# SPARK

By
**Niranjan Hegde** 1BM19IS103
**Prashanth Jaganathan** 1BM19IS115
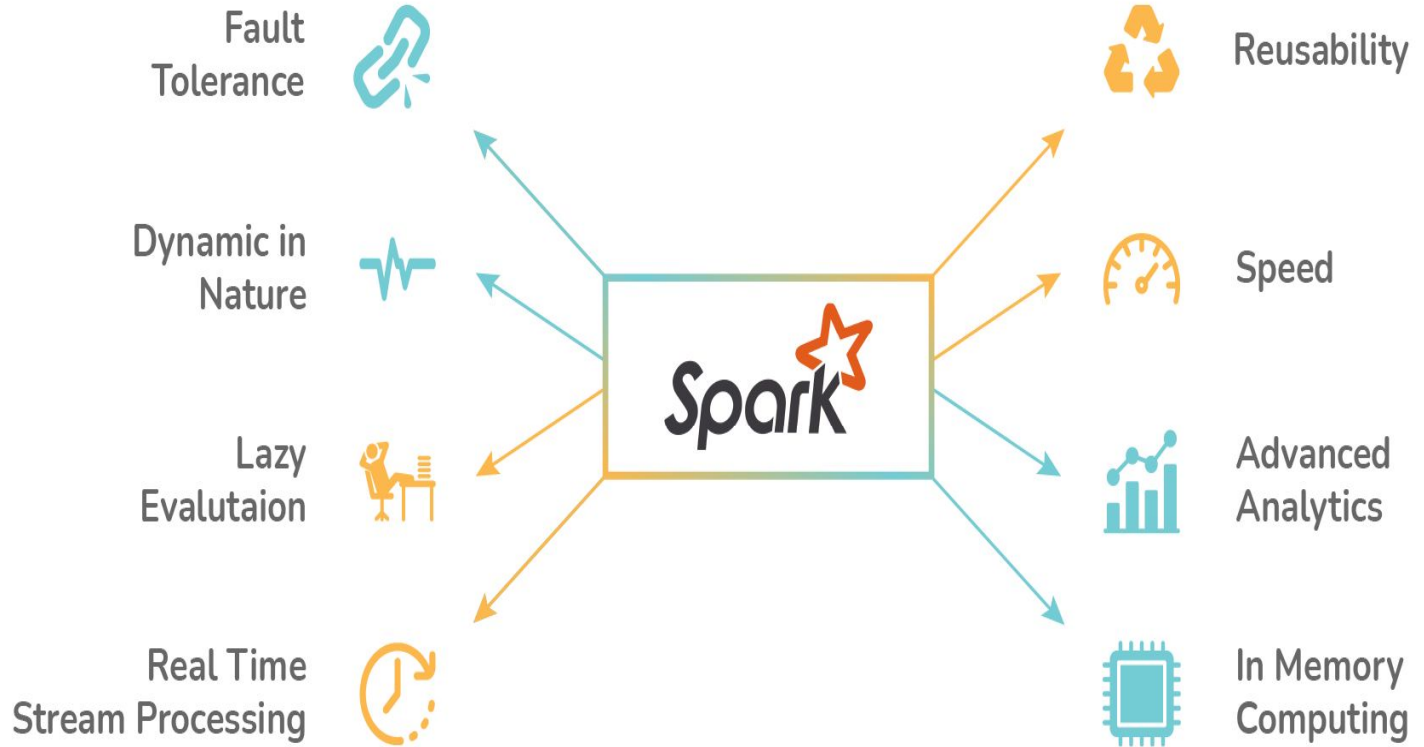**Prateek Gummaraju** 1BM19IS117
**Samartha S** 1BM19IS219

Apache Spark is an open source, distributed processing system used for Big Data workloads.

Developed in 2009 in UC Berkeley's AMPLab

The main feature of Spark is its **in-memory cluster computing technology** that increases the processing speed of an application

# What is Apache Spark?

# Performance

**Hadoop**

Hadoop is generally **slow** as it performs operations on the disk and **cannot deliver** near **real-time analytics** from the data

No real-time analytics

**APACHE Spark**

Spark runs **100 times faster** in-memory, and **10 times faster** on disk. If Spark runs on YARN with other resources demanding services, there could be major degradation
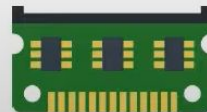
Faster in-memory processing

# Fault Tolerance



Hadoop is **highly fault-tolerant** because it was designed to **replicate data** across many nodes. Each file is split into blocks and replicated numerous times across many machines



Spark uses **Resilient Distributed Datasets** (RDDs), which are fault-tolerant collections of elements that can be operated on in parallel
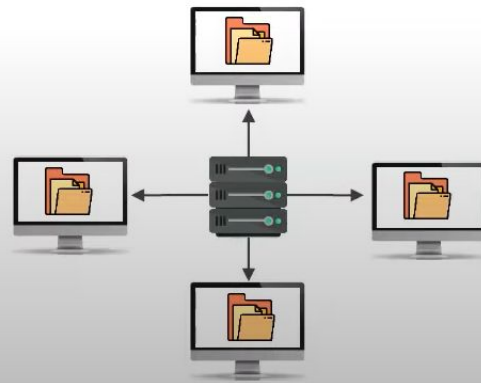
# Data Processing

## Hadoop

Hadoop **processes** data in **batches**. MapReduce operates in **sequential steps** by reading data from the cluster, performing its operations on the data, writing the results back to the cluster

Batches of input data → Hadoop → Output data

## Spark

Sparks performs **batch**, **real-time**, and **graph processing** of data. It reads data from the cluster, performs its operation on the data, and then writes it back to the cluster

Batch →
Real-time →
Graph →
Spark

# Ease of Use



Hadoop's MapReduce has **no interactive mode** and is complex. It needs to handle low-level APIs to process the data, which requires lots of coding

Spark supports **user-friendly APIs** for different languages. It has an **interactive mode** and provides intermediate feedback for queries and actions

# Language Support

## hadoop

Hadoop framework is developed in Java programming language. While, MapReduce applications can be written in Python, R and C++

**Java**

MapReduce supports programming languages

## Apache Spark

Apache Spark is developed in Scala language and supports other programming languages like Python, R, and Java

**Scala**

Spark supports other programming languages

# Scalability

**Hadoop** is **highly scalable** as we can add n number of nodes in the cluster. Yahoo reportedly used a **42,000** node Hadoop cluster

The largest known Spark cluster has **8,000** nodes. But as big data grows, it's expected that cluster sizes will increase to maintain throughput expectations.
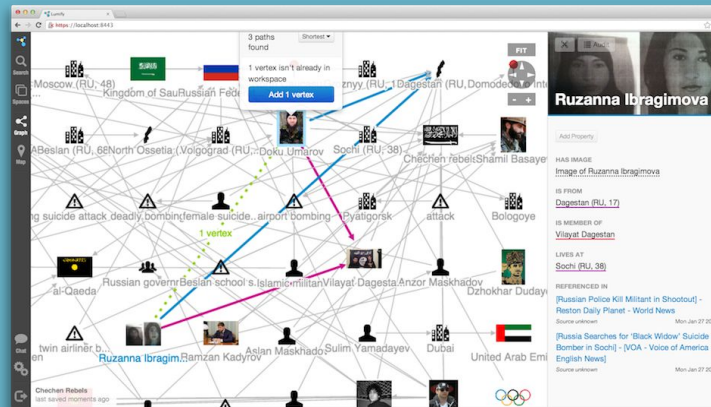
Pyspark Demo

- https://colab.research.google.com/drive/1dOV2TuRV5EIjfxII2jHnNgml8I4_Dn8O
- https://colab.research.google.com/drive/1hoX7JLNGtZxUJSn3gT6msQerKWUE2Juq?ts=62beb310

**Big Data Analytics and Visualization using LUMIFY**

- Lumify is a big data fusion, analysis, and visualization platform. Like all big data analytics tools, it too enables you to understand connections and explore the relationship between your data.

- Lumify is considered as a good big data analytics tool because it facilitates its users to get a set of analytics options that include graph visualizations, full-text faceted search, dynamic histograms, interactive geospatial views, and collaborative workspaces that can be shared in real-time.

- Lumify offers both 2D and 3D graph visualizations with automatic layouts. It also provides a plethora of options to analyze the links between different entities in a graph.

- Lumify comes with specific ingest processing and interface elements for textual content, images, and videos. The platform allows you to organize your work in different workspaces.

- The platform is built on proven, scalable big data technologies. It is secure, scalable, and backed by a motivated full-time development team.

- Lumify enables users to discover complex connections and explore diverse relationships in their data through a suite of analytic options, including graph visualizations, full-text faceted search, dynamic histograms, interactive geospatial views, and collaborative workspaces shared in real-time.

- It works well in cloud environments, especially AWS.

- Datawrapper is a free, intuitive and interactive tool that does not require any coding or design knowledge in order to visualize data.

- It lets you plot data as insightful maps, charts and tables. The map, chart or table can be downloaded as PNG, PDFs or they can be embedded directly onto your website.

- Let me take you through a quick demo.