**Program: MSc of Data Science**
**Module: Big Data Tools and Techniques**

Week 4 – Part 1

Apache, Spark, Databricks?

2025

# Ground Rules

1. Choose a quiet place to attend the class and please concentrate during the lecture

2. Put your questions in Padlet (not Teams' chat box) and I will review them in the due time (Padlet link is in the Bb, week 4, Lecture)

3. Turn off your mic during the lecture

4. We will have 5 mins break after the first hour of the lecture (please remind me)

5. Jisc code will be shared during the break time

SCHOOL OF
SCIENCE, ENGINEERING
& ENVIRONMENT

# What is Apache?

SCHOOL OF
SCIENCE, ENGINEERING
& ENVIRONMENT

# Apache is a name for a multidimensional project. The dimensions like:

**Apache** Spark, **Apache** Hadoop, **Apache** Hive, **Apache** Avro, **Apache** Ambari, **Apache** Kafka, …

# So, what is Apache exactly?
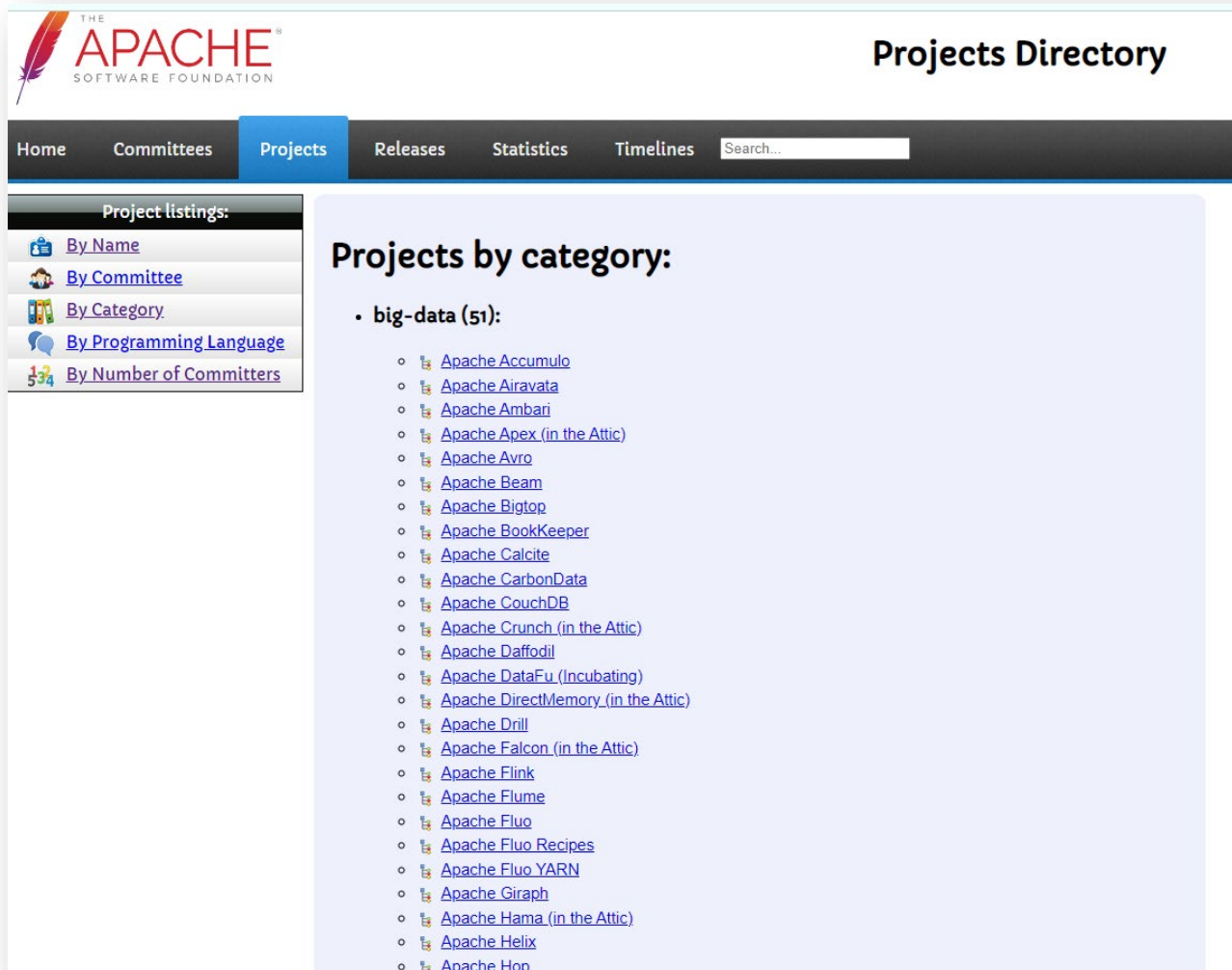
SCHOOL OF
SCIENCE, ENGINEERING
& ENVIRONMENT

# Apache project

Founded in 1999, the ASF functions as a charitable organisation in the United States, relying on contributions from individuals and corporate sponsors for funding. An all-volunteer board manages over 300 prominent Open-Source projects, among them the <span style="color:red">Apache HTTP Server</span>, recognized as the foremost web server software globally.



https://www.apache.org/

# Apache Projects







https://projects.apache.org/projects.html?category

# Why has been named "Apache"

A group of people calling themselves the "Apache Group" created the foundation in 1999. They had come together several years earlier, **to continue to support and maintain the HTTPD web server** written by the National Center for Supercomputing Applications (NCSA) at the University of Illinois.

That server was freely available, came with its source code, its license allowed very open modification and redistribution, but the original developers lost interest in that project and moved onto something else, leaving users with no support for the application.

Some of those users started to exchange fixes (called "patches") and information on how to prevent problems and improve the existing software. Brian Behlendorf created a mailing list on which those users could collaborate to fix, maintain and improve that software.



https://www.apache.org/foundation/how-it-works.html#what

SCHOOL OF
SCIENCE, ENGINEERING
& ENVIRONMENT

# Why has been named "Apache"

Co-founder <u>Brian Behlendorf</u> states how the name Apache was chosen:

"I had just seen a documentary about <span style="color:red">Geronimo</span> and the <span style="color:red">last days of a Native American tribe called Apaches</span>, right, who succumbed to the invasion from the West, from the United States, and **they were the <span style="color:red">last tribe to give up</span> their territory and for me that almost romantically represented what I felt we were doing with this <span style="color:red">web-server project</span>** …"

SCHOOL OF SCIENCE, ENGINEERING & ENVIRONMENT

# What is Apache Spark?

# Apache Spark?

Imagine you have a giant **puzzle with millions of pieces**, each representing a piece of information. Analysing all that data by hand would be impossible! That's where Spark comes in.



**Think of Spark as a superpowered puzzle solver**. It's a tool that can:

➢ Work on a single computer or many computers working together (like a team, depending on the size of the puzzle.)

➢ Hold some of the data in its memory (like remembering where you saw a specific colour or shape, making it much faster than searching through the whole puzzle each time.)

➢ Break the data down into smaller parts, allowing it to analyse different sections simultaneously (like having multiple people working on different areas of the puzzle at once.)

➢ Do different things with the data pieces (like sorting them by colour, shape, or size,
➢ or even using them to build new things, like a picture or a story.)

SCHOOL OF
**SCIENCE, ENGINEERING
& ENVIRONMENT**

# Apache Spark?

**Is Spark a <u>programming language</u>? or a <u>platform</u>? or an <u>architecture</u>? what is it exactly?**

SCHOOL OF
**SCIENCE, ENGINEERING
& ENVIRONMENT**

# Apache Spark?

**Spark isn't exactly one of those categories on its own.** It's more like a combination of them, serving as a powerful tool in the big data world:

**It's not a programming language:** While you can interact with Spark using various languages like Python, Scala, and Java, Spark itself isn't a language you directly write code in.

**It's not just a platform:** Although platforms like Databricks offer ways to use and manage Spark, Spark exists independently and can be implemented in various environments.

**It's not just an architecture:** While Spark utilizes a distributed architecture across clusters or single machines, the core is more than just its structure.

SCHOOL OF
SCIENCE, ENGINEERING
& ENVIRONMENT

# Apache Spark?

**Spark is a unified analytics engine**

It's a software framework specifically designed to perform large-scale data processing efficiently across various domains like data engineering, data science, and machine learning.

SCHOOL OF
SCIENCE, ENGINEERING
& ENVIRONMENT

# Apache Spark



**What is Apache Spark?**
Apache Spark is a multi-language engine for executing data engineering, data science, and machine learning even on single-node machines or clusters.

**RDD is the main concept of the Apache Spark**

The primary abstraction of **the Spark is the concept of RDD, which Spark uses to achieve faster and efficient MapReduce operations**.

Resilient Distributed Dataset (**RDD**) is the fundamental data structure of Spark.

https://spark.apache.org/

SCHOOL OF
SCIENCE, ENGINEERING
& ENVIRONMENT

# Apache Spark Structure

# What is Databricks?

SCHOOL OF
SCIENCE, ENGINEERING
& ENVIRONMENT

# Comparing Apache Spark™ and Databricks



Apache Spark capabilities provide speed, ease of use and breadth of use benefits and include APIs supporting a range of use cases:

- Data integration and ETL
- Interactive analytics
- Machine learning and advanced analytics
- Real-time data processing

Databricks builds on top of Spark and adds:

- Highly reliable and performant data pipelines
- Productive data science at scale



Want to learn more? Visit our **platform page**.

SCHOOL OF
**SCIENCE, ENGINEERING
& ENVIRONMENT**

# Why use Databricks not plain Apache Spark‼



Comparison table with Apache Spark and databricks columns:

| Feature | Apache Spark | databricks |
|---|:---:|:---:|
| **DATABRICKS RUNTIME** — *Built on Apache Spark and optimized for performance* Learn More | ✕ | ✓ |
| **MANAGED DELTA LAKE** — *Reliable and Performant Data Lakes* | ✕ | ✓ |
| **INTEGRATED WORKSPACE** — *Interactive Data Science and Collaboration* | ✕ | ✓ |
| **PRODUCTION JOBS AND WORKFLOWS** — *Data Pipelines and Workflow Automation* | ✕ | ✓ |
| **ENTERPRISE SECURITY** — *End-to-End Data Security and Compliance* Learn More | ✕ | ✓ |
| **INTEGRATIONS** — *Compatible with Common Tools in the Ecosystem* | ✕ | ✓ |
| **EXPERT SUPPORT** — *Unparalled Support by the Leading Committers of Apache Spark* | ✕ | ✓ |

https://www.databricks.com/spark/comparing-databricks-to-apache-spark

SCHOOL OF SCIENCE, ENGINEERING & ENVIRONMENT

# What is Databricks?

Databricks is an enterprise software company that provides Data Engineering tools for **Processing** and **Transforming** huge volumes of data to build machine learning models. Traditional Big Data processes are not only sluggish to accomplish tasks but also consume more time to set up clusters using Hadoop. However, Databricks is built on top of distributed Cloud computing environments like **Azure, AWS,** or **Google Cloud** that facilitate running applications on CPUs or GPUs based on analysis requirements. Databricks platform is said to be **100 times faster** than **Apache Spark.**

https://hevodata.com/learn/what-is-databricks/

SCHOOL OF
SCIENCE, ENGINEERING
& ENVIRONMENT

# What is Databricks?

Let us start by answering this main question of What is Databricks. Databricks, developed by the creators of **Apache Spark**, is a Web-based platform, which is also a **one-stop product** for all Data requirements, like Storage and Analysis. It can derive insights using SparkSQL, provide active connections to visualization tools such as **Power BI, Qlikview, and Tableau,** and build Predictive Models using **SparkML**. Databricks also can create interactive **displays, text,** and **code** tangibly. Databricks is an alternative to the **MapReduce** system.

SCHOOL OF
SCIENCE, ENGINEERING
& ENVIRONMENT

# What have you learnt so far?

SCHOOL OF
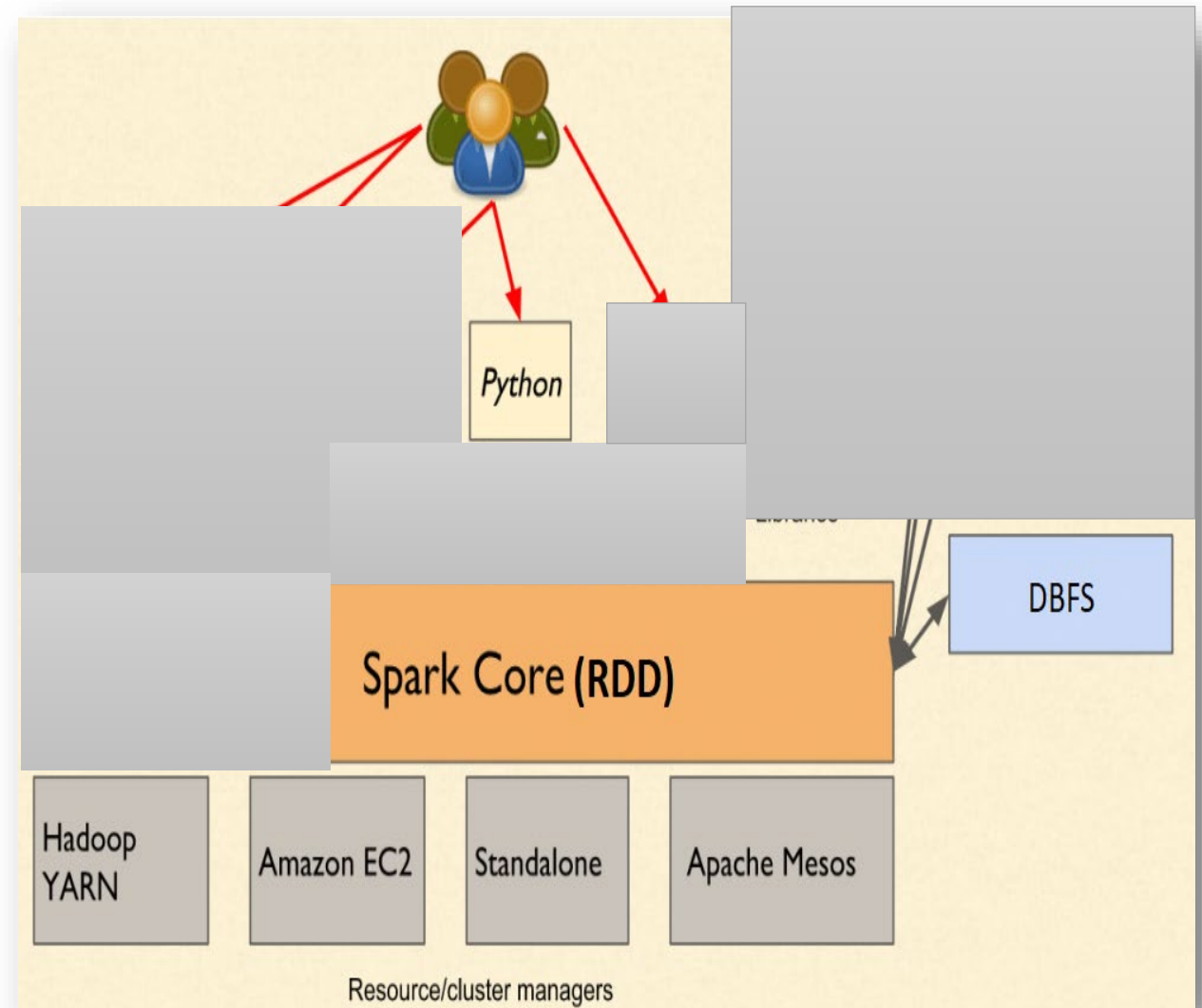SCIENCE, ENGINEERING
& ENVIRONMENT

# What you have learnt so far?

**What you have learnt so far?**

You have used **Python** to communicate with **Spark Core** inside the **Databricks** platform.

You have uploaded Data sets into **DBFS** and then you created many **RDDs**.

You have done some **Big Data analysis in its early stages** like extracting RDDs containing parts of the data sets that matter to you.

SCHOOL OF
SCIENCE, ENGINEERING
& ENVIRONMENT

# What will you learn after this?

# What will you learn?

**What you will learn in BDTT module.**

You will use **Python** and **SQL** to work with Spark SQL, Spark Streaming and MLib.

You will create and analyse data frames (instead of RDDs)

You will execute SQL statements against data sets.

You will upload batch data into **DBFS** or bring **stream data** into DBFS and then doing advance **Big Data analysis.**