

KELAS BIMBINGAN PEMBANGUNAN KERANGKA DATA RAYA

**10 – 12 DIS 2024 (SELASA – KHAMIS)
8:30 PAGI HINGGA 4:30 PETANG
OrenG ACADEMY**





Day #1

Introduction to Big Data and Types of Platforms

BIG DATA



BIG DATA

refers to extremely **large and complex datasets** that are challenging to process, store, and analyze using traditional methods due to their size, speed of generation, and variety of formats.

Cost Saving

- Optimize resources and reduce inefficiencies, lowering operational costs

Time Saving

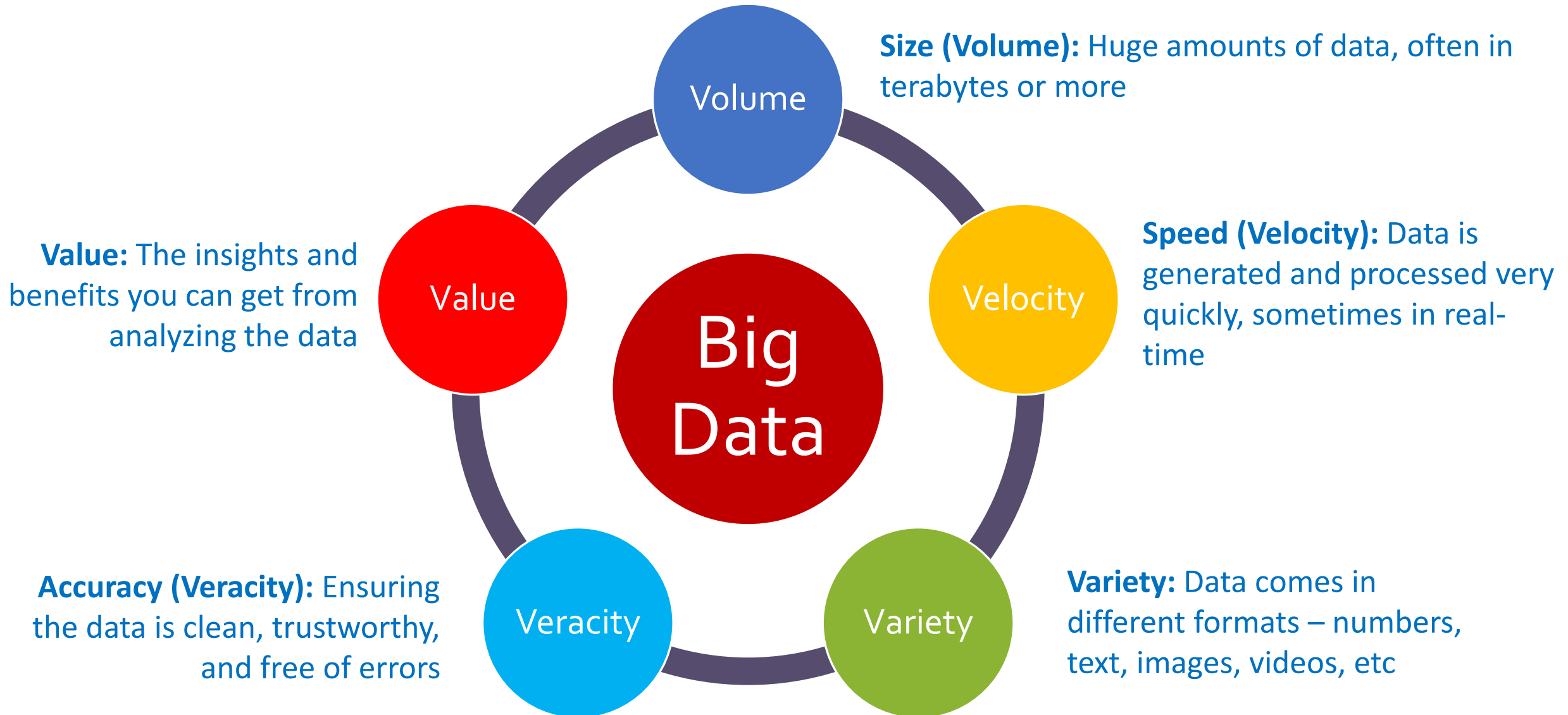
- Automated data processing and real-time analytics speed up tasks that would take hours or days manually

Faster & Better Decision Making

- Data-driven insights enable more accurate and timely decisions, improving outcomes

Data-Driven Policy Making

- In sectors like government, Big Data provides empirical evidence for crafting effective policies





Big Data Platform

What Data
(Represent)?

How Data is
Used?

SIMPLIFIED EXPLANATION



How Data is
Stored?

Data Life
Cycle/Stage

What Data / What Data Represents?

Data is **information** that can be collected, stored, and analyzed

How Data is Stored?

Structured Data:

Highly organized and easily searchable within relational databases

Unstructured Data:

Data that lacks a predefined model or format

Semi-structured Data:

Data that has some organizational properties but does not fit neatly into a structured database.

How Data is Used?

Operational Data:

Used to run day-to-day tasks

Analytical Data:

Used to find trends and make decisions

Real-Time Data:

Used for live updates

Historical Data:

Stored for later use

Data Life Cycle/Stage

Raw Data:

Like ingredients in the kitchen—collected but unprocessed

Cleaned Data:

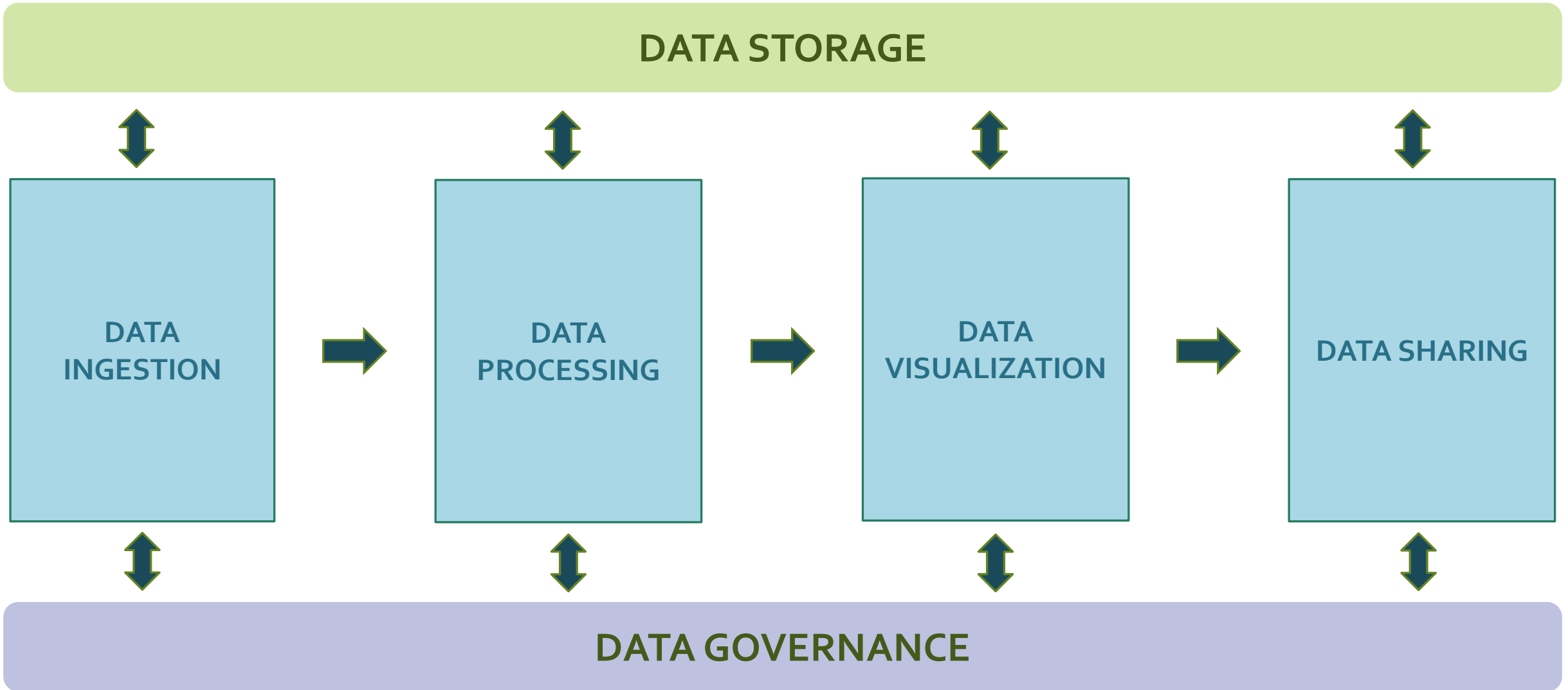
Prepped and ready to use—errors removed

Analyzed Data:

Insights drawn—helps make decisions

Archived Data:

Stored away for future reference

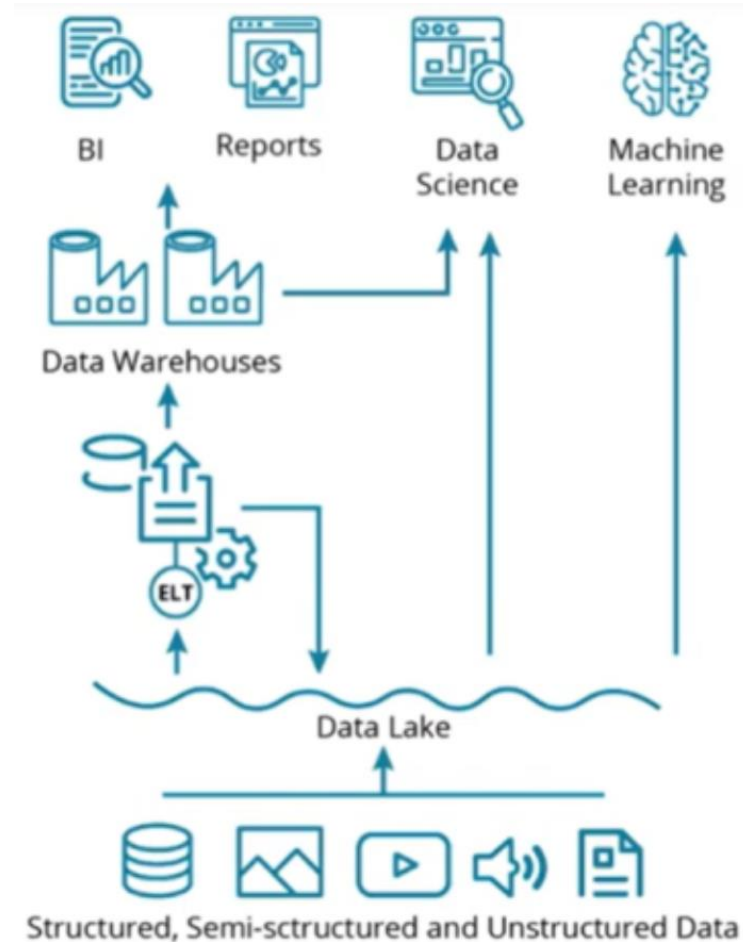


Types of Big Data Platforms

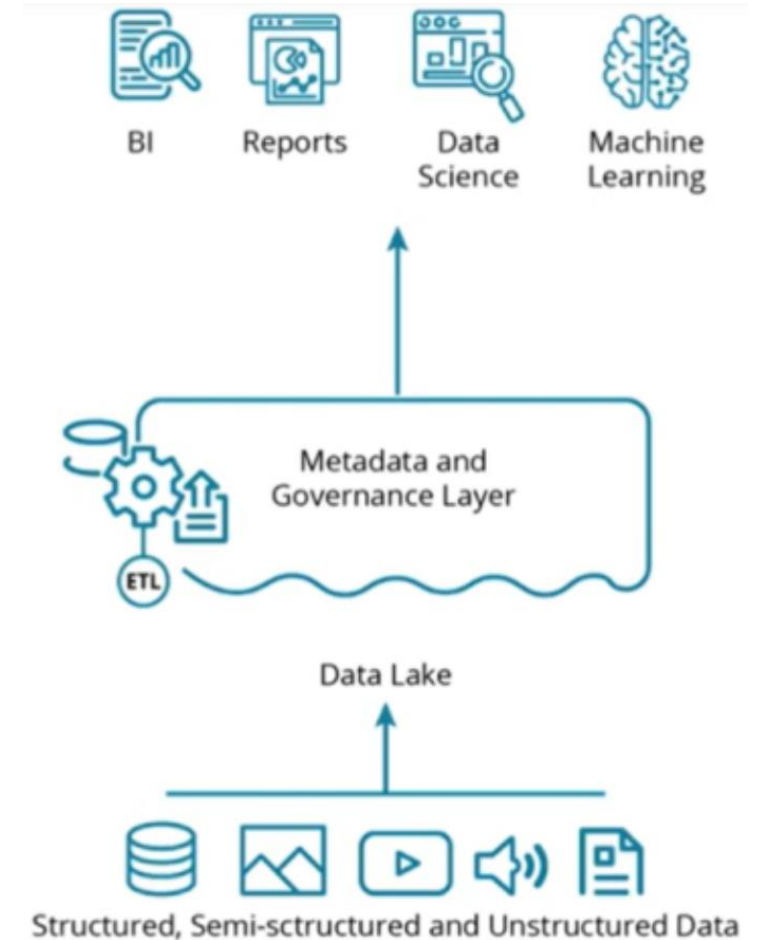
Data Warehouse



Data Lake



Data Lakehouse



Comparison of Data Management Systems

Feature	Data Warehouse	Data Lake	Data Lakehouse
Data Type	Structured	Structured, Semi-structured, Unstructured	Structured, Semi-structured, Unstructured
Use Case	Analytical, Reporting	Storage, Exploration, Analytics	Hybrid Storage and Analytics
Schema	Defined	No strict rules; you can dump data in any format	A hybrid approach, combining both strict and flexible organization
Query Complexity	Complex	Varied	Complex
ACID Compliance	Limited	Typically None	Provided by technologies like Delta Lake
Data Processing Tools	Business Intelligence Tools (e.g., Tableau)	Big Data Tools (e.g., Spark, Hadoop)	Hybrid Approach with Big Data Tools
Scalability	Scalable for analytics	Highly Scalable	Scalable, but may require a robust architecture
Example Tools	Snowflake, Redshift, BigQuery	Amazon S3, Azure Data Lake Storage, Hadoop	Delta Lake, Databricks, AWS Glue
Data Integration	ETL and ELT (Extract, Load, Transform) processes	Often uses ETL/ELT, supports data from various sources	ETL and ELT for structured and raw data
Common Use Cases	Historical Sales Analysis, Reporting	Raw data storage, Sensor Data, User-generated Content	Healthcare Data, IoT Data, Financial Data
Storage Efficiency	Optimized for query performance	Low-cost storage for diverse data types	Storage efficiency can vary based on architecture

**** ACID is principles make sure your data processes are reliable, accurate, and protected, even when working with massive amounts of information or during unexpected problems**

Database

Like a well-organized filing cabinet

Data Warehouse

Like a library where books (data) are categorized for easy access

Data Lake

Like a giant storage room where you toss everything without sorting

Data Lakehouse

Like a library that also has a storage room for unsorted items

Relationship Between Big Data And Data Analysis



BIG DATA

refers to extremely **large and complex datasets** that are challenging to process, store, and analyze using traditional methods due to their size, speed of generation, and variety of formats.

DATA ANALYTICS

is the **process** of examining, interpreting, and analyzing Big Data to uncover patterns, trends, correlations, and insights.

Aspect	Big Data	Data Analytics
What it is	Large amounts of data	Process of analyzing that data
Focus	It emphasizes the storage, processing, and management of large, complex datasets	It focuses on extracting value and actionable insights from Big Data
Purpose	Serves as the raw material. It provides the information but doesn't derive meaning from it	Serves as the tool or technique used to interpret and make use of the raw data
Example	Collecting data from millions of users about their online shopping habits	Analyzing shopping habit data to recommend products to users or improve sales strategies

Data Sources

Data Processing

Data Cleaning

Analysis & Modelling

Visualization & Presentation

Objective:

Define the problem or question and determine the required data

Actions:

- Identify internal and external sources
- Classify data as structured or unstructured

Objective:

Gather raw data from identified sources

• Actions:

- Web scraping
- Surveys
- Database queries
- Upload Module
- API
- IoT sensors

Objective:

Ensure the data is accurate, consistent, and usable

Actions:

- Handle missing data (impute, remove, or flag)
- Remove duplicates data
- Standardize data formats

Objective:

Apply analytical techniques to derive insights

Actions:

- Summarize data (mean, median, standard deviation).
- Perform statistical analysis (e.g., hypothesis testing, regression)

Objective: Present insights visually for clarity and impact

Actions:

- Create dashboards
- Create Charts
- Create Reports

Descriptive Analytics

- Summarizes past data to understand what has happened

Diagnostic Analytics

- Investigates data to determine why something happened

Predictive Analytics

- Uses statistical models and machine learning to forecast future trends

Prescriptive Analytics

- Suggests actions based on predictive insights to achieve desired outcomes

Data Quality Issues

- Missing, incomplete, or unreliable data

Volume and Variety

- Handling large and diverse datasets

Skill Gap

- Need for expertise in programming, statistics, and domain knowledge

Privacy and Security

- Ethical considerations and legal compliance in data handling

Data Ingestion

Data
Storage

Data
Processing

Data Sharing

Visualization &
Presentation



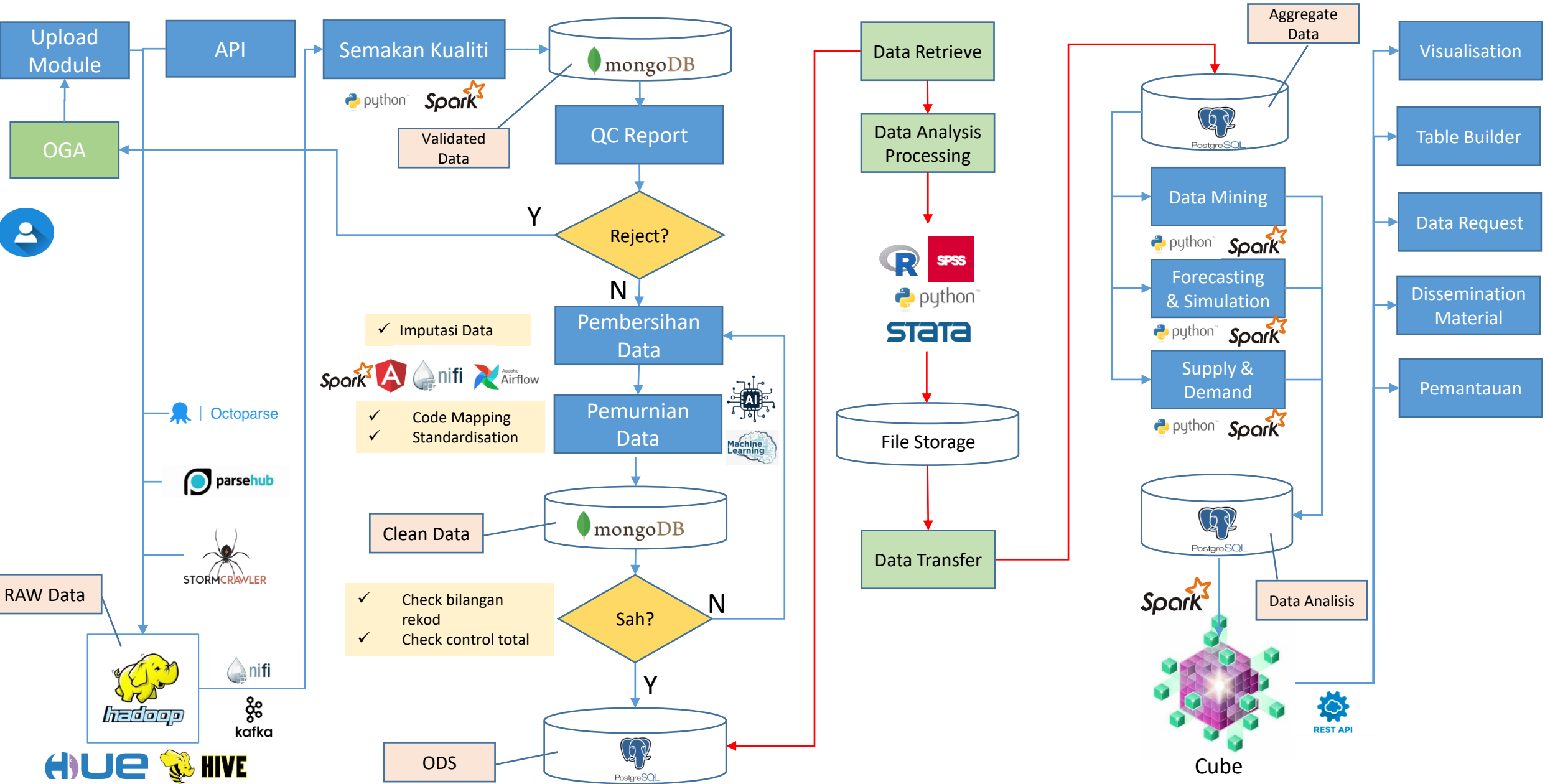
Data Collection

Transformation & Processing (ETL/ELT)

Computation & Analysis Cycle

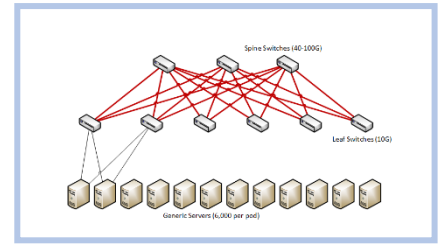
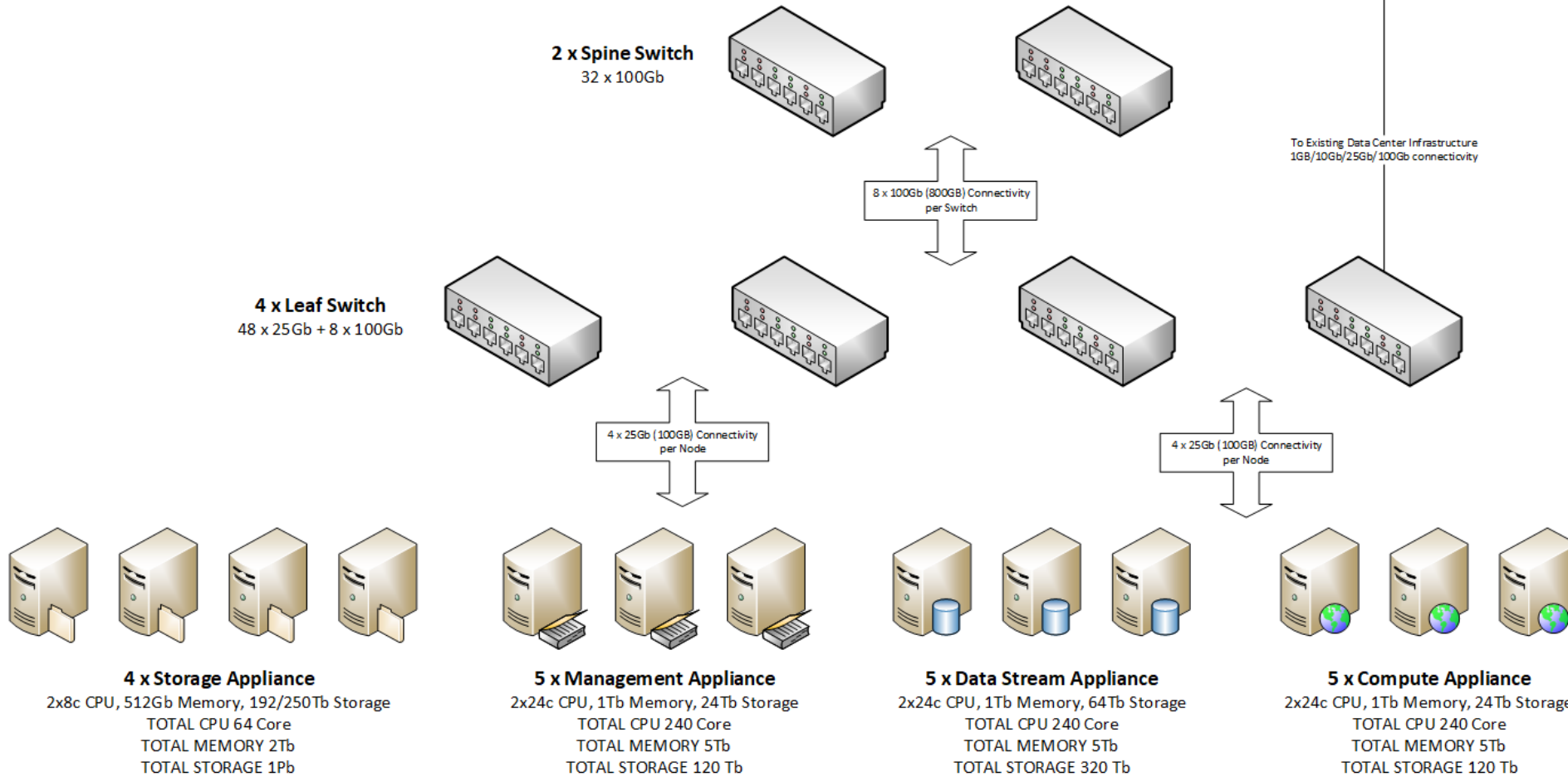
Analytics & Modelling

Output



Massive Scale Data Analytic Platform

Spine-Leaf network architecture with all software running in containers



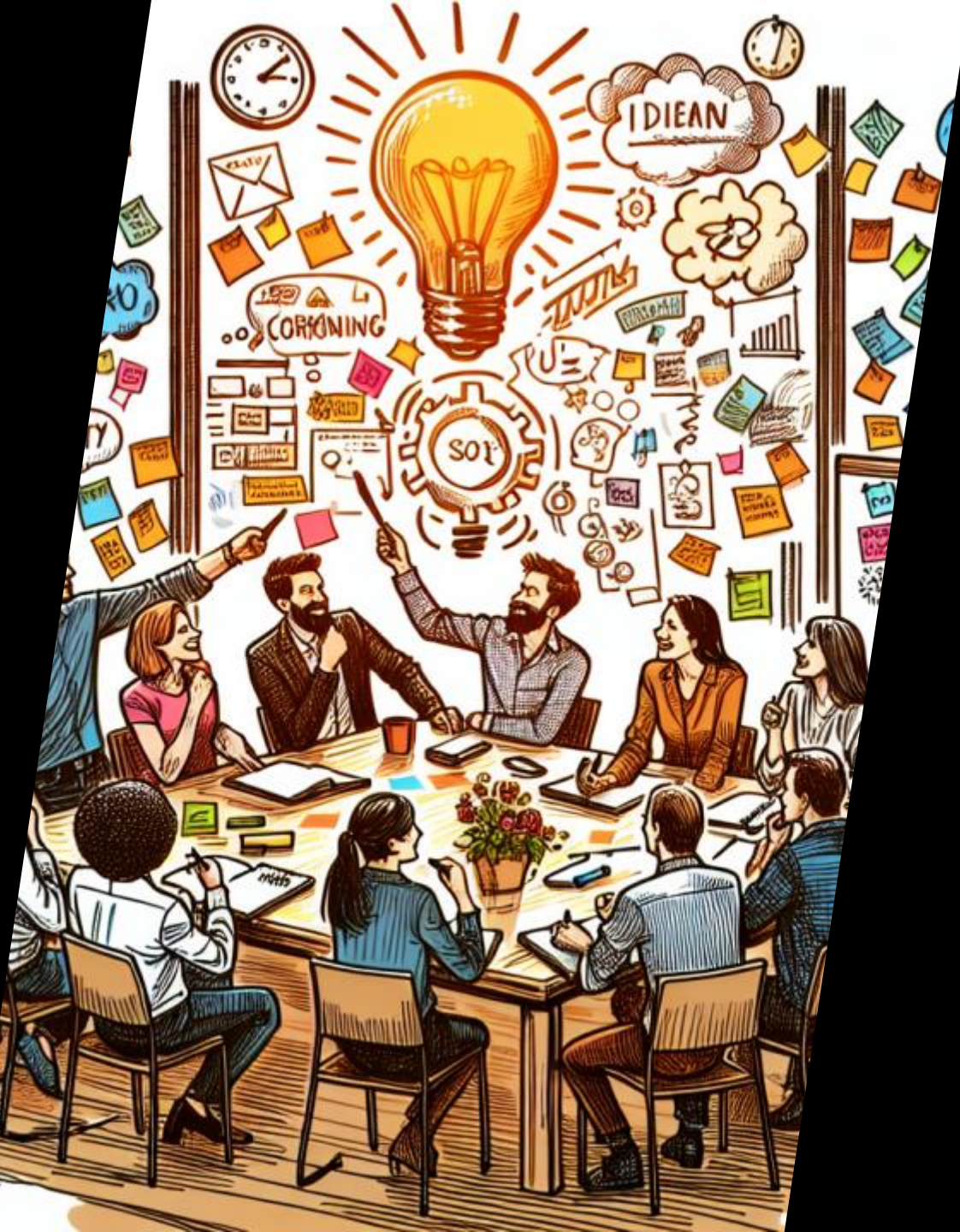
Typical Cluster Size

TOTAL CPU 1,040 Core
TOTAL MEMORY 27.8Tb
TOTAL STORAGE 1,704Tb
TOTAL APPLIANCE 24 Nodes

796,480 Nodes Max

Horizontally scale up to 760 Appliance per Pod, 1,048 pods per Cluster, for maximum of 796,480 Nodes of Data Analytic Nodes per Data Center.

Providing long-term, low cost, on-premise, private, hyperscale technology for large scale data analytic



Brainstorming Activity

