

## Survey on Various Tools in Big Data

Afsha Akkalkot<sup>1</sup>, Aakanksha Barakade<sup>2</sup>, Avisha Mulchandani<sup>2</sup>, Tanmayi Yere<sup>2</sup>, Sonali Dalvi<sup>2</sup>, Vinayak Yadav<sup>2</sup>

<sup>1</sup> Assistant Professor, Department of Computer Engineering, Savitribai Phule Pune University, Pune, India.

<sup>2</sup> Student, Zeal College of Engineering and Research (ZCOER), Pune, Maharashtra, India.

\*\*\*

**Abstract** - In recent years, big data has attracted a lot of interest. It's a fairly typical necessity today to analyse large amounts of data, and these requirements can become nightmares when done with a large data source like Twitter tweets. It can be quite difficult to analyse a large number of tweets in order to obtain timely, relevant information with a variety of patterns. This paper will examine the idea of big data analysis and identify some important data from a sample big data source, like Twitter tweets, using one of the recently developed tools in the market, called Spark by Apache Spark, Tableau, Jupyter Notebook, Microsoft Power BI.

**Key Words:** Apache Spark, Tableau, Jupyter Notebook, Microsoft Power BI

### 1. INTRODUCTION

Bigdata is a vast collection of organised, unorganised, and semi-organized information. Numerous scientific disciplines, including chemistry, zoology, computer, physics, industry, etc., produce a vast amount of organised and unorganised data. Large data sets make it extremely difficult to perform traditional procedures. Big data analytics is the challenging process of examining vast amounts of data and various data sets to find patterns, connections, and relationships that may go undetected and provide valuable insight to the user. Big data platforms combine numerous tools and methods into a single, well-rounded product that aids in managing and analysing data. These data operate on a variety of platforms and industries, including cloud computing, both Mango DB and Apache Spark.

### 2 Tools for Big Data Processing

#### 1 Tableau

A business intelligence tool for visual data analysis is Tableau. Users can design a dashboard that is interactive and shared and uses graphs and charts to show the trends, variances, and density of the data. To gather and process data, Tableau can connect to files, relational databases,

and Big Data sources. The programme is highly distinctive since it supports real-time collaboration and data merging. For the study of visual data, it is utilised by corporations, academic researchers, and many government agencies. It is also positioned as a leader in the Gartner Magic Quadrant for Business Intelligence and Analytics Platform.

#### 1.1 Features

- **Data Blending** : The key function of Tableau is data blending. When combining relevant data from many data sources that you wish to analyse in a single view and portray as a graph, it is employed.
- **Real-time Analysis** : When the Velocity is high and real-time data analysis is challenging, real-time analysis enables users to swiftly grasp and analyse dynamic data. With interactive analytics, Tableau can assist in obtaining valuable information from rapidly changing data.
- **The collaboration of data** : Data analysis is not a lonely endeavour. Tableau is designed for collaboration because of this. Members of the team can distribute data, conduct follow-up research, and send simple visualisations to those who could benefit from the information. Success depends on ensuring that everyone can grasp the data and make informed decisions.

#### 1.2 Tools Of Tableau

Data analytics in Tableau is classified into two parts:-

**1. Developer Tools** : Developer's tools refer to the Tableau tools that are used for development, including the production of charts, dashboards, report generating, and visualisation. Examples of this kind include Tableau Desktop and Tableau Public.

**2.Sharing Tools :** These tools' purpose is to distribute the dashboards, reports, and visualisations that were produced with the developer tools. Tableau Server, Tableau Online, and Tableau Reader are examples of Tableau tools that fit under this category.

### 1.2.1Tableau Desktop

We can code and customise reports using Tableau Desktop, which has a robust feature set. The entire process is completed in Tableau Desktop, starting with the creation of the reports and charts and concluding with their mixing to generate a dashboard. Tableau Desktop creates communication between the Data Warehouse and other different sorts of files for real-time data processing. The worksheets and dashboards created here can be shared privately or publicly. Based on the connectivity to the publishing option and data sources, Tableau Desktop is also classified into two parts-

**Tableau Desktop Personal:-** The worksheet is kept private and access is restricted in the personal edition of the Tableau desktop. There is no way to publish the workbooks online. Therefore, it ought to be shared either offline or through Tableau Public.

**Tableau Desktop Professional:-** Comparable to Tableau desktop, it. The primary distinction is the ability to publish workbooks made in Tableau Desktop online or on Tableau Server. All types of datatypes are fully accessible in the professional edition. For those who wish to publish their workbook in Tableau server, it is the ideal option.

### 1.2.2Tableau Public

This Tableau version was created with budget-conscious consumers in mind. The phrase "Public" indicates that the generated workbooks cannot be locally saved. They ought to be stored on Tableau's public cloud, which anyone can access and observe. The files stored in the cloud have no privacy, therefore anyone can view and download the same information. For people who want to study Tableau and for those who wish to publish their data with the world, this version is optimal.

### 1.2.3Tableau online

Although its functionality is comparable to that of the tableau server, data is kept on servers hosted in the cloud that are managed by the Tableau group. The data that is made available via Tableau Online can be stored indefinitely. Over 40 cloud-hosted data sources, including Hive, MySQL, Spark SQL, Amazon Aurora, and many more, are directly connected via Tableau Online. The workbooks produced by Tableau Desktop must be published in order for Tableau Server and Tableau Online to function. Google Analytics and Salesforce.com are two web programmes that Tableau Server and Tableau Online may access data from.

### 1.2.4Tableau Server

The software is properly utilized to distribute workbooks and visualizations produced by the Tableau Desktop application around the company. You must publish your worksheet in Tableau Desktop before sharing dashboards on the Tableau Server. Only the authorized users will have access to the worksheet once it has been uploaded to the server. Authorized users don't always need to have Tableau Server installed on their computers. They merely need the login information in order to examine reports using a web browser. Tableau Server's high level of security is advantageous for efficient and speedy data exchange. The organization's administrator has complete control over the server. Both the software and the hardware are maintained by the organization.

### 1.2.5 Tableau Reader

We may view the visualizations and workbooks made with Tableau Desktop or Tableau Public using the free utility Tableau Reader. Filtering the data is possible, but changes and editing are limited. Tableau Reader has no security because anyone may use it to read workbooks. The recipient of the dashboards you build themselves must have Tableau Reader in order to read the file.

## 1.3 Architecture

The multiple data levels depicted in Figure 1 can be connected by Tableau Server thanks to its design. It is able to link clients from desktop, mobile, and online platforms. A powerful tool for data visualization is

Tableau desktop. It is really secure and available. Both real and virtual machines are capable of running it. It is a system with several users, processes, and threads. You may access databases, work together on projects, create reports, and more with its assistance. It also makes COVID-19 data analytics simple. The fact that Tableau can connect to several sources simultaneously is its strongest feature. In order to provide you with accurate results, it combines the data it receives from those sources.

### 1.3.1 Tableau Data Server Process

The data sources that Tableau Architecture may link to are its main component. It carries out numerous related crucial responsibilities, including storing data in the repository, protecting user data, and many others. Multiple data sources can be connected to by Tableau.

These data sources may be situated locally or remotely. It may simultaneously connect to a web application, an excel file, and a database. Tableau is capable of integrating data from several settings. All of these can be used simultaneously by Tableau. With our rapid in-memory data engine, Tableau offers straightforward alternatives to updating your data quickly and responsively. It has the ability to combine data from many data sources. Additionally, it can improve the link between different types of data sources.

### 1.3.2 Tableau Data Connectors Process

The data connectors offer a means of integrating Tableau Data Server with external data sources. There is an ODBC/SQL connector built into Tableau. Without using the native connector for each database, this ODBC connector can connect to any database. Tableau offers the choice of current data or data that has been extracted. Tableau offers a variety of database connectors, like as. In addition to many more, these include Microsoft Excel, SQL Server, Oracle, Teradata, Vertica, and Cloudera Hadoop. You have two choices for where to store this transferred data from Tableau. First off, real-time data is transferred using this method, which also uses data that is obtained directly from an external source. For data transfer, Tableau sends SQL statements and multi-dimensional expressions. The alternative two Using this strategy, we can retrieve data from a specific source in addition to relying on a current data source. You can make a local copy of the data as an extract file using Tableau. Millions of records can be extracted by

Tableau's data extraction tool from a data source. The user-friendly interface makes sure that data extraction isn't too challenging for you.

### 1.3.3 Tableau Application Server Process

The authentications and authorizations are offered by the application server. It takes care of online and mobile interface permissions and management. By logging each session id on the Tableau server, it ensures security. The server's default session timeout can be set by the administrator [30]. The REST API and web application calls are handled by the application server. Searching and browsing are also supported by the application server. Configure instances on each node in the Tableau server cluster to provide high availability of the application server. Additionally, it manages VizQL server-related tasks unrelated to data visualisation.

### 1.3.4 Tableau VizQL Server Process

The data source's queries are transformed into visualisations using the VizQL server. After the VizQL process receives the client request, it sends the query directly to the data source and retrieves the information as images [33]. The user is shown this visualisation or image. To speed up loading, Tableau server stores visualisations in a cache. Many people that have the right to view the visualisation can share the cache. Configure one or more instances to run on several nodes in order to ensure high availability for the VizQL server process.

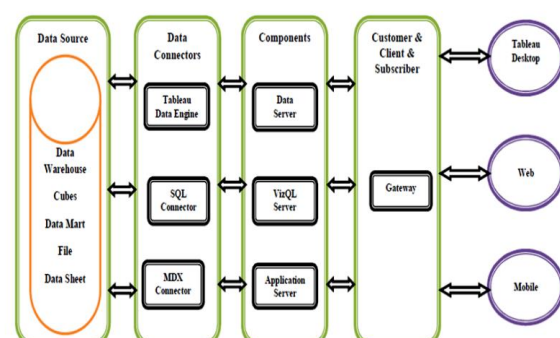


Figure 1 Architecture of Tableau

### 1.3.5 Tableau Gateway Process

A gateway is a type of web server that facilitates HTTP or https communications between clients and the server. Requests from users are routed through the gateway to Tableau components. When a client submits a request, the external load balancer receives it and processes it. The gateway acts as a conduit for procedures to other parts. Gateway also serves as a load balancer when an external load balancer is not present. One primary server or gateway controls all processes in a single server architecture. One physical system serves as the primary server in systems with several servers, while others are employed as worker servers. A part of the Apache web server, httpd.exe, is used by the Tableau server gateway process. Its job is to manage all client requests sent to the server by means of proxy servers, load balancers, mobile devices, desktops from Tableau, etc. If the system lacks a load balancer, the gateway can also serve in that capacity.

### 1.3.6 Tableau Data Engine Process

For greater effectiveness, the data engine speeds up analytical procedures. It produces refreshes or runs extract queries. If you employ data sources with many connections, it can also aid with cross-database joins. Data engine technology designed to handle analytical queries on big or complicated data sets quickly. When producing, updating, or requesting extracts, the data engine is utilized. In order to handle federated data sources with several connections, it is also utilized for cross-database joins. The data engine is intended to use all CPU and memory that are available on the device to deliver the quickest response times.

### 1.3.7 Tableau Backgrounder Process

Backgrounder is a crucial multi-process that handles information refresh schedules and makes sure the Tableau server and data engine are running correctly. Server tasks including extract refreshes, subscriptions, flow runs, and data-driven alerts are all carried out by the backgrounder process. Jobs can be started manually or automatically by using the 'Run Now', REST API, or tab cmd commands. It also aids in synchronising directory groups, checking for available disc space, and rebuilding search indexes.

### 1.3.8 Tableau Repository Process

The repositories on the Tableau server store server metadata pertaining to users, permissions, assignments, groups, and projects, as well as extract metadata and refreshed data. It also maintains performance information for auditing together with visualisations in flat files (TWS, TDS) and metadata. The repository is used to supply metadata whenever a server service or component requests it. Additionally, it keeps the visualisations in flat files. Additionally, performance data can be saved for upcoming audits. Information is sent to the application server for login verification using a partnership with the active directory.

### 1.3.9 Tableau Search and License Process

As its name implies, the license component handles the server's licensing requirements. On the other side, you can search the repository's index for our needs using the search area. These elements can appear straightforward, yet they are crucial to the server's smooth operation. On the main server of the Tableau server cluster, both of these services are active.

### 1.3.10 Tableau Server File Store Process

The extract storage is managed by the Tableau server file store process. Unless the node already has a data engine instance, when the file store is installed, a data engine instance is also installed. File stores, however, can be used both locally and outside of the Tableau server. In scenarios with high availability, the file store makes sure that extracts are synchronised to other file store nodes so they are still accessible if one file store node fails to function.

### 1.3.11 Tableau Server Administration Controller Process

The TSM REST API is hosted by the Administration Controller process and allows you to configure and manage your Tableau server setup. In the entire cluster, there can only be one instance of the administration controller.



### 1.3.12 Tableau Clients

The clients are the end users of Tableau who access it via the web, mobile devices, on-premises, in the cloud, or through a command-line interface for development. For accessing workbooks or visualisations, these end users primarily communicate with the Tableau server. Through the dashboards of Tableau online, you may change the contents of your visualisation using web browsers like Safari, Google Chrome, and Mozilla Firefox. Using the data you obtained from sources, Tableau desktop assists you in constructing the dashboard, workbooks, and visualisations. Additionally, we are uploading the outcomes to the server for later use. Additionally, you may use this tool to design your dashboards in accordance with tablets, phones, and PCs.

## 2 Jupyter Notebook

### 2.1 TOOLS

An open-source web tool called Jupyter Notebook enables you to create and share documents with real-time code, equations, visuals, and text. The notes, which are referred to as notebooks, can be used to:

**Run code:** Python, R, Julia, and Scala are just a few of the programming languages that can be used in Jupyter Notebooks.

**Create visualizations:** Making interactive visualisations with well-known tools like Matplotlib, Seaborn, and Plotly is possible with Jupyter Notebooks.

**Write narrative text:** You can use Jupyter Notebooks to create narrative writing that explains your code and findings.

**Share your work:** By posting them to a number of online repositories, such as GitHub and Binder, Jupyter Notebooks can be shared with others.

### 2.2 INFORMATION

Jupyter notebooks are executable documents that may be viewed online. The components of the notebook are composed of human-written contextual factors, computer code, and the output produced by the computer once the computer code has been executed. Tables and graphs can be among these outputs.

The components of the notebook may be used interactively, and the entire notebook may be saved, loaded again and executed, or converted to read-only forms like HTML, LaTeX, and PDF. Jupyter notebooks can enhance the efficiency of computational and data exploration, documenting, communication, reproducibility, and reusability of scientific research findings by taking advantage of these features.

Jupyter is a notebook interface that supports a number of different programming languages thanks to backends called kernels that may be installed separately. Through a web browser's interface, Jupyter note-books can be browsed and modified.

### 2.3 ARCHITECTURE

The Jupyter Notebook architecture is composed of four main components:

**1. The Kernel:** The Jupyter Notebook's kernel serves as its foundation. It is in charge of running the code contained in the notepad cells. Any programming language, including Python, R, Julia, or Scala, can implement the kernel.

**2. The Notebook Server:** The kernel is managed by the notebook server, which also gives users a web interface through which to communicate with the notebook. Either a local workstation or a distant server can be used to run the notebook server.

**3. The Notebook Client:** The web browser that users use to communicate with the notebook server is known as the notebook client. Any web browser that supports HTML5 and JavaScript can be the notebook client.

**4. The Notebook File Format:** Jupyter Notebooks are stored in a file format called a "notebook" that is based on JSON. Users can share their notebooks with others via the notebook file format.

#### Architecture description about Jupyter Notebook:

**The browser:** The browser is the client-side application that users interact with to view and edit Jupyter Notebooks. The browser communicates with the Jupyter server over a web socket connection. The browser sends requests to the Jupyter server to view, edit, and execute Jupyter Notebooks. The Jupyter server sends responses to

the browser that contain the HTML, CSS, and JavaScript that the browser needs to render the Jupyter Notebook.

**The Jupyter server:** The Jupyter server is a web server that hosts the Jupyter Notebook application and provides access to Jupyter Notebooks over the network. The Jupyter server is implemented in Python and uses the Tornado web framework. The Jupyter server listens for requests on a specified port, and when a request is received, it parses the request and determines what action to take. For example, if the request is to view a Jupyter Notebook, the Jupyter server will load the Notebook file and render it in the browser.

**The kernel:** The kernel is a process that executes the code in a Jupyter Notebook. The kernel is implemented in a variety of languages, including Python, Julia, and R. The kernel communicates with the Jupyter server over a ZeroMQ socket connection. The kernel sends requests to the Jupyter server to execute code. The Jupyter server sends responses to the kernel that contain the results of the executed code.

**The Jupyter notebook file:** The Jupyter notebook file is a JSON file that stores the code, output, and other metadata for a Jupyter Notebook. The Jupyter notebook file is saved on the file system, and it can be opened in the Jupyter Notebook application.

The Jupyter Notebook application will load the Jupyter notebook file and render it in the browser.v

## 2.4 FEATURE

Some of the specific features of Jupyter Notebook that make it a powerful tool for data analysis, scientific computing, and machine learning:

**Code cells:** Code cells are where you can write and run code. This is a great way to experiment with different

programming concepts and to see the results of your work immediately.

**Markdown cells:** Markdown cells are where you can write narrative text. This is a great way to explain your code and results, and to create documentation for your projects.

**Equation cells:** Equation cells are where you can write mathematical equations using LaTeX. This is a great way to add mathematical notation to your notebooks.

**Images:** You can insert images into your notebooks. This is a great way to illustrate your work and to make your notebooks more visually appealing.

**Charts and graphs:** You can create charts and graphs using popular libraries like Matplotlib, Seaborn, and Plotly. This is a great way to visualize your data and to communicate your results to others.

**Cell magics:** Cell magics are commands that can be used to run code or change the behavior of a cell. This can be used to automate tasks and to add new features to your notebooks.

**Line magics:** Line magics are commands that can be used to run code or change the behavior of the entire notebook. This can be used to set up your environment and to run tasks before and after each cell.

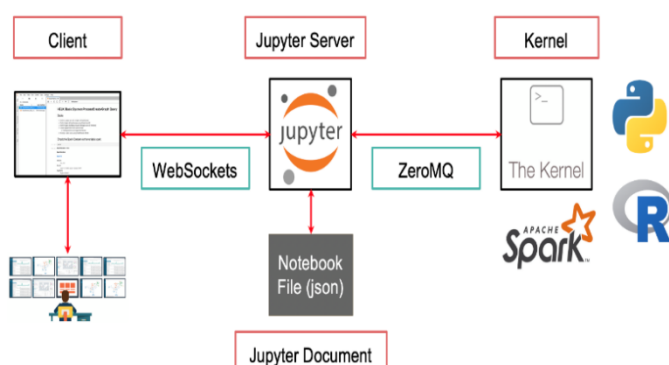
**Notebook extensions:** There are a variety of extensions available for Jupyter Notebook that can add new features and functionality. This can be used to add new visualization tools, to improve the formatting of your notebooks, and to collaborate with others.

## 2.5 ADVANTAGES

**Interactive:** Since Jupyter Notebooks are interactive, you can run code and see the results right away. This makes it simple to experiment with various concepts and troubleshoot code.

**Extensible:** Jupyter Notebooks are expandable, so by installing extensions, you can increase their functionality and add new features.

Figure 2 Jupyter Notebook Architecture



**Portable:** Jupyter Notebooks are portable, allowing you to use them on any computer equipped with a web browser.

**Collaborative:** Jupyter Notebooks allow for collaboration, so you can share your work and work on projects with other people.

### 3 Apache Spark

#### 3.1 Introduction

Apache Spark is a powerful and dependable general-purpose cluster computing engine. This system offers APIs for application programming in programming languages, including Java, Python, and Scala. Spark is an Apache cluster computing system that has incubator status. This tool is designed to accelerate data processing, and it is fairly quick at both running programmes and writing data. When compared to disk-based engines like Hadoop, Spark's in-memory processing capabilities allow it to query data significantly more quickly. It also provides a universal execution architecture that can optimise any operator graph.

#### 3.2 Tools

The software components of the Spark framework known as spark tools are utilised for scalable and effective data processing for big data analytics. The Apache licence governs the open-source distribution of the Spark framework. It includes GraphX, MLlib, Spark Streaming, Spark SQL, and Spark Core, 5 crucial data processing tools. The programme used to handle and perform graph data analysis is called GraphX. For applying machine learning to distributed datasets, utilise the MLlib Spark utility. While Spark SQL is the main tool for structured data analysis, users simultaneously use Spark Streaming for stream data processing. RDD, or resilient data distribution, is managed by the Spark Core tool.

#### 3.2.1.Tool GraphX

This is the Spark API for graphs and computations that take advantage of parallel graphs. Resilient Distributed Property Graph, a Spark RDD extension, is offered by GraphX. The form has an expanding library of graph builders and algorithms to make graph analytics tasks easier.

This essential tool creates and modifies graph data to conduct comparative analytics. The former quickly transforms and integrates structured data while using the least amount of time and resources. Use the friendly Graphical User Interface to select an algorithm from a rapidly expanding list. Even your own unique algorithms can be created to track ETL insights. You can run graph operations on data frames using the GraphFrames module. For graph queries, this involves making use of the Catalyst optimizer. This important tool has a number of distributed algorithms at its disposal. The latter's function is to process network structures that incorporate a PageRank algorithm implementation from Google. To model critical data, these specialised algorithms make use of Spark Core's RDD methodology.

#### 3.2.2 MLlib Tool

A library called MLlib contains fundamental machine learning services. The library provides many types of machine learning algorithms that enable numerous actions on data with the goal of gaining insightful information. The Spark platform combines libraries to use machine learning and graph analysis methods on data at scale.

A framework for creating machine learning pipelines is provided by the MLlib programme, making it easy to execute transformations, feature extraction, and selections on any given structured dataset. Basic machine learning techniques like filtering, regression, classification, and clustering are part of the first category. However, there are no facilities for modelling or training deep neural networks. Machine learning libraries that power business intelligence may be built and maintained quickly and with the help of reliable algorithms thanks to MLlib.

### 3.2.3. Spark Streaming Tool

Processing live data streams is the aim of this technology. Data created by various sources is processed in real-time. Examples of this type of data include log files, messages with visitor-posted status updates, and others. Additionally, this tool uses Spark Core's quick scheduling ability to implement streaming analytics. Mini-batches of data are ingested, and then these mini-batches are subjected to RDD (Resilient Distributed Dataset) modifications. High throughput of real-time data streams and fault-tolerant stream processing are made possible by Spark Streaming. DStream is the primary stream unit.

In the latter, real-time data processing is carried out by a collection of resilient distributed datasets. This useful utility expanded the batch processing paradigm of Apache Spark to streaming. The stream was divided into several micro-batches and processed using the Apache Spark API. The backbone of robust applications that require real-time data is Spark Streaming. The former is highly desirable for development since it has the dependable fault tolerance of the big data platform. Interactive analytics are now available for real-time data from virtually any common repository source thanks to Spark Streaming.

### 3.2.4. Spark SQL Tool

This recently added Spark module mixes relational processing and the platform's functional programming interface. Both Standard SQL and the Hive Query Language are supported for data queries.

Spark SQL consists of 4 libraries:

- SQL Service
- Interpreter and Optimizer
- Data Frame API
- Data Source API

The purpose of this tool is to handle structured data. Integrated access to the most popular data sources is provided by the former. This comprises JDBC, JSON, Hive, Avro, and other technologies. The program organizes data into labelled columns and rows that are ideal for distributing the outcomes of quick queries. Spark SQL seamlessly interacts with both new and current Spark programs, resulting in low computational costs and

excellent performance. In order to develop an effective query plan for computation and data locality, Apache Spark uses a query optimizer called Catalyst. The required computations will be carried out by the plan throughout the cluster. The current recommendation is to use the Spark SQL interface of datasets and data frames for development purposes.

### 3.2.5. Spark Core Tool

This serves as the platform's fundamental building piece. It includes parts for carrying out memory operations, scheduling tasks, and other things. The API containing RDD is hosted by Core. The former, GraphX, offers APIs for creating and modifying RDD data. Dispatching distributed tasks and basic I/O capabilities are also offered by the core. When compared to components of Apache Hadoop, the Spark Application Programming Interface is very straightforward and user-friendly for developers. Behind comparatively straightforward method calls, the API hides a significant portion of the complexity involved in a distributed processing engine. By combining a driver core process, which divides a specific Spark application into various jobs and distributes them among numerous ways that do the work, Spark functions in a distributed manner. Depending on the needs of the application, these specific executions could be scaled up or down.

The technologies that make up the Spark ecosystem all work well together and with little overhead. This makes Spark a tremendously powerful platform that is also extremely scalable. The tools are still being worked on to make them more efficient and convenient to use.



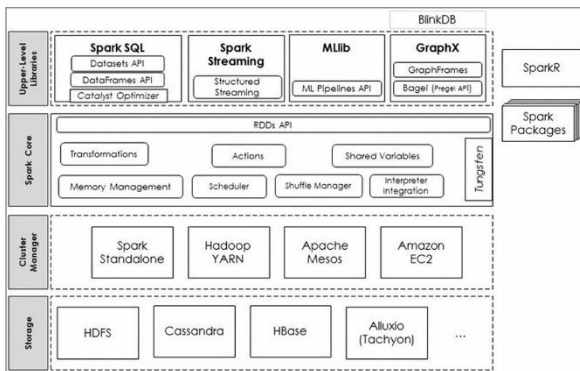


Figure 3 High Level Architecture of Apache Spark Stack

### 3.3 Architecture

In the Apache Spark architecture, when the driver program runs, it calls the actual application program and generates a SparkContext. The fundamental operations are all included in Spark Context. A DAG Scheduler, Task Scheduler, Backend Scheduler, and Block Manager are just a few of the many parts that make up the Spark Driver and are all in charge of converting user-written code into jobs that are actually run on the cluster.

The Cluster Manager controls how various jobs are carried out within the cluster. The Cluster Manager and Spark Driver collaborate to manage the execution of numerous additional activities. The responsibility of distributing resources for the project is carried out by the cluster manager.

SparkDriver will manage the execution after the job has been divided into smaller jobs and deployed to worker nodes. An RDD produced in the SparkContext can be processed by a large number of worker nodes, and the outcomes can also be cached. Information about tasks is enqueued on worker nodes by the Spark Context after being received from the Cluster Manager.

These responsibilities fall under the executor's purview. Executors have the same lifespan as the Spark Application. If we wish to boost the system's performance, we can hire more employees. By doing so, we may divide jobs into more logical sections.

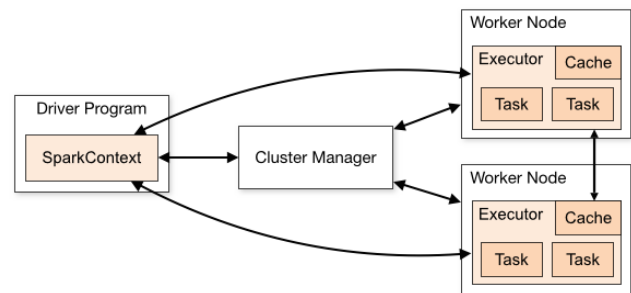


Figure 4 Apache Spark Architecture Diagram

A high-level view of the architecture of the Apache Spark application is as follows:

#### 1.The Spark driver

In a driver process, the master node (process) oversees the tasks and coordinates the workers. Spark is divided into jobs and scheduled to run on clusters of executors. The driver creates Spark contexts (gateways) to connect to a Spark cluster and monitor the jobs running in a particular cluster. In the diagram, the driver programmes connect to a Spark cluster, call the main application, and construct a spark context (which serves as a gateway) to monitor the jobs running in the cluster. The spark context is employed for all operations.

The Spark context contains a record for each Spark session. Cluster managers and additional components for executing jobs in clusters are included in Spark drivers.

#### 2.The Spark Executors

An executor is in charge of starting a job execution and putting data in a cache. At the start, executors sign up with the driver scheme. These executors can run the program concurrently during a variety of time windows. When the executor has loaded data and is not in idle state, the task is executed. When data is loaded and unloaded while the tasks are being completed, the executor is running in the Java process. During task execution, the executors are dynamically assigned and are frequently added and withdrawn. The executors are observed in action by a driving program. Tasks from users are carried out via the Java process.

### 3. Cluster Manager

Data is cached and controlled by a driver program, which also manages the execution of jobs. Executors initially register with the drivers. This executor has a number of time slots available for concurrently running the program. Executors fulfil client requests as well as read and write external data. When the executor has loaded data and has been deleted while in the idle state, a job is executed. Depending on how long it will be used, the executor is dynamically allocated and is constantly added and removed. Executors are watched over by a driver program while they complete duties for users. When an executor completes a user's task, code is run in the Java process.

### 3.4 Features

The goal of the development of Apache Spark, a well-known cluster computing platform, was to speed up data processing applications. Popular open source framework Spark uses in-memory cluster computing to speed up application performance. A cluster is a group of nodes that communicate and share information. Spark may be used to fulfil a variety of sequential and interactive processing requirements because of implicit data parallelism and fault tolerance.

- **Speed:** When processing huge amounts of data, Spark is up to 100 times faster than MapReduce. Additionally, it has the ability to break the data into manageable bits.
- **Powerful Caching:** A straightforward programming layer provides powerful caching and disc persistence capabilities.
- **Deployment options** include Mesos, Hadoop via YARN, and Spark's own cluster management.
- **Real-Time:** It provides real-time calculation and reduced latency due to its in-memory processing.
- **Polyglot:** Spark supports all four of these languages in addition to Java, Scala, Python, and R. Any one of these languages can be used to create Spark code. Additionally, Python and Scala command-line interfaces are offered by Spark.

### 4 Power BI

The bigger Power BI platform includes Power BI Desktop, Pro, Premium, Mobile, Embedded, and Report Server. Even while some of these tools are free to use, paid memberships are required to access the pro and premium versions, which provide more sophisticated insights.

Power BI is a component of Microsoft's Power Platform, which also includes Power Apps, Power Pages, Power Automate, and Power Virtual Agents. These applications, created as "low-code tools," help businesses with data analysis and visualisation, business solution design, process automation, and chatbot development.

### 4.1 Tools

**4.1.1 PowerBI.tips - Business Ops** a simple deployment mechanism for Power BI Desktop extensions of external apps. The objective of Business Ops is to offer a single location for installing all recent versions of external tools.

**4.1.2 Tabular Editor** Tabular models can be readily built, maintained, and managed by model creators utilising an easy-to-use editor. In a hierarchical view, which supports multi-select property editing and DAX syntax highlighting, all of the items in your tabular model are arranged by display folders.

**4.1.3 DAX Studio** a powerful tool with lots of features for DAX analysis, performance tuning, and diagnosis. Among the features are object browsing, integrated tracing, breakdowns of how queries were executed with specific statistics, and formatting and highlighting of the DAX syntax.

**4.1.4 ALM Toolkit** A tool for Power BI models and datasets schema comparison that is employed in ALM scenarios. Simple deployment between environments is possible, and you can keep incremental refresh historical data. Metadata files, branches, and repositories can all be diffed and merged. Additionally, you can apply shared definitions across datasets.

**4.1.5 Metadata Translator** simplifies the process of localising Power BI models and datasets. The programme can automatically translate table, column, measure, and hierarchy captions, descriptions, and display folder names. The tool translates utilising Azure Cognitive Services' machine translation technology. Additionally, Comma Separated Values (.csv) files can be used to export and import translations for convenient mass editing in Excel or a localization programme.

## 4.2 Architecture

### Components

#### Power BI Mobile Apps'

The mobile apps of Power BI keep you connected with the data no matter where you are. You can see live reports and dashboards on your iOS and Android smartphones and make better market decisions on the go. Only pro Power BI architecture provides the feature of Mobile reports and dashboards.

#### Power BI Query

Power Query allows users to connect distinct information from multiple sources and convert them to satisfy their business requirements. Power Query is included in the Power Query Editor of Power BI Desktop.

#### Power Map

Power BI queries offer a 3D visualization tool, Power Map, that shows differences in your datasets with shadings ranging from dark to light.

#### Power View

For a quick and effective visualization in your Excel workbooks, you can try Power View's drag-n-drop feature and save your time. It's an important part of MS Power BI architecture that enables the user to quickly visualize the data in a few clicks.

#### Working of Power BI Architecture

The architecture is mainly divided into two parts: On-cloud and on-premises services.

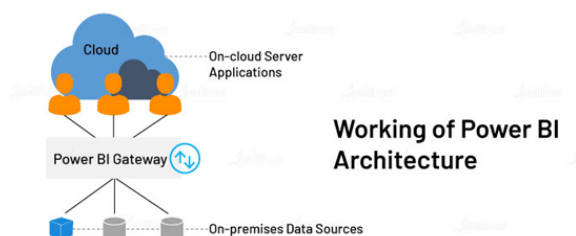


Figure 5 Power BI architecture

#### On-premises

The end-user receives all types of reports that have been made available on the Power BI Report Server here. The user can publish Excel worksheets to the Power BI

Report Server using Power Publisher. You may construct datasets, paginated reports, mobile reports, and more with the use of report server and publisher tools.

#### On-cloud

The BI gateway serves as a bridge in the Power BI Gateway architecture, allowing data to be moved from on-premises data sources to on-cloud servers or apps. The cloud is made up of a variety of parts, including reports, dashboards, datasets, Power BI Embedded, and more. The Power BI tools are linked to these on-cloud data sources.

## 4.3 Advantages

### 4.3.1 Easy to Use

You can create dashboards with ease using Power BI's incredibly user-friendly interface. This tool has built-in intelligence that suggests the best reporting component based on your selections. Power BI, for instance, can instantly recognise the map chart if you select sales and location as your data sources. Just how clever it is.

### 4.3.2 Low Learning Curve

Using and mastering the Power BI dashboard doesn't involve any coding. Further reducing the learning curve when it comes to building Power BI dashboards is the fact that Power BI was built on the Microsoft Excel platform. Utilising the Power BI dashboard won't present any problems if you are familiar with EXCEL.

### 4.3.3 Customizable Dashboards

When it comes to designing and sharing dashboards, Power BI provides incredible customization. To make the HR process simpler, you can develop a Power BI HR Analytics Dashboard, a Power BI for Banking Dashboard for financial analysis, or a Power BI Marketing Dashboard for campaign success analysis.

### 3. CONCLUSIONS

We have reviewed Apache Spark, Tableau, Jupyter Notebook & Microsoft Power BI features for large data analytics in this paper. Apache Spark is a multi-use cluster computing framework that supports advanced execution DAGs and APIs in Java, Scala, Python, and R. It also has an engine that is optimised for this purpose. MLlib in Spark. Tableau can connect to files, relational databases, and Big Data sources to gather and process data. Because it offers real-time cooperation and data merging, the programme is very unique. You may create and share documents with real-time code, equations, graphics, and text with Jupyter Notebook. A platform for data visualisation with a focus on business intelligence is called Microsoft Power BI. With varied levels of data expertise, business professionals can use this tool.

### REFERENCES

- 1.Drs. Yusuf P., Nikhat A., Nazia T., Dr. Asif P., and 2020. Tableau data analysis and visualisation are helpful for the COVID19 (Coronavirus) virus. *Advances in Global Engineering and Technology*, 3(2), 28–50
- 2."IPython: a system for interactive scientific computing," F. Pérez and B. E. Granger, *Computing in Science and Engineering*, vol. 9, no. 3, pp. 21-29, May 2007.
- 3.Ruslan Dautov, Xiaojun Chen, Patrick Xiaogang Peng, and Joshua Zhexue Huang. Salman Salloum. *Analytics for big data using Apache Spark*. Pages 145–164 of Volume 1 of the *International Journal of Data Science and Analytics* (2016)
- 4.International Conference on Robots & Intelligent Systems (ICRIS), 10.1109/ICRIS.2016.30, Hu Shuijing, "Big Data Analytics: Key Technologies and Challenges," 2016.
- 5.International Conference on Networking and Network Applications, IEEE DOI 10.1109/NaNA, 2017; Anku Jaiswal and Purrushottam Bagale, "A Survey on Big Data in the Financial Sector."
- 6."Big data: Issues, Challenges, Tools and Good Practises," IEEE "Contemporary Computing (IC3), 2013 Sixth International Conference, pp. 404–409. Authors: Katal, A.; Wazid, M.; Goudar, R.H.

7."Big data for C4I systems: goals, applications, challenges, and tools," Fazal-e-Amin, A.S. Alghamdi, I. Ahmad, and T. Hussain