

Big Data Architecture & Ecosystem Overview

Mehal Sakthi M S

Analyst SDE

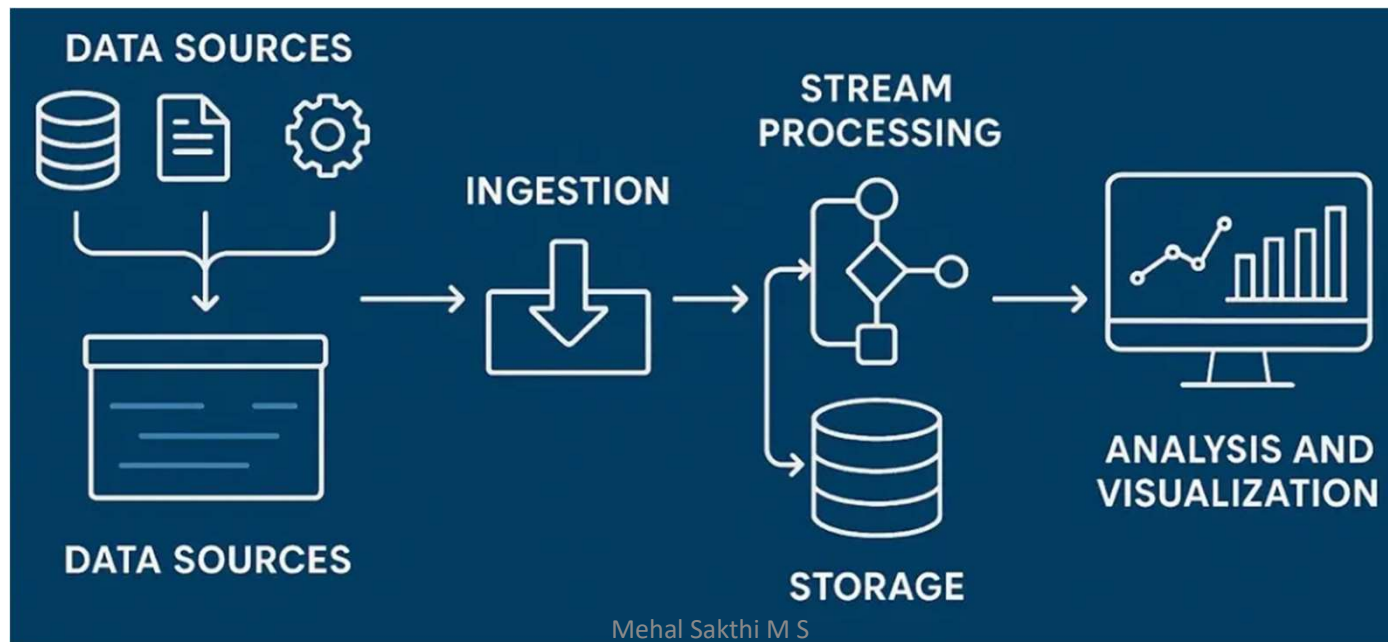
What is Big Data Architecture

Where is data stored?

How does it move?

How is it processed?

How does it become insights?



Data Sources Layer

Web
applications

Mobile apps

IoT sensors

Payment
systems

Logs

Social media
feeds

Databases

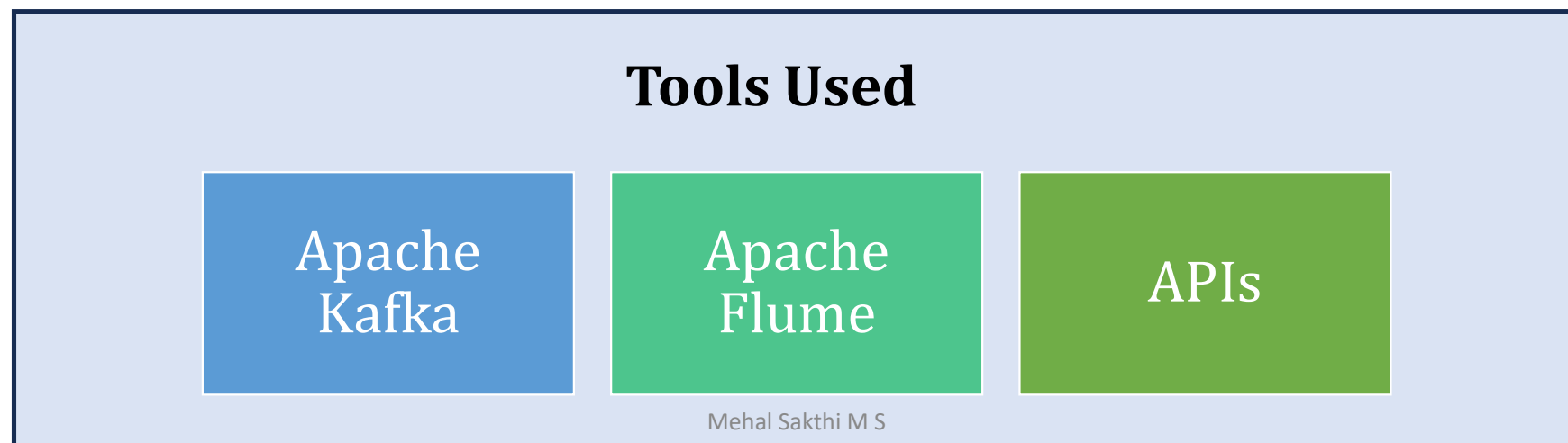
Structured

Semi-
structured

Unstructured
data

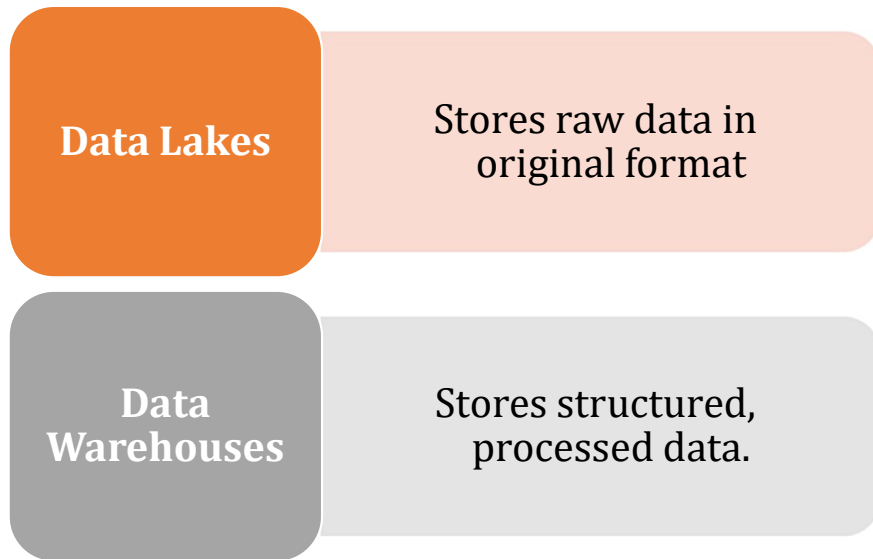
Data Ingestion Layer

The process of collecting and moving data into storage systems.



Storage Layer

Where is Big Data Stored?



Distributed File Systems

A system that stores data across multiple machines (nodes) but appears as a single storage system.

HDFS (Hadoop Distributed File System)

Splits large files into blocks

Stores blocks across different machines

Replicates data for fault tolerance

Automatically handles failures

*Instead of moving data to computation
Move computation to where data exists.*

Storage Layer

Data Lake

Raw data

All formats

Flexible

Schema-on-read

Data Warehouse

Processed data

Structured

Optimized

Schema-on-write



Data Warehouse



Data Lake



Lake House

Data Lakehouse

Data Lake flexibility
+
Data Warehouse performance

No duplication of data

Unified analytics + AI

Better governance

ACID transactions on data lakes

Processing Layer

Batch Processing

Processes data in bulk

Scheduled jobs (hourly/daily)

High latency
(minutes–hours)

Suitable for
historical analysis

Simpler architecture

Example: Monthly sales report

Stream Processing

Processes data continuously

Real-time / near real-time

Low latency
(milliseconds–seconds)

Suitable for instant decision-
making

More complex distributed
systems

Example: Fraud detection alert

Processing Layer



Apache Spark

An open-source distributed processing engine designed for large-scale data processing.

Key Characteristics

Distributed
Processing
Engine

In-Memory
Processing

Fault
Tolerant

Lazy
Evaluation

Unified
Analytics
Engine

Spark Ecosystem Components

Spark Core

Spark SQL

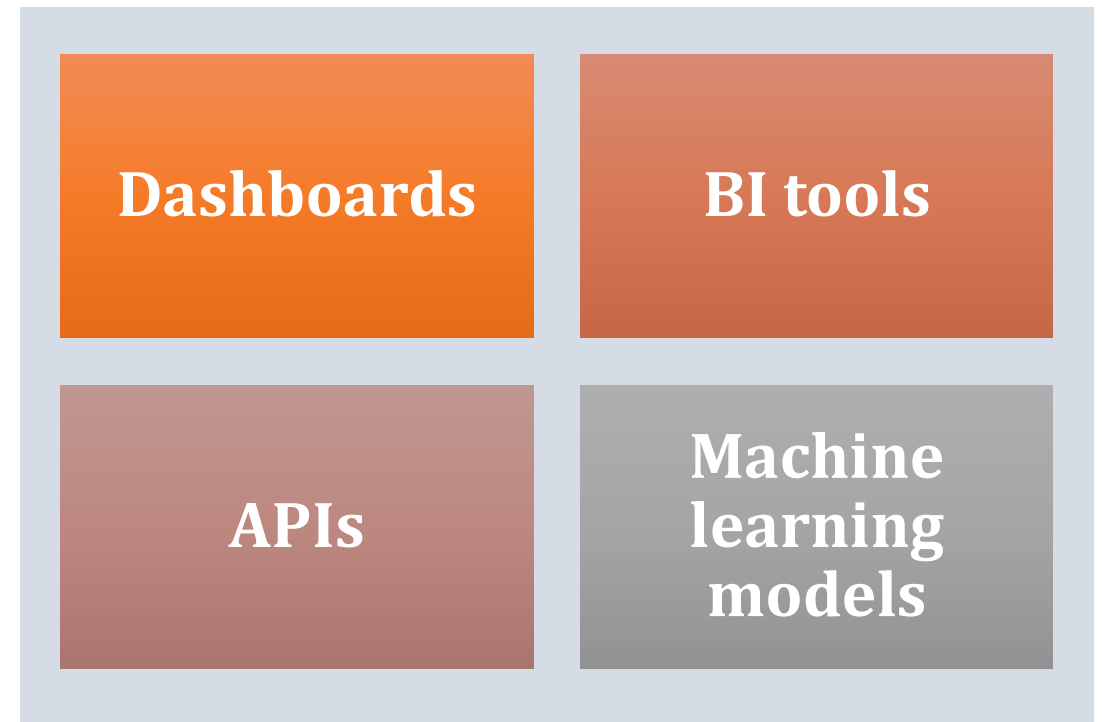
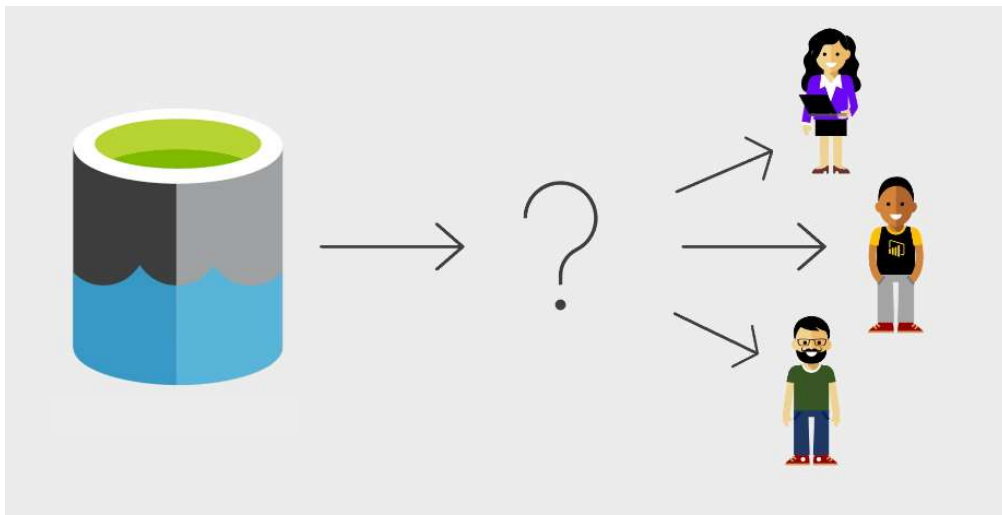
Spark Streaming

Spark MLlib

GraphX

Serving Layer

This is where we make data usable



Analytics & AI Layer

Reports are created

Insights are generated

ML models are trained

Predictions are made



Use Case Discussion

Thank You