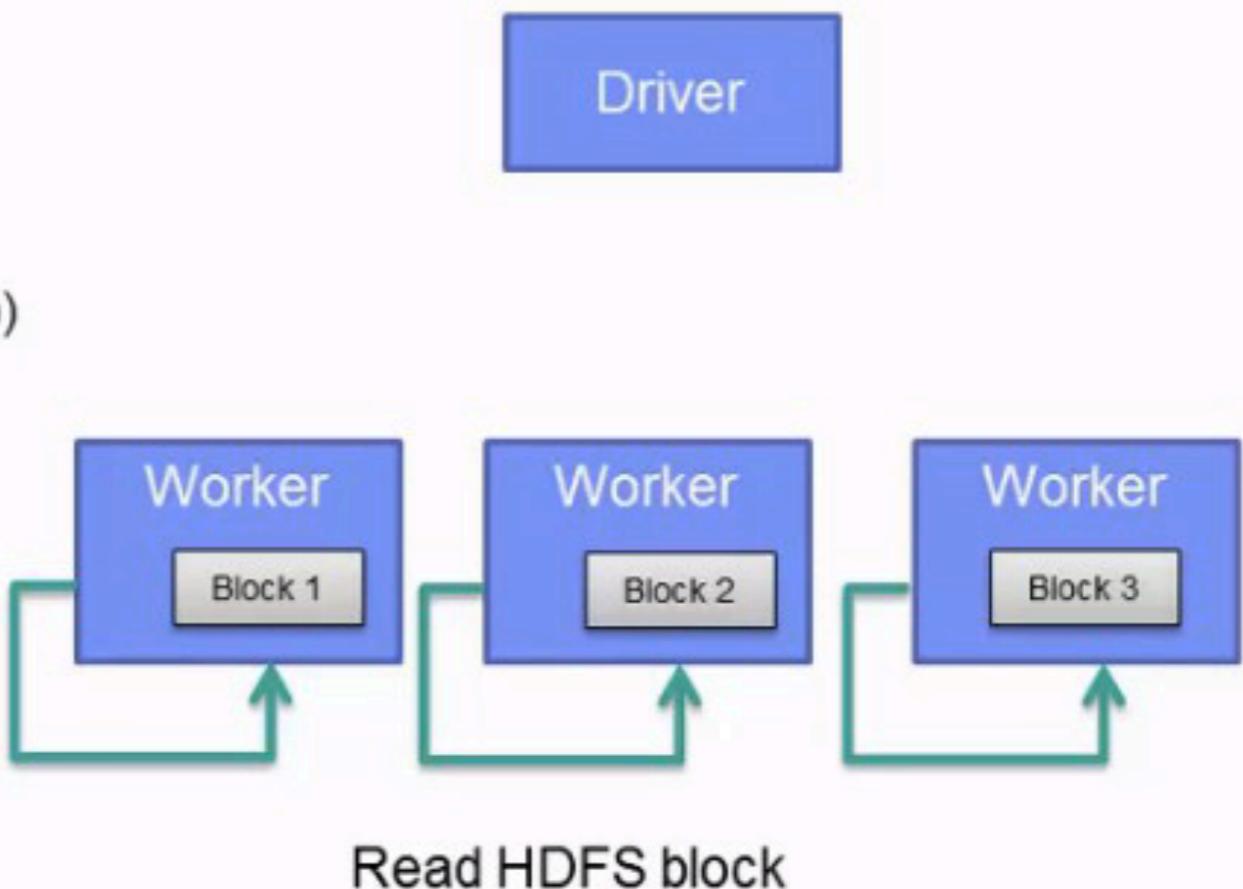
```
// Creating the RDD  
  
val logFile = sc.textFile("hdfs://...")  
  
// Transformations  
  
val errors = logFile.filter(_.startsWith("ERROR"))  
  
val messages = errors.map(_.split("\t")).map(r => r(1))  
  
// Cache  
  
messages.cache()  
  
// Actions  
  
messages.filter(_.contains("mysql")).count()  
  
messages.filter(_.contains("php")).count()
```



Driver sends the code to be
executed on each block

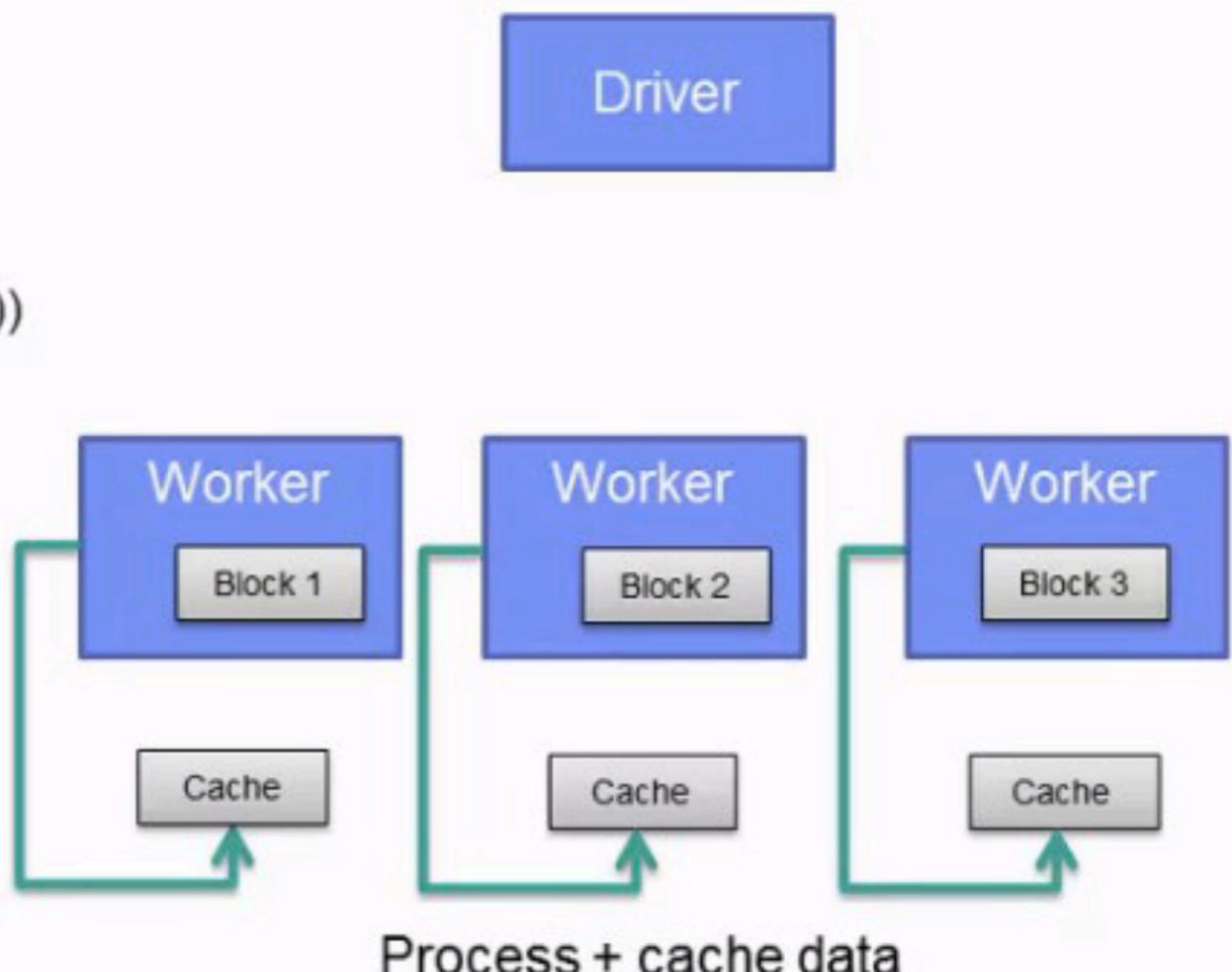
What happens when an action is executed

```
// Creating the RDD  
  
val logFile = sc.textFile("hdfs://...")  
  
// Transformations  
  
val errors = logFile.filter(_.startsWith("ERROR"))  
  
val messages = errors.map(_.split("\t")).map(r => r(1))  
  
//Caching  
  
messages.cache()  
  
// Actions  
  
messages.filter(_.contains("mysql")).count()  
messages.filter(_.contains("php")).count()
```



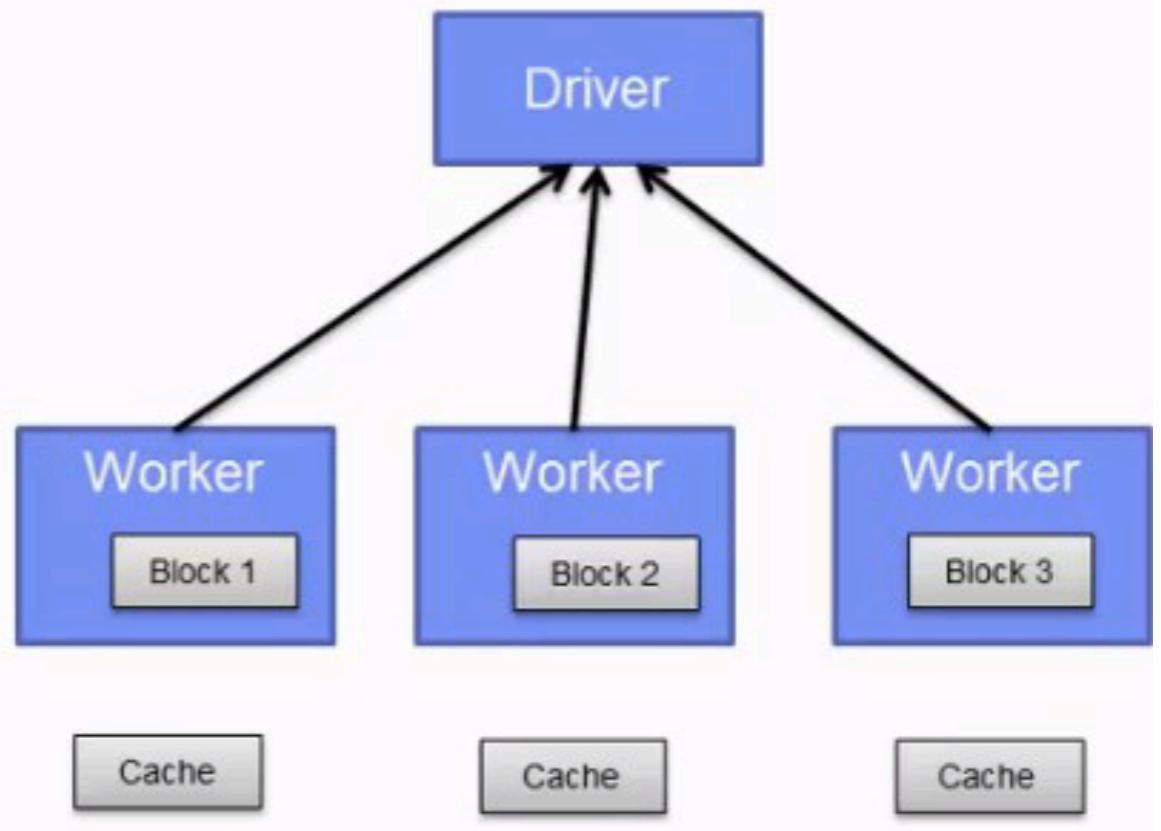
What happens when an action is executed

```
// Creating the RDD  
  
val logFile = sc.textFile("hdfs://...")  
  
// Transformations  
  
val errors = logFile.filter(_.startsWith("ERROR"))  
  
val messages = errors.map(_.split("\t")).map(r => r(1))  
  
//Caching  
  
messages.cache()  
  
// Actions  
  
messages.filter(_.contains("mysql")).count()  
messages.filter(_.contains("php")).count()
```



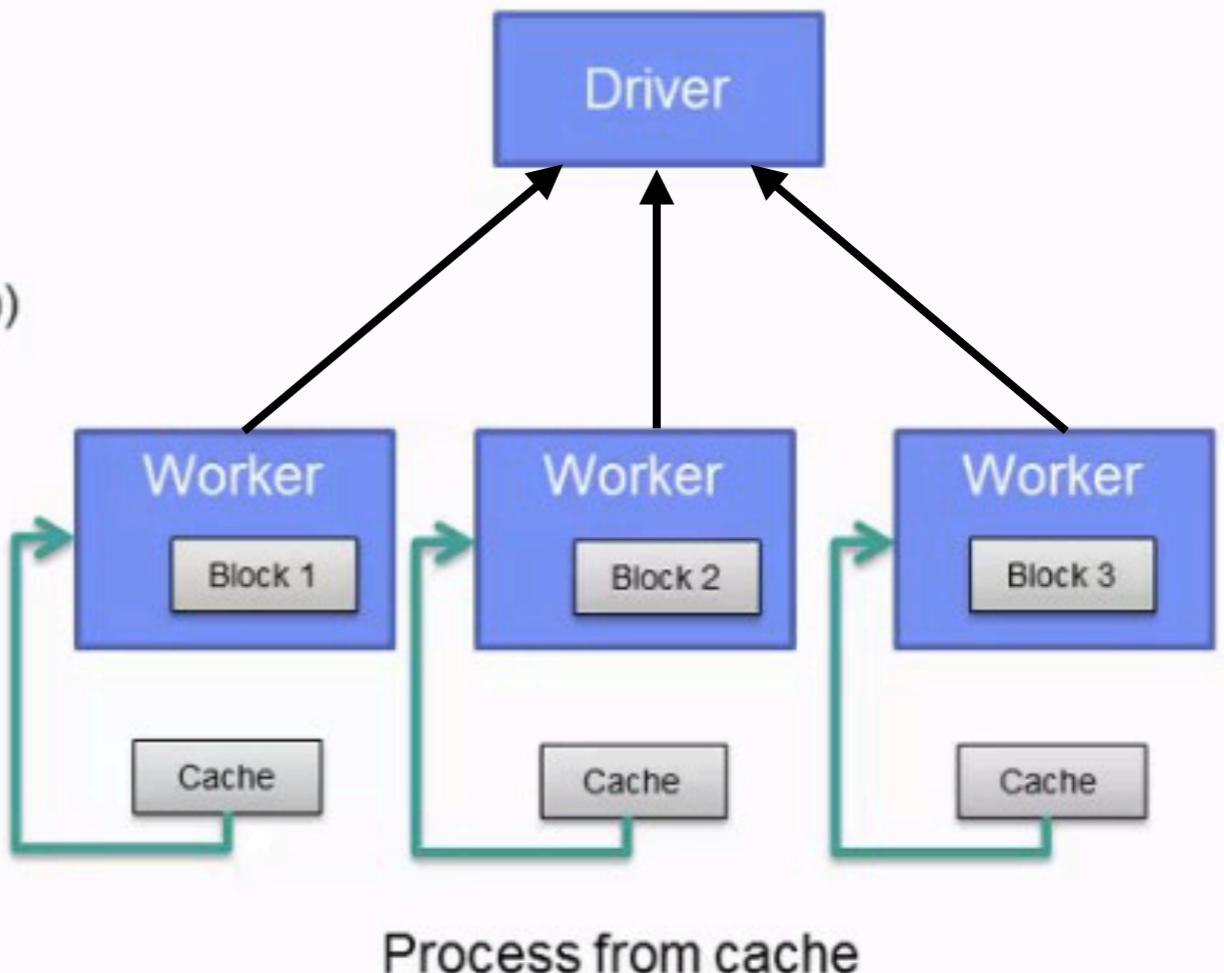
What happens when an action is executed

```
// Creating the RDD
val logFile = sc.textFile("hdfs://...")
// Transformations
val errors = logFile.filter(_.startsWith("ERROR"))
val messages = errors.map(_.split("\t")).map(r => r(1))
//Caching
messages.cache()
// Actions
messages.filter(_.contains("mysql")).count()
messages.filter(_.contains("php")).count()
```



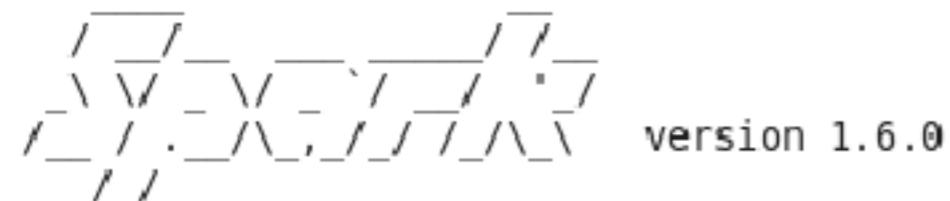
What happens when an action is executed

```
// Creating the RDD
val logFile = sc.textFile("hdfs://...")
// Transformations
val errors = logFile.filter(_.startsWith("ERROR"))
val messages = errors.map(_.split("\t")).map(r => r(1))
//Caching
messages.cache()
// Actions
messages.filter(_.contains("mysql")).count()
messages.filter(_.contains("php")).count()
```



What happens when an action is executed

```
16/11/01 01:03:54 INFO storage.BlockManagerMaster: Registered BlockManager  
Welcome to
```



```
Using Python version 2.6.6 (r266:84292, Jul 23 2015 15:22:56)  
SparkContext available as sc, HiveContext available as sqlContext.
```

```
>>> rdd1 = sc.parallelize(range(1, 11)) [1,2,3,4,5,6,7,8,9,10]  
>>> rdd1 = rdd1.filter(lambda n : n % 2 == 0) [2,4,6,8,10]  
>>> rdd1 = rdd1.map(lambda n : n ** 2) [4, 16, 36, 64, 100]  
>>> rdd1 = rdd1.map(lambda n : n / 2) [2, 8, 18, 32, 50]  
>>> ans = rdd1.reduce(lambda a, b : a + b)
```

```
16/11/01 01:04:52 INFO scheduler.DAGScheduler: ResultStage 0 (reduce at <stdin>:  
1) finished in 0.360 s  
16/11/01 01:04:52 INFO scheduler.TaskSetManager: Finished task 0.0 in stage 0.0  
(TID 0) in 329 ms on localhost (1/1)  
16/11/01 01:04:52 INFO scheduler.TaskSchedulerImpl: Removed TaskSet 0.0, whose t  
asks have all completed, from pool  
16/11/01 01:04:52 INFO scheduler.DAGScheduler: Job 0 finished: reduce at <stdin>  
:1, took 0.990547 s  
>>> print(ans)  
110
```

Spark: Transformation

<i>transformation</i>	<i>description</i>
map(func)	return a new distributed dataset formed by passing each element of the source through a function <i>func</i>
filter(func)	return a new dataset formed by selecting those elements of the source on which <i>func</i> returns true
flatMap(func)	similar to map, but each input item can be mapped to 0 or more output items (so <i>func</i> should return a Seq rather than a single item)
sample(withReplacement, fraction, seed)	sample a fraction <i>fraction</i> of the data, with or without replacement, using a given random number generator <i>seed</i>
union(otherDataset)	return a new dataset that contains the union of the elements in the source dataset and the argument
distinct([numTasks])	return a new dataset that contains the distinct elements of the source dataset

Spark: Transformation

transformation	description
groupByKey([numTasks])	when called on a dataset of (K, V) pairs, returns a dataset of (K, Seq[V]) pairs
reduceByKey(func, [numTasks])	when called on a dataset of (K, V) pairs, returns a dataset of (K, V) pairs where the values for each key are aggregated using the given reduce function
sortByKey([ascending], [numTasks])	when called on a dataset of (K, V) pairs where K implements ordered, returns a dataset of (K, V) pairs sorted by keys in ascending or descending order, as specified in the boolean ascending argument
join(otherDataset, [numTasks])	when called on datasets of type (K, V) and (K, W), returns a dataset of (K, (V, W)) pairs with all pairs of elements for each key
cogroup(otherDataset, [numTasks])	when called on datasets of type (K, V) and (K, W), returns a dataset of (K, Seq[V], Seq[W]) tuples – also called groupWith
cartesian(otherDataset)	when called on datasets of types T and U, returns a dataset of (T, U) pairs (all pairs of elements)

Single RDD Transformation

filter females to analyze female buying patterns

male1, male2, female1 -> female1

map squared values

2, 5, 6 -> 4, 25, 36

flatMap to break up a sentence into words

my name is ray -> my, name, is, ray

find the **distinct** values in a dataset

apple, apple, banana -> apple, banana

sample two values at random

apple, banana, guava -> banana, apple

Multiple RDD Transformation

union

apple, orange, banana, guava,
banana, pear

intersection

banana

subtract anything shown in Dataset B
from Dataset A

apple, orange

cartesian (every possible pair combo)

(apple, guava), (apple, banana), ...

Dataset A

apple
orange
banana

Dataset B

guava
banana
pear

Pair RDD Transformation

- reduceByKey
- groupByKey
- combineByKey
- mapValues
- flatMapValues
- keys
- values
- subtractByKey
- join
- rightOuterJoin
- leftOuterJoin
- cogroup
- sortByKey

Spark:Actions

action	description
reduce(func)	aggregate the elements of the dataset using a function <i>func</i> (which takes two arguments and returns one), and should also be commutative and associative so that it can be computed correctly in parallel
collect()	return all the elements of the dataset as an array at the driver program – usually useful after a filter or other operation that returns a sufficiently small subset of the data
count()	return the number of elements in the dataset
first()	return the first element of the dataset – similar to <i>take(1)</i>
take(n)	return an array with the first <i>n</i> elements of the dataset – currently not executed in parallel, instead the driver program computes all the elements
takeSample(withReplacement, fraction, seed)	return an array with a random sample of <i>num</i> elements of the dataset, with or without replacement, using the given random number generator <i>seed</i>

Spark:Actions

action	description
saveAsTextFile(path)	write the elements of the dataset as a text file (or set of text files) in a given directory in the local filesystem, HDFS or any other Hadoop-supported file system. Spark will call <code>toString</code> on each element to convert it to a line of text in the file
saveAsSequenceFile(path)	write the elements of the dataset as a Hadoop SequenceFile in a given path in the local filesystem, HDFS or any other Hadoop-supported file system. Only available on RDDs of key-value pairs that either implement Hadoop's <code>Writable</code> interface or are implicitly convertible to <code>Writable</code> (Spark includes conversions for basic types like <code>Int</code> , <code>Double</code> , <code>String</code> , etc).
countByKey()	only available on RDDs of type <code>(K, V)</code> . Returns a 'Map' of <code>(K, Int)</code> pairs with the count of each key
foreach(func)	run a function <code>func</code> on each element of the dataset – usually done for side effects such as updating an accumulator variable or interacting with external storage systems

Transformations

```
>>> nums = sc.parallelize([1,2,3])  
>>> squared = nums.map(lambda x : x*x)  
>>> even = squared.filter(lambda x: x%2 == 0)  
>>> evens = nums.flatMap(lambda x: range(x))
```

Actions

```
>>> nums = sc.parallelize([1,2,3])  
>>> nums.collect()  
>>> nums.take(2)  
>>> nums.count()  
>>> nums.reduce(lambda:x, y:x+y)  
>>> nums.saveAsTextFile("hdfs://user/cloudera/output/test")
```

Spark Programming

Functional tools in Python

map

filter

reduce

lambda

IterTools

Chain, flatmap

map

```
>>> a= [1,2,3]
>>> def add1(x) : return x+1
>>> map(add1, a)
Result: [2,3,4]
```

```
val input = sc.parallelize(List(1,2,3,4))
val result = input.map(x => x*x)
println(result.collect().mkString(","))

Result: [1,4,9,16]
```

Filter

```
>>> a= [1,2,3,4]  
>>> def isOdd(x) : return x%2==1  
>>> filter(isOdd, a)
```

Result: [1,3]

Reduce

```
>>> a= [1,2,3,4]  
>>> def add(x,y) : return x+y  
>>> reduce(add, a)  
Result: 10
```

lambda

```
>>> (lambda x: x + 1)(3)
```

Result: 4

```
>>> map((lambda x: x + 1), [1,2,3])
```

Result: [2,3,4]

Exercises

- `(lambda x: 2*x)(3) => ?`
- `map(lambda x: 2*x, [1,2,3]) =>`
- `map(lambda t: t[0], [(1,2), (3,4), (5,6)]) =>`
- `reduce(lambda x,y: x+y, [1,2,3]) =>`
- `reduce(lambda x,y: x+y, map(lambda t: t[0], [(1,2), (3,4), (5,6)]))=>`

Spark Platform

Spark SQL

Spark
Streaming

PySpark
SparkR

MLlib
spark.ml

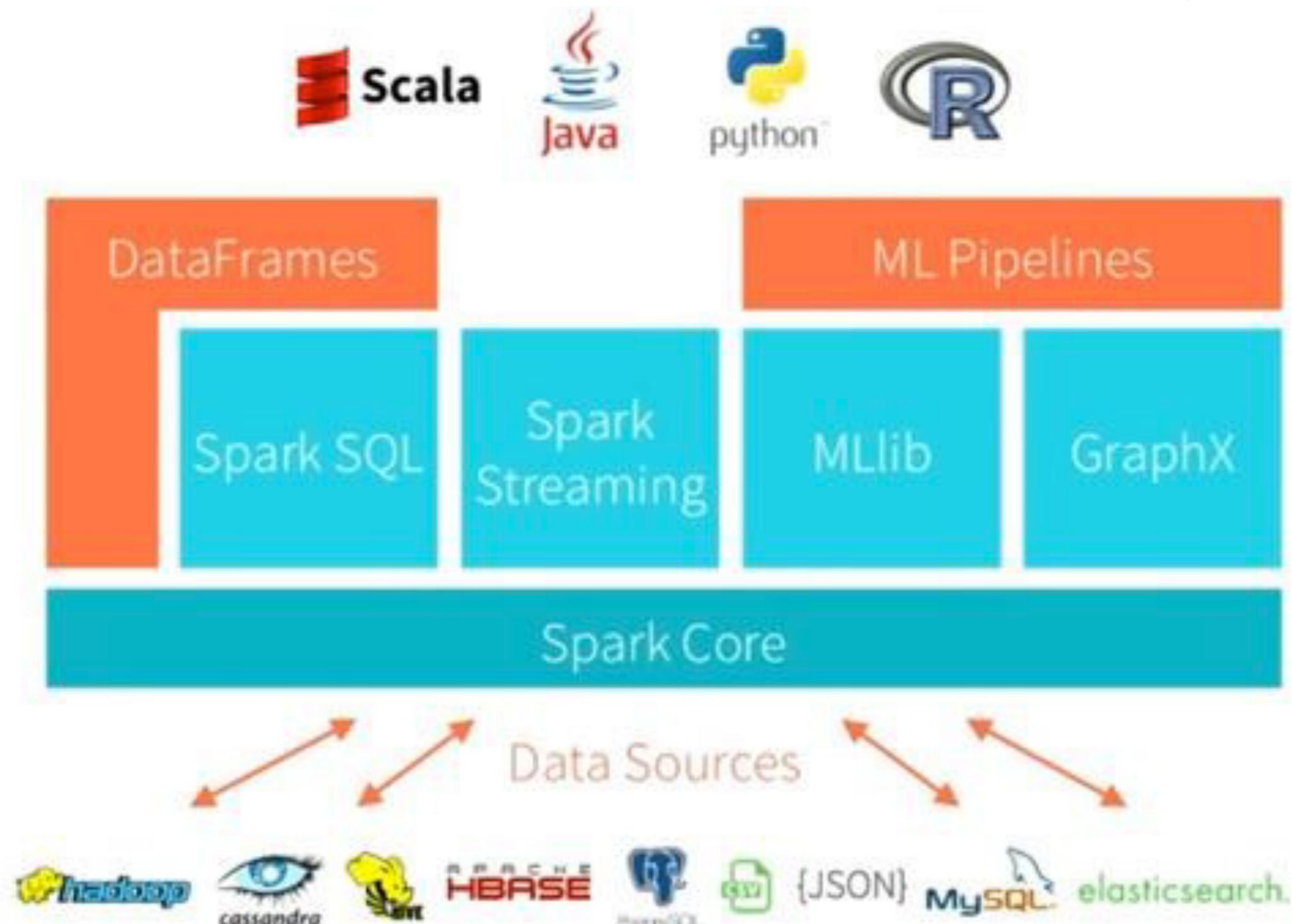
GraphX

Spark Core

Standalone Scheduler

YARN

Mesos



Spark Platform

Spark SQL

- Structured Data
- Querying with SQL/HQL
- DataFrames

Spark Streaming

- Processing of live streams
- Micro-batching

MLlib

- Machine Learning
- Multiple types of ML algorithms

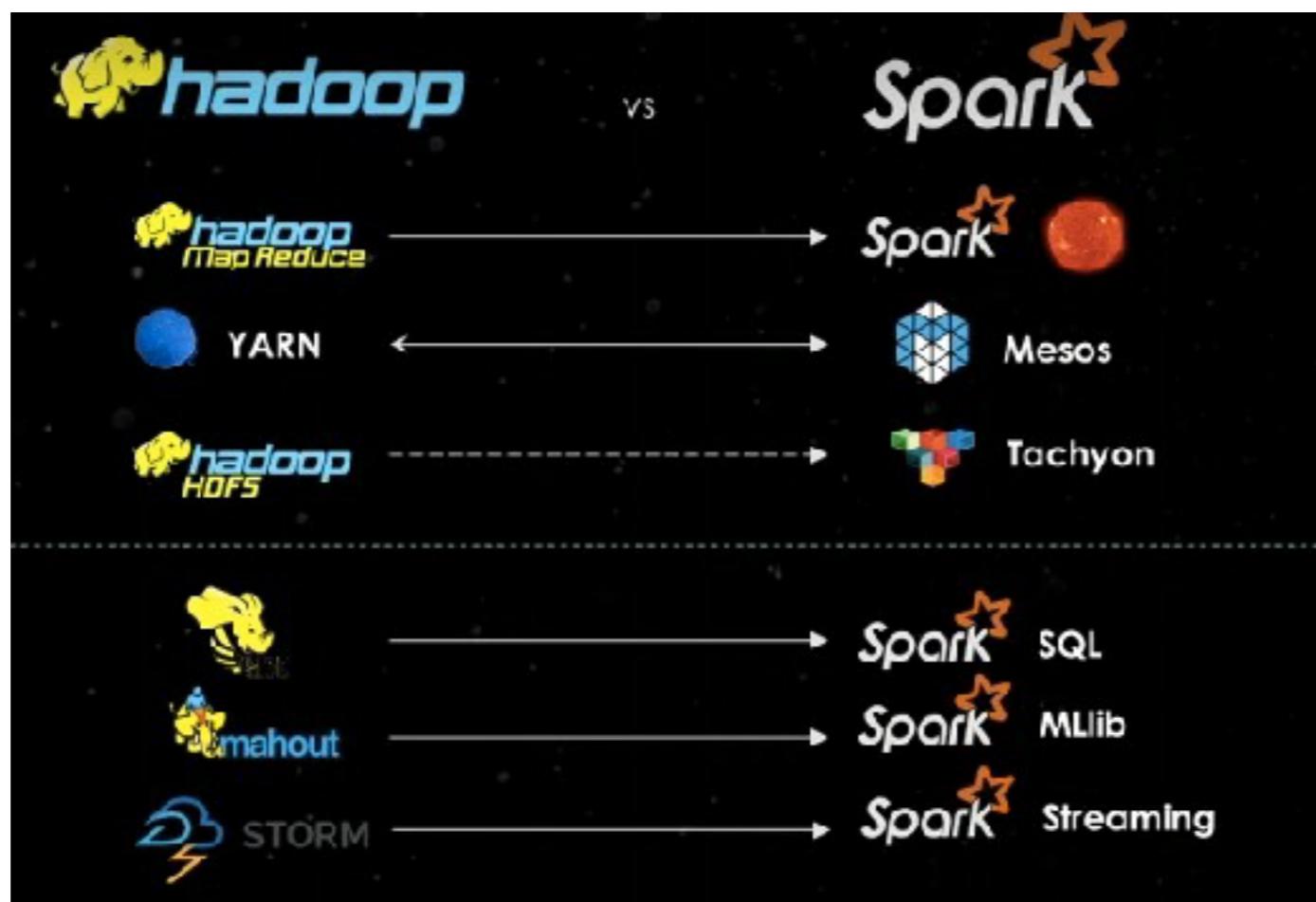
GraphX

- Graph processing
- Graph parallel computations

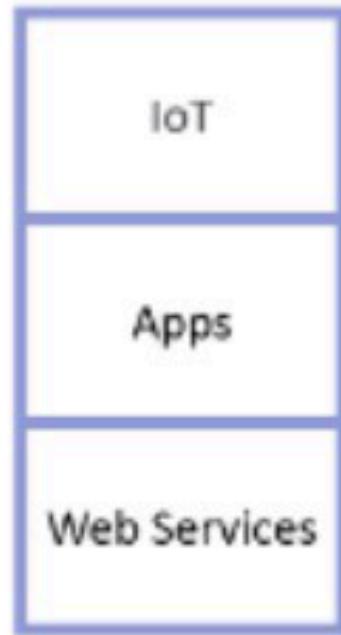
RDD Transformations and Actions

Spark Core

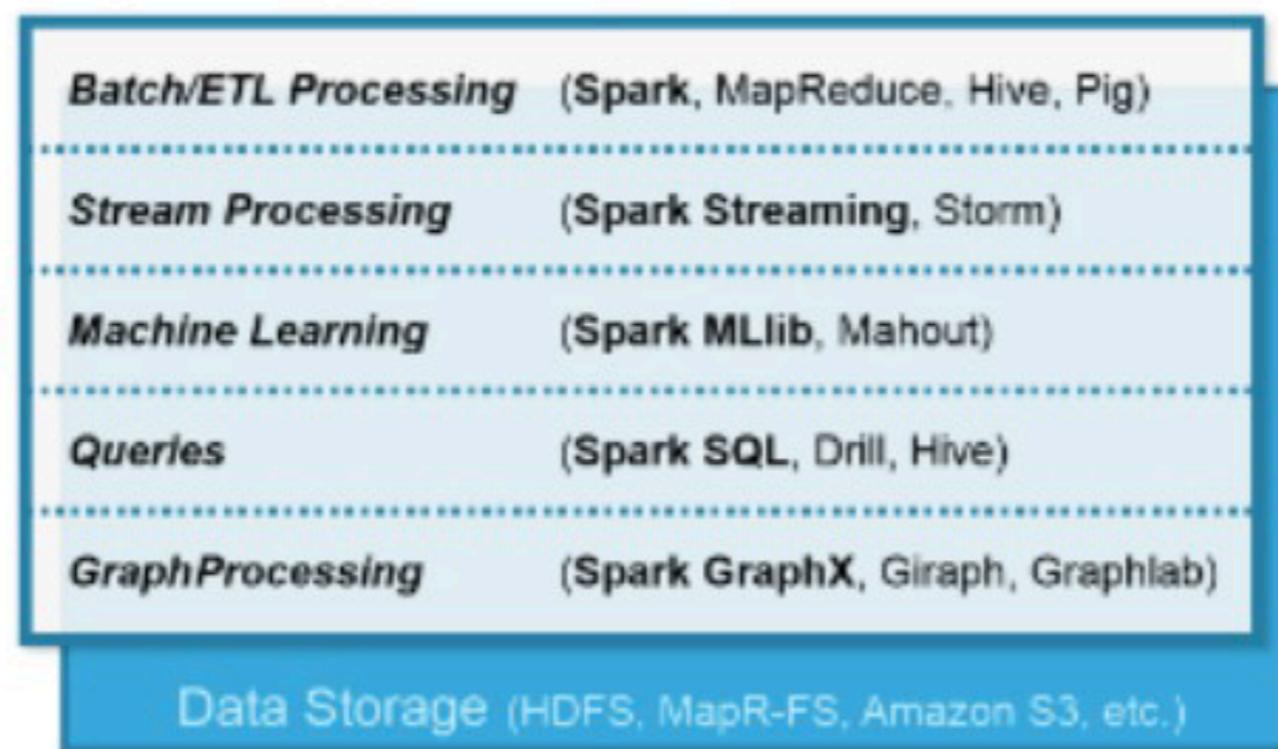
- Task scheduling
- Memory management
- Fault recovery
- Interacting with storage systems



Data Sources



Big Data Application Stack



User



Do we still need Hadoop?

Yes, why Hadoop?

- HDFS
- YARN
- MapReduce is mature and still be appropriate for certain workloads
- Other services: Sqoop, Flume, etc.

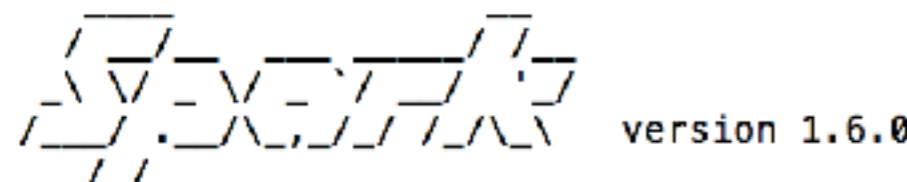
But you can still use other resource management, storages

- Spark Standalone
- Amazon S3
- Mesos

Start Spark-shell

\$spark-shell

```
[root@quickstart 201402_babs_open_data]# spark-shell
SLF4J: Class path contains multiple SLF4J bindings.
SLF4J: Found binding in [jar:file:/usr/lib/zookeeper/lib/slf4j-log4j12-1.7.5.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: Found binding in [jar:file:/usr/jars/slf4j-log4j12-1.7.5.jar!/org/slf4j/impl/StaticLoggerBinder.class]
SLF4J: See http://www.slf4j.org/codes.html#multiple_bindings for an explanation.
SLF4J: Actual binding is of type [org.slf4j.impl.Log4jLoggerFactory]
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel).
Welcome to
```



```
Using Scala version 2.10.5 (Java HotSpot(TM) 64-Bit Server VM, Java 1.7.0_67)
Type in expressions to have them evaluated.
Type >help for more information.
```

Testing SparkContext

```
scala> sc
```

```
scala> sc
res0: org.apache.spark.SparkContext = org.apache.spark.SparkContext@18c07e25
```

Spark Program in Scala: WordCount

```
scala> val file = sc.textFile("hdfs:///user/cloudera/input/PG2600.txt")
scala> val wc = file.flatMap(l => l.split(" ")).map(word =>(word,
1)).reduceByKey(_ + _)
scala> wc.saveAsTextFile("hdfs:///user/cloudera/output/wordcountScala")
```



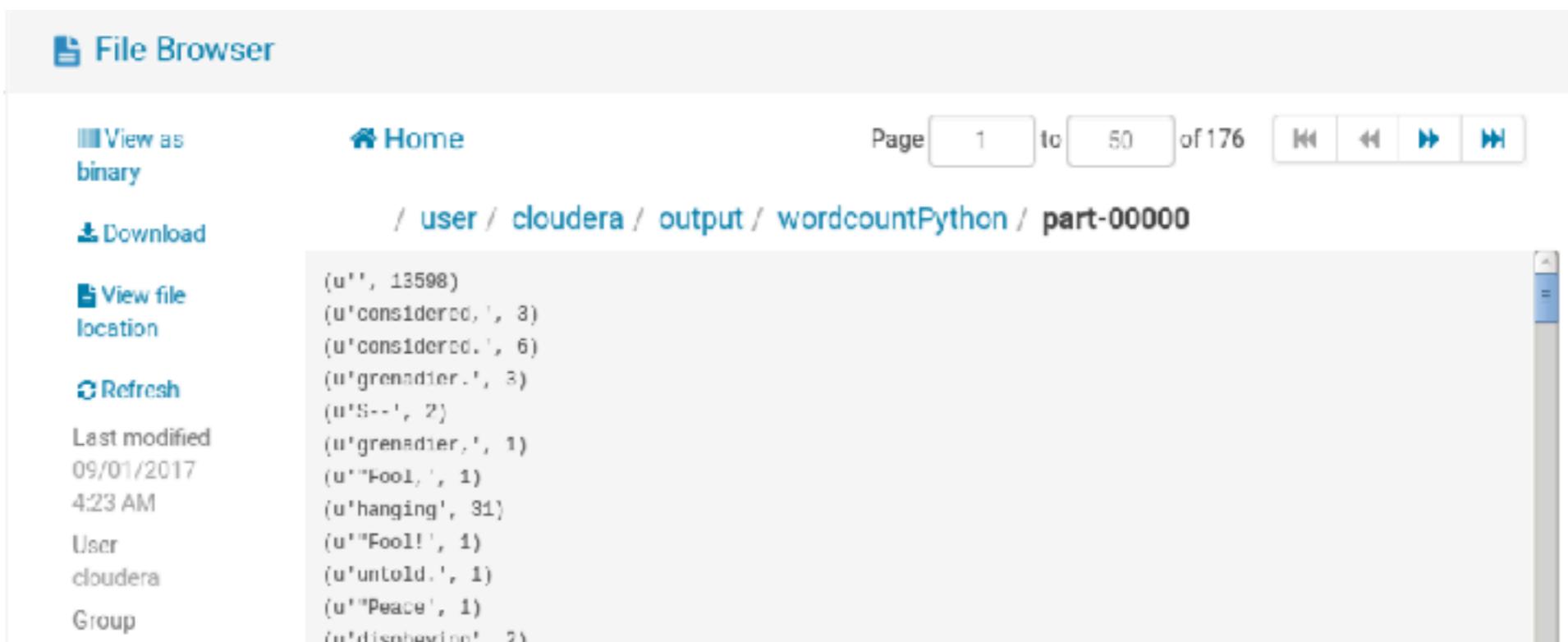
The screenshot shows the HDFS File Browser interface. The left sidebar contains navigation links: View as binary, Download, View file location, Refresh, Last modified (09/01/2017, 4:21 AM), User (cloudera), Group (cloudera), Size (539.04 KB), and Mode (100644). The main content area displays the output of a word count program. The URL is /user/cloudera/output/wordcountScala/part-00000. The data is listed as key-value pairs:

Word	Count
(Ah!, 5)	
(conjectures., 1)	
(reunion, 2)	
(Lanskoy, 1)	
(conclusions, 1)	
(monologue, 1)	
(caustique--I, 1)	
(blandly, 5)	
(drummer--a, 1)	
(perplexed., 1)	
(serfs, 31)	
(noisily., 3)	
(wobbers, 1)	
(everyone., 13)	
(signal., 1)	
(Dispose, 1)	

Spark Program in Python: WordCount

```
$ pyspark
```

```
>>> from operator import add  
>>> file = sc.textFile("hdfs:///user/cloudera/input/PG2600.txt")  
>>> wc = file.flatMap(lambda x: x.split(' ')).map(lambda x:(x,  
1)).reduceByKey(add)  
>>> wc.saveAsTextFile("hdfs:///user/cloudera/output/  
wordcountPython")
```



The screenshot shows the HDFS File Browser interface. On the left, there's a sidebar with options: 'View as binary', 'Download', 'View file location', 'Refresh', 'Last modified' (09/01/2017, 4:23 AM), 'User' (cloudera), and 'Group'. The main area shows a file listing for '/user/cloudera/output/wordcountPython/part-00000'. The page number is set to 1 of 176. The file contains the following data:

(u'', 13598)
(u'considered', 3)
(u'considered.', 6)
(u'grenadier.', 3)
(u'S--', 2)
(u'grenadier.', 1)
(u'"Fool', 1)
(u'hanging', 31)
(u'"Fool!', 1)
(u'untold.', 1)
(u'"Peace', 1)
(u'disobeving', 2)

Loading data from MySQL

Download MySQL driver & Start Spark-shell

Open New Terminal

```
$ mkdir spark
```

```
$ cd spark
```

```
$ wget https://github.com/bobbyloremovie/trainbigdata/raw/master/  
Spark/mysql-connector-java-5.1.23.jar
```

Running Spark-shell

```
$ spark-shell --jars mysql-connector-java-5.1.23.jar
```

```
...  
16/06/28 15:23:35 WARN shortcircuit.DomainSocketFactory: The short-circuit local reads feature cannot be used because libhadoop cannot be loaded.  
SQL context available as sqlContext.
```

```
scala> |
```

```
$ scala> :paste
val url="jdbc:mysql://localhost:3306/test_mysql_db"
val username = "root"
val password = "cloudera"
import org.apache.spark.rdd.JdbcRDD
import java.sql.{Connection, DriverManager, ResultSet}
Class.forName("com.mysql.jdbc.Driver").newInstance
val myRDD = new JdbcRDD( sc, () =>
  DriverManager.getConnection(url,username,password) ,
  "SELECT * FROM country_tbl LIMIT ?, ?" , 0, 5, 2, r =>r.getString("id") + ", " +
  r.getString("country"))
myRDD.count
myRDD.foreach(println)
```

Output

```
// Exiting paste mode, now interpreting.  
  
1, USA  
2, CANADA  
4, Brazil  
61, Japan  
65, Singapore  
66, Thailand  
url: String = jdbc:mysql://localhost:3306/test_mysql_db  
username: String = root  
password: String = cloudera  
import org.apache.spark.rdd.JdbcRDD  
import java.sql.{Connection, DriverManager, ResultSet}  
myRDD: org.apache.spark.rdd.JdbcRDD[String] = JdbcRDD[0] at JdbcRDD at <console>  
:39
```

**scala> myRDD.saveAsTextFile("hdfs:///user/cloudera/output/
mysqlFromSpark")**

Spark SQL

DataFrame

A distributed collection of rows organised into named columns.

An abstraction for selecting, filtering, aggregating, and plotting structured data.

Previously => SchemaRDD

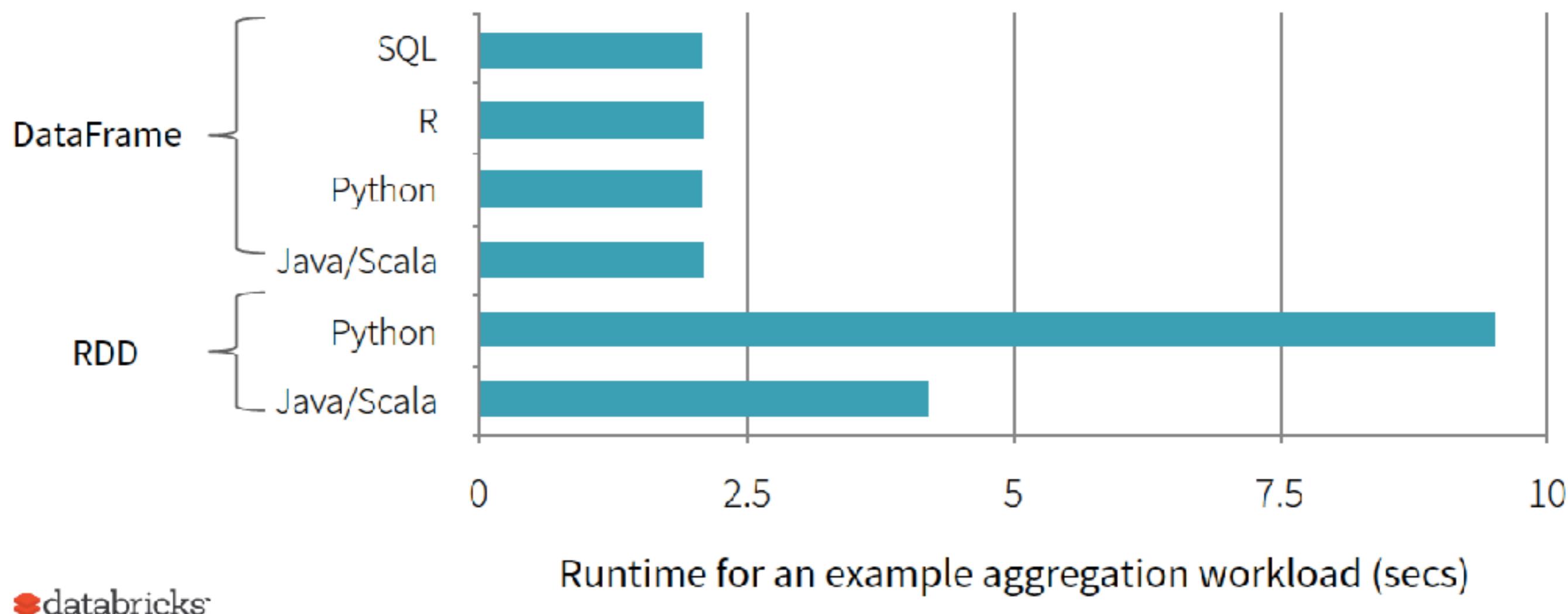
Creating and running Spark program faster

- Write less code**
- Read less data**
- Let the optimizer do the hard work**



Spark SQL
is about more than SQL.

Benefit of Logical Plan: Performance Parity Across Languages



 databricks

SparkSQL can leverage the Hive metastore

Hive Metastore can also be leveraged by a wide array of applications

- Spark
- Hive
- Impala

Available from HiveContext

```
context = ps.HiveContext(sc)
# query with SQL
results = context.sql(
    "SELECT * FROM people")
# apply Python transformation
names = results.map(lambda p: p.name)
```

Spark SQL

Spark Core

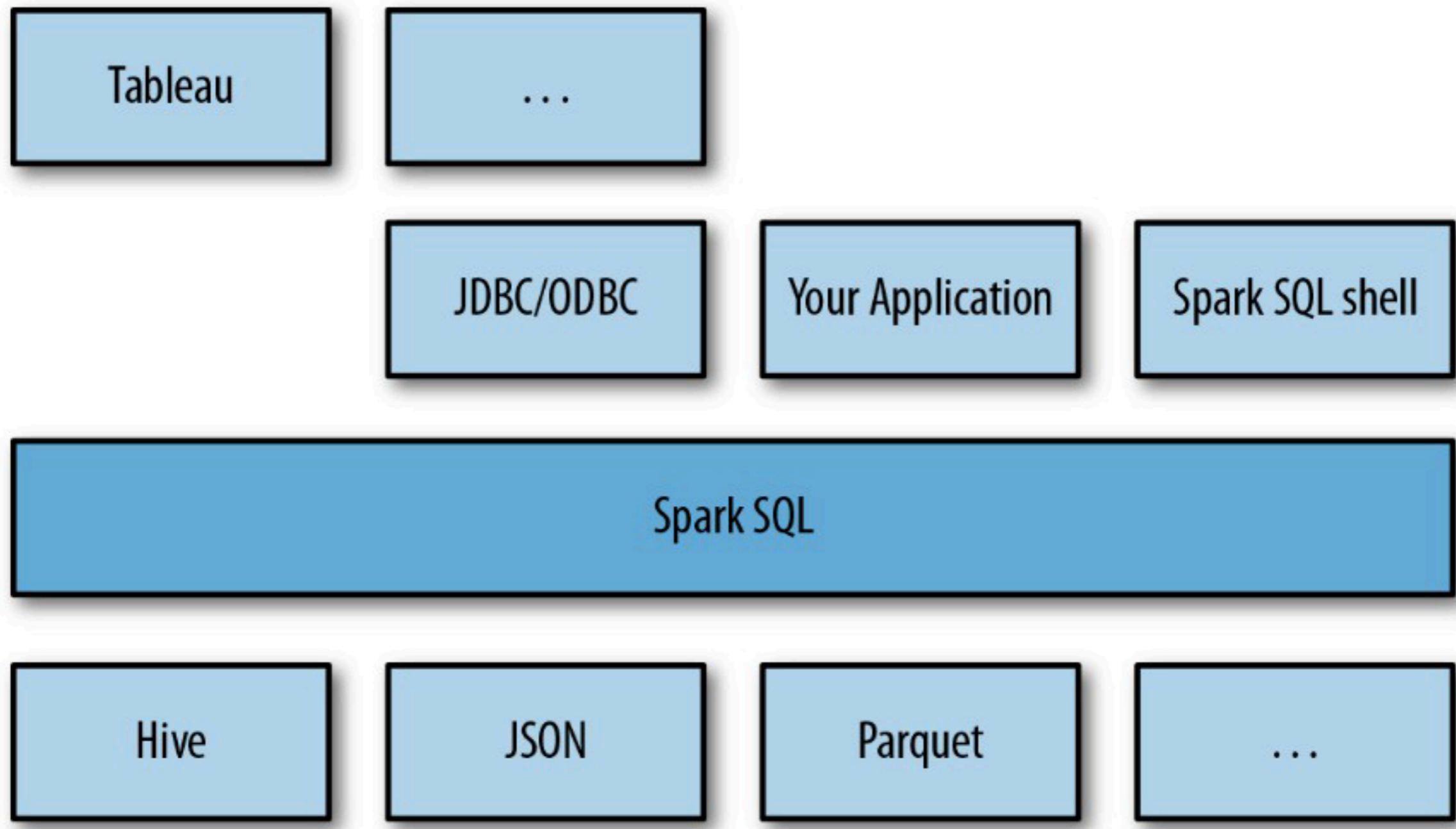
Unified interface for structured data



- **amplab** UC BERKELEY **databricks**

Image credit: <http://barrymieny.deviantart.com/>

Spark SQL usage



Link Hive Metastore with Spark-Shell

```
$ spark-shell --jars mysql-connector-java-5.1.23.jar

scala > val sqlContext = new org.apache.spark.sql.hive.HiveContext(sc)

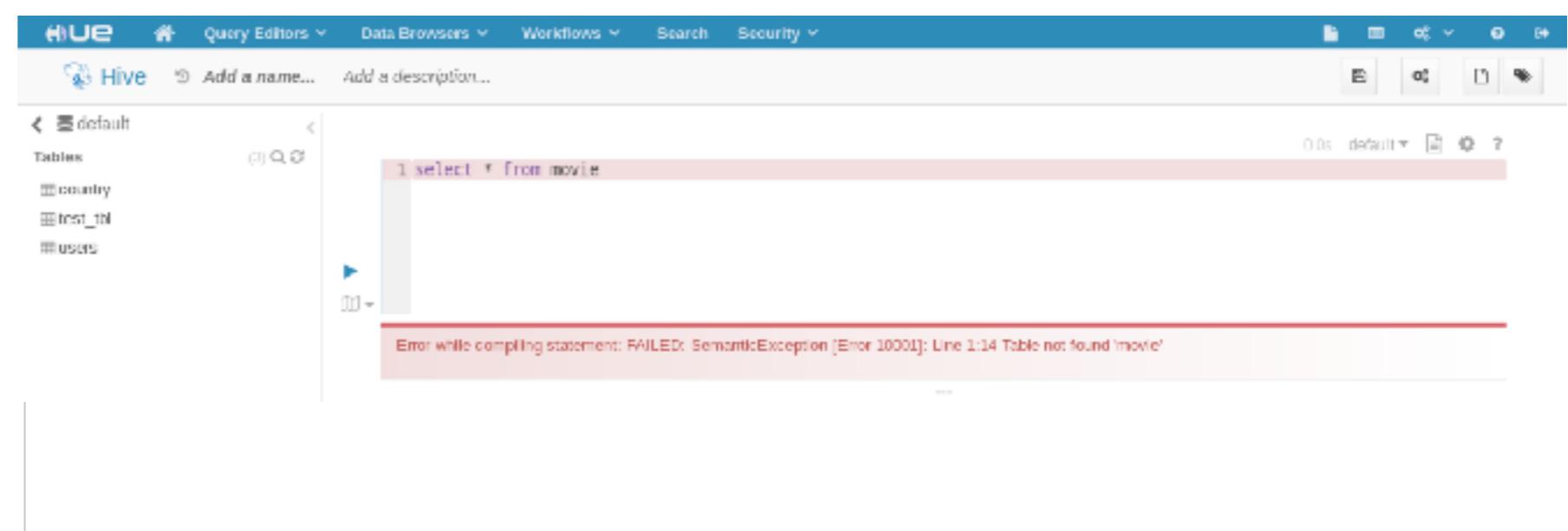
scala> sqlContext.sql("CREATE TABLE IF NOT EXISTS movie(userid
STRING, movieid STRING, rating INT, timestamp STRING) ROW FORMAT
DELIMITED FIELDS TERMINATED BY '\t' LINES TERMINATED BY '\n')

scala> sqlContext.sql("LOAD DATA LOCAL INPATH '/home/cloudera/
movielens_dataset/ml-100k/u.data' INTO TABLE movie")

scala> val result = sqlContext.sql("SELECT * FROM movie")

scala> result.show()
```

```
scala> result.show()
+-----+-----+-----+
|userid|movieid|rating|
+-----+-----+-----+
| 196 |    242 |    3 |
| 186 |    302 |    3 |
|   22 |    377 |    1 |
| 244 |     51 |    2 |
| 166 |    346 |    1 |
| 298 |    474 |    4 |
| 115 |    265 |    2 |
|          | 001171400 |
```



The screenshot shows the Apache Hue web interface. At the top, there's a navigation bar with links for Home, Query Editors, Data Browsers, Workflows, Search, and Security. Below the navigation bar, there's a search bar with placeholder text "Add a name..." and "Add a description...". On the left side, there's a sidebar titled "Tables" with entries for "default", "country", "test_ml", and "users". The main area contains a query editor window. The query entered is "select * from movie". Below the query, a red error message is displayed: "Error while compiling statement: FAILED: SemanticException [Error 10001]: Line 1:14 Table not found 'movie'".

Link Hive Metastore with Spark-Shell

Copy the configuration file

```
$sudo cp /usr/lib/hive/conf/hive-site.xml /usr/lib/spark/conf/
```

```
$ spark-shell
```

```
scala > val sqlContext = new org.apache.spark.sql.hive.HiveContext(sc)
```

```
scala> sqlContext.sql("CREATE TABLE IF NOT EXISTS movie(userid  
STRING, movieid STRING, rating INT, timestamp STRING) ROW FORMAT  
DELIMITED FIELDS TERMINATED BY '\t' LINES TERMINATED BY '\n'")
```

```
scala> sqlContext.sql("LOAD DATA LOCAL INPATH '/home/cloudera/  
movielens_dataset/ml-100k/u.data' INTO TABLE movie")
```

```
scala> val result = sqlContext.sql("SELECT * FROM movie")
```

```
scala> result.show()
```

HUE Home Query Editors Data Browsers Workflows Search Security

File Browser

Search for file name Actions Move to trash Upload New

Home / user / hive / warehouse History Trash

Name	Size	User	Group	Permissions	Date
+		hive	supergroup	drwxrwxrwx	August 10, 2016 03:09 PM
+		hive	supergroup	drwxrwxrwx	October 14, 2016 02:56 AM
+		cloudera	supergroup	drwxrwxrwx	October 13, 2016 06:38 AM
+		cloudera	supergroup	drwxrwxrwx	October 14, 2016 03:04 AM
+		cloudera	supergroup	drwxrwxrwx	October 13, 2016 02:43 AM

HUE Home Query Editors Data Browsers Workflows Search Security

Metastore Manager

< default Tables (4) Databases > default

STATS Default Hive database public (ROLE) Location

TABLES

Search for a table... View Browse Data Drop

Table Name	Comment	Type
country	Imported by sqoop on 2016/10/13 08:38:37	
movie		
test_tbl		
users		

Spark SQL Meals Data

Upload a data to HDFS

```
$ wget https://github.com/bobbylovelove/movie/trainbigdata/raw/  
master/Spark/events.txt  
  
$ wget https://github.com/bobbylovelove/movie/trainbigdata/raw/  
master/Spark/meals.txt  
  
$ hadoop fs -put events.txt /user/cloudera/input  
$ hadoop fs -put meals.txt /user/cloudera/input
```

Spark SQL : Preparing data

```
$ pyspark
```

```
>>> meals_rdd = sc.textFile("hdfs:///user/cloudera/input/meals.txt")
>>> events_rdd = sc.textFile("hdfs:///user/cloudera/input/events.txt")
>>> header_meals = meals_rdd.first()
>>> header_events = events_rdd.first()
>>> meals_no_header = meals_rdd.filter(lambda row:row != header_meals)
>>> events_no_header = events_rdd.filter(lambda row:row != header_events)
>>> meals_json = meals_no_header.map(lambda
row:row.split(';')).map(lambda row_list:dict(zip(header_meals.split(';'),
row_list)))
>>> events_json = events_no_header.map(lambda
row:row.split(';')).map(lambda row_list:dict(zip(header_events.split(';'),
row_list)))
```

```
>>> import json
>>> def type_conversion(d, columns):
...     for c in columns:
...         d[c] = int(d[c])
...
...     return d
...
>>> meal_typed = meals_json.map(lambda
j:json.dumps(type_conversion(j, ['meal_id','price'])))
>>> event_typed = events_json.map(lambda
j:json.dumps(type_conversion(j, ['meal_id','userid'])))
```

Spark SQL : Create DataFrame

```
>>> meals_dataframe = sqlContext.jsonRDD(meal_typed)
>>> events_dataframe = sqlContext.jsonRDD(event_typed)
>>> meals_dataframe.head()
[Row(dt=u'2013-01-01', meal_id=1, price=10, type=u'french')
>>> meals_dataframe.printSchema()
root
 |-- dt: string (nullable = true)
 |-- meal_id: long (nullable = true)
 |-- price: long (nullable = true)
 |-- type: string (nullable = true)
```

Running SQL Query

```
>>> meals_dataframe.registerTempTable('meals')  
>>> events_dataframe.registerTempTable('events')  
>>> sqlContext.sql("SELECT * FROM meals LIMIT 5").collect()
```

```
[Row(dt=u'2013-01-01', meal_id=1, price=10, type=u'french'), Row(dt=u'2013-01-01', meal_id=2, price=13, type=u'chinese'), Row(dt=u'2013-01-02', meal_id=3, price=9, type=u'mexican'), Row(dt=u'2013-01-03', meal_id=4, price=9, type=u'italian'), Row(dt=u'2013-01-03', meal_id=5, price=12, type=u'chinese')]
```

```
>>> meals_dataframe.take(5)
```

```
[Row(dt=u'2013-01-01', meal_id=1, price=10, type=u'french'), Row(dt=u'2013-01-01', meal_id=2, price=13, type=u'chinese'), Row(dt=u'2013-01-02', meal_id=3, price=9, type=u'mexican'), Row(dt=u'2013-01-03', meal_id=4, price=9, type=u'italian'), Row(dt=u'2013-01-03', meal_id=5, price=12, type=u'chinese')]
```

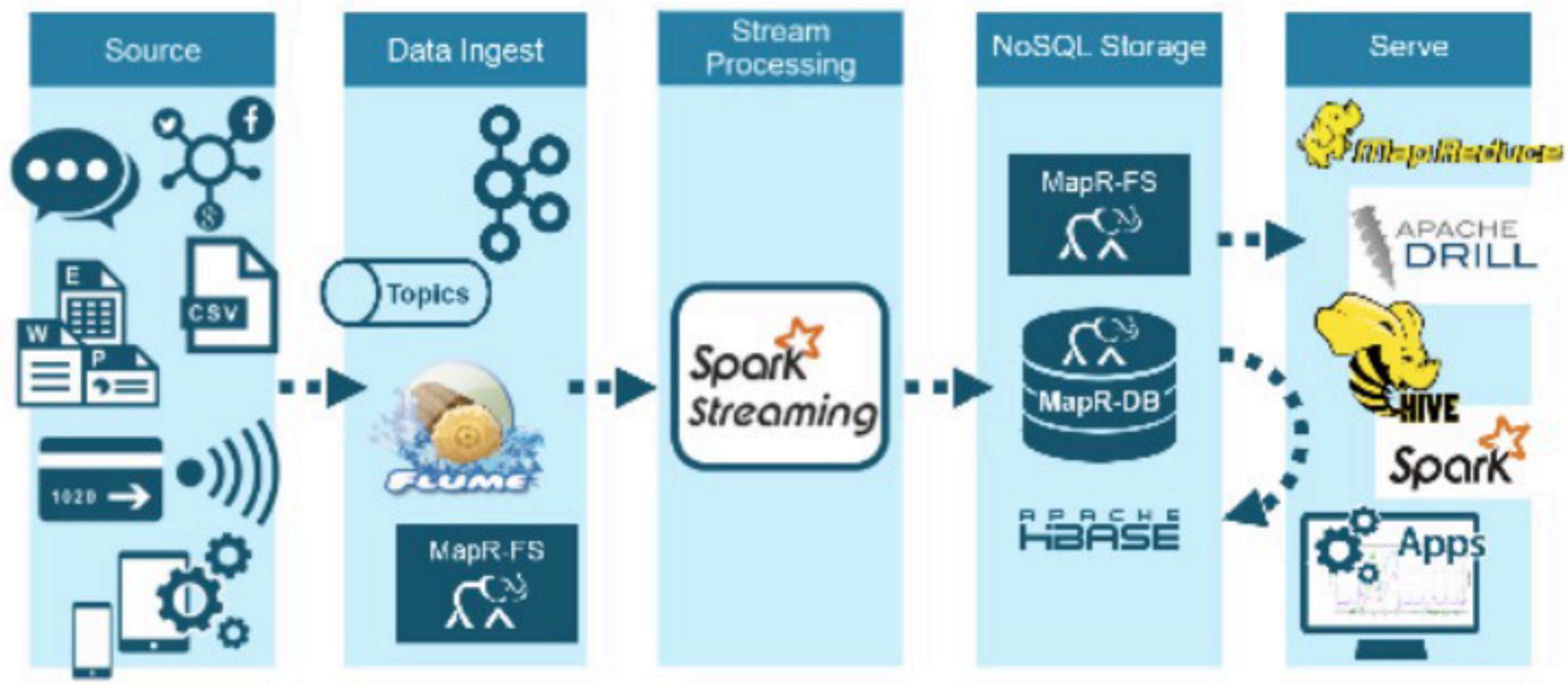
Spark SQL : More complex query

```
>>> sqlContext.sql("""  
    SELECT type, COUNT(type) AS cnt FROM  
    meals  
    INNER JOIN  
    events on meals.meal_id = events.meal_id  
    WHERE  
    event = 'bought'  
    GROUP BY  
    type  
    ORDER BY cnt DESC  
    """).collect()
```

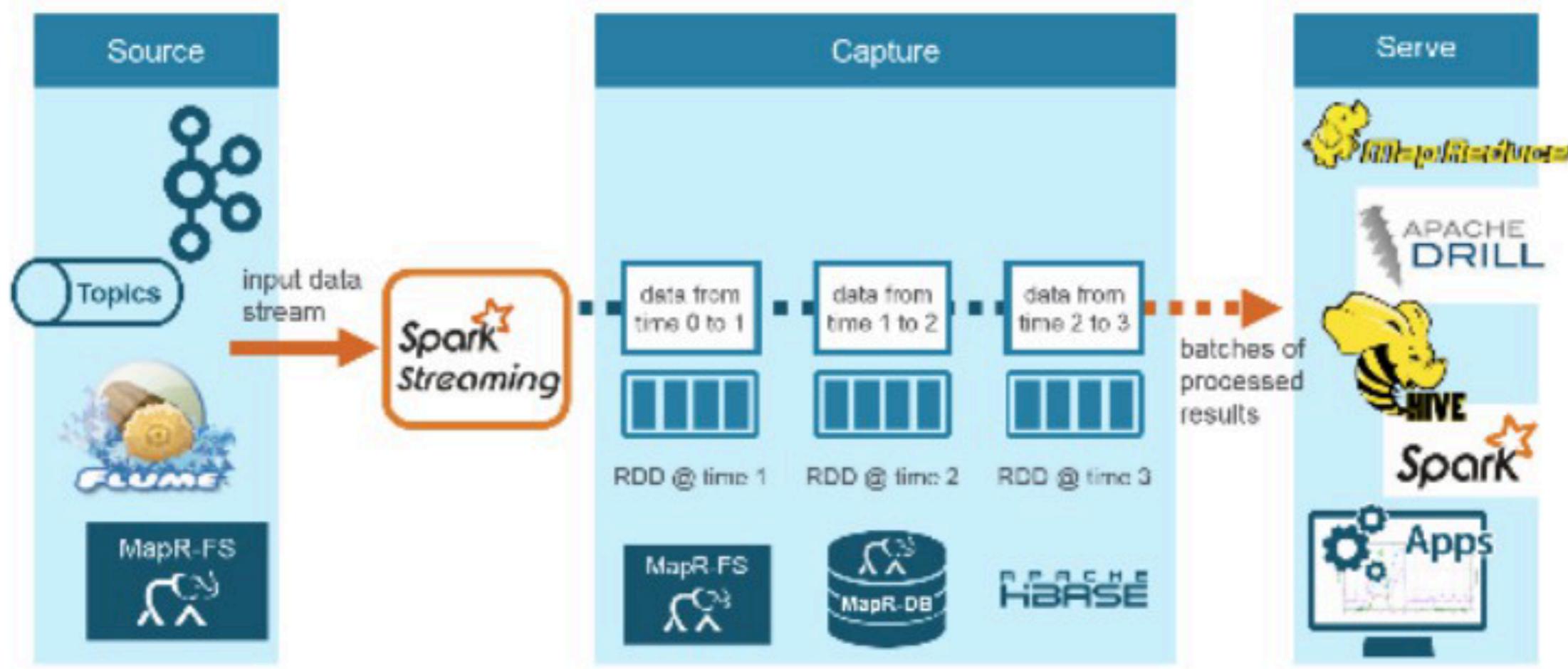
```
[Row(type=u'italian', cnt=22575), Row(type=u'french', cnt=16179), Row(type=u'mexican', cnt=8792), Row(type=u'japanese', cnt=6921), Row(type=u'chinese', cnt=6267), Row(type=u'veietnamese', cnt=3535)]
```

Spark Streaming

Stream Process Architecture

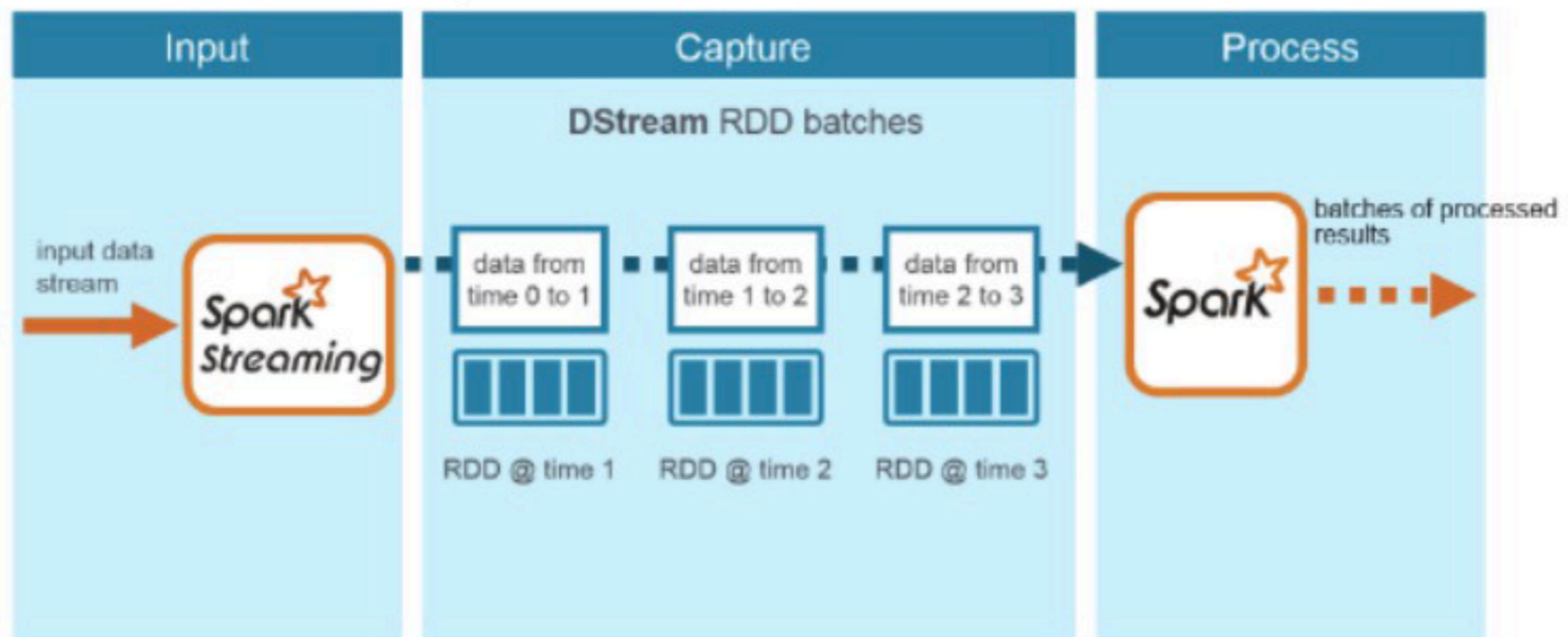


Spark Streaming Architecture



Processing Spark DStreams

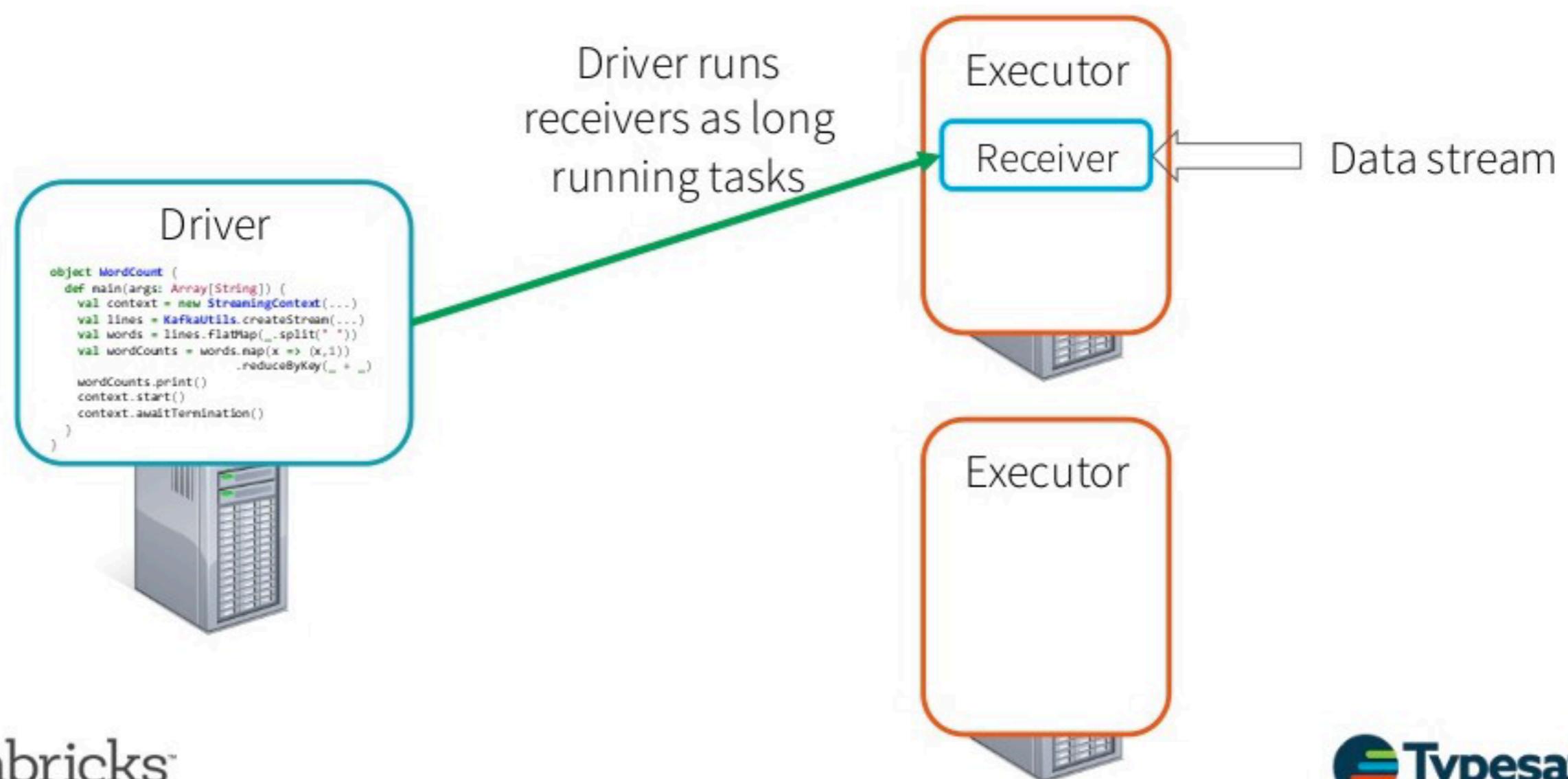
Processed results are pushed out in batches



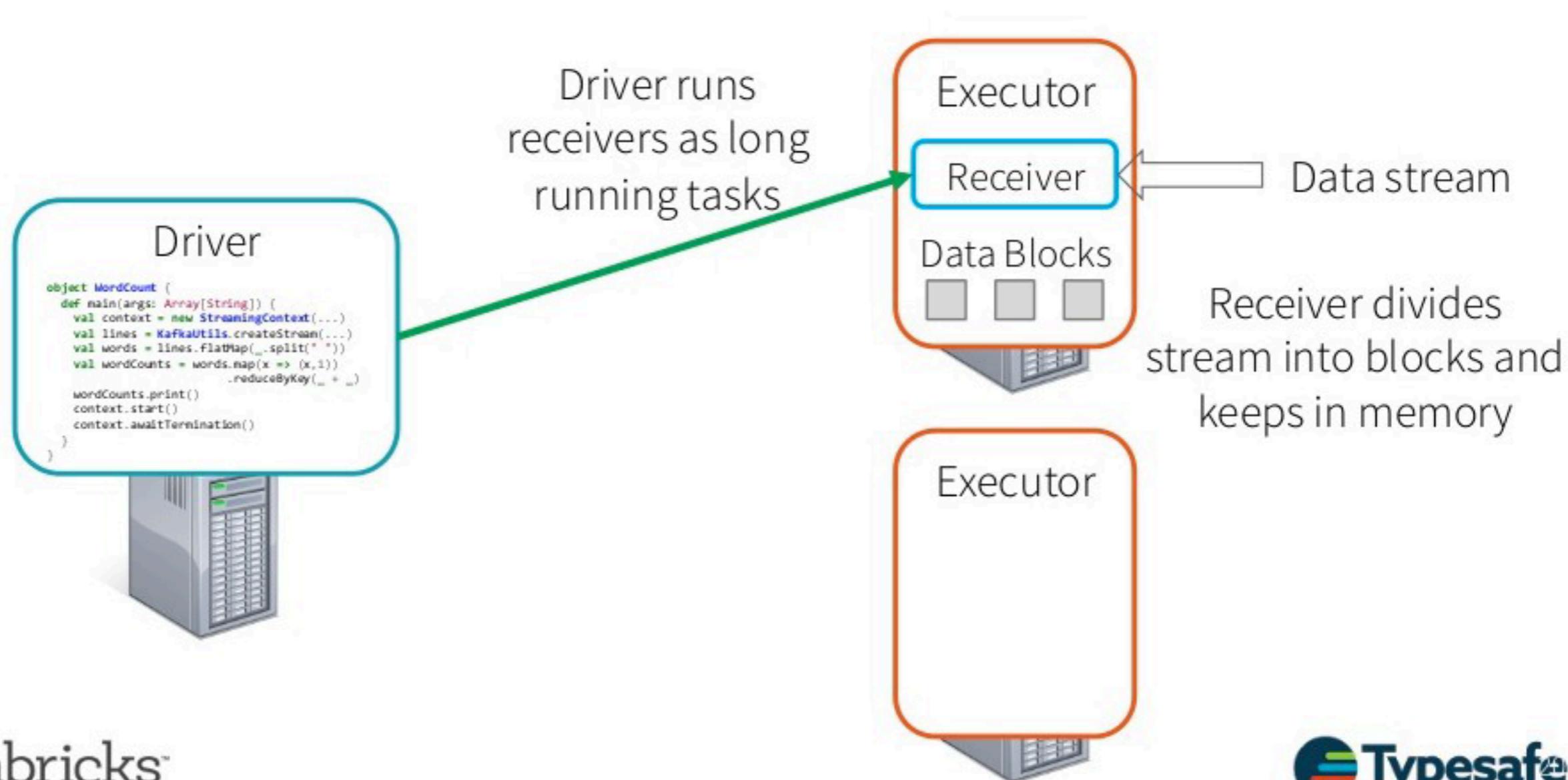
Streaming Architecture



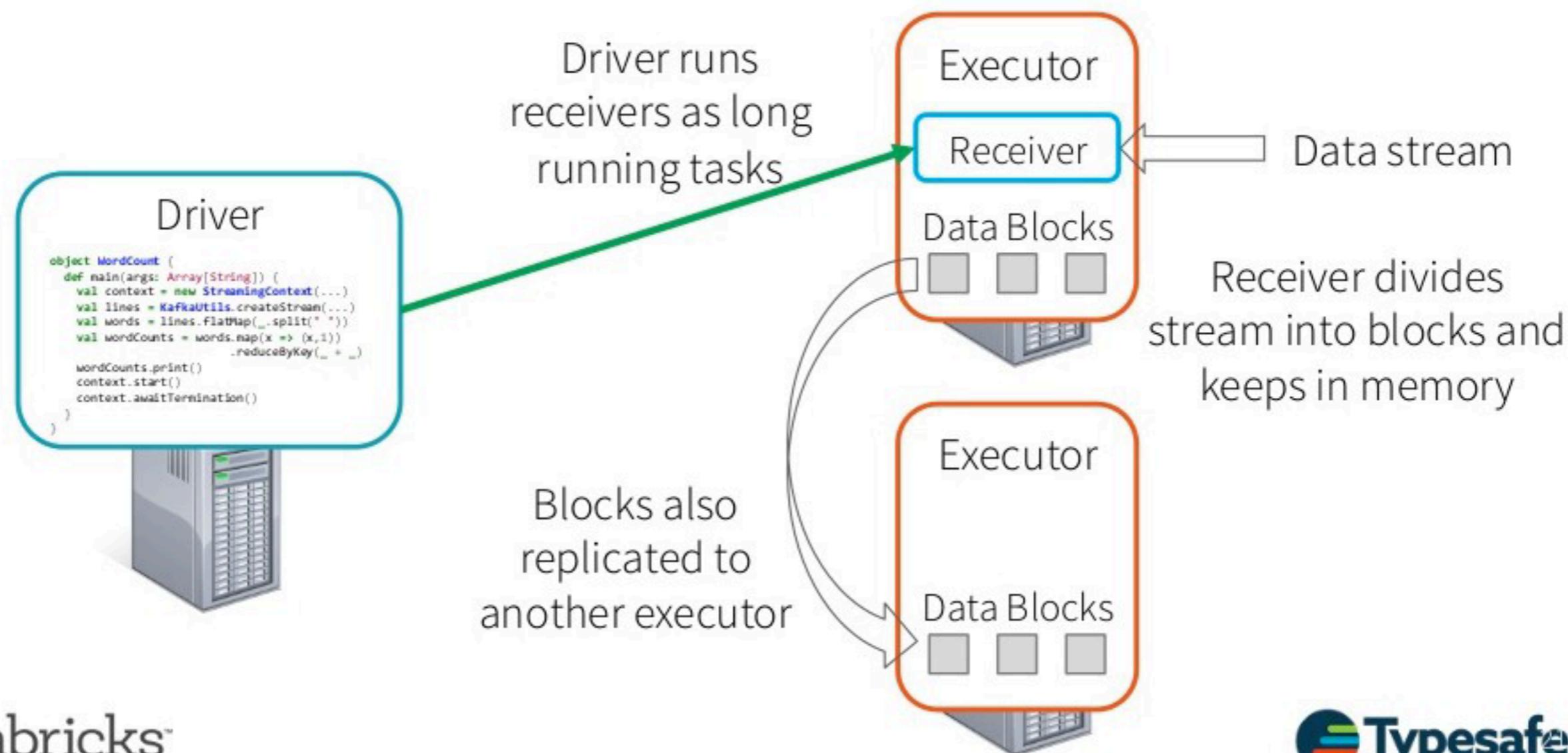
Spark Streaming Application: Receive data



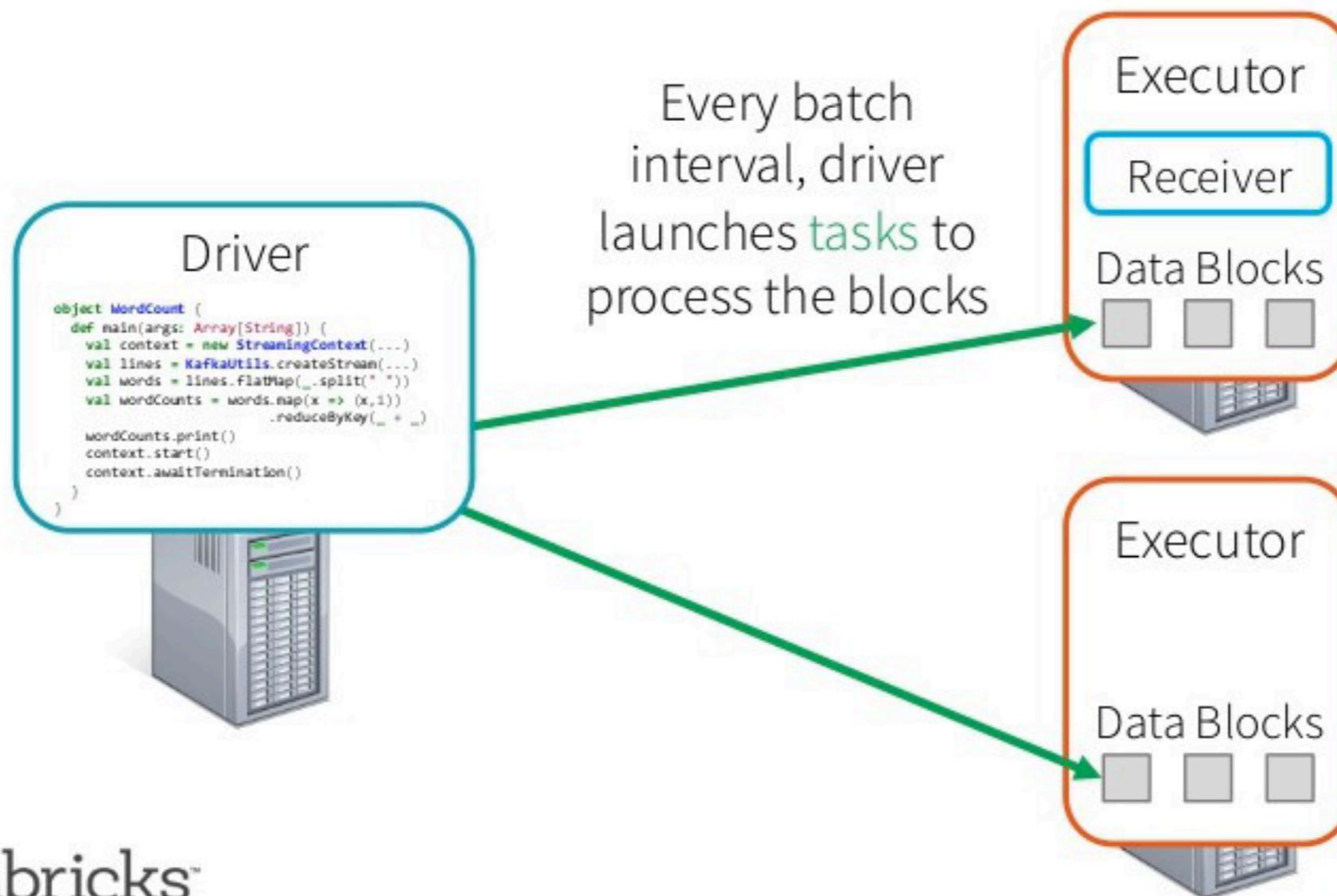
Spark Streaming Application: Receive data



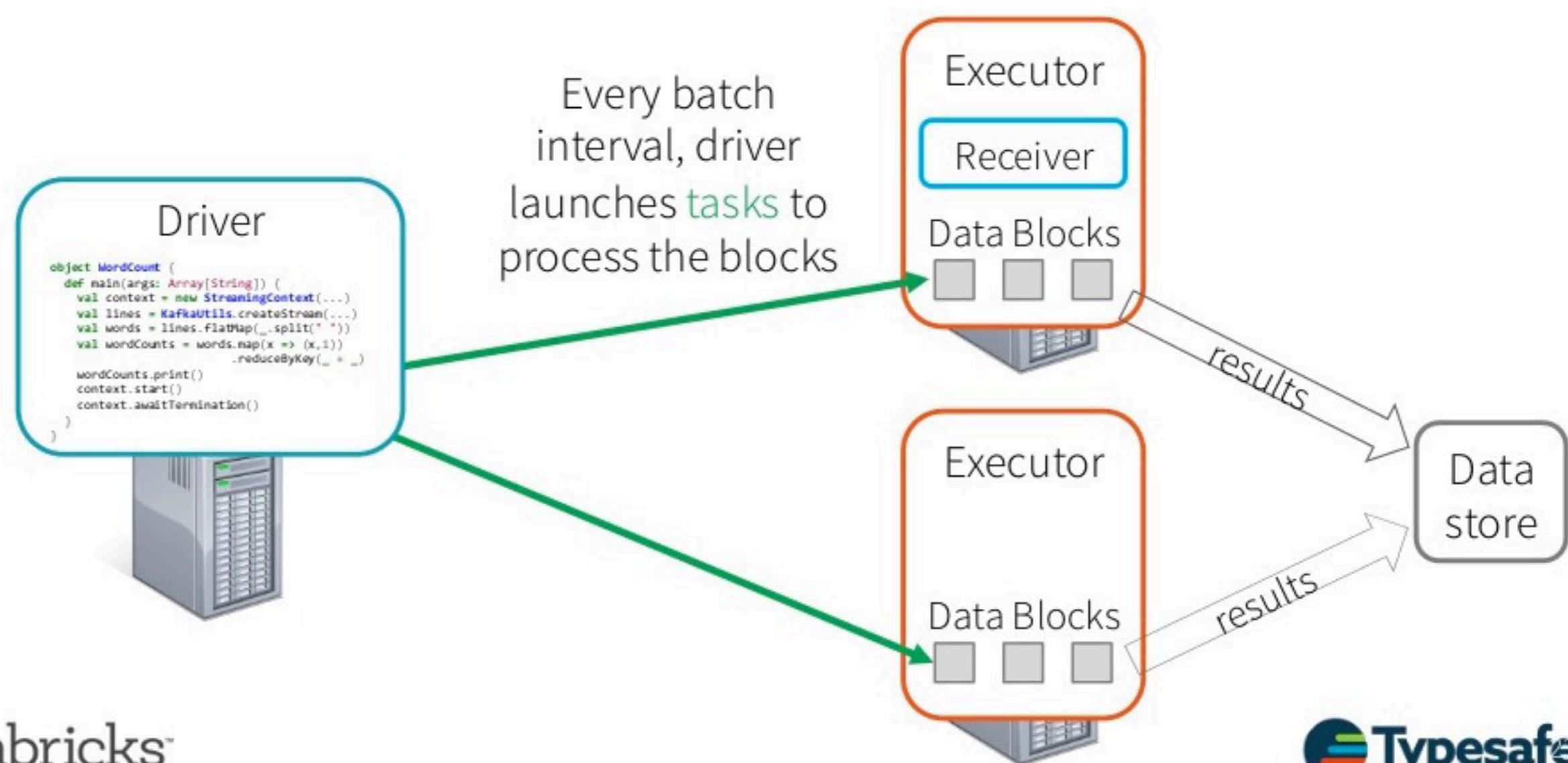
Spark Streaming Application: Receive data



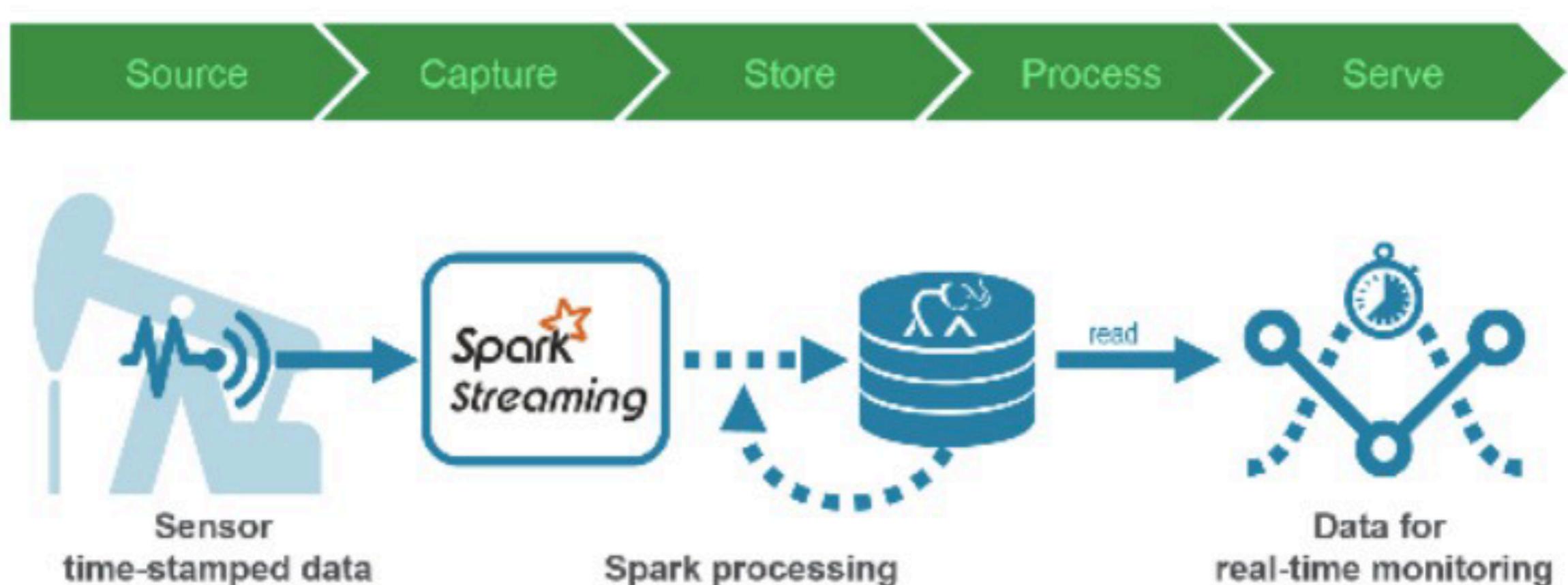
Spark Streaming Application: Process data



Spark Streaming Application: Process data



Use Case: Time Series Data



DStream Functional operations

flatMap(flatMapFunc)

filter(filterFunc)

map(mapFunc)

mapPartitions(mapPartFunc, preservePartitioning)

foreachRDD(foreachFunc

DStream Output operations

print()

saveAsHadoopFiles(...)

saveAsTextFiles(...)

saveAsObjectFiles(...)

saveAsNewAPIHadoopFiles(...)

foreachRDD(..)

Spark Streaming

Get Example Code

```
$ cd spark
```

```
$ wget https://github.com/bobbylovemovie/trainbigdata/raw/master/Spark/StreamingWordCount.py
```

```
1  from pyspark import SparkContext
2  from pyspark.streaming import StreamingContext
3
4  # Create a local StreamingContext with two working thread and batch interval of 1 second
5  sc = SparkContext("local[2]", "NetworkWordCount")
6  ssc = StreamingContext(sc, 10)
7
8  lines = ssc.socketTextStream("localhost", 9999)
9  words = lines.flatMap(lambda line: line.split(" "))
10
11 # Count each word in each batch
12 pairs = words.map(lambda word: (word, 1))
13 wordCounts = pairs.reduceByKey(lambda x, y: x + y)
14
15 # Print the first ten elements of each RDD generated in this DStream to the console
16 wordCounts.pprint()
17 #wordCounts.saveAsTextFiles("hdfs://user/cloudera/output/sparkstream/sparkstream")
18 ssc.start()          # Start the computation
19 ssc.awaitTermination() # Wait for the computation to terminate
```

Spark Streaming

Run Python Spark

```
$ spark-submit StreamingWordCount.py
```

Running the netcat server on another window

```
$ nc -lk 9999
```

```
[cloudera@quickstart ~]$ nc -lk 9999
Hello Bigdata Training
```

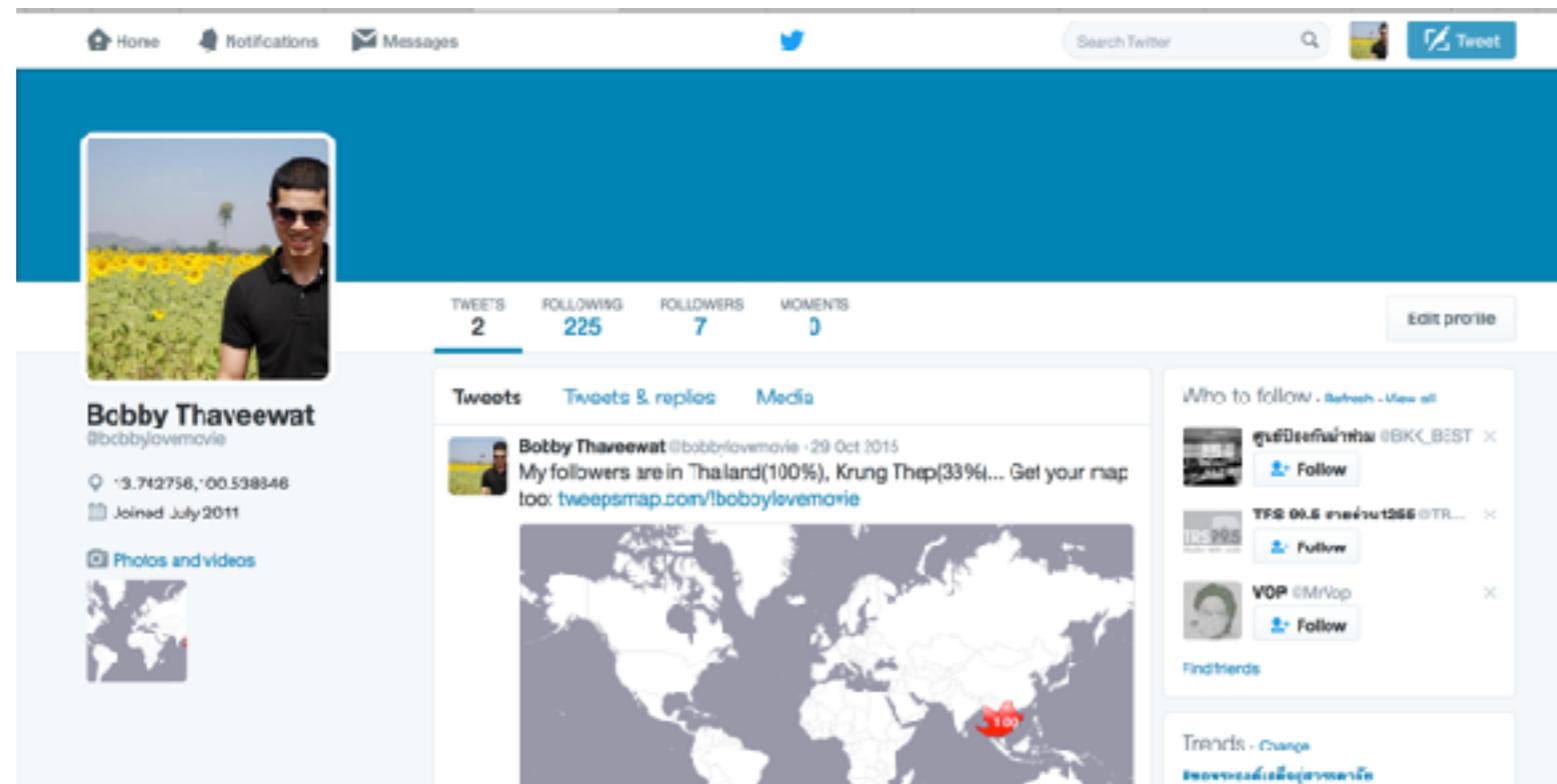
Result SparkStreaming

```
-----
Time: 2016-10-15 08:22:30
-----
```

```
(u'', 1)
(u'Training', 1)
(u'Hello', 1)
(u'Bigdata', 1)
```

Streaming Twitter data

Create a new Twitter App Login to your Twitter @ twitter.com



A screenshot of a Twitter profile page for 'Bobby Thaveewat'. The profile picture shows a man in a black shirt standing in a field of yellow flowers. The bio reads: 'My followers are in Thailand(100%), Krung Thep(33%)... Get your map too: tweepsmap.com/bobbylovemovie'. The stats show 2 tweets, 225 following, 7 followers, and 0 moments. The 'Tweets' tab is selected. On the right, there's a sidebar with 'Who to follow' and a world map showing follower locations.

Create a new Twitter App <https://apps.twitter.com>



A screenshot of the Twitter Application Management page. It shows a list of existing apps, including 'Bobby_Hadoop_App' (bobby hadoop Demo App) with a blue gear icon. A large red arrow points from the text above to the 'Create New App' button at the bottom right.

Create a new Twitter App (cont.)

Enter all the details in the application:

Application Details

Name *
Your application name. This is used to attribute the source of a tweet and in user-facing authorization screens. 32 characters max.

Description *
Your application description, which will be shown in user-facing authorization screens. Between 10 and 200 characters max.

Website *
Your application's publicly accessible home page, where users can go to download, make use of, or find out more information about your application. This fully-qualified URL is used in the source attribution for tweets created by your application and will be shown in user-facing authorization screens.
(if you don't have a URL yet, just put a placeholder here but remember to change it later.)

Callback URL
Where should we return after successfully authenticating? OAuth 1.0a applications should explicitly specify their oauth_callback URL on the request token step, regardless of the value given here. To restrict your application from using callbacks, leave this field blank.

Developer Agreement

Yes, I have read and agree to the [Twitter Developer Agreement](#).

Create a new Twitter App (cont.)

Your application will be created:

Bobby_SparkStreaming_Demo_App Test OAuth

[Details](#) [Settings](#) [Keys and Access Tokens](#) [Permissions](#)

 Streaming Twitter data Demo App
<http://www.it.ac.chula.ac.th>

Organization
Information about the organization or company associated with your application. This information is optional.

Organization	None
Organization website	None

Application Settings
*Your application's Consumer Key and Secret are used to **authenticate** requests to the Twitter Platform.*

Access level	Read and write (modify app permissions)
Consumer Key (API Key)	eSevkYKyO94uGtFxxHoGwXDpX (manage keys and access tokens)
Callback URL	None
Callback URL Locked	No
Sign in with Twitter	Yes
App-only authentication	https://api.twitter.com/oauth2/token
Request token URL	https://api.twitter.com/oauth/request_token

Create a new Twitter App (cont.)

Click on Keys and Access Tokens:

Bobby_SparkStreaming_Demo_App

Details Settings **Keys and Access Tokens** Permissions

Application Settings

Keep the "Consumer Secret" a secret. This key should never be human-readable in your application.

Consumer Key (API Key)	eSevkYKyO94uGtFxxHoGwXDpX
Consumer Secret (API Secret)	OvjqlF8VmNScaeMGZjz9e1flfq0TxehmwkJ454wHCCUG0AvED
Access Level	Read and write (modify app permissions)
Owner	bobbylovemovie
Owner ID	344100790

Create a new Twitter App (cont.)

Click on Keys and Access Tokens:

Application Actions

[Regenerate Consumer Key and Secret](#) [Change App Permissions](#)

Your Access Token

You haven't authorized this application for your own account yet.

By creating your access token here, you will have everything you need to make API calls right away. The access token generated will be assigned your application's current permission level.

Token Actions 

[Create my access token](#)

Your Access Token

This access token can be used to make API requests on your own account's behalf. Do not share your access token secret with anyone.

Access Token 344100790-rOcjGTw9wdmD5FVwacOaR7ZriGPRILBZiha43Tyg

Access Token Secret 2t41euCCQiwmdukihxkBWAS2rOTfTJpTAHUy27fOAZILWE

Access Level Read and write

Owner bobbylovemovie

Owner ID 344100790

Download the third-party libraries

```
$ wget https://github.com/bobbylovemovie/trainbigdata/raw/master/Spark/  
twitter4j-core-4.0.2.jar  
$ wget https://github.com/bobbylovemovie/trainbigdata/raw/master/Spark/  
twitter4j-stream-4.0.2.jar  
$ wget https://github.com/bobbylovemovie/trainbigdata/raw/master/Spark/  
spark-streaming-twitter_2.10-1.2.0.jar
```

Run Spark-shell

```
$ spark-shell --jars spark-streaming-twitter_2.10-1.2.0.jar,twitter4j-  
stream-4.0.2.jar,twitter4j-core-4.0.2.jar
```

Running Spark commands

```
$ scala> :paste
// Entering paste mode (ctrl-D to finish)
import org.apache.spark.streaming.twitter._
import twitter4j.auth._
import twitter4j.conf._
import org.apache.spark.streaming.{Seconds, StreamingContext}
import org.apache.spark._
import org.apache.spark.streaming._
import org.apache.spark.streaming.StreamingContext._
val ssc = new StreamingContext(sc, Seconds(10))
val cb = new ConfigurationBuilder
```

Running Spark commands(cont.)

```
cb.setDebugEnabled(true).setOAuthConsumerKey("eSevkYKyO94uGtFxxHoG  
wXDpX").setOAuthConsumerSecret("OvjqlF8VmNScaeMGZjz9e1flfq0Txehmw  
kMJ454wHCCUG0AvED").setOAuthAccessToken("344100790-  
rOcjGTw9wdmD5FVwacOaR7ZriGPRILBZiha43Tyg").setOAuthAccessTokenSe  
cret("2t41euCCQiwmdukihxkBWAS2rOTfTJpTAHUy27fOAZILWE")  
  
val auth = new OAuthAuthorization(cb.build)  
  
val tweets = TwitterUtils.createStream(ssc,Some(auth))  
  
val status = tweets.map(status => status.getText)  
  
status.print  
  
ssc.checkpoint("hdfs:///user/cloudera/data/tweets")  
  
ssc.start  
  
ssc.awaitTermination
```

HUE Home Query Editors ▾ Data Browsers ▾ Workflows ▾ Search Security ▾ ≡ ☰ ⌂ ⌂ ▾ ⌂ ⌂ ⌂ ⌂

File Browser

Search for file name Actions ▾ Move to trash ▾ Upload ▾ New ▾

Home / user / cloudera / data / tweets History Trash

<input type="checkbox"/>	Name	Size	User	Group	Permissions	Date
<input type="checkbox"/>	..		cloudera	cloudera	drwxr-xr-x	October 15, 2016 09:07 AM
<input type="checkbox"/>	.		cloudera	cloudera	drwxr-xr-x	October 15, 2016 09:09 AM
<input type="checkbox"/>	5f370383-6967-4d7f-8e04-20526f5fce01		cloudera	cloudera	drwxr-xr-x	October 15, 2016 09:07 AM
<input type="checkbox"/>	checkpoint-1476547700000	5.0 KB	cloudera	cloudera	-rw-r--r--	October 15, 2016 09:08 AM
<input type="checkbox"/>	checkpoint-1476547710000	5.0 KB	cloudera	cloudera	-rw-r--r--	October 15, 2016 09:08 AM
<input type="checkbox"/>	checkpoint-1476547720000	5.0 KB	cloudera	cloudera	-rw-r--r--	October 15, 2016 09:08 AM
<input type="checkbox"/>	checkpoint-1476547730000	5.0 KB	cloudera	cloudera	-rw-r--r--	October 15, 2016 09:08 AM
<input type="checkbox"/>	checkpoint-1476547740000	5.0 KB	cloudera	cloudera	-rw-r--r--	October 15, 2016 09:09 AM
<input type="checkbox"/>	checkpoint-1476547750000	5.0 KB	cloudera	cloudera	-rw-r--r--	October 15, 2016 09:09 AM
<input type="checkbox"/>	checkpoint-1476547760000	5.0 KB	cloudera	cloudera	-rw-r--r--	October 15, 2016 09:09 AM