



Machine Reasoning

Day 4

Course Manager: GU Zhan (Sam)
zhan.gu@nus.edu.sg

Machine Reasoning

Day 4

4.1 Knowledge Discovery by Machine Learning (Big Data)

- 4.1.1 Big Data Overview
- 4.1.2 NoSQL (Not Only SQL)
- 4.1.3 Exercise

4.2 Contemporary Reasoning Systems (Big Data)

- 4.2.1 Big Data Applications
- 4.2.2 Development Eco-systems (Hadoop & Spark)
- 4.2.3 Exercise

4.3 Building Machine Reasoning System [Workshop]

- 4.3.1 Big Data Machine Learning & Reasoning System
- 4.3.2 Workshop Submission

4.4 In-class Assessment

- 4.4.1 In-class Assessment

4.1 Knowledge Discovery by Machine Learning (Big Data)

4.1.1 Big Data Overview

4.1.2 NoSQL (Not Only SQL)

4.1.3 Exercise

4.1 Knowledge Discovery by Machine Learning (Big Data)

4.1.1 Big Data Overview

4.1.2 NoSQL (Not Only SQL)

4.1.3 Exercise

Big Data

- Big data is a field that treats ways to analyze, systematically extract information from, or otherwise deal with data sets that are too large or complex to be dealt with by traditional data-processing application software.
- It is also a term used to describe large volumes of data which is hard to process using traditional database management systems (DBMS).

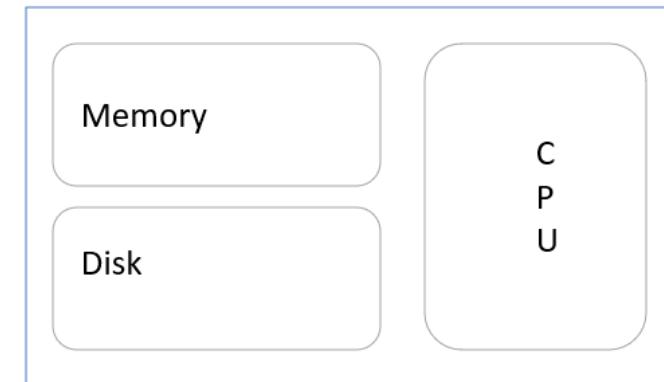
Big Data 5 V's

- 1. Volume** : refers to the huge amount of data produced;
- 2. Velocity** : refers to the speed at which data is generated;
- 3. Variety** : The variety of data available/generated, which includes (structured/semi-structured/unstructured data);
- 4. Veracity** : It refers to the uncertainties in data produced, i.e., real world data is messy and can be hard to analyze;
- 5. Value**: Insights generated must be based on accurate data so that the knowledge discovered is of value to the company;

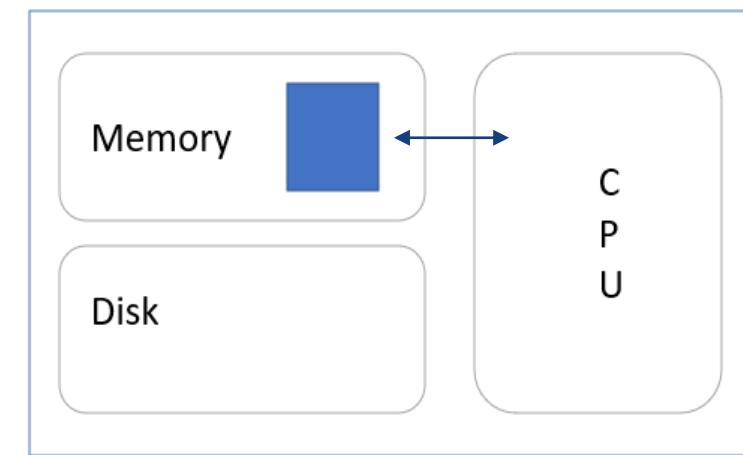
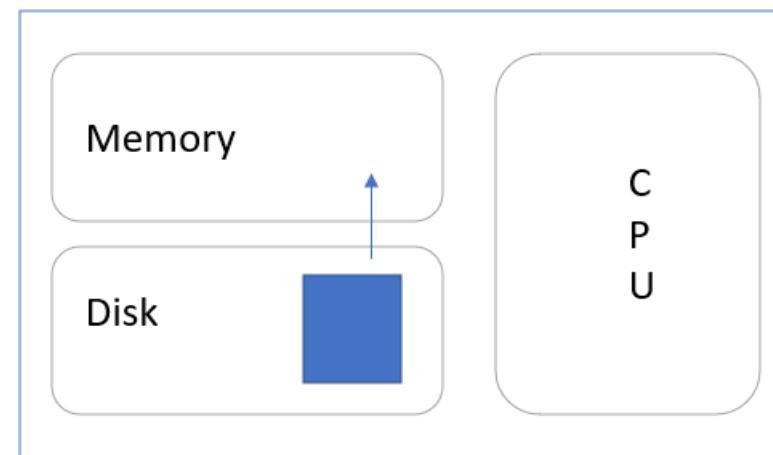
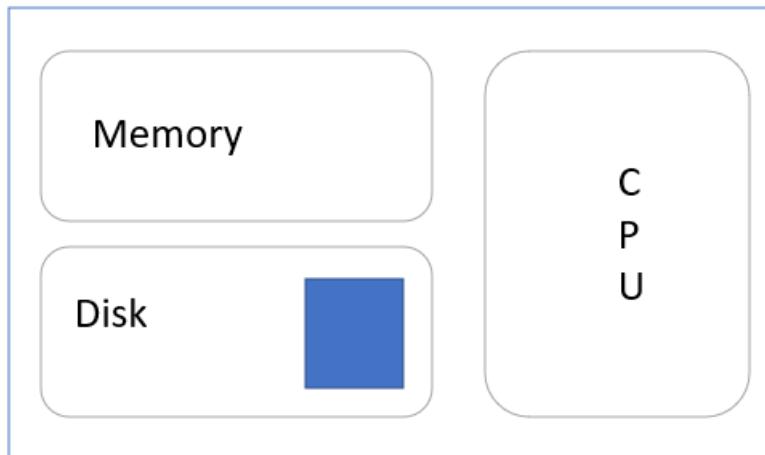
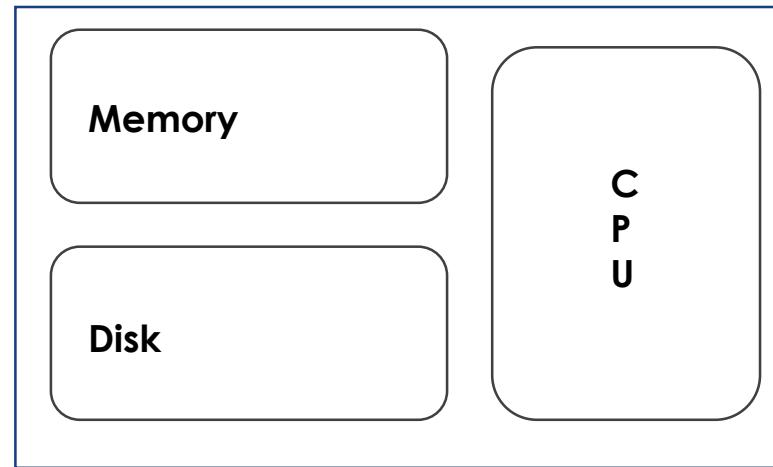
Traditional DBMS Systems

In a traditional database management system (DBMS):

- Data resides on disk and loaded into memory when needed, we do processing by spending CPU cycles on the server and then the data is presented to the user for analysis, which in turn spends/uses User CPU and memory.
- Resources (**memory size**) are limited and its **slow (I/O: inputs/outputs)** to load data from disk. Big data or extremely large data sets that could be analyzed computationally to reveal insights like patterns, trends and associations is unlikely to fit in to memory **all at once**.



Traditional DBMS Systems

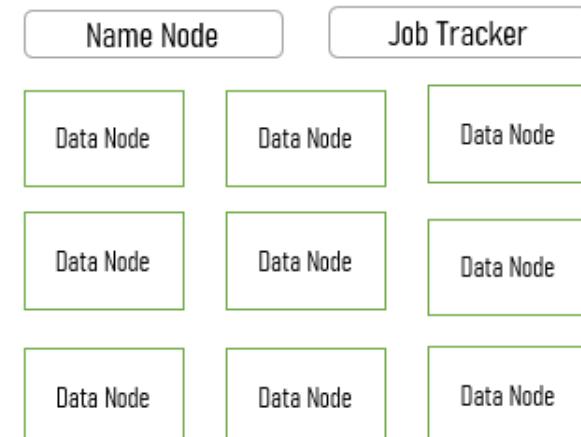


How to Handle Big Data?

- The solution is to scale out and use ***distributed computing*** by connecting many servers to perform jobs.
- A distributed system is a system whose components are located on ***different networked computers***, which communicate and coordinate their actions by passing messages to one another, via Remote Procedure Calls.

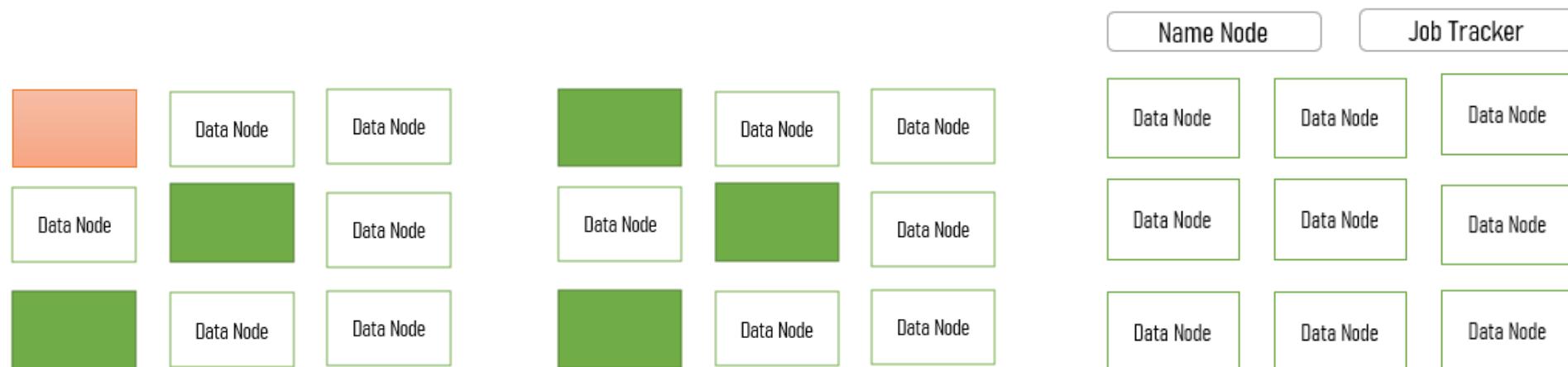
Distributed Computing

- Recover from failures
- Shared nothing architecture
- Hadoop distributed file system (HDFS)
- Map-Reduce (Divide-Conquer method)



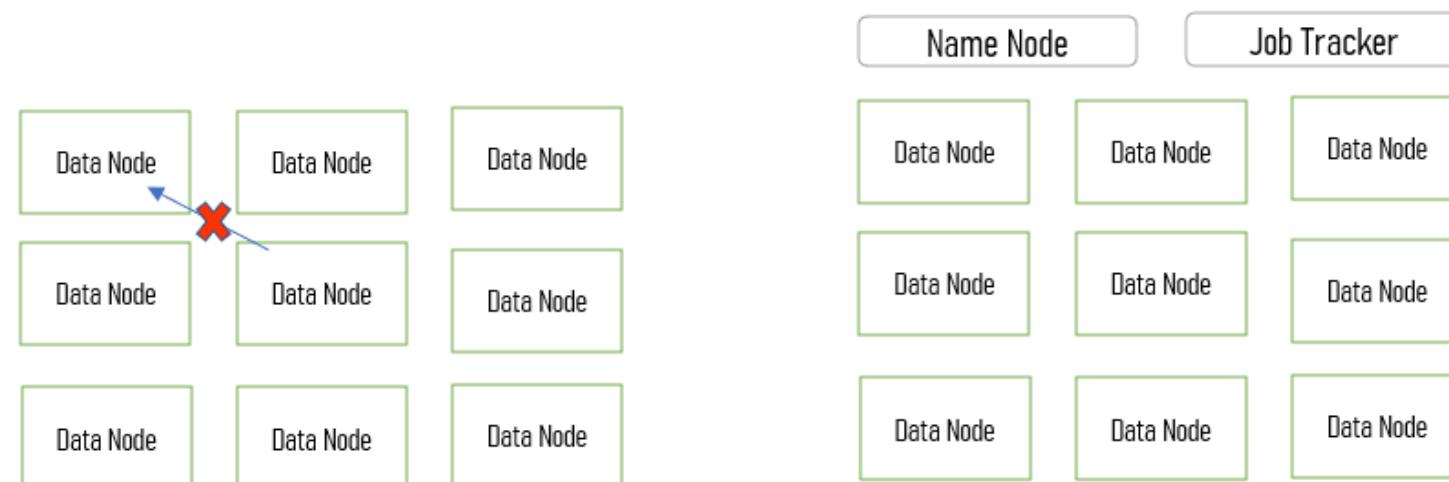
Distributed Computing

- Recover from failures: With more resources things are bound to fail more frequently, because of that they must be able to recover from failure in a scalable manner;



Distributed Computing

- Shared Nothing Architecture: In order to scale out efficiently, we do not burden the master servers / name nodes, like coordination and communication with other data nodes: reduced networking communication overheads;
- Each **data node** is unaware of the other nodes in a cluster;
- The Coordination task is left to the master servers: Name Node and Job Tracker;



Distributed Computing

- Each Data node is unaware of the other nodes in a cluster;
- Data Nodes be like ... I GUESS I'M THE ONLY ONE HERE



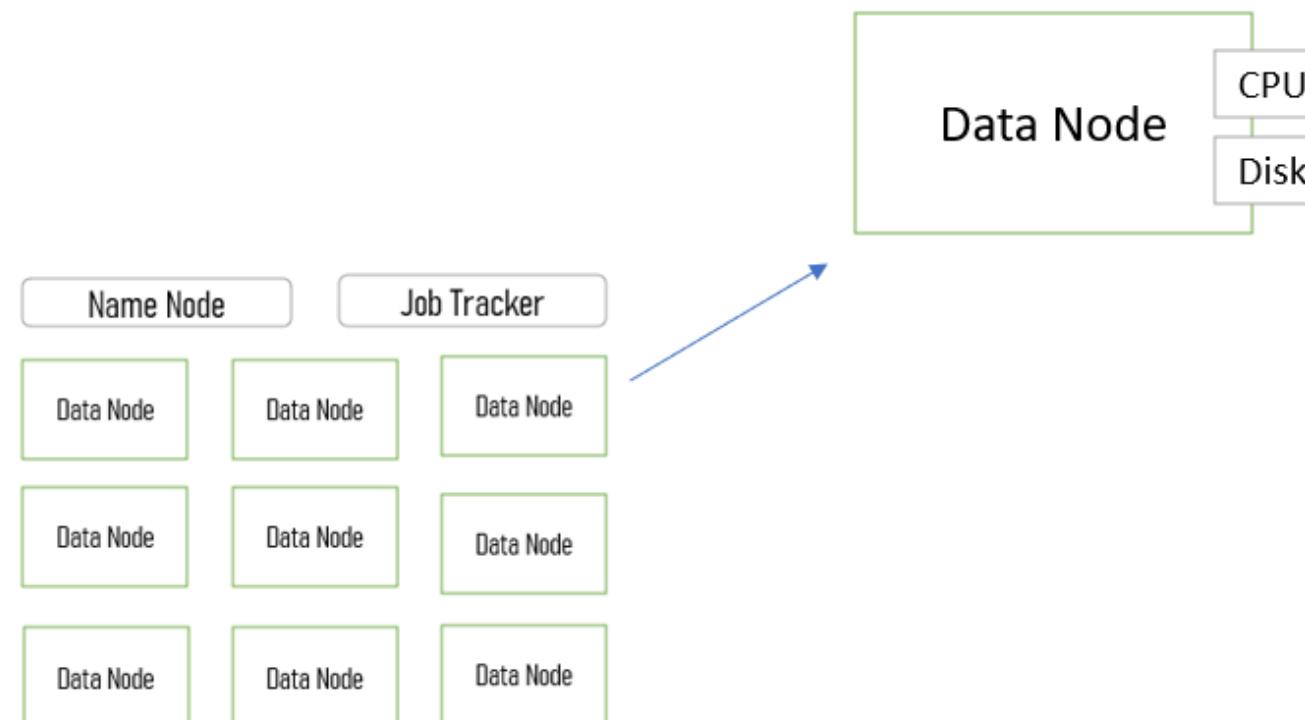
Distributed Computing

HDFS (Hadoop Distributed File System)

- Hadoop is an opensource implementation of the Hadoop distributed file System. It masks the complexities when working with distributed systems;
- Provides fault tolerant distributed file system on top of these nodes; (servers/storages) carrying forward all the characteristics as well as an API that makes it easier to work on Hadoop;

Distributed Computing

- Map-Reduce: Takes compute against ***data stored on disk*** instead of ***data loaded in memory***,



4.1 Knowledge Discovery by Machine Learning (Big Data)

4.1.1 Big Data Overview

4.1.2 NoSQL (Not Only SQL)

4.1.3 Exercise

What is NoSQL?

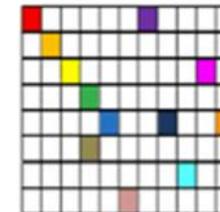
- **NoSQL** databases (aka. "not only SQL") are non tabular, and store data differently than relational tables;
- NoSQL databases come in a variety of types based on their data model / db schema;
- The main types are: column-family/wide-column, graph, document, key-value;

Types of NoSQL Databases

- **Column-family / Wide-column stores** store data in tables, rows, and dynamic columns;
- **Graph databases** store data in nodes and edges;
- **Document databases** store data in documents similar to JSON (**JavaScript Object Notation**) objects. Each document contains pairs of (embedded) fields and values;
- **Key-value databases** are a simpler type of database where each item contains keys and values;

NoSQL Database

Column-Family



SQL Database

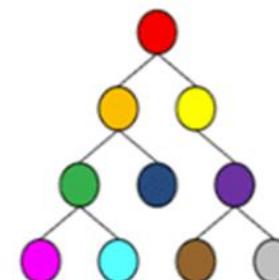
Graph



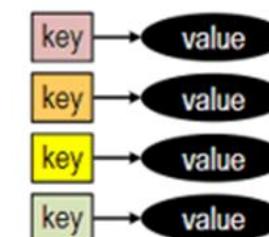
Relational



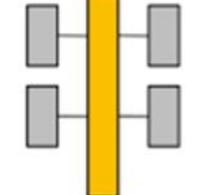
Document



Key-Value



Analytical (OLAP)



SQL Databases

Data Storage Model

Tables with fixed rows and columns

Development History

Developed in the 1970s with a focus on reducing data duplication

Examples

Oracle, MySQL, Microsoft SQL Server, and PostgreSQL

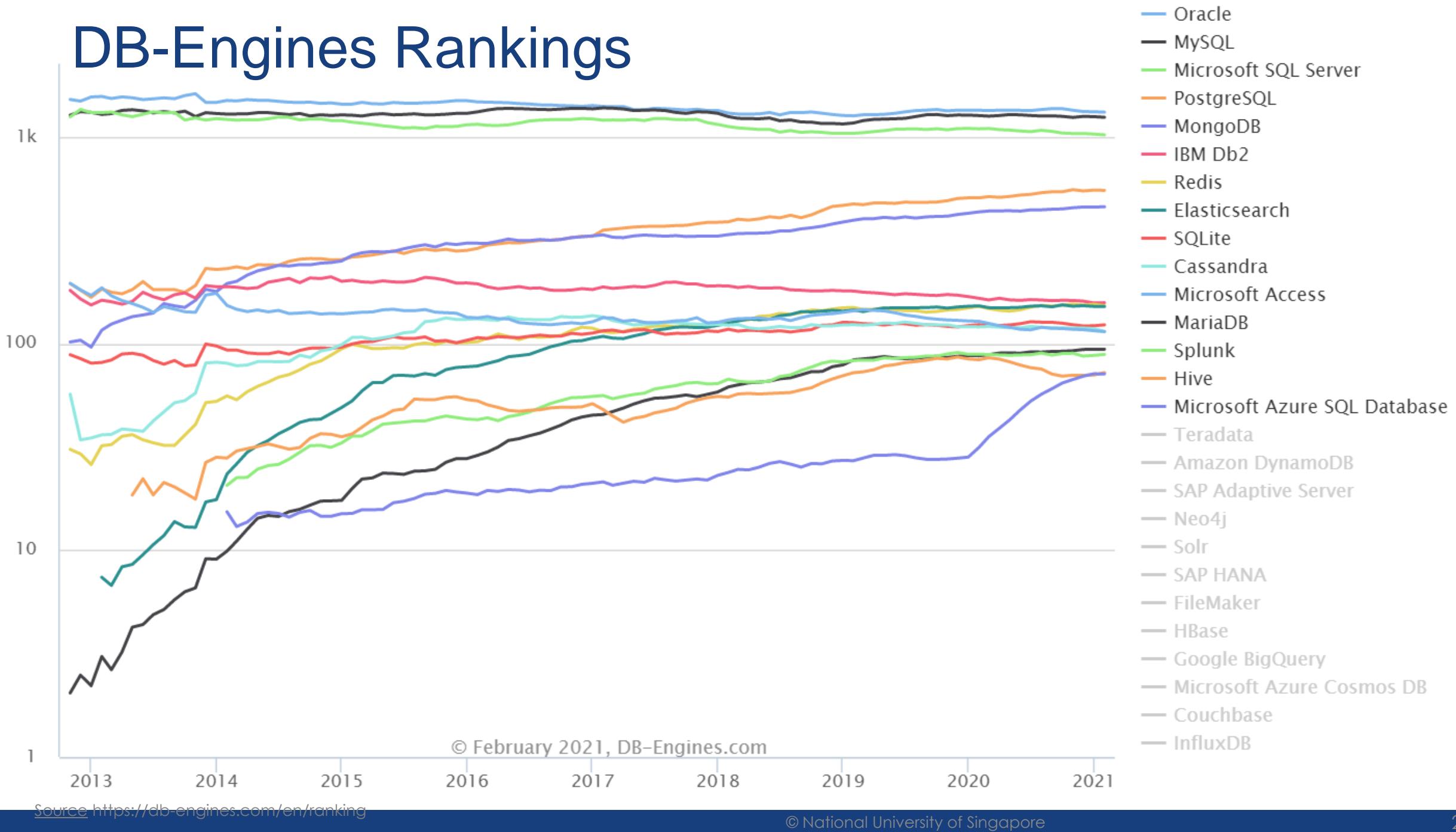
NoSQL Databases

Document: JSON documents, Key-value: key-value pairs, Wide-column: tables with rows and dynamic columns, Graph: nodes and edges

Developed in the late 2000s with a focus on scaling and allowing for rapid application change driven by agile and DevOps practices.

Document: MongoDB and CouchDB, Key-value: Redis and DynamoDB, Wide-column: Cassandra and HBase, Graph: Neo4j and Amazon Neptune

DB-Engines Rankings



Data Structure in SQL DB

- Example: Employee DB
- Each row in each table is a single value(employee), and each column has fixed attributes:

Emp_ID	Emp_firstname	Emp_lastname	Emp_phone
1	Geet	Jethwani	81815656

Data Structure in NoSQL DB

- Example: document database – JSON;
- Each of the SQL column attributes would be fields;
- The details of an employee's record would be the data values associated with each field;
- For Example: Emp_firstname: “Geet”, Emp_lastname: “Jethwani”

Data Structure: JSON

- Experience data represented in semi-structured format;
- We have looked at the structured (tabular) formats;
- However, we must understand JSON representation as in real world big data application, one may come across this very frequently;
- When setting up a data store, your first task is to answer the question: “What data would I like to store and how do the fields relate to each other”?

name	quantity	size	status	tags	rating
journal	25	14x21,cm	A	brown, lined	9
notebook	50	8.5x11,in	A	college-ruled,perforated	8

Data Structure: JSON

- JSON is formatted as name-value pairs;
- In JSON documents, fieldnames and values are separated by a colon, fieldname and value pairs are separated by commas, and sets of fields are encapsulated in “curly brackets” ({});
- Example:

```
{  
    "name": "notebook",  
    "quantity": 50,  
    "size": { "length": 11, "width": 8.5, "unit": "in" },  
    "status": "A",  
    "tags": [ "college-ruled", "perforated" ],  
    "rating": 8  
}
```

4.1 Knowledge Discovery by Machine Learning (Big Data)

4.1.1 Big Data Overview

4.1.2 NoSQL (Not Only SQL)

4.1.3 Exercise

[Exercise] Knowledge Discovery by Machine Learning (Big Data)

- **Setup Databricks account use for free: Community Edition.**
- **Structure data in JSON and flatten it into a table.**

What is DataBricks?

- Unified big data analytics platform;
- Run all knowledge discovery & analytics in one place;
- Powering dashboards, running reports;
- Machine learning and streaming jobs;
- Databricks is more optimized than open-source apache spark;
- Easier installation and no hassles of managing clusters;

- <https://databricks.com/>

All your data,
analytics and AI on one
unified data platform

TRY FOR FREE LEARN MORE

Databricks is
the data and AI
company



Try Databricks for free

An open and unified data analytics platform for data engineering, data science, machine learning, and analytics.

From the original creators of Apache Spark™, Delta lake, MLflow, and Koalas.



Databricks trial:

- Collaborative environment for data teams to build solutions together.
- Interactive notebooks to use Apache Spark™, SQL, Python, Scala, Delta Lake, MLflow, TensorFlow, Keras, Scikit-learn and more.
- Available as a 14-day full trial in your own cloud, or as a lightweight trial hosted by Databricks.

Used by:



Please tell us about yourself

First Name: *

Last Name: *

Company *

Company Email *

Title *

Phone Number

Keep me informed with occasional updates about Databricks and related open source products

By Clicking "Get Started For Free", you agree to the [Privacy Policy](#).

GET STARTED FOR FREE

I'm not a robot



reCAPTCHA

[Privacy](#) · [Terms](#)

Choose a cloud provider

 Amazon Web Services

 Microsoft Azure

 Google Cloud Platform

Get started

By clicking "Get started", you agree to the [Privacy Policy](#) and [Terms of Service](#)

Don't have a cloud account?

Community Edition is a limited Databricks environment for personal use and training.

[Get started with Community Edition](#)

By clicking "Get started with Community Edition", you agree to the [Privacy Policy](#) and [Community Edition Terms of Service](#)

Select
Community Edition



Sign In to Databricks
Community Edition

 Email / Username

 Password

[Forgot Password?](#)

Sign In

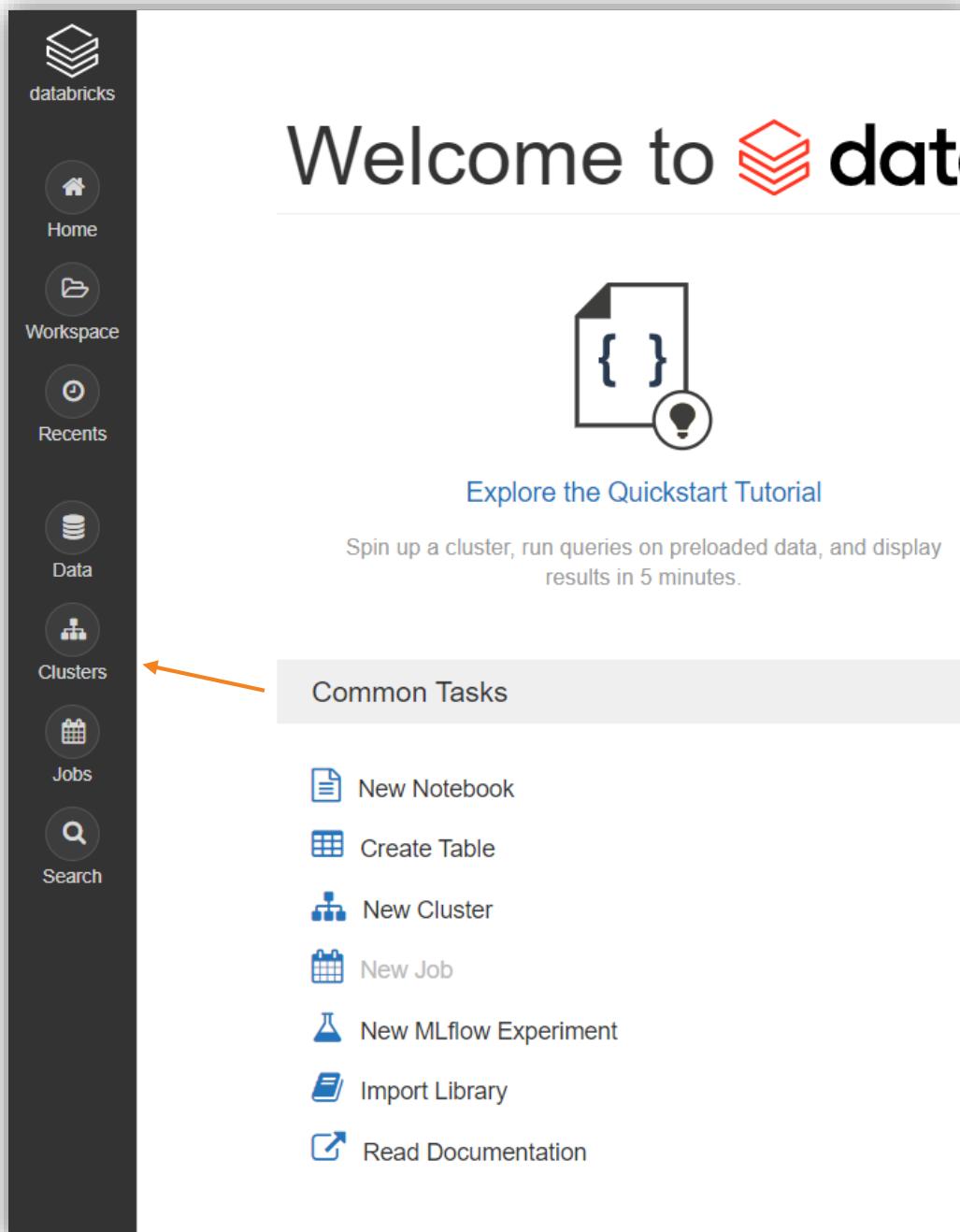
New to Databricks? [Sign Up.](#)

[Privacy Policy](#) | [Terms of Use](#)

Login to your account;

Community Edition log in portal:

<https://community.cloud.databricks.com/login.html>

A screenshot of the Databricks web interface. On the left is a dark sidebar with white icons and text for Home, Workspace, Recents, Data, Clusters (which has an orange arrow pointing to it), Jobs, and Search. The main area is titled "Welcome to databricks" and features a "Quickstart Tutorial" section with a code editor icon and a "Explore the Quickstart Tutorial" link. Below this is a callout about spinning up a cluster in 5 minutes. A "Common Tasks" section is listed on the right, containing links for New Notebook, Create Table, New Cluster, New Job, New MLflow Experiment, Import Library, and Read Documentation.

- Now create clusters;
- To get your cluster up and running;
- Click the clusters button as shown;

Clusters

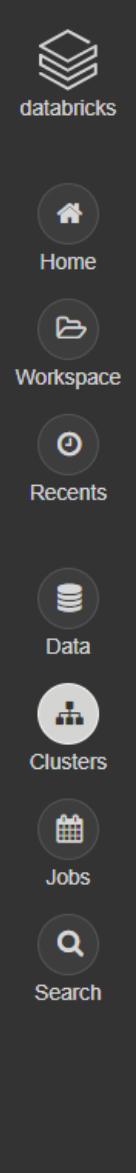
+ Create Cluster 

▼ All-Purpose Clusters

No clusters found

▼ Job Clusters

No clusters found



- Home
- Workspace
- Recents
- Data
- Clusters
- Jobs
- Search

Create Cluster

New Cluster

Cancel Create Cluster

0 Workers: 0.0 GB Memory, 0 Cores, 0 DBU
1 Driver: 15.3 GB Memory, 2 Cores, 1 DBU 

Cluster Name: ISDemoCluster

Databricks Runtime Version: Runtime: 6.5 ML (Scala 2.11, Spark 2.4.5) 

Instance: Free 15GB Memory: As a Community Edition user, your cluster will automatically terminate after an idle period of two hours. [For more configuration options, please upgrade your Databricks subscription.](#)

Instances 

Availability Zone: us-west-2c 

 databricks

Home Workspace Recents Data Clusters Jobs Search

Clusters /

ISDemoCluster

[Edit](#) [Clone](#) [Restart](#) [Terminate](#) [Delete](#)

Configuration Notebooks Libraries Event Log Spark UI Driver Logs Metrics Apps Spark Cluster UI - Master ▾

Databricks Runtime Version

7.5 (includes Apache Spark 3.0.1, Scala 2.12)

New This Runtime version supports only Python 3.

Driver Type

Community Optimized 15.3 GB Memory, 2 Cores, 1 DBU

Instance

Free 15GB Memory: As a Community Edition user, your cluster will automatically terminate after an idle period of two hours.
For [more configuration options](#), please [upgrade your Databricks subscription](#).

Instances Spark JDBC/ODBC Permissions

Availability Zone ?

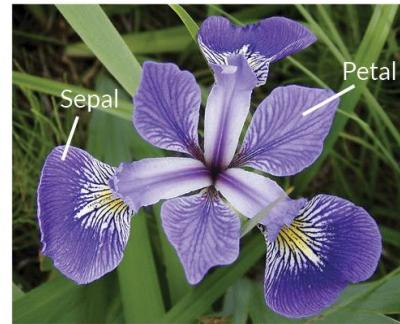
us-west-2c

Data Representation in JSON

- In our Machine Memory module we used SQLite DB to represent and store data a structured tabular format;
- However, in other scenarios we may obtain data in different formats which may not be as simple to interpret as the relational data format (tables);
- Let us now structure our data in JSON format: <https://en.wikipedia.org/wiki/JSON>.

Data Representation in JSON Example

- Lets take an example of Iris flower dataset.
- When represented in the tabular format:



Iris Versicolor



Iris Setosa



Iris Virginica

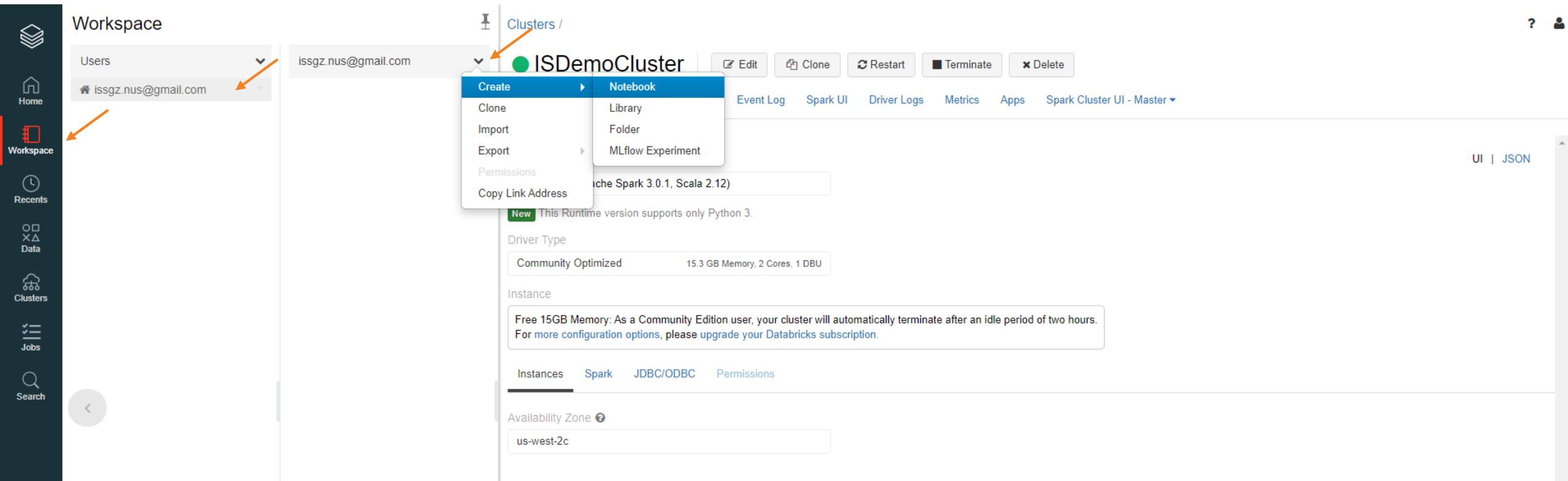
SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Species
5.1	3.5	1.4	0.2	Iris-setosa
4.9	3	1.4	0.2	Iris-setosa
4.7	3.2	1.3	0.2	Iris-setosa
4.6	3.1	1.5	0.2	Iris-setosa
5	3.6	1.4	0.2	Iris-setosa
5.4	3.9	1.7	0.4	Iris-setosa
4.6	3.4	1.4	0.3	Iris-setosa
5	3.4	1.5	0.2	Iris-setosa
4.4	2.9	1.4	0.2	Iris-setosa

Data Representation in JSON Example

- The same content can be structured in JSON as:

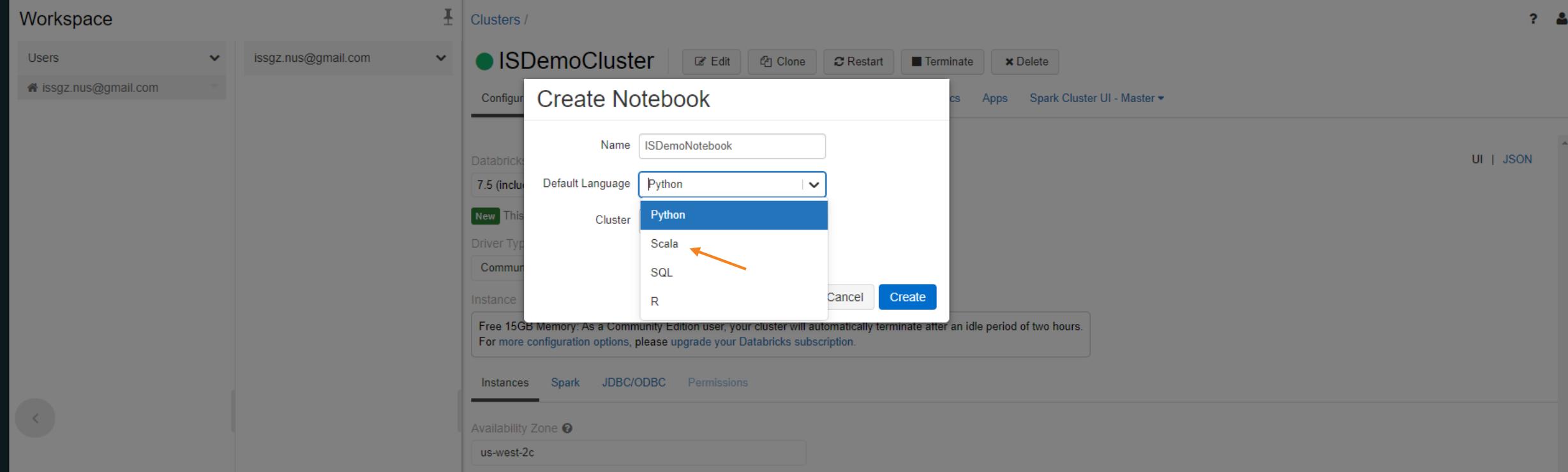
```
//code is in scala %scala
dbutils.fs.put("/tmp/test.json", """
{"SepalLengthCm": 5.1, "SepalWidthCm": 3.5, "PetalLengthCm": 1.4, "PetalWidthCm": 0.2, "Species": "setosa"},
 {"SepalLengthCm": 4.9, "SepalWidthCm": 3.0, "PetalLengthCm": 1.4, "PetalWidthCm": 0.2, "Species": "setosa"},
 {"SepalLengthCm": 4.7, "SepalWidthCm": 3.2, "PetalLengthCm": 1.3, "PetalWidthCm": 0.2, "Species": "setosa"},
 {"SepalLengthCm": 4.6, "SepalWidthCm": 3.1, "PetalLengthCm": 1.5, "PetalWidthCm": 0.2, "Species": "setosa"},
 {"SepalLengthCm": 5.0, "SepalWidthCm": 3.6, "PetalLengthCm": 1.4, "PetalWidthCm": 0.2, "Species": "setosa"},
 {"SepalLengthCm": 5.4, "SepalWidthCm": 3.9, "PetalLengthCm": 1.7, "PetalWidthCm": 0.4, "Species": "setosa"},
 {"SepalLengthCm": 4.6, "SepalWidthCm": 3.4, "PetalLengthCm": 1.4, "PetalWidthCm": 0.3, "Species": "setosa"},
 {"SepalLengthCm": 5.0, "SepalWidthCm": 3.4, "PetalLengthCm": 1.5, "PetalWidthCm": 0.2, "Species": "setosa"},
 {"SepalLengthCm": 4.4, "SepalWidthCm": 2.9, "PetalLengthCm": 1.4, "PetalWidthCm": 0.2, "Species": "setosa"},
"""
, true)
// COMMAND -----
val testJsonData = spark.read.json("/tmp/test.json")
display(testJsonData)
```

Your email id should appear in the workspace, create notebooks here



The screenshot shows the Databricks workspace interface. On the left, there's a sidebar with icons for Home, Workspace (which is selected), Recents, Data, Clusters, Jobs, and Search. The main area is titled 'Clusters /' and shows a cluster named 'ISDemoCluster'. A dropdown menu is open over the cluster name, with 'Notebook' selected. Other options in the menu include Create (Clone, Import, Export, Permissions, Copy Link Address), Notebook (Library, Folder, MLflow Experiment), and Cluster actions (Edit, Clone, Restart, Terminate, Delete). Below the cluster name, there's a note: 'New This Runtime version supports only Python 3.' Under 'Driver Type', it says 'Community Optimized' with '15.3 GB Memory, 2 Cores, 1 DBU'. In the 'Instance' section, there's a message: 'Free 15GB Memory: As a Community Edition user, your cluster will automatically terminate after an idle period of two hours. For more configuration options, please upgrade your Databricks subscription.' At the bottom, tabs for 'Instances', 'Spark', 'JDBC/ODBC', and 'Permissions' are visible, with 'Instances' being the active tab. The availability zone is set to 'us-west-2c'.

Choose: Scala



The screenshot shows the Databricks workspace interface. A modal dialog box titled "Create Notebook" is open in the foreground. Inside the dialog, the "Name" field contains "ISDemoNotebook". The "Default Language" dropdown is set to "Python", with "Scala" highlighted by an orange arrow. Other options in the dropdown are "SQL" and "R". At the bottom right of the dialog are "Cancel" and "Create" buttons. Below the dialog, a message states: "Free 15GB Memory: As a Community Edition user, your cluster will automatically terminate after an idle period of two hours. For more configuration options, please upgrade your Databricks subscription." The background shows the "Clusters" page with a cluster named "ISDemoCluster".

Workspace

 Users
 issgz.nus@gmail.com

 issgz.nus@gmail.com
 ISDemoNotebook

ISDemoNotebook (Scala)



Detached File View: Standard Permissions Run All Clear Publish Comments Experiment Revision history

```

1 // Databricks notebook source
2 // code is in scala %scala
3 dbutils.fs.put("/tmp/test.json", """
4 {"sepallength": 5.1, "sepalWidth": 3.5, "petalLength": 1.4, "petalWidth": 0.2, "species": "setosa"},
5 {"sepallength": 4.9, "sepalWidth": 3.0, "petalLength": 1.4, "petalWidth": 0.2, "species": "setosa"},
6 {"sepallength": 4.7, "sepalWidth": 3.2, "petalLength": 1.3, "petalWidth": 0.2, "species": "setosa"},
7 {"sepallength": 4.6, "sepalWidth": 3.1, "petalLength": 1.5, "petalWidth": 0.2, "species": "setosa"},
8 {"sepallength": 6.7, "sepalWidth": 3.1, "petalLength": 5.6, "petalWidth": 2.4, "species": "virginica"},
9 {"sepallength": 6.9, "sepalWidth": 3.1, "petalLength": 5.1, "petalWidth": 2.3, "species": "virginica"},
10 {"sepallength": 5.8, "sepalWidth": 2.7, "petalLength": 5.1, "petalWidth": 1.9, "species": "virginica"}
11 """), true)
12
13 // COMMAND -----
14
15 val testJsonData = spark.read.json("/tmp/test.json")
16
17 display(testJsonData)
18

```

(2) Spark Jobs

testJsonData: org.apache.spark.sql.DataFrame = [petalLength: double, petalWidth: double ... 3 more fields]

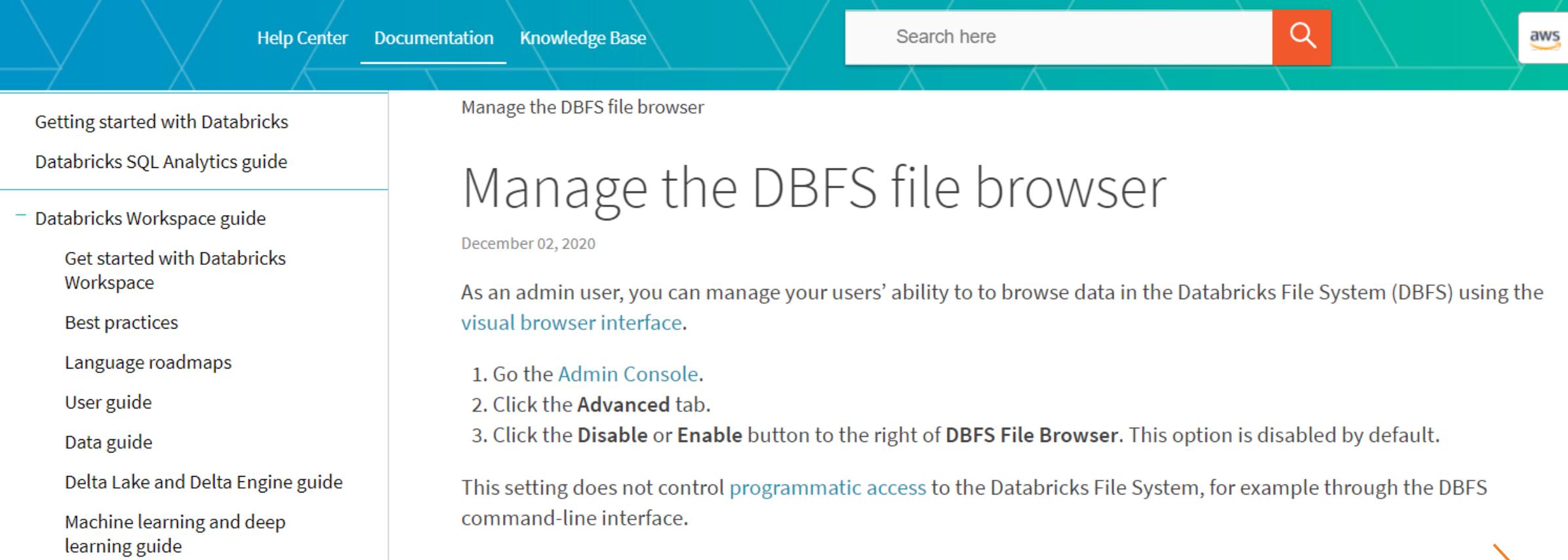
	petalLength	petalWidth	sepallength	sepalWidth	species
1	1.4	0.2	5.1	3.5	setosa
2	1.4	0.2	4.9	3	setosa
3	1.3	0.2	4.7	3.2	setosa
4	1.5	0.2	4.6	3.1	setosa
5	5.6	2.4	6.7	3.1	virginica
6	5.1	2.3	6.9	3.1	virginica
7	5.1	1.9	5.8	2.7	virginica

Showing all 7 rows.

Command took 2.89 seconds -- by issgz.nus@gmail.com at 1/27/2021, 6:08:02 PM on ISDemoCluster

Shift+Enter to run [shortcuts](#)

How to enable DBFS file browser?



The screenshot shows the Databricks Documentation website. The top navigation bar includes links for Help Center, Documentation (which is underlined), Knowledge Base, a search bar with a magnifying glass icon, and an AWS logo. On the left, a sidebar menu lists various guides: Getting started with Databricks, Databricks SQL Analytics guide, Databricks Workspace guide (with sub-links for Get started with Databricks, Workspace, Best practices, Language roadmaps, User guide, Data guide, Delta Lake and Delta Engine guide, and Machine learning and deep learning guide). The main content area is titled "Manage the DBFS file browser" and was last updated on December 02, 2020. It contains instructions for enabling the DBFS File Browser in the Admin Console. Below the instructions, a note states that this setting does not control programmatic access. An orange arrow points from the bottom right towards the "Admin Console" link in the instructions.

Getting started with Databricks

Databricks SQL Analytics guide

- Databricks Workspace guide

- Get started with Databricks
- Workspace
- Best practices
- Language roadmaps
- User guide
- Data guide
- Delta Lake and Delta Engine guide
- Machine learning and deep learning guide

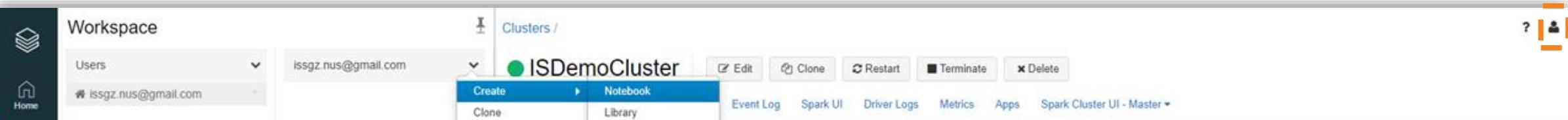
Manage the DBFS file browser

December 02, 2020

As an admin user, you can manage your users' ability to browse data in the Databricks File System (DBFS) using the visual browser interface.

1. Go to the [Admin Console](#).
2. Click the **Advanced** tab.
3. Click the **Disable** or **Enable** button to the right of **DBFS File Browser**. This option is disabled by default.

This setting does not control [programmatic access](#) to the Databricks File System, for example through the DBFS command-line interface.



The screenshot shows the Databricks Cluster UI. The top navigation bar includes a workspace switcher, user dropdown, and cluster dropdown. The cluster dropdown is set to "ISDemoCluster". The main interface shows the cluster status (green dot) and controls for Edit, Clone, Restart, Terminate, and Delete. Below the cluster controls, there are tabs for Create, Notebook, Event Log, Spark UI, Driver Logs, Metrics, Apps, and Spark Cluster UI - Master. A red arrow points from the bottom right towards the "Edit" button.

[Exercise]

- Refer to Iris flower example, structure all (or first 10) records of the toddler Autism data (csv) into JSON format;
- Use the created JSON format to import it into databrick cluster using notebook;
- **[Optional]** Visualize using the different graphs. Choose appropriate graphs to visualize your data;

4.2 Contemporary Reasoning Systems (Big Data)

4.2.1 Big Data Applications

4.2.2 Development Eco-systems (Hadoop & Spark)

4.2.3 Exercise

4.2 Contemporary Reasoning Systems (Big Data)

4.2.1 Big Data Applications

4.2.2 Development Eco-systems (Hadoop & Spark)

4.2.3 Exercise

Big Data Applications (Hadoop)



- “One of the factors driving transformation for Western Union right now is the massive amount of transactional information that we accumulate as we serve our customers,” noted Sanjay Saraf, the company’s chief technology officer and senior vice president. “We’ve built an enterprise data hub on Cloudera to drive actionable insights that help the company create products and services that are relevant to our customers and help differentiate Western Union in a competitive marketplace.”
- Solution: Western Union’s enterprise data hub (EDH) feeds in structured data from multiple data warehouses as well as unstructured data including click streams, behavioral data, logs, and sentiment data collected by tools such as transactional, marketing, and other outreach systems. Western Union uses a combination of Apache Flume, Apache Sqoop, and Informatica Big Data Edition (BDE) to collect data from the various sources.



- Tesla is using a Hadoop cluster to collect the increasing amount of data being generated by its connected cars.
- CIO Jay Vijayan said: “We are working on a big data platform... The car is connected, but it does not really talk to the network every minute because we want to keep it as smart and efficient as possible. It alerts us if the car is not functioning properly so service teams can take action.”

NETFLIX

- Netflix has been determined to be able to predict what exactly its customers will enjoy watching with Big Data. As such, Big Data analytics is the fuel that fires the ‘recommendation engine’ designed to serve this purpose.
- More recently, Netflix started positioning itself as a content creator, not just a distribution method. Unsurprisingly, this strategy has been firmly driven by data. Netflix’s recommendation engines and new content decisions are fed by data points such as what titles customers watch, how often playback stopped, ratings are given, etc. The company’s data structure includes Hadoop, Hive and Pig with much other traditional business intelligence.

The Secret Ingredients to Netflix's Success

1

Targeted Use Case

What are the business use cases upon which your big data and data science initiative should focus?

Why are these use cases important to the business?

How: Netflix's success rests on their ability to increase customer engagement through their **recommendation engine**.

2

Create Analytic Profiles

Viewer preferences are integrated with external data sources (like social media) in an **Analytic Profile**. These profiles capture analytic assets in a way that can be utilized across multiple use cases.



3

Capture Show Characteristics and Viewing Patterns

Build **Analytic Profiles** for each individual product.
Netflix uses 100+ data points to tag each title.

5

Management Fortitude to Become "Netflix Intelligent"

The management team must have a willingness to learn how to properly value and use the organization's data and analytics.

4

Mastering Machine Learning

Data is fed into **Machine Learning** algorithms to create critical scores. The key to success here is to have a data science exploration and learning process to discover better predictors of customer behavior.



- A big technical challenge for eBay as a data-intensive business to exploit a system that can rapidly analyze and act on data as it arrives (streaming data).
- There are many rapidly evolving methods to support streaming data analysis. eBay is working with several tools including Apache Spark, Storm, Kafka. It allows the company's data analysts to search for information tags that have been associated with the data (metadata) and make it consumable to as many people as possible with the right level of security and permissions (data governance).
- The company has been at the forefront of using big data solutions and actively contributes its knowledge back to the open-source community.

Apache Spark at ebay



eBay uses Apache Spark to provide targeted offers, enhance customer experience



Apache Spark is leveraged at eBay through Hadoop YARN. YARN manages all the cluster resources to run generic tasks



EBay spark users leverage the Hadoop clusters in the range of 2000 nodes, 20,000 cores and 100TB of RAM through YARN

4.2 Contemporary Reasoning Systems (Big Data)

4.2.1 Big Data Applications

4.2.2 Development Eco-systems (Hadoop & Spark)

4.2.3 Exercise



is a distributed hard-disk computation platform.

- Hadoop distributed file systems (HDFS)
- Divide and conquer approach: Map-Reduce approach



is a distributed in-memory computation platform.

- Resilient Distributed Dataset (RDD)
- Spark RDD & ML library (decision tree)

Spark vs Hadoop MapReduce

Factors

Speed

Written In

Data Processing

Ease of Use

Caching

Spark

100x times than MapReduce

Scala

Batch / real-time / iterative /
interactive / graph

Compact & easier than Hadoop

Caches the data in-memory &
enhances the system performance

Hadoop MapReduce

Faster than traditional system

Java

Batch processing

Complex & lengthy

Doesn't support caching of data

4.2 Contemporary Reasoning Systems (Big Data)

4.2.1 Big Data Applications

4.2.2 Development Eco-systems (Hadoop & Spark)

4.2.3 Exercise

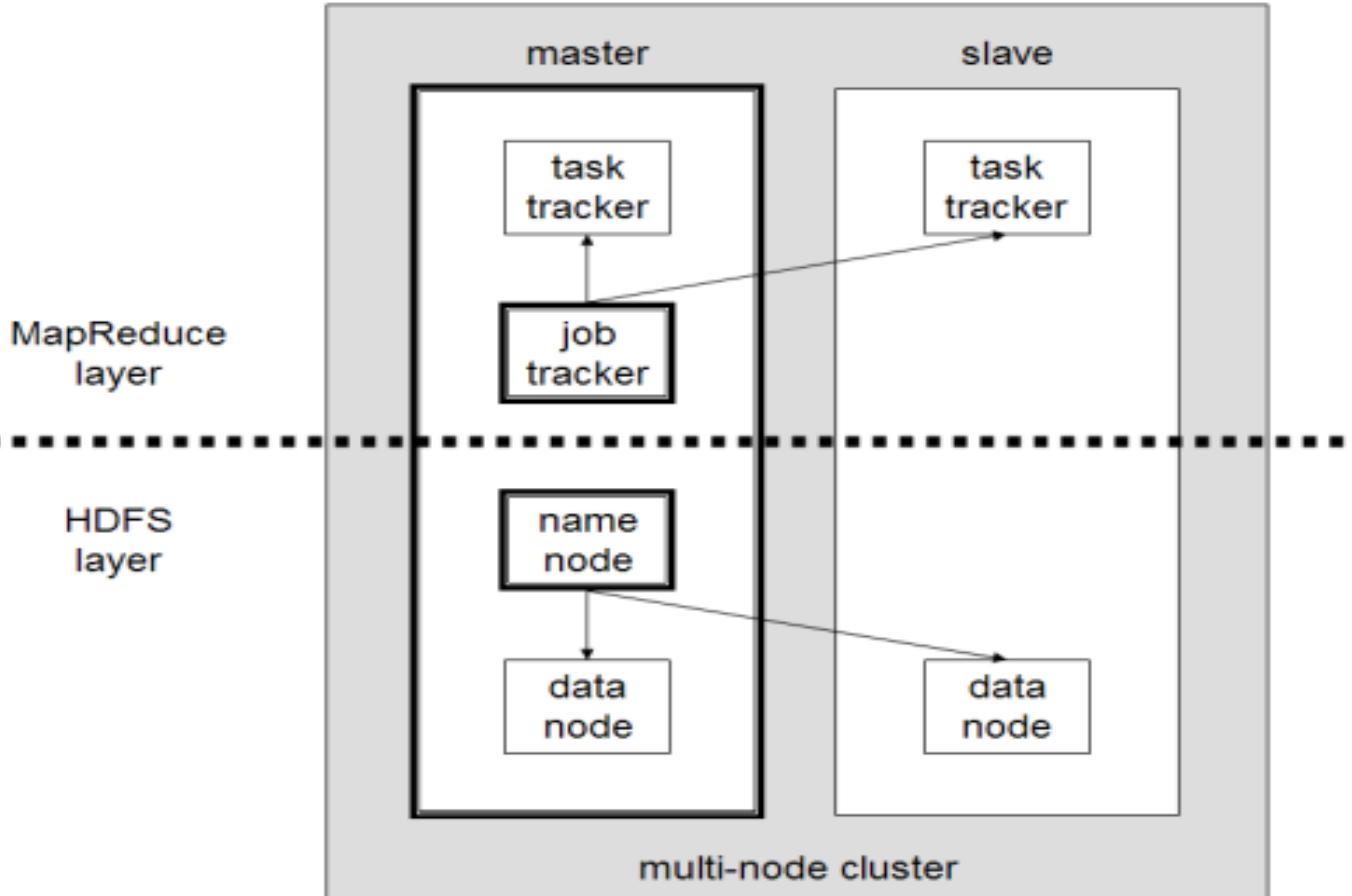
Introduction to Hadoop Eco-systems

- Hadoop is the ***open source*** implementation of the google file system and Map-Reduce ***distributed computation***;
- Hadoop masks the complexities that are involved in working with distributed systems and provides a ***fault tolerant*** distributed file system (HDFS) as well as an API that makes programming in Hadoop a lot easier.
- HDFS: Hadoop Distributed File System is based upon Google File System.
- Hadoop is divided into ***HDFS*** and ***Map-Reduce***. HDFS is used for storing the data and Map-Reduce is used for processing data.
- Hadoop deals with ***files***; It is ***not a database***.

Hadoop HDFS Architecture

- The Hadoop distributed file system (HDFS) is a distributed, scalable, and portable file system written in Java for the Hadoop framework.
- HDFS has five services as follows:
 1. Name Node
 2. Data node
 3. Job tracker
 4. Secondary Name Node
 5. Task Tracker

Hadoop HDFS Architecture: Services



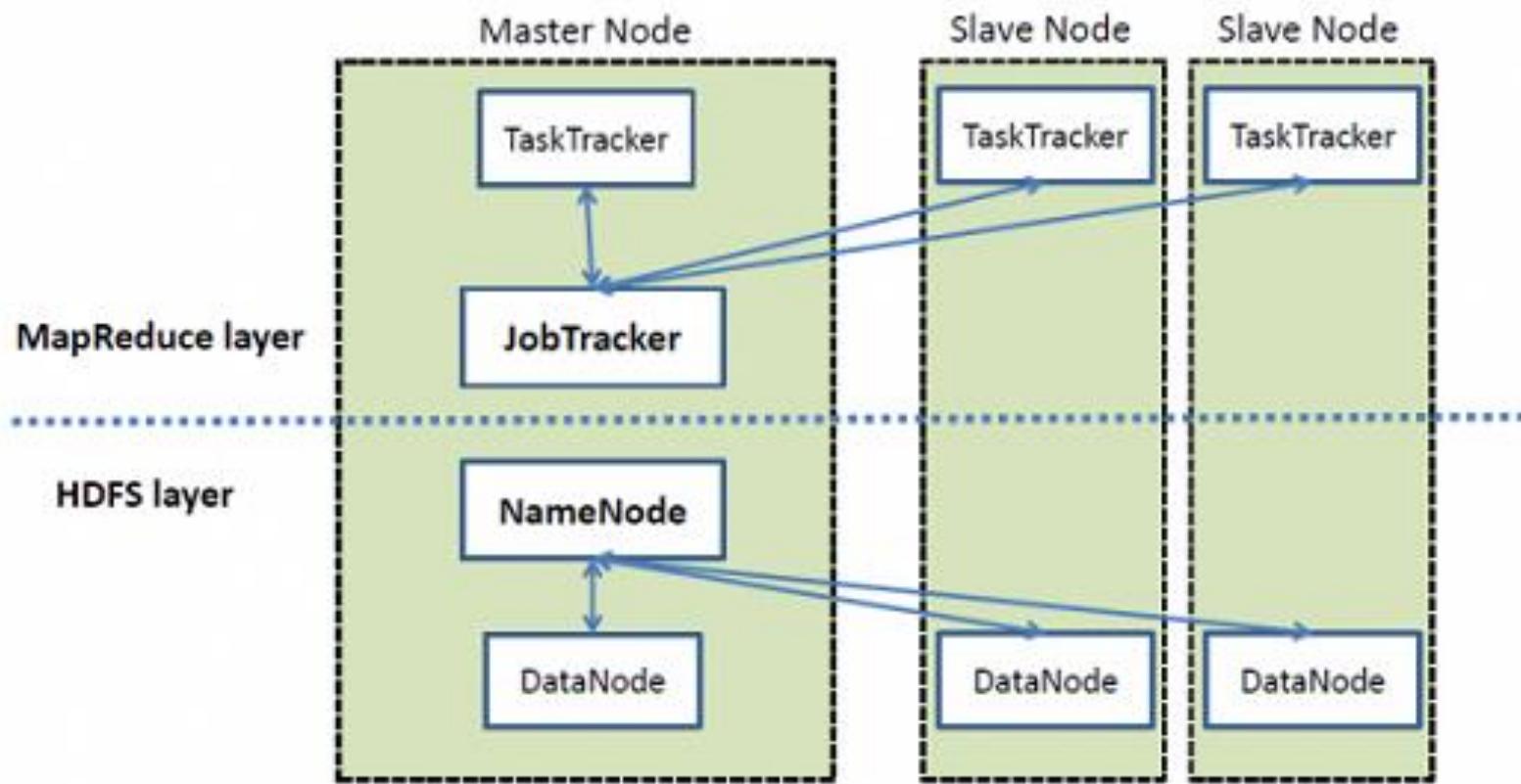
- 1. Name Node:** HDFS consists of only one Name Node that is called the Master Node. The master node can track files, manage the file system and has the metadata of all of the stored data within it.
- 2. Data Node:** This is also known as the slave node and it stores the actual data into HDFS.
- 3. Job tracker:** Job tracker talks to the Name Node to know about the location of the data that will be used in processing.
- 4. Secondary Name Node:** This is also known as the checkpoint Node. It is the helper Node for the Name Node.
- 5. Task Tracker:** It is the Slave Node for the Job Tracker and it takes the task from the Job Tracker.

Source

https://s3.amazonaws.com/files.dezyre.com/images/blog/Hadoop+Architecture+Explained-What+it+is+and+why+it+matters/Hadoop+Architecture_OpenSource.png

Hadoop HDFS Architecture: Server Cluster

- Name Node is a master node and Data node is its corresponding Slave node and can talk with each other
- Enables ***multiple*** machines to perform computation of data;



Source

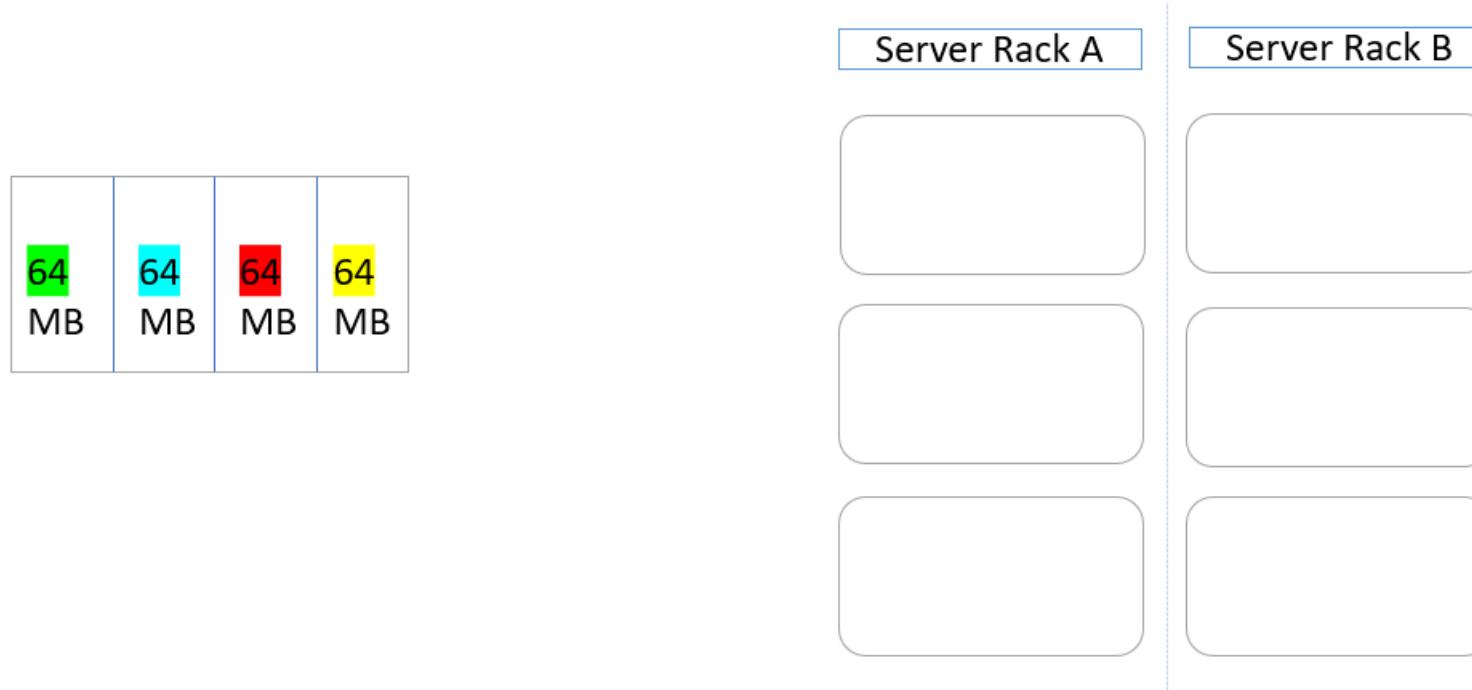
https://s3.amazonaws.com/files.dezyre.com/images/blog/Hadoop+Architecture+Explained-What+it+is+and+why+it+matters/Hadoop+Architecture_OpenSource.png

Hadoop HDFS Example

- Hadoop does two main things when a file is to be saved on HDFS:
 - It splits the file into chunks or blocks typically 64 to 128 mb data sizes;
 - Replicate each block of data (x3) and place them on a different data node;
- With the replication factor of 3, the cluster has the power to choose among 3 servers, which one to use for computation.
- In case any node goes down, there is another copy of the data that is available, achieving ***fault tolerance***;

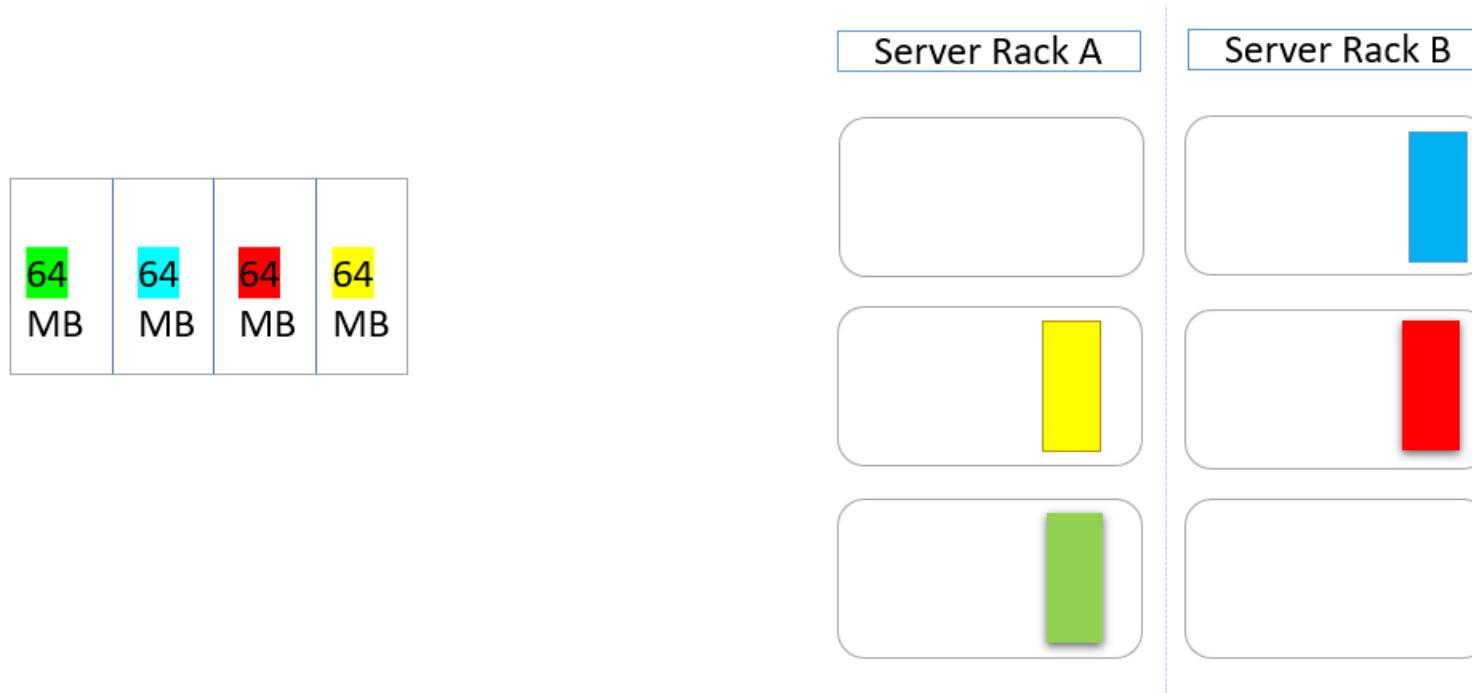
Hadoop HDFS Example

- It splits the file into chunks or blocks typically 64 to 128 MB in size



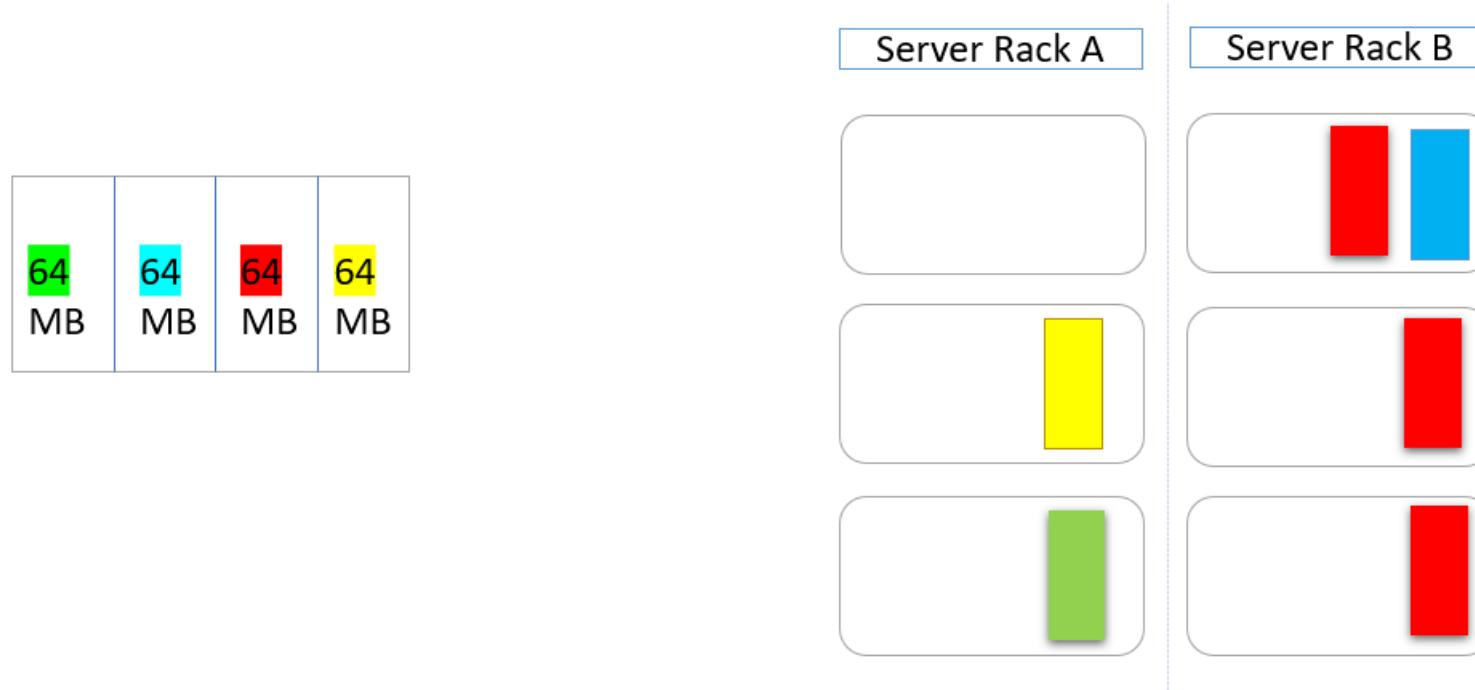
Hadoop HDFS Example

- Placing each block of data on a different data node



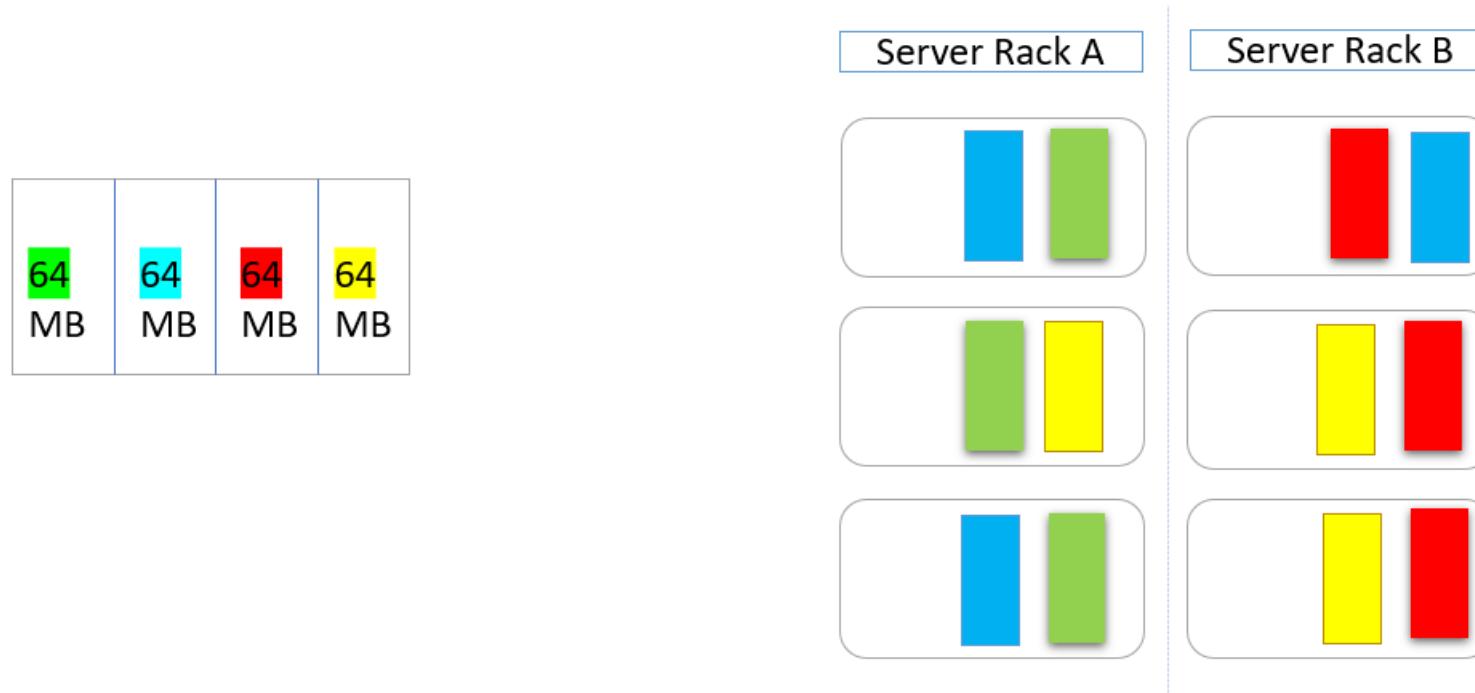
Hadoop HDFS Example

- Replicates each block to 3 nodes (data servers) by default



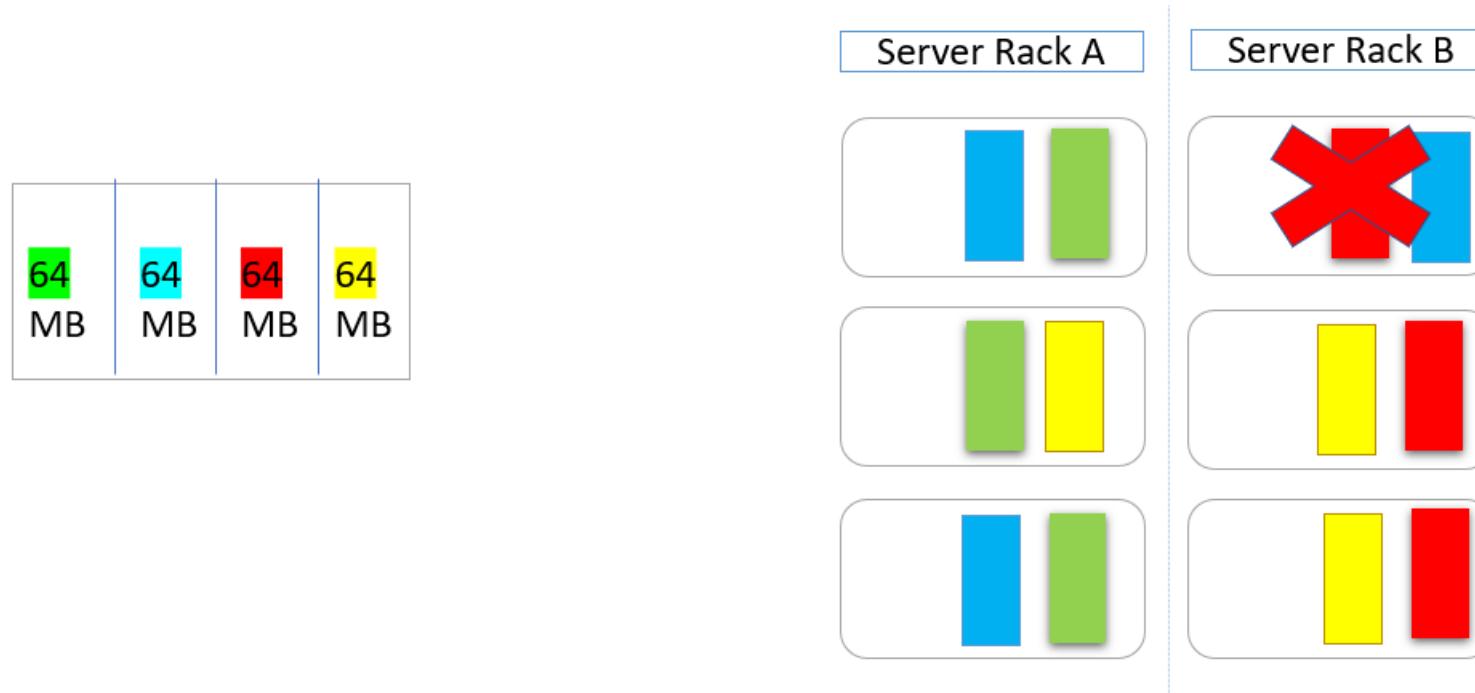
Hadoop HDFS Example

- Enables multiple machines to perform computation of data



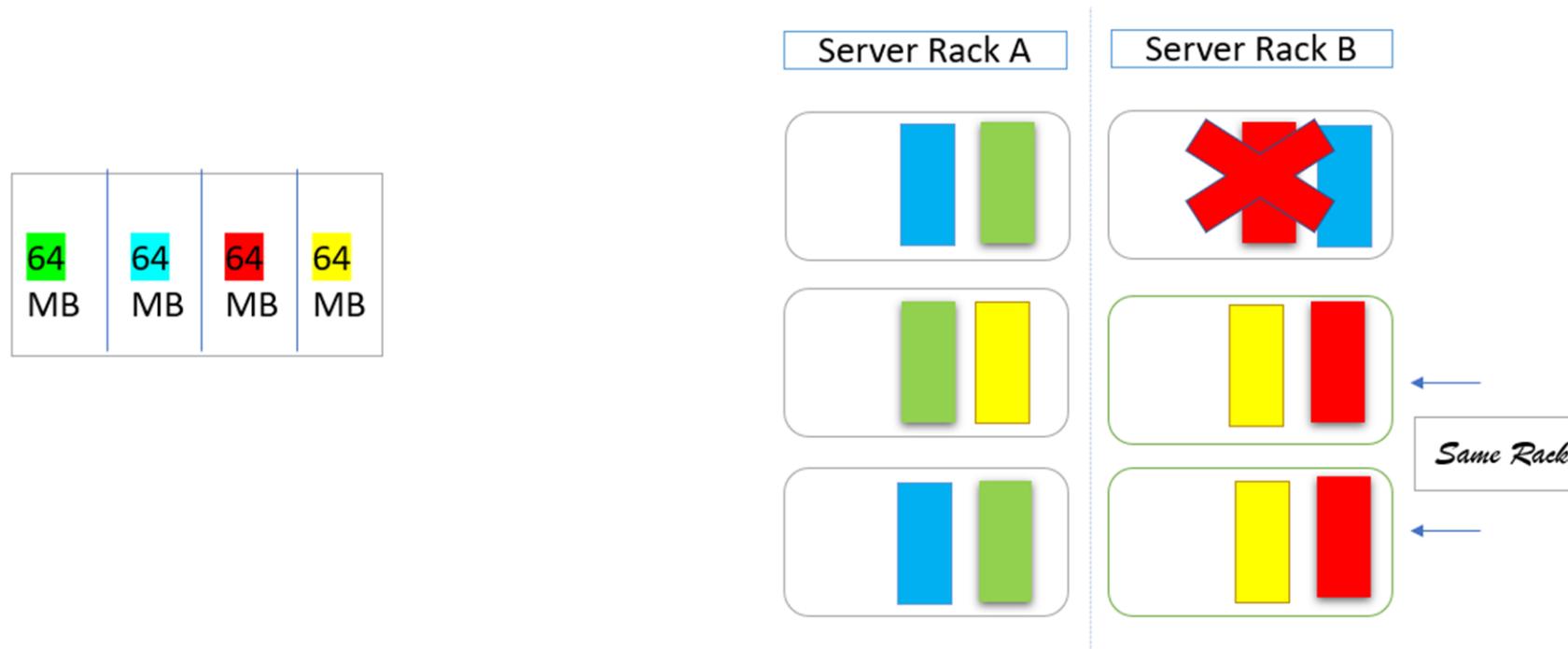
Hadoop HDFS Example

- Fault recovery: in case any node goes down, there's another copy of the data that is available, achieving fault tolerance



Hadoop HDFS Example

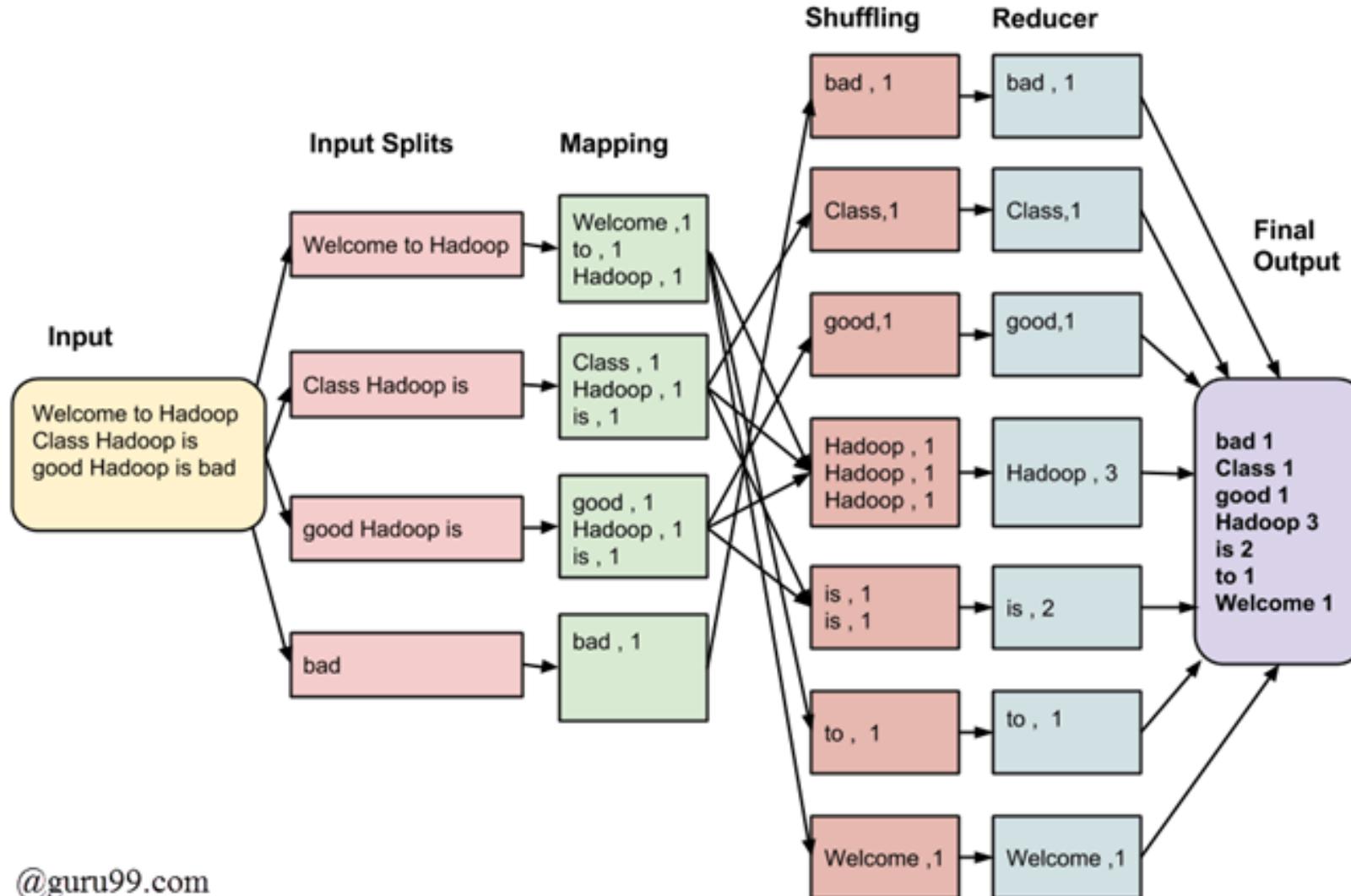
- It tries to serve data nodes on the same rack, to avoid network traffic between nodes, thus called: rack aware



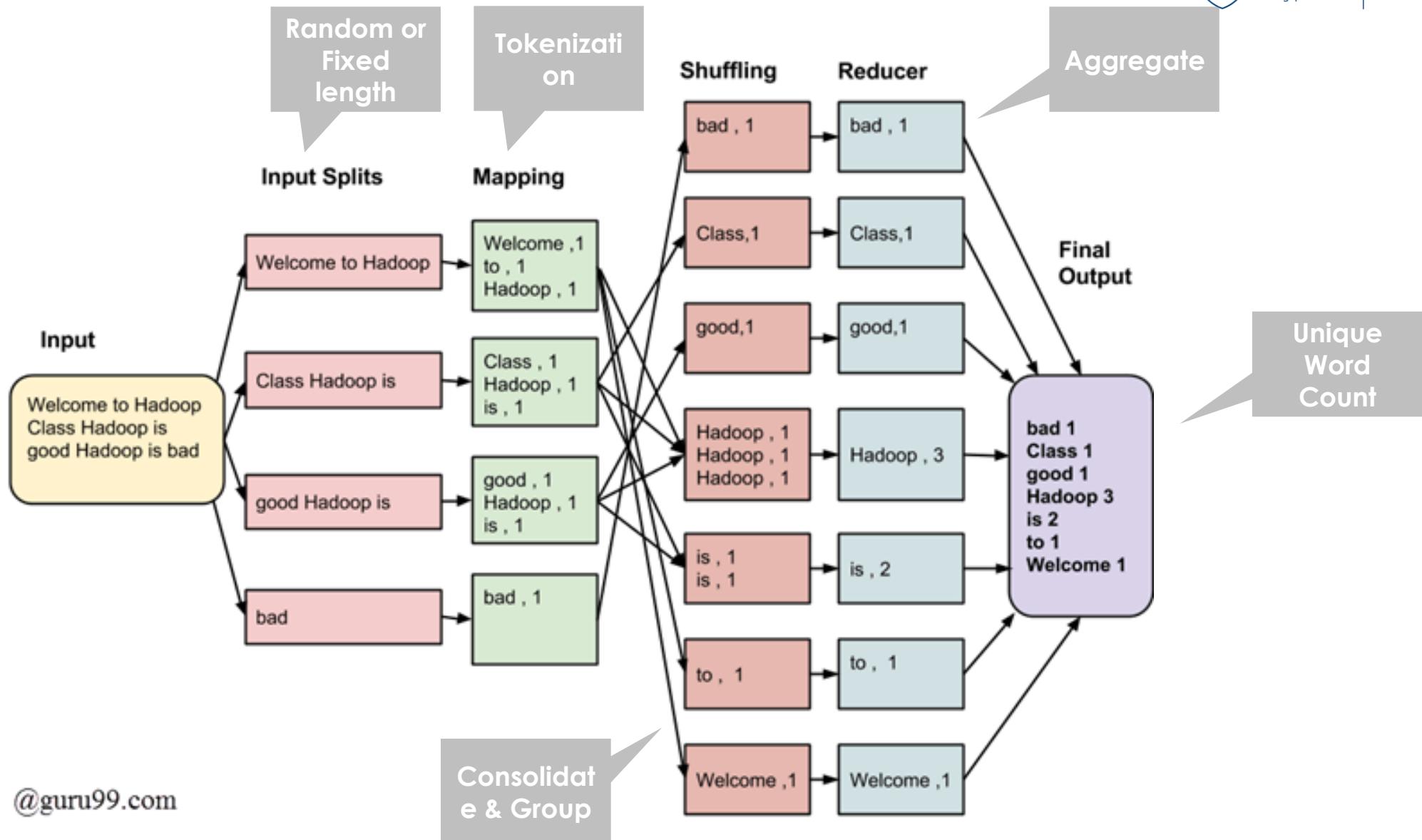
Map-Reduce Computation

- Map-Reduce is a software framework used for ***parallelly*** computing or processing data that we store in HDFS;
- Framework where (Hive Sql) queries are interpreted;
- A Map-Reduce job has three phases: (***divide & conquer method***)
 1. Map
 2. Shuffle & Sort
 3. Reduce

Map-Reduce Computation: Word Count

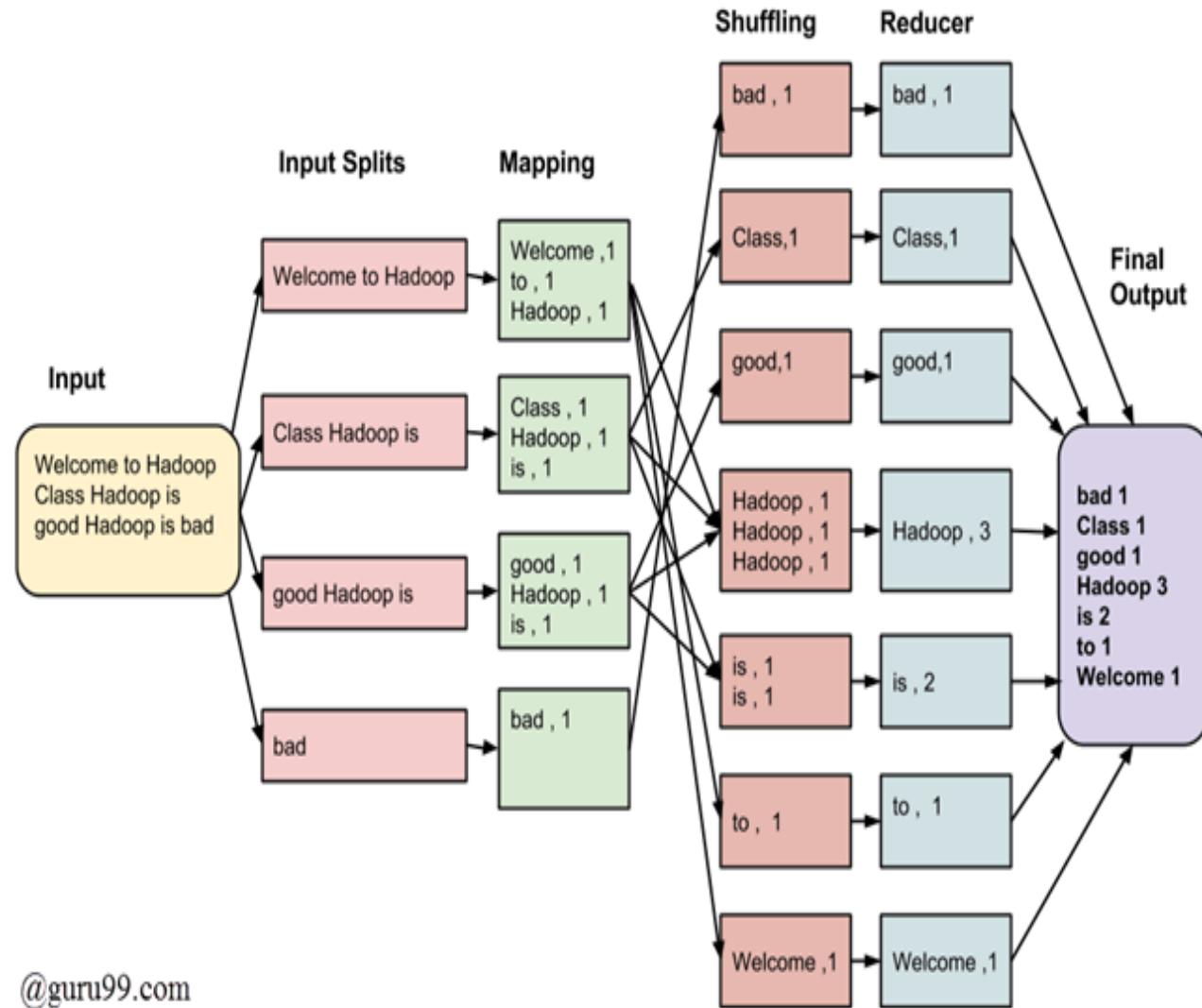


Map-Reduce Computation: Word Count



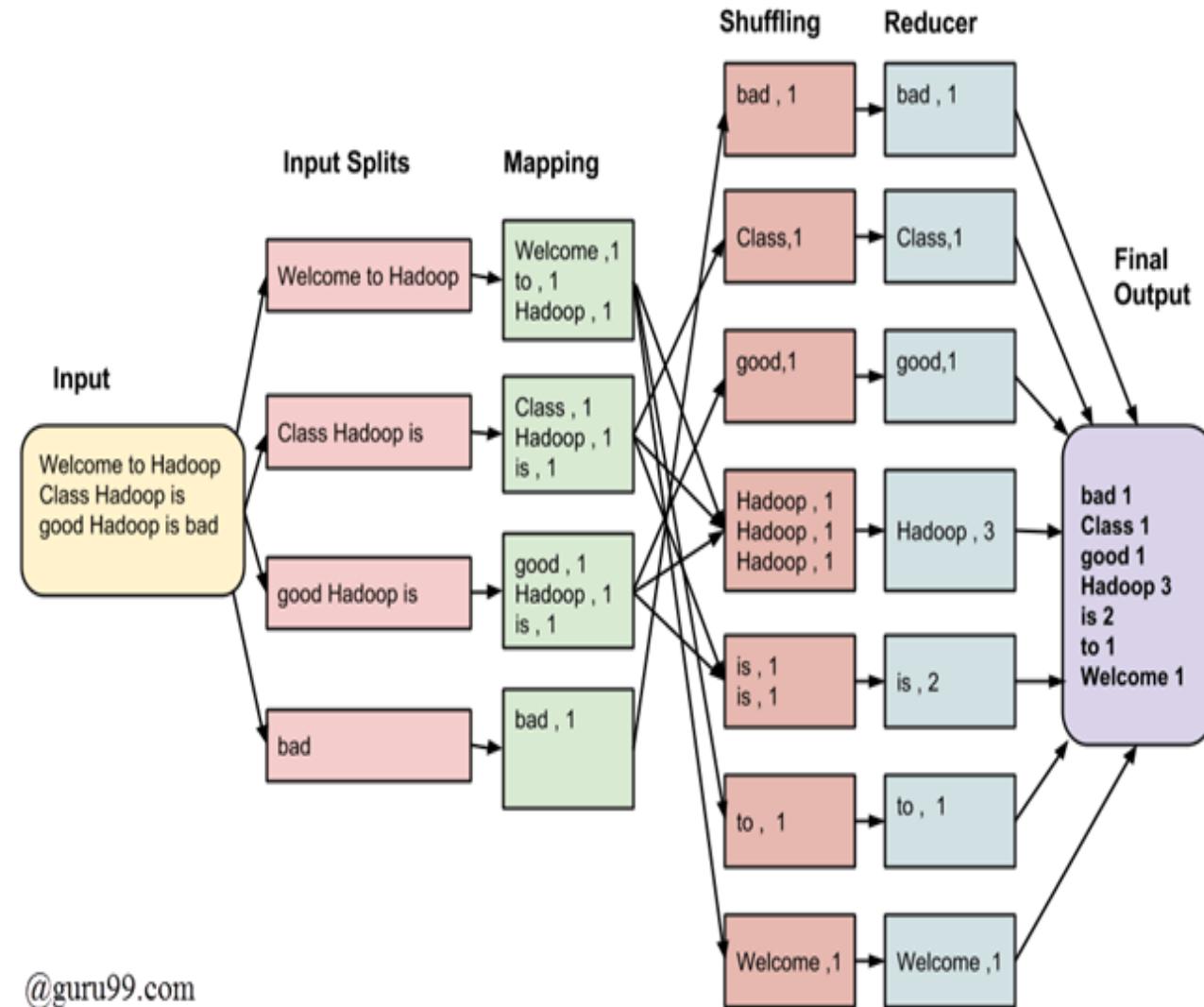
Map-Reduce Computation: Mappers

- A **map-reduce job** can use one or more mapper depending on the number of blocks the file spreads across.
- A **mapper** is assigned to each block. Here parallel distributed processing takes place given a file is split into blocks across multiple servers.
- Mappers take elements as a **key** and **value** and process them one at a time.



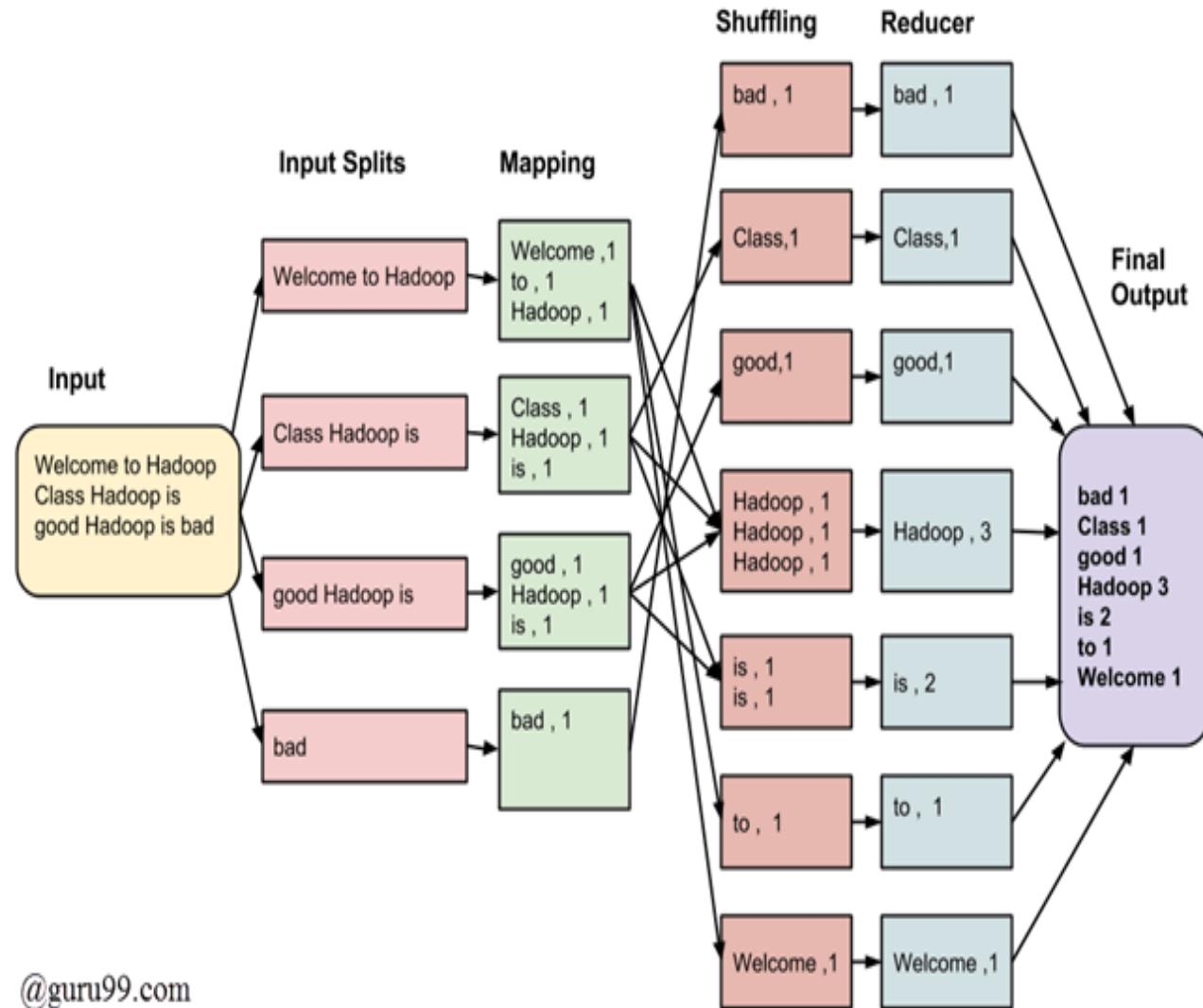
Map-Reduce Computation: Shuffle and Sort

- Maps are associated with the number of blocks of data required to read the input.
- Shuffle and sort** has the task of sorting the key-value pairs.
- It uses the hash value of the key and splits keys into buckets according to the number of reducers.



Map-Reduce Computation: Reducer

- Mapper makes sense of data (calculation) while **reducers** get the key value pair as an output from mappers and **aggregate** the results together.
- No reducers are not dependent on the number of mappers. It can be configurated in the map-reduce job.
- Output of reducer is sent to a file directory in HDFS.



4.2 Contemporary Reasoning Systems (Big Data)

4.2.1 Big Data Applications

4.2.2 Development Eco-systems (Hadoop & Spark)

4.2.3 Exercise

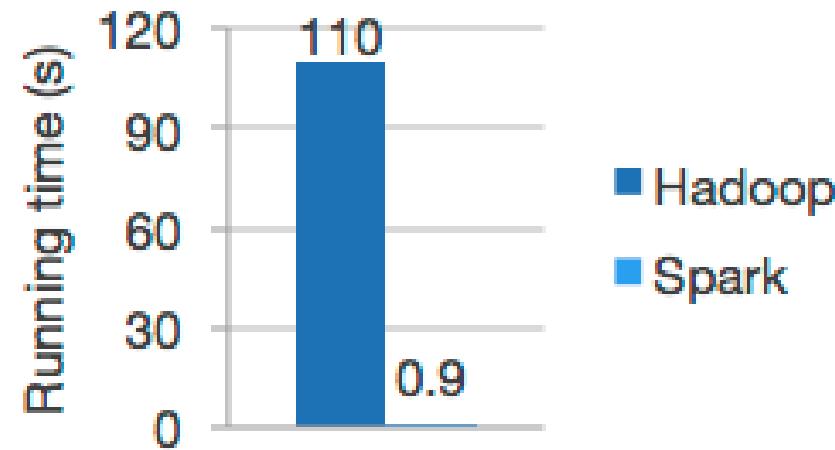
Spark Introduction

- **Apache Spark** is a unified analytics engine for large-scale data processing.
- It was the first unified analytics engine
- Spark simplifies working with data by supporting different languages (e.g. SQL, Java, Python, etc.)
- Spark brings data processing and analytics to one platform.



Why Spark?

- Speed: Spark runs workloads 100x faster than Hadoop.
- Apache Spark achieves high performance for both batch and streaming data, using ***in-memory computation*** compared with Hadoop's ***compute-on-disk***.



Spark Ease of Use

- Write applications quickly in Java, Scala, Python, R, and SQL;
- Spark's Python DataFrame API Read JSON files with automatic schema inference;

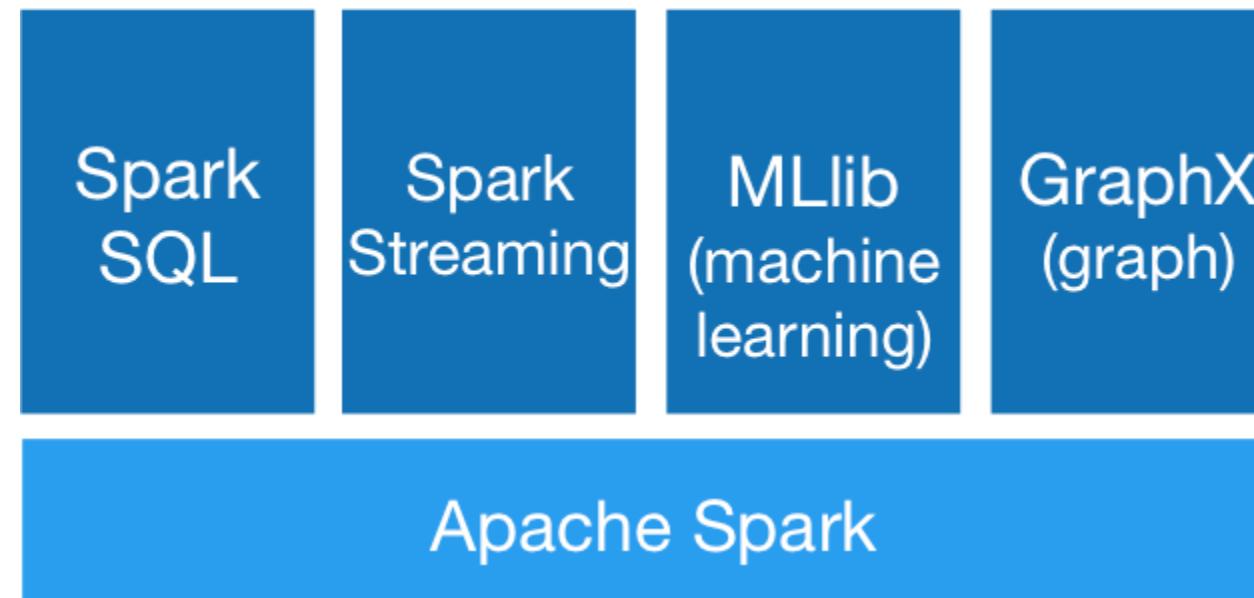
<https://www.codegrepper.com/code-examples/python/spark+read+json+schema>

```
# Read JSON files with automatic schema inference
```

```
df = spark.read.json("logs.json")
df.where("age > 21").select("name.first").show()
```

Spark Generality

- Spark powers a stack of libraries including SQL and DataFrames, MLlib for machine learning, GraphX, Spark Streaming, interactive dashboarding, and advanced analytics.



Spark Runs Everywhere

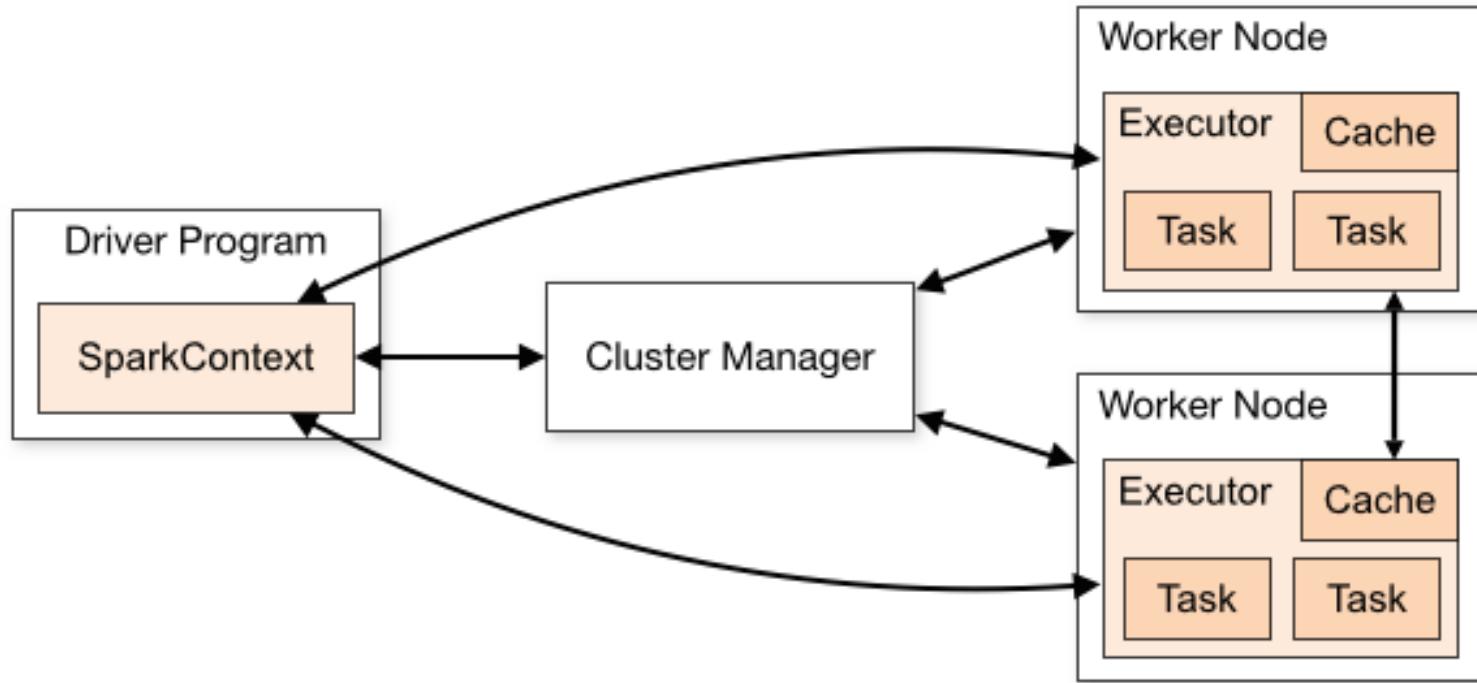
- Spark can run on Hadoop HDFS, Apache Mesos, Kubernetes, standalone, or in the cloud.
- It can access diverse data sources.



kubernetes

Spark Architecture

- Spark applications run as independent sets of processes on a cluster, coordinated by the **SparkContext** object in your main program (called the driver program).
- Specifically, to run on a cluster, the SparkContext can connect to several types of cluster managers (either Spark's own standalone cluster manager, Mesos or YARN), which allocate resources across applications.



1. **Once** connected, Spark acquires **executors** on worker nodes/servers in the cluster, which are processes that run computations and store data for your application.
2. **Next**, it sends your application code (Java JAR, Scala or Python passing via SparkContext) to the executors.
3. **Finally**, SparkContext sends tasks to the executors to run.

Spark Computation: Color Ball Count

- Imagine if you have to count the number of balls of each color in this bag;
- You are given the responsibility of performing this task, but you can call your friends (**executors**) over for help;
- How many friends will you call?
- Why?



Resilient Distributed Datasets:

- **Definition**: a collection of elements partitioned across the nodes of the cluster that can be operated on in parallel.
- Example: consider the previous example , if a friend leaves before the completion of your task?
- How would you compute the count of balls: Recompute the whole set or only the set that was taken up by that left individual?
- RDD makes this re-work easier, when fault happens.

Spark Summary

- RDD is ***distributed*** and stored in various clusters on your system;
- Each ***executor (data node server)*** works on their part of the data, the results are then aggregated and send back to the ***driver (name node server)***;
- Dataframe in Spark inherits RDD properties (***resilient and distributed***) and metadata;
- ***SparkSQL*** commands execute against dataframes (table alike);
- Spark is not a dataframe, it's a ***in-memory compute engine*** that can read from databases;
- Data is ***ephemeral***: you never lose your data even when spark is down;
- Dataframe aggregation takes ***shorter*** time;

4.2 Contemporary Reasoning Systems (Big Data)

4.2.1 Big Data Applications

4.2.2 Development Eco-systems (Hadoop & Spark)

4.2.3 Exercise

Analyze the following case study: Google Flu Trends (GFT)

In 2008, researchers from Google explored this potential, claiming that they could “nowcast” the flu based on people’s searches. The essential idea, published in a paper in *Nature*, was that when people are sick with the flu, many search for flu-related information on Google, providing almost instant signals of overall flu prevalence. The paper demonstrated that search data, if properly tuned to the flu tracking information from the Centers for Disease Control and Prevention, could produce accurate estimates of flu prevalence two weeks earlier than the CDC’s data — turning the digital refuse of people’s searches into potentially life-saving insights.

Google Flu Trends

- And then, GFT failed — and failed spectacularly — missing at the peak of the 2013 flu season by 140 percent. When Google quietly euthanized the GFT service, it turned the poster child of big data into the poster child of the foibles of big data.
- Google Flu Trends (GFT): <https://youtu.be/lEDt89eQ64o>

Answer questions:

- Why did GFT fail?
- What could Google have done differently?
- What does this teach us about analyzing Big data?

4.3 Building Machine Reasoning System [Workshop]

4.3.1 Big Data Machine Learning & Reasoning System

Special thanks to Geet Jethwani (A0215395B) for his contribution.

4.3.2 Workshop Submission

4.3 Building Machine Reasoning System [Workshop]

4.3.1 Big Data Machine Learning & Reasoning System

Special thanks to Geet Jethwani (A0215395B) for his contribution.

4.3.2 Workshop Submission

Learning Outcome

- Import data into big data platform.
- Apply SQL knowledge and RDBMS concepts over Spark.
- Apply Data Analysis techniques and Spark SQL queries on autism dataset
-
- Finally, we are to build an ASD diagnosis solution using Spark ML big data algorithm library.

Issues with Local Machine (read 2GB csv in Jupyter)

The screenshot shows a Jupyter Notebook interface with the following details:

- Title Bar:** A1234567X Donald Duck - Decision Tree ASD v001.ipynb
- Toolbar:** File, Edit, View, Insert, Cell, Kernel, Widgets, Help, Dead kernel (highlighted in red), Trusted, Python 3.
- Section Header:** 2. Import ASD Data
- In [5]:** !ls -l

```
total 24219692
-rwxrwx--- 1 root vboxsf 24423 Jan 27 11:05 'A1234567X Donald Duck - Decision Tree ASD v001.ipynb'
-rwxrwx--- 1 root vboxsf 20960909651 Feb 3 14:17 Toddler_Autism_dataset_BigData20GB.csv
-rwxrwx--- 1 root vboxsf 1795350682 Feb 3 15:43 Toddler_Autism_dataset_BigData20GB.zip
-rwxrwx--- 1 root vboxsf 2044598281 Feb 3 14:40 Toddler_Autism_dataset_BigData2GB.csv
-rwxrwx--- 1 root vboxsf 71201 Jan 18 09:26 'Toddler_Autism_dataset.csv'
```

- In []:** # Loading the dataset
ASD_data = pd.read_csv('./data.csv')
ASD_data = pd.read_csv('./Toddler_Autism_dataset_BigData2GB.csv')

Print the first 5 rows of the dataframe.
ASD_data.head()
- Kernel Status:** Dead kernel
- Kernel Restarting Dialog:** The kernel appears to have died. It will restart automatically. OK

- Databricks (import data/files);
- Spark SQL on Databricks (toddler autism case);
- SparkML (machine learning) library to build decision tree to detect autism;
- Anaconda with pyspark installed (if needed);

Additional Setup (if needed):

- Java 8 required on system;
- Set up JAVA_HOME and HADOOP_HOME System Variable;

[Workshop] Big Data Query

- Log in Community Edition log in portal:
<https://community.cloud.databricks.com/login.html>
- Import & Query toddler ASD data;

The screenshot shows the Databricks interface with the 'Data' tab selected in the sidebar. The main area displays the 'Create New Table' wizard. The 'Data source' tab is selected, showing options for 'Upload File', 'S3', 'DBFS', 'Other Data Sources', and 'Partner Integrations'. The 'DBFS Target Directory' field contains '/FileStore/tables/' with '(optional)' and a 'Select' button. A note states: 'Files uploaded to DBFS are accessible by everyone who has access to this workspace. [Learn more](#)'. Below this is a 'Files' section with a box labeled 'Drop files to upload, or browse.' On the left, the sidebar lists 'Databases' (with a warning about no cluster specified), 'Tables' (with a similar warning), 'Clusters', 'Jobs', and 'Search'. A red arrow points from the bottom right towards the 'Data' tab in the sidebar.

Import Toddler Autism dataset.csv via “Upload File”

© National University of Singapore

Data

Databases

Tables

You need to create a cluster to access tables

Create Table

Create New Table

Data source [?](#)

Upload File S3 DBFS Other Data Sources Partner Integrations

DBFS Target Directory [?](#)
/FileStore/tables/ (optional) [Select](#)

Files uploaded to DBFS are accessible by everyone who has access to this workspace. [Learn more](#)

Files [?](#)

Toddler Autism dataset BigData.csv
63.9 MB Remove file

✓ File uploaded to /FileStore/tables/Toddler_Autism_dataset_BigData.csv

[Create Table with UI](#) [Create Table in Notebook](#) [?](#)

File uploaded to
/FileStore/tables/Toddler_Autism_dataset_BigData.csv

Option 1: Create table with UI

© National University of Singapore

93

The screenshot shows the Databricks Data interface. On the left, there's a sidebar with icons for Home, Workspace, Recents, Data (selected), Clusters, Jobs, and Search. The main area has tabs for 'Databases' and 'Tables'. A 'Create Table' button is at the top. The 'Create New Table' wizard is open, showing a file upload step. A file named 'Toddler Autism dataset BigData.csv' (63.9 MB) has been uploaded to the DBFS target directory '/FileStore/tables/'. Below the file list, a success message says '✓ File uploaded to /FileStore/tables/Toddler_Autism_dataset_BigData.csv'. There are two buttons: 'Create Table with UI' (highlighted in blue) and 'Create Table in Notebook'. A dashed orange box highlights the 'Cluster' dropdown, which shows 'My Cluster' selected. The text 'Select a Cluster to Preview the Table' and 'Choose a cluster with which you will read and preview the data.' is displayed above the dropdown.

**Start a new cluster if needed.
Or use Option 2 to create table**

© National University of Singapore

94

Databases



Databases

Tables

⚠ You need to create a cluster to access tables

Create New Table

Files uploaded to DBFS are accessible by everyone who has access to this workspace. [Learn more](#)

Files

Toddler Autism dataset
BigData.csv

63.9 MB
[Remove file](#)

✓ File uploaded to /FileStore/tables/Toddler_Autism_dataset_BigData.csv

[Create Table with UI](#)[Create Table in Notebook](#)

Select a Cluster to Preview the Table

Choose a cluster with which you will read and preview the data.

Cluster

My Cluster

[Preview Table](#)

Specify Table Attributes

Specify the Table Name, Database and Schema to add this to the data UI for other users to access

Table Name

toddler_autism_dataset_big

Create in Database

default

File Type

CSV

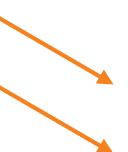
Column Delimiter

,

 First row is header Infer schema Multi-line[Create Table](#)[Create Table in Notebook](#)

Table Preview

A1	A2	A3	A4	A5	A6
INT	INT	INT	INT	INT	INT
0	0	0	0	0	0
1	1	0	0	0	1
1	0	0	0	0	0
1	1	1	1	1	1
1	1	0	1	1	1
1	1	0	0	1	1



The screenshot shows the Databricks Data interface. On the left, there's a sidebar with icons for Home, Workspace, Recents, Data (which is selected), Clusters, and Search. The main area is titled 'Create New Table'. It shows a file named 'Toddler Autism dataset BigData.csv' has been uploaded to the DBFS target directory '/FileStore/tables/'. There are two buttons at the bottom: 'Create Table with UI' (blue) and 'Create Table in Notebook' (gray). An orange arrow points to the 'Create Table in Notebook' button.

Option 2: Create table in Notebook

Home
Workspace
Recents
Data
Clusters
Jobs
Search

Users
issgz.nus@gmail.com
issgz.nus@gmail.com

2021-02-02 - DBFS Example
ISDemoNotebook
ISDemoNotebookASD

2021-02-02 - DBFS Example (Python)

Detached File Edit View: Standard Permissions Run All Clear Publish Comments Experiment Revision history

Cmd 1

Overview

This notebook will show you how to create and query a table or DataFrame that you uploaded to DBFS. DBFS is a Databricks File System that allows you to store data for querying inside of Databricks. This notebook assumes that you have a file already inside of DBFS that you would like to read from.

This notebook is written in Python so the default cell type is Python. However, you can use different languages by using the %LANGUAGE syntax. Python, Scala, SQL, and R are all supported.

Cmd 2

```
1 # File location and type
2 file_location = "/FileStore/tables/Toddler_Autism_dataset_BigData.csv"
3 file_type = "csv"
4
5 # CSV options
6 infer_schema = "true"
7 first_row_is_header = "true"
8 delimiter = ","
9
10 # The applied options are for CSV files. For other file types, these will be ignored.
11 df = spark.read.format(file_type) \
12 .option("inferSchema", infer_schema) \
13 .option("header", first_row_is_header) \
14 .option("sep", delimiter) \
15 .load(file_location)
16
17 display(df)
```

(3) Spark Jobs

df: pyspark.sql.dataframe.DataFrame = [A1: integer, A2: integer ... 15 more fields]

	A1	A2	A3	A4	A5	A6	A7	A8	A9	A10	Age_Mons	Sex
1	0	0	0	0	0	0	1	1	0	1	28	f
2	1	1	0	0	0	1	1	0	0	0	36	m
3	1	0	0	0	0	0	1	1	1	1	36	m
4	1	1	1	1	1	1	1	1	1	1	24	m
5	1	1	0	1	1	1	1	1	1	1	20	f
6	1	1	0	0	1	1	1	1	1	1	21	m
7	1	0	0	1	1	1	0	0	1	0	33	m
8	0	1	0	0	1	0	1	1	1	1	33	m

Showing the first 1000 rows.



Command took 19.18 seconds -- by issgz.nus@gmail.com at 2/2/2021, 10:12:40 AM on My Cluster

© National University of Singapore

Update:

CSV options

From 'false' to 'true'

Then execute Notebook

Workspace

Users issgz.nus@gmail.com

2021-02-02 - DBFS Example (Python)

Detached File Edit View Standard Permissions Run All Clear Publish Comments Experiment Revision history

Cmd 4

```
1 %sql
2
3 /* Query the created temp table in a SQL cell */
4
5 select * from `Toddler_Autism_dataset_BigData_csv`
```

(1) Spark Jobs

	A1	A2	A3	A4	A5	A6	A7	A8	A9	A10	Age_Mons	Sex
1	0	0	0	0	0	0	1	1	0	1	28	f
2	1	1	0	0	0	1	1	0	0	0	36	m
3	1	0	0	0	0	0	1	1	0	1	36	m
4	1	1	1	1	1	1	1	1	1	1	24	m
5	1	1	0	1	1	1	1	1	1	1	20	f
6	1	1	0	0	1	1	1	1	1	1	21	m
7	1	0	0	1	1	1	0	0	1	0	33	m
8	n	1	n	n	1	n	1	1	1	1	23	m

Showing the first 1000 rows.

Command took 0.59 seconds -- by issgz.nus@gmail.com at 2/2/2021, 10:27:57 AM on My Cluster

Cmd 5

```
1 %sql
2 select Ethnicity, count(*) from `Toddler_Autism_dataset_BigData_csv` Group by Ethnicity;
```

(2) Spark Jobs

Ethnicity	count(1)
south asian	59252
Native Indian	2960
mixed	7906
asian	295295
Latino	25675
Pacifica	7901
Others	34570
black	52339

Showing all 11 rows.

Command took 4.62 seconds -- by issgz.nus@gmail.com at 2/2/2021, 10:28:00 AM on My Cluster

Construct new queries

© National University of Singapore

Home
Workspace
Recents
Data
Clusters
Jobs
Search

Users
issgz.nus@gmail.com

issgz.nus@gmail.com
2021-02-02 - DBFS Example
ISDemoNotebook
ISDemoNotebookASD

2021-02-02 - DBFS Example (Python)

Detached File Edit View Standard Permissions Run All Clear Publish Comments Experiment Revision history

Command took 0.59 seconds -- by issgz.nus@gmail.com at 2/2/2021, 10:27:57 AM on My Cluster

Cmd 5

```
1 %sql
2 select Ethnicity, count(*) from `Toddler_Autism_dataset_BigData_csv` Group by Ethnicity;
```

(2) Spark Jobs

Ethnicity	Count (approx.)
south asian	60k
Native Indian	5k
mixed	5k
asian	290k
Latino	25k
Pacifica	5k
Others	40k
black	55k
middle eastern	185k
Hispanic	45k
White European	320k

Plot Options...

Command took 4.62 seconds -- by issgz.nus@gmail.com at 2/2/2021, 10:28:00 AM on My Cluster

Cmd 6

```
1 # With this registered as a temp view, it will only be available to this particular notebook. If you'd like other users to be able to query this table,
2 # you can also create a table from the DataFrame.
3 # Once saved, this table will persist across cluster restarts as well as allow various users across different notebooks to query this data.
4 #
5 permanent_table_name = "Toddler_Autism_dataset_BigData_csv"
6 #
7 # df.write.format("parquet").saveAsTable(permanent_table_name)
```

Shift+Enter to run [shortcuts](#)

The screenshot shows the Databricks Data browser interface. On the left, there's a sidebar with icons for Home, Workspace, Recents, Data (which is selected), Clusters, Jobs, and Search. The main area has tabs for 'Databases' and 'Tables'. Under 'Tables', it says 'No tables'. A 'Create Table' button is at the top. Below it, a 'Create New Table' dialog is open, showing 'DBFS' as the data source. It displays a file tree under 'FileStore': 'import-stage' and 'tables'. Inside 'tables' are two files: 'Toddler_Autism_dataset_BigData....csv' and 'Toddler_Autism_dataset_BigData...'. An orange arrow points from the text 'Location of uploaded files' down to the 'tables' folder in the file tree. At the bottom of the dialog are buttons for 'Create Table with UI' (highlighted in blue) and 'Create Table in Notebook'.

Location of uploaded files

You can delete file in a notebook cell by:
dbutils.fs.rm("/FileStore/tables/FileName.csv")

© National University of Singapore 100

[Workshop] ASD Diagnosis in PySpark

- Open the **A1234567X Donald Duck – Spark DTTree** PySpark notebook file and complete the workshop;
- Execute notebook to be familiar with PySpark scripts;
- Interpret decision tree results, e.g. which feature is most important on toddler ASD detection: different ethnicity? age?

[Optional] Enhancements:

- Evaluate decision tree using confusion matrix;
- Dummy / One-hot encoding of categorical features;
- Feature engineering and selection, e.g. whether A1 – A10 is useful?

```
1 # Some functions that convert our CSV input data into numerical
2 # features for each job candidate
3 def binary(YN): # To be completed by Student
4     if (YN == 'yes' or YN == 'Yes'):
5         return 1
6     else:
7         return 0
8
9 def mapSex(degree): # To be completed by Student
10    if (degree == 'f'):
11        return 1
12    else:
13        return 0
14
15 def mapEthnicity(ethnic): # To be completed by Student
16    if(ethnic == 'asian'):
17        return 1
18    else:
19        return 0
20
21 def mapAssessor(assessor): # To be completed by Student
22    if(assessor == 'Health Care Professional'):
23        return 1
24    else:
25        return 0
26
```

Workspace

Users issgz.nus@gmail.com

ISDemoNotebook
ISDemoNotebookASD

Create
Clone
Import
Export
Permissions
Copy Link Address

Welcome to  databricks

Drop files or [click to browse](#)

Create a Blank Notebook

Create a notebook to start querying, visualizing, and modeling your data.

Import Notebooks

Import from: File URL

A1234567X ✓
Donald Duck –
Spark DTree
v001.ipynb
5.9 KB

Accepted formats: .dbc, .scala, .py, .sql, .r, .ipynb, .Rmd, .html
(To import a library, such as a jar or egg, [click here](#))

Cancel Import

What's new in v3.37

[View latest release notes](#)

Open/Import notebook

© National University of Singapore

103

Workspace

Users issgz.nus@gmail.com

Home issgz.nus@gmail.com

Recent 2021-02-02 - DBFS Example

A1234567X Donald Duck – Sp...

ISDemoNotebook

ISDemoNotebookASD

Clusters

Jobs

Search

A1234567X Donald Duck – Spark DTree v001 (Python)

My Cluster

Attached cluster:

- My Cluster
15.25 GB | 2 Cores | DBR 7.5 | Spark 3.0.1 | Scala 2.12
Detach | Restart Cluster | Detach & Re-attach | Spark UI | Driver Logs | ⚙

Detach & Attach:

- My Cluster
DBR 7.5 | Spark 3.0.1 | Scala 2.12
- My Cluster
DBR 7.5 | Spark 3.0.1 | Scala 2.12

Command took 0.68 seconds -- by issgz.nus@gmail.com at 2/2/2021, 10:37:59 AM on My Cluster

Cmd 3

```
1 import matplotlib.pyplot as plt
2 import numpy as np
3 import pandas as pd
4 from numpy import array
5 from sklearn.datasets import load_iris
6 from sklearn.datasets import load_breast_cancer
7 from sklearn.tree import DecisionTreeClassifier
8 from sklearn.ensemble import RandomForestClassifier
9 from sklearn.model_selection import train_test_split
10 from sklearn import tree
```

1

Command took 1.73 seconds -- by issgz.nus@gmail.com at 2/2/2021, 10:46:28 AM on My Cluster

Cmd 4

```
1 from pyspark.mllib.regression import LabeledPoint
2 from pyspark.mllib.tree import DecisionTree
3 from pyspark import SparkConf, SparkContext
4 from pyspark.ml.feature import StringIndexer
5 from pyspark.ml import Pipeline
6 from pyspark.ml.classification import DecisionTreeClassifier
```

Command took 0.28 seconds -- by issgz.nus@gmail.com at 2/2/2021, 10:46:32 AM on My Cluster

Cmd 5

```
1 # Some functions that convert our CSV input data into numerical
2 # features for each job candidate
3 def binary(YN): # To be completed by Student
4     if (YN == 'yes' or YN == 'Yes'):
5         return 1
6     else:
7         return 0
8
9 def mapSex(degree): # To be completed by Student
10    if (degree == 'f'):
```

Connect to active cluster

© National University of Singapore 104

Home
Workspace
Recents
Data
Clusters
Jobs
Search

Users

issgz.nus@gmail.com

2021-02-02 - DBFS Example

A1234567X Donald Duck – Sp...

ISDemoNotebook

ISDemoNotebookASD

Detached

File Edit View: Standard Permissions Run All Clear Publish Comments Experiment Revision history

```
1 # Create a test candidate
2 # testCandidates = [ array([36, 0, 2, 1, 0]) ] # test different candidates
3 testCandidates = [ array([1, 1, 0, 0, 0, 1, 1, 0, 0, 36, 0, 2, 1, 0, 1]) ] # test different candidates
4 testData = sc.parallelize(testCandidates)
```

Command took 0.06 seconds -- by issgz.nus@gmail.com at 2/2/2021, 10:59:08 AM on My Cluster

Cmd 12

```
1 # Now get predictions (Note, you could separate
2 # the source data into a training set and a test set while tuning
3 # parameters and measure accuracy as you go!)
4 predictions = model.predict(testData)
5 print('ASD prediction:')
6 results = predictions.collect()
7 for result in results:
8     print(result)
```

(1) Spark Jobs

ASD prediction:

1.0

Command took 0.30 seconds -- by issgz.nus@gmail.com at 2/2/2021, 10:59:10 AM on My Cluster

Cmd 13

```
1 # We can also print out the decision tree itself:
2 print('Learned classification tree model:')
3 print(model.toDebugString())
```

```
Learned classification tree model:
DecisionTreeModel classifier of depth 5 with 35 nodes
If (feature 8 <= 0.5)
    If (feature 6 <= 0.5)
        If (feature 7 <= 0.5)
            If (feature 3 <= 0.5)
                Predict: 0.0
            Else (feature 3 > 0.5)
                If (feature 9 <= 0.5)
                    Predict: 0.0
                Else (feature 9 > 0.5)
                    Predict: 1.0
            Else (feature 7 > 0.5)
                If (feature 5 <= 0.5)
                    If (feature 1 <= 0.5)
                        Predict: 0.0
                    Else (feature 1 > 0.5)
                        Predict: 1.0
                Else (feature 5 > 0.5)
                    Predict: 1.0
            Else (feature 6 > 0.5)
                Predict: 1.0
        Else (feature 6 > 0.5)
            Predict: 0.0
```

← Knowledge in decision tree

- **Naming convention: StudentID YourFullName, e.g.
A1234567X Donald Duck – sln – Spark DTree.ipynb/py/zip**
- **Use zip to a single file, then rename, if you plan to submit multiple files.**

4.3 Building Machine Reasoning System [Workshop]

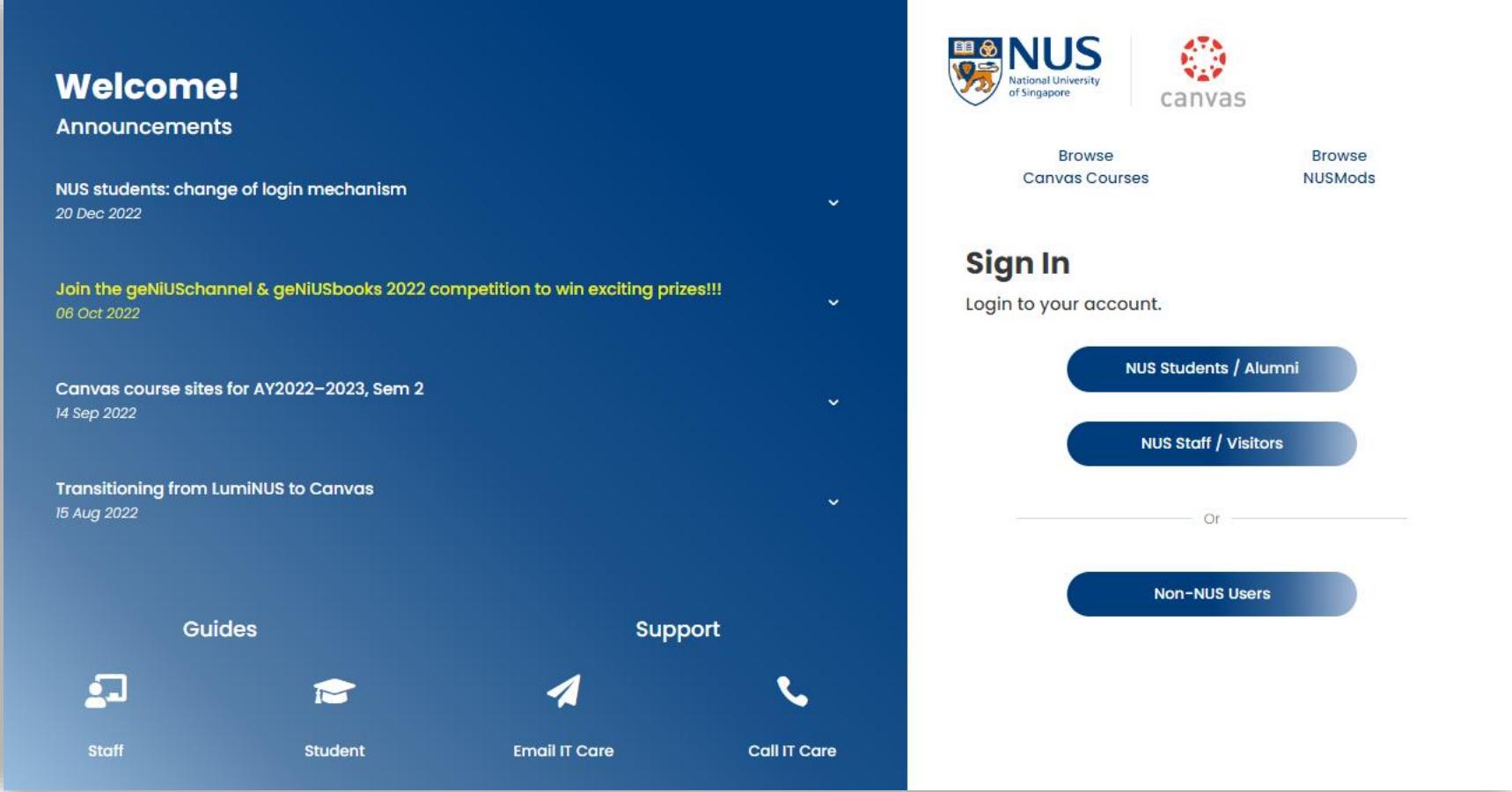
4.3.1 Big Data Machine Learning & Reasoning System

Special thanks to Geet Jethwani (A0215395B) for his contribution.

4.3.2 Workshop Submission

Workshop Submission

- **Naming convention: StudentID YourFullName**
- **Use zip to a single file, then rename, if you plan to submit multiple files.**



The image shows two side-by-side screenshots. The left screenshot is the NUS Canvas LMS homepage, featuring a dark blue header with "Welcome!" and "Announcements". It lists several announcements, including one about a competition and another about transitioning from LumiNUS to Canvas. The right screenshot is the NUS IT Support page, showing the "Sign In" section with three login options: "NUS Students / Alumni", "NUS Staff / Visitors", and "Non-NUS Users".

Welcome!
Announcements

NUS students: change of login mechanism
20 Dec 2022

Join the geNiUSchannel & geNiUSbooks 2022 competition to win exciting prizes!!!
06 Oct 2022

Canvas course sites for AY2022–2023, Sem 2
14 Sep 2022

Transitioning from LumiNUS to Canvas
15 Aug 2022

Guides

Staff Student

Support

Email IT Care Call IT Care

NUS
National University
of Singapore

canvas

Browse Canvas Courses Browse NUSMods

Sign In
Login to your account.

NUS Students / Alumni
NUS Staff / Visitors
Non-NUS Users
Or

4.4 In-class Assessment

4.4.1 In-class Assessment

Machine Reasoning Summary

What we have learnt:

Day 1

- 1.1 Machine Reasoning Overview
- 1.2 Reasoning System Architectures
- 1.3 Rule/Process Reasoning System **Workshop**

Day 2

- 2.1 Machine Reasoning Enabler: Knowledge Representation
- 2.2 Machine Reasoning Enabler: Knowledge Acquisition
- 2.3 Knowledge Representation and Acquisition/Discovery **Workshop**

Day 3

- 3.1 Deductive Reasoning by Logical Inference
- 3.2 Reasoning under Uncertainty
- 3.3 Deductive Reasoning (under Uncertainty) **Workshop**

Day 4

- 4.1 Knowledge Discovery by Machine Learning (Big Data)
- 4.2 Contemporary Reasoning Systems (Big Data)
- 4.3 Building Machine Reasoning System **Workshop**

Start

How can I apply AI to my business problem?

Define & Quantify the business problem/opportunity {Goal, KPI}, e.g. {Increase customer satisfaction, Higher rating}; {Get rich fast, More \$ Less time}

Anyone (on Earth) knows how to solve?
Or known existing/manual approaches?



Is **Reasoning/Inference Task**:
Use existing knowledge to solve problem at scale and speed.
(decision automation; optimization)

End

Is **Learning Task**:
Acquire new knowledge to solve the problem.

Has relevant data collection?



Is a special learning task:
Reinforcement Learning Task:
Use generate & test approach.

Can simulate /generate relevant data?
Or conduct experiment?



Use reinforcement learning algorithms; simulation techniques to **acquire** knowledge, e.g. best series of actions to achieve goal.

Currently not solvable by AI (even human)

Future Work (Homework)

Homework Topic	Priority	Example Tools	Reference
From Human Intelligence to Machine Intelligence	1		https://youtu.be/HQUxSi52Ujk
Computer Programming	2	Python Anaconda	https://www.anaconda.com/
Machine Memory: Database & SQL	3	SQLite; MySQL	https://github.com/agarcialeon/awesome-database
Text Processing	4	Python lib: SpaCy	https://github.com/keon/awesome-nlp
Optimization	5	Google OR-Tools; KIE-OptaPlanner	https://developers.google.com/optimization
Cloud Computing	6	GCP; AWS	https://github.com/tmrts/awesome-cloud-computing
Machine Reasoning	7	Semantic Reasoner; KIE-Drools; KIE-jBPM;	https://github.com/semantalytics/awesome-semantic-web
Speech Virtual Assistant	8	MyCrost; RASA;	https://mycroft.ai/
Information Retrieval / Search Engine	9	Chatter-Bot; Lucene; Lemur;	https://github.com/harpribot/awesome-information-retrieval
Knowledge Graph & GraphDB	10	protégé; grakn.ai	https://github.com/totogo/awesome-knowledge-graph
Recommender	11		https://github.com/jihoo-kim/awesome-RecSys
Big Data	12	Spark; Hadoop;	https://github.com/onurakpolat/awesome-bigdata
Machine Learning	13	Scikit-Learn; Orange3	https://github.com/josephmisiti/awesome-machine-learning
Object/Face Detection	14		https://github.com/amusi/awesome-object-detection
Time Series	15	Facebook Prophet	https://facebook.github.io/prophet/
Full Stack	16		https://github.com/kevindeasis/awesome-fullstack

END OF NOTES

APPENDICES



INTRODUCING Neo4j Aura Enterprise

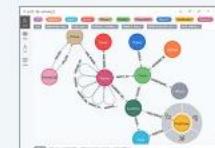
Neo4j's fully managed cloud service – the zero-admin, always-on graph database – now for the global enterprise

[Learn More](#)

Free O'Reilly Ebook

[Graph Algorithms: Examples in Spark and Neo4j](#)

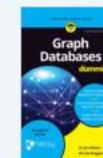
Sample code and tips for over 20 practical algorithms. Find vulnerabilities, detect communities, improve machine learning.

[Read Now](#)

No Download Required

[Get Started Quickly with the Neo4j Sandbox](#)

Start using Neo4j in seconds, with built-in guides and datasets for popular use cases. No experience necessary.

[Try It Now](#)

Neo4j Special Edition
[Graph Databases for Dummies](#)

Learn the basics of graph database technology, from building a data model to deploying a graph-powered application.

[Download Now](#)



Book

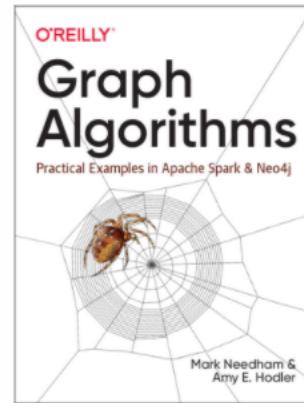
Graph Algorithms: Practical Examples in Apache Spark and Neo4j

By Mark Needham & Amy Hodler

Published by O'Reilly Media

Print Length: 300 pages

Available Formats: PDF - EN US, iBooks, Kindle



Register now for your copy of the O'Reilly book, *Graph Algorithms: Practical Examples in Apache Spark and Neo4j* by Mark Needham and Amy E. Hodler. You'll receive a link to download an electronic version as soon as it's available this spring.

Whether you are trying to build dynamic network models or forecast real-world behavior, this book demonstrates how graph algorithms deliver value – from finding vulnerabilities and bottlenecks to detecting communities and improving machine learning predictions.

We walk you through **hands-on examples** of how to use graph algorithms in Apache Spark and Neo4j. We include sample code and tips for over 20 practical graph algorithms that cover importance through centrality, community detection and optimal pathfinding. Read this book to:

- Learn how graph analytics vary from conventional statistical analysis
- Understand how classic graph algorithms work and how they are applied
- Dive into popular algorithms like PageRank, Label Propagation and Louvain to find out how subtle parameters impact results

Register to Download O'Reilly's Graph Algorithms for Free!

First Name

Last Name

Business Email

Company Name

Country

Get My Free Copy

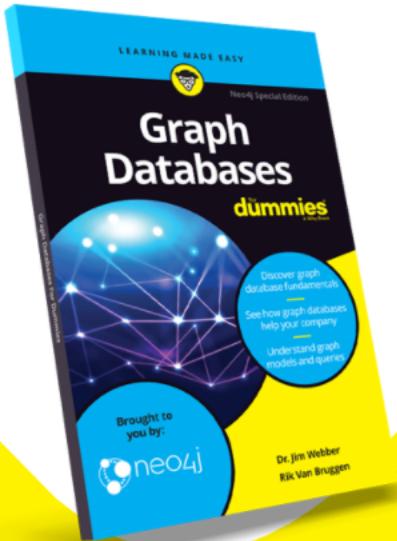
The information you provide will be used in accordance with the terms of our [privacy policy](#).



FREE BOOK

Graph Databases For Dummies

Graph Databases For Dummies, a Neo4j Special Edition, introduces you to the basics of graph database technology from building a rich graph data model to deploying your first graph-powered application.



Download Your Free Copy

First Name Last Name

Business Email

Company Name

Country

The information you provide will be used in accordance with the terms of our [privacy policy](#).



Learn graph database
fundamentals



Identify when to use graph
technology



Deploy your first graph
database

END OF APPENDICES