

UNIT 5: Data Mining and Warehousing

Introduction

- In today's digital era, data has become one of the most valuable assets for organizations.
- Enterprises generate vast amounts of data daily through systems like On-Line Transaction Processing (OLTP), Point-of-Service (PoS), financial ATMs, and web platforms.
- However, the challenge lies not just in collecting this data, but in managing it effectively—especially when it exists in multiple formats and needs to be transformed and loaded into a centralized system for analysis.
- To address this, two powerful technologies have emerged: Data Warehousing and Data Mining.
- Data Warehousing is the process of collecting, integrating, and storing large volumes of historical data from various sources into a unified repository.
- This data is structured in a way that supports complex queries and decision-making processes.
- The concept was first introduced by Bill Inmon in 1990, who defined a data warehouse as a subject-oriented, integrated, time-variant, and non-volatile collection of data.
- It enables knowledge workers—such as executives, managers, and analysts—to make informed decisions quickly and accurately.
- Data Mining, on the other hand, is the process of discovering meaningful patterns, trends, and relationships within large datasets.

- It uses advanced tools and techniques from statistics, artificial intelligence, machine learning, and mathematics to extract hidden insights.
- These insights can be used for predictive analytics, customer segmentation, fraud detection, and more.

Data Warehouse (DW)

- A Data Warehouse is a centralized repository that stores large volumes of structured data from multiple sources. It's designed specifically for querying, reporting, and analysis, rather than day-to-day transaction processing.

Data Mining

- Data Mining is the process of discovering meaningful patterns, trends, and relationships in large datasets using techniques from statistics, machine learning, and artificial intelligence. It's a key step in the broader process known as Knowledge Discovery in Databases (KDD).

Data Warehousing

- Data Warehousing is the process of collecting, integrating, storing, and managing large volumes of data from multiple sources into a centralized repository called a Data Warehouse (DW). This system is optimized for querying, reporting, and analysis, rather than day-to-day transaction processing.

Key Characteristics of Data Warehouse (DW)

1. **Subject-Oriented:** Focuses on specific business areas (like sales, finance, or inventory) rather than day-to-day operations. This helps in analyzing data by topic rather than by transaction.

2. **Integrated:** Combines data from multiple sources into a consistent format. This ensures uniform naming conventions, measurements, and encoding structures across the organization.
3. **Time-Variant:** Stores historical data over time, allowing for trend analysis and forecasting. Every data entry is time-stamped to track changes and support long-term decision-making.
4. **Non-Volatile:** Once data is entered into the warehouse, it is not updated or deleted. This stability ensures reliable reporting and analysis without interference from operational changes.

Goals of Data Warehouse (DW)

1. **Support Decision-Making:** Provide clean, consistent, and timely access to data for strategic and operational decisions.
2. **Centralized Data Storage:** Consolidate data from multiple sources into a single, unified repository.
3. **Historical Data Analysis:** Store time-variant data to enable trend analysis, forecasting, and performance tracking.
4. **Improve Data Quality and Consistency:** Ensure standardized formats, naming conventions, and definitions across the organization.
5. **Enable Fast Query Performance:** Optimize data structures for quick retrieval, even with massive datasets.
6. **Enhance Business Intelligence (BI):** Feed BI tools with reliable data for dashboards, reports, and analytics.

Advantages of Data Warehousing

1. **Potential High Returns on Investment:** Strategic insights from data can lead to smarter decisions and increased profitability.

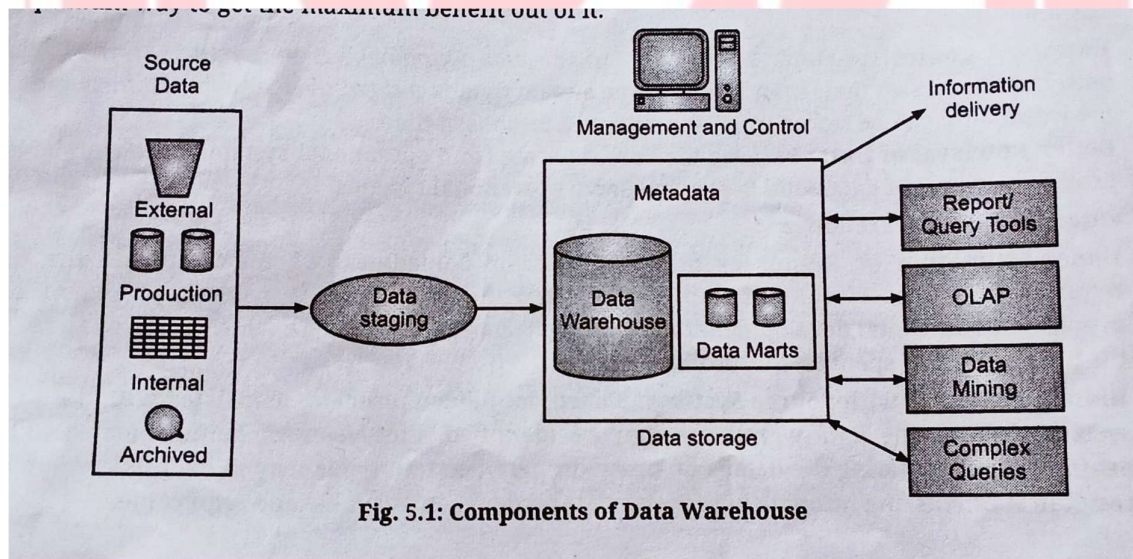
2. **Competitive Advantages:** Timely access to analytics helps businesses stay ahead of market trends and rivals.
3. **Increased Productivity of Corporate Decision Makers:** Decision-makers spend less time gathering data and more time acting on insights.
4. **More Cost-Efficient Decision-Making:** Centralized data reduces duplication and streamlines analysis, saving resources.
5. **Better Enterprise Intelligence:** Integrated data enables deeper understanding of operations, customers, and markets.
6. **Improved Control of Data:** A single source of truth enhances governance, accuracy, and data consistency.
7. **Better Retrieval of Data:** Optimized queries and structured storage allow faster access to relevant information.

Disadvantages of Data Warehousing

1. **Under-estimation of Resources for Data Loading:** Data extraction and transformation often require more time and computing power than initially planned.
2. **Hidden Problems with Source Systems:** Legacy systems may contain inconsistent, incomplete, or incompatible data that complicates integration.
3. **Required Data Not Captured:** Some critical data may be missing or unavailable, limiting the effectiveness of analysis.
4. **Increased End-User Demands:** As users see the potential, they may request more features, reports, and data access than anticipated.
5. **Data Homogenization:** Standardizing data can oversimplify or strip away valuable context from original sources.

6. **High Demand for Resources:** Warehousing requires significant hardware, software, and skilled personnel to maintain performance.
7. **Data Ownership:** Conflicts may arise over who controls and governs the data, especially across departments.
8. **High Maintenance:** Regular updates, monitoring, and troubleshooting are needed to keep the warehouse reliable and secure.
9. **Long Duration Projects:** Building and deploying a full-scale data warehouse can take months or even years.
10. **Complexity of Integration:** Merging data from diverse systems with varying formats and structures is technically challenging.

Components of Data Warehouse



1. Source Data

This is the raw input collected from various origins:

- **Production Data:** Real-time operational data from transactional systems.

- **Internal Data:** Data generated within the organization (e.g., HR, finance).
- **Archived Data:** Historical records stored for long-term analysis.
- **External Data:** Data from outside sources like market feeds or third-party vendors.

2. Data Staging Component

Prepares data before it enters the warehouse:

- **Data Extraction:** Pulls data from source systems.
- **Data Transformation:** Cleanses and formats data to match warehouse standards.
- **Data Loading:** Inserts transformed data into the warehouse.

3. Data Storage

- Final storage area for transformed data.
- Components:
 - **Data Warehouse** – Central repository using schemas (star, snowflake, galaxy).
 - **Data Mart** – Subset of warehouse for specific departments (e.g., sales, finance).
 - **Metadata** – Info about data (source, extraction time, frequency, etc.).
 - **Data Loading** – Periodic or one-time loading from staging area.
 - **Indexing** – Partitioning and indexing for fast access.

1. Information Delivery Component

- Tools and methods to present data to users.
- Includes:
 - **Query and reporting tools**

- **Analysis and visualization tools**
- **Data mining tools**

2. Metadata Component

- Stores data about data which includes:
 - Logical structures
 - File paths
 - Indexes
 - Source and transformation details

3. Data Management & Control Component

- Oversees warehouse operations and services.
- Tracks data movement and interactions.
- Coordinates updates, refreshes, and purges.
- Manages metadata and controls access paths.

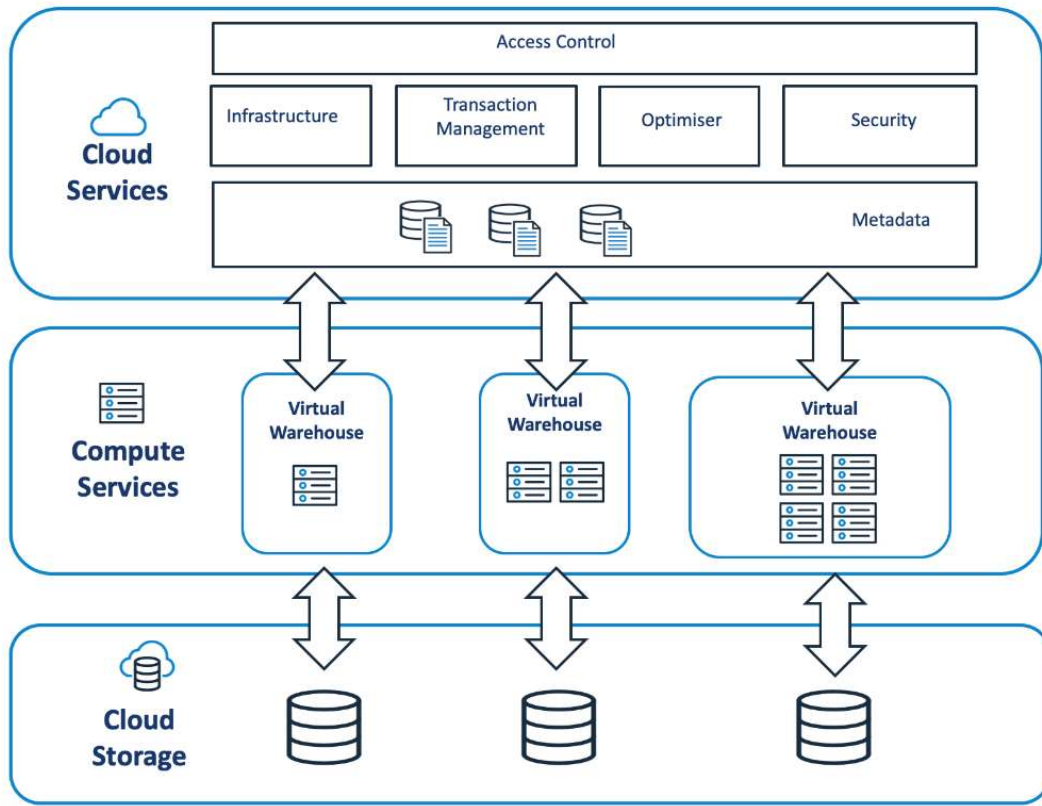
Virtual warehouse

Virtual Warehouse

- A virtual warehouse is a set of views over operational databases. For efficient query processing, only some of the possible summary views may be materialized.
- The view over an operational data warehouse is known as a virtual warehouse.
- A virtual warehouse is a digital representation of inventory, enabling businesses to manage stock across multiple locations or channels through software, providing real-time visibility and control, even if the physical inventory is stored in different places.
- Virtual warehousing leverages technology to track and manage inventory, providing a single, holistic view of stock levels, regardless of where it's physically stored.

Architecture of Virtual warehouse

- Data warehouse architecture is typically organized into layered components that streamline the flow of data from source systems to end-user analysis.
- It begins with source data—including internal, external, and archived information—which is processed through the data staging area using ETL (Extract, Transform, Load) operations.
- The cleaned and transformed data is then stored in the data storage layer, optimized for querying and analysis.
- Above this, the information delivery layer provides tools for reporting, dashboards, and OLAP operations.
- Supporting all layers are the metadata component, which describes the data, and the management and control component, which oversees performance, security, and scheduling.
- In modern cloud-based platforms, virtual warehouses act as scalable compute engines that interact with cloud storage and services to execute queries efficiently.



Layer	Interaction with Virtual Warehouse
Cloud Services	Manages query optimization, transaction control, security, access control, and metadata. It directs the VW on what tasks to perform.
Virtual Warehouse	Executes queries, transforms data, and handles compute-intensive operations. Multiple VWs can run in parallel for scalability.
Cloud Storage	Stores the actual data. VWs retrieve data from here for processing and analysis.

Benefits of Virtual warehouse

1. **Automation:** Virtual warehousing automates various inventory management tasks, such as stock replenishment and picking. This reduces the need for manual labor and eliminates human error.
2. **Enhanced Efficiency:** Optimized inventory management and streamlined supply chain operations.
3. **Improved Inventory Control:** Real-time visibility and tracking of inventory levels across all locations.
4. **Better Decision-Making:** Data-driven insights for informed decisions about inventory levels and replenishment times.
5. **Reduced Costs:** Improved inventory accuracy and reduced waste through better planning and management.
6. **Improved Data Analytics:** Monitoring the performance of multiple physical warehouses from a centralized view can provide valuable insights into the overall inventory health, allowing businesses to analyze trends and optimize stock accordingly.
7. **Increased Visibility:** In virtual warehousing systems, all operations are monitored digitally from one source, giving teams extensive visibility over their inventories in real-time. This helps increase accuracy and drive efficiency in warehouse operations.
8. **Reduced Storage Costs:** By managing stocks more efficiently through an automated system, companies can reduce storage costs associated with traditional warehousing.
9. **Increased Agility:** Virtual warehousing enables businesses to easily adapt to changing customer needs, enabling them to deliver products quickly and efficiently while providing better customer service.

Challenges of Virtual warehouse

1. **Data Security:** sensitive data need more protection
2. **Integration:** Businesses may find it difficult to integrate the necessary systems into their operations, and many companies need to hire skilled personnel to guide them through the process
3. **Initial Setup Costs:** Setting up a virtual warehouse require initial investment of time, resources and money which may be costly for some businesses.

Three-Tier Architecture of Data Warehouse

1. Bottom Tier

- Relational database system (warehouse server).
- Uses ETL tools for data extraction, cleaning, transformation, loading, and refreshing.
- Gateways like ODBC, JDBC, OLEDB enable SQL-based data access.
- Contains metadata repository for warehouse information.

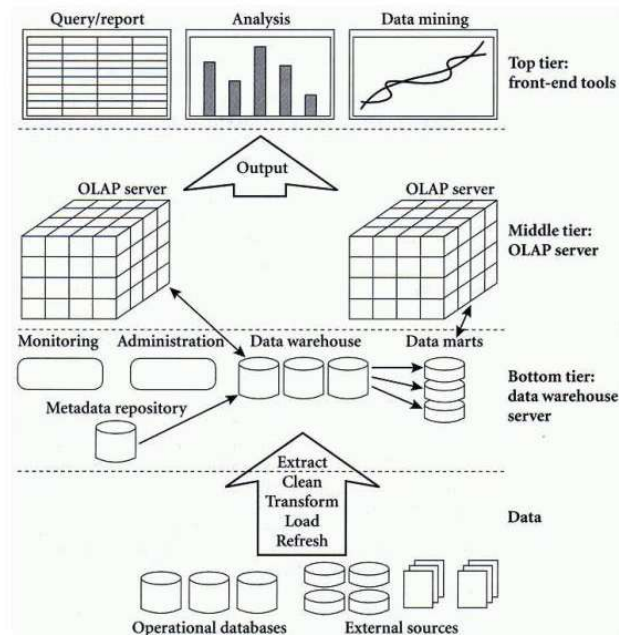
2. Middle Tier

- OLAP server for analytical processing.
- ROLAP: Extended relational DBMS.
- MOLAP: Multidimensional data engine.

3. Top Tier

- Front-end tools for query, reporting, analysis, and data mining.
- Supports trend analysis, prediction, and decision-making.

Three-Tier DW Architecture



One of the key tools used in data warehousing is ETL (Extract, Transform, Load) tools.

The ETL Process Explained



Extract-Transform-Load function involves a 5-step workflow.

1. **Data Extraction** - Collection of raw data from different sources.

2. **Data Cleaning** - Sanitisation of raw data collected from unstructured sources, flat files; removing anomalies.
3. **Data Transformation** - The raw data is cleaned, formatted, and transformed into a usable structure.
4. **Load** - The processed data is loaded into a target system, such as a data warehouse or database.
5. **Refresh** - Replaces all existing data in the target system with the latest data.

Data extraction

- It is the process of retrieving relevant information from various sources such as databases, websites, APIs, or documents.
- It serves as the first step in data processing workflows like ETL (Extract, Transform, Load), enabling organizations to consolidate raw data into a structured format for analysis, reporting, and decision-making.
- Efficient data extraction helps improve data accessibility, accuracy, and integration across systems.

Data Extraction Tools

1. **Batch processing tool:** These tools automate the execution of jobs in bulk, often scheduled or triggered by events.
Ex. Apache Hadoop, IBM DataStage
2. **Open-source tools:** Free and community-supported tools for flexible data workflows.
Ex. Apache NiFi
3. **Cloud-based tools:** Scalable tools hosted on cloud platforms, ideal for modern data ecosystems.
Ex. Azure Data Factory, Google Cloud Dataflow

Data Mining

- In practical applications, data mining is also used to mine the past and predict the future.
- Data mining is a deeper level of active design and analysis of business data.
- It is a process of using knowledge discovery tools to mine previously unknown and potentially useful knowledge.
- It is an active method of automatic discovery.
- Data mining techniques can be broadly classified **into two categories based on their purpose: Predictive and Descriptive.**

1. Prediction data mining

Predictive data mining aims to forecast future outcomes using historical data.

Techniques: Classification, regression, and time-series analysis.

Goal: To build models that predict unknown or future values.

Examples: Predicting customer churn, loan default risk, or product demand.

2. Descriptive data mining

Descriptive data mining focuses on identifying patterns and relationships within existing data.

Techniques: Clustering, association rule mining, anomaly detection.

Goal: To summarize and understand the structure of data.

Examples: Segmenting customers, discovering frequent item sets, detecting unusual transactions.

Advantages of data mining

1. **Predict Future Trends:** Helps forecast customer behavior, market shifts, and business outcomes using historical data.

2. **Decision Making:** Provides actionable insights that support strategic and operational decisions.
3. **Cost Reduction:** Identifies inefficiencies and optimizes processes to reduce operational expenses.
4. **Market-Based Analysis:** Reveals patterns in consumer preferences and buying behavior for targeted marketing.
5. **Fast and Feasible Decision:** Speeds up decision-making by automating data analysis and reducing guesswork.
6. **Signifies Customer Habits:** Uncovers hidden patterns in customer interactions to personalize services and improve retention.

Disadvantages of data mining

1. **Violates User Privacy:** Mining personal data without consent can lead to breaches of privacy and ethical concerns.
2. **Additional Irrelevant Information:** Large datasets may produce patterns that are statistically valid but practically useless, cluttering analysis.
3. **Misuse of Information:** Sensitive insights can be exploited for manipulation, discrimination, or unauthorized surveillance.
4. **Accuracy of Data (Misleading Data):** Poor-quality or biased data can result in incorrect conclusions, affecting decisions and trustworthiness.

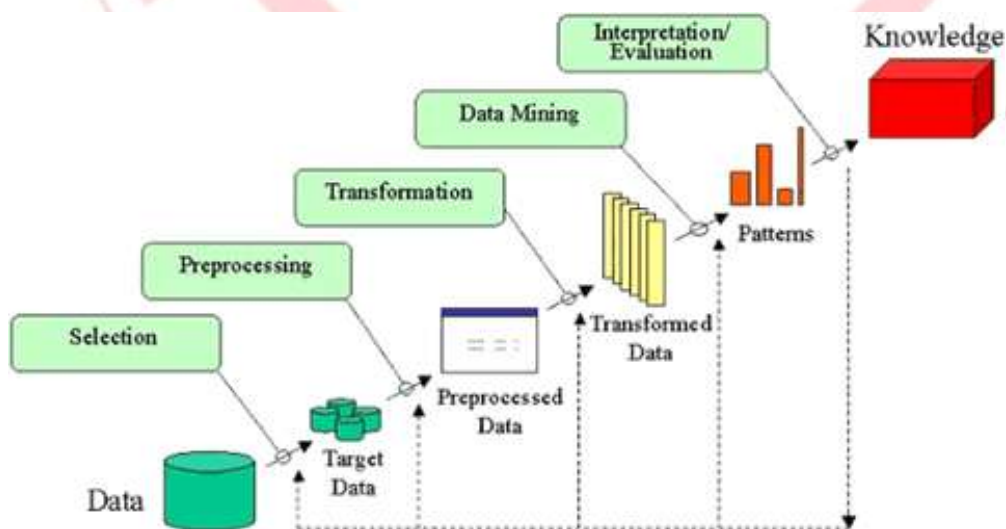
Applications of data mining

1. **Biomedical and DNA Data Analysis:** Identifies genetic patterns and disease markers to support precision medicine and drug discovery.
2. **Image Processing:** Enhances image recognition, classification, and pattern detection in fields like medical imaging and security.

3. **Financial Data Analysis:** Detects fraud, predicts market trends, and assesses credit risk using transactional and behavioral data.
4. **Manufacturing Sectors:** Optimizes production processes, predicts equipment failures, and improves quality control through sensor data analysis.
5. **Telecommunication Industry:** Analyses call records and usage patterns to reduce churn, personalize services, and detect anomalies.
6. **Government:** Supports policy-making, crime detection, and resource allocation by mining census, tax, and public service data.
7. **Website Optimization:** Tracks user behaviour to improve site structure, content relevance, and conversion rates.
8. **Data Mining in E-commerce:** Recommends products, personalizes shopping experiences, and analyses customer purchase patterns for targeted marketing.

Data mining process

The Knowledge Discovery in Databases (KDD)



1. **Data Cleaning:** Removes noise, duplicates, and inconsistencies to ensure high-quality input data.
2. **Data Integration:** Combines data from multiple sources (e.g., databases, files, APIs) into a unified dataset.
3. **Data Selection:** Filters and selects relevant data based on the mining objectives or domain requirements.
4. **Data Transformation:** Converts data into suitable formats (e.g., normalization, aggregation) for mining algorithms.
5. **Data Mining:** Applies techniques like classification, clustering, or association to discover patterns and insights.
6. **Pattern Evaluation:** Assesses the interestingness, validity, and usefulness of discovered patterns.
7. **Knowledge Presentation:** Visualizes and communicates results through charts, reports, or dashboards for decision-making.

This process is iterative and often revisited to refine results and improve accuracy.

Architecture of data mining process

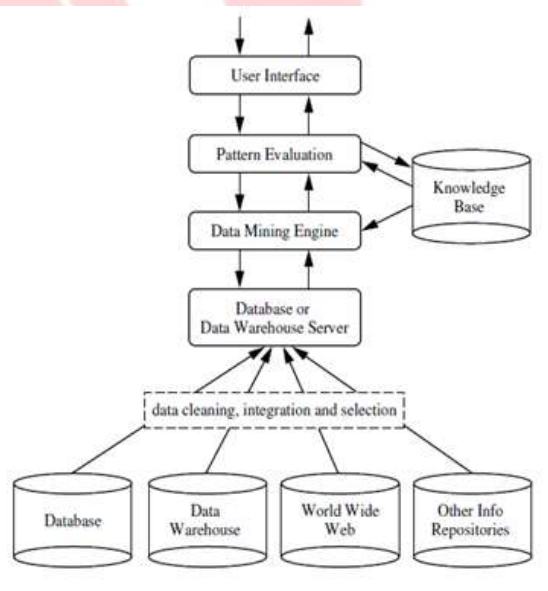
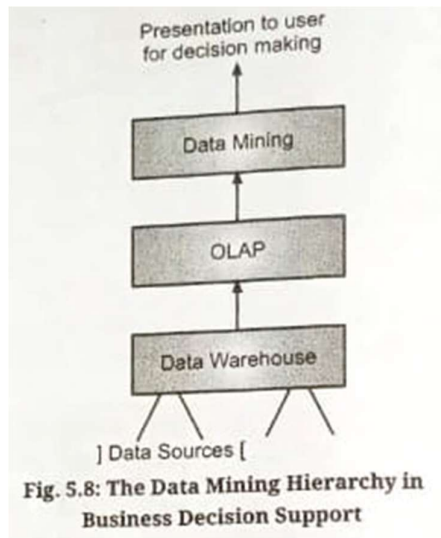


Figure 1.5 Architecture of a typical data mining system.

1. **Data Sources:** Database, World Wide Web (WWW), and data warehouse are parts of data sources. The data in these sources may be in the form of plain text, spreadsheets, or other forms of media like photos or videos. WWW is one of the biggest sources of data.
2. **Database Server:** The database server contains the actual data ready to be processed. It performs the task of handling data retrieval as per the request of the user.
3. **Data Mining Engine:** It is one of the core components of the data mining architecture that performs all kinds of data mining techniques like association, classification, characterization, clustering, prediction, etc.
4. **Pattern Evaluation Modules:** They are responsible for finding interesting patterns in the data and sometimes they also interact with the database servers for producing the result of the user requests.
5. **Knowledge Base:** Knowledge Base is an important part of the data mining engine that is quite beneficial in guiding the search for the result patterns. Data mining engines may also sometimes get inputs from the knowledge base. This knowledge base may contain data from user experiences. The objective of the knowledge base is to make the result more accurate and reliable.
6. **User Interface:** Since the user cannot fully understand the complexity of the data mining process so graphical user interface helps the user to communicate effectively with the data mining system.

Data mining technology and its relation to Data warehousing

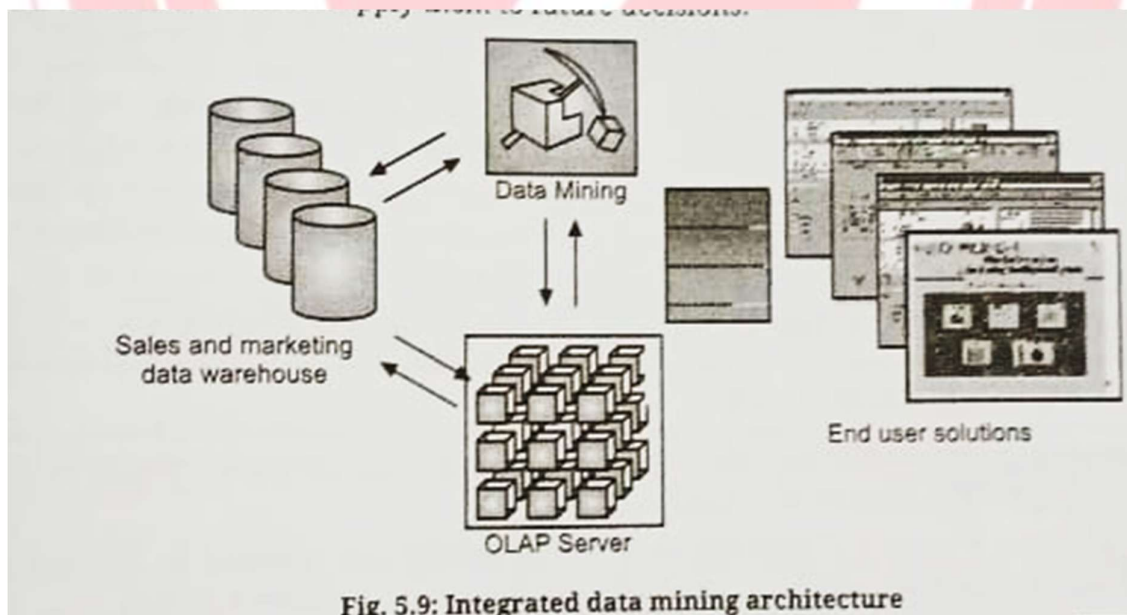


- Data mining depends on the data warehouse as its primary data source. The warehouse consolidates data, making it ready for mining.
- OLAP serves as a complementary tool—it helps verify and explore the patterns discovered by data mining.
- The hierarchy shows that while OLAP provides retrospective analysis, data mining offers predictive insights.
- This layered approach allows businesses to move from raw data → structured analysis → actionable insights.

Integrated Data mining architecture with data warehouse

- Database/Data Warehouse is the base.
- OLAP (On-line Analytical Processing) and Data Mining Engines are built on top.
- Pattern Evaluation Module and GUI help interpret and present findings.

- The tight-coupling architecture integrates all these layers seamlessly, enabling efficient data flow and analysis.
- The figure visually reinforces the idea that data mining is not standalone—it thrives when embedded within a well-structured data warehouse system.
- The **data mining server** is directly connected to the **data warehouse**.
- Queries and analysis are performed seamlessly without intermediate file storage.
- This setup allows **real-time pattern discovery** and **business rule generation**.
- How data flows from raw sources → warehouse → mining → business analysis.
- The role of mining in uncovering insights like customer behavior, fraud detection, or marketing trends.
- That data mining is not a business solution itself, but a technology enabling smarter decisions.



State key differences between data mining and data warehouse

Comparison between Data Mining and Data Warehouse:

[S-23, W-23, W-24]

Sr. No.	Data Mining	Data Warehouse
1.	Data mining is the process of analyzing unknown patterns of data.	A data warehouse is used to integrate data from multiple sources and then combine it into a single database.
2.	Data mining techniques are applied on data warehouse in order to discover useful patterns.	It provides the organization a mechanism to store huge amount of data.
3.	Data mining is the process of analyzing unknown patterns of data.	Data warehousing is the process of pooling all relevant data together.
4.	Data mining is a broad set of activities used to uncover patterns, and give meaning to this data.	Data warehouse is the repository to store data.
5.	As Facebook stores all the data in central aggregate database, now users can extract meaningful data patterns from it. That is, users can see the ads, get friends suggestions relevant to them, etc. This implies the data mining phase.	For example, Facebook gathers (collects) all user's data such as his/her friends, likes and messages, notifications. Etc., and then stores them into a central repository.

6.	This process is always carried out after data warehousing process because it needs compiled data in order to extract useful patterns.	This process must take place before the data mining process because it compiles and organizes the data into a common database.
7.	Low maintenance.	Data warehouses has high maintenance.
8.	After successful initial queries, users may ask more complicated queries which would increase the workload.	Data Warehouse (DW) is complicated to implement and maintain.
9.	The data mining methods are cost-effective and efficient compares to other statistical data applications.	Data warehouse's responsibility is to simplify every type of business data. Most of the work that will be done on user's part is inputting the raw data.
10.	Data mining helps to generate actionable strategies built on data insights.	Once we input any information into Data warehouse system, you will unlikely to lose track of this data again. You need to conduct a quick search, helps you to find the right statistic information.
11.	Data is analyzed regularly in data mining.	Data is stored periodically in DW.

Data mining concept:

- AS it transforms raw data into meaningful information.
- Common techniques include classification, clustering, regression, and association rule mining.

Mob No : [9326050669](tel:9326050669) / [9372072139](tel:9372072139) | Youtube : [@v2vedtechllp](https://www.youtube.com/@v2vedtechllp)

Insta : [v2vedtech](https://www.instagram.com/v2vedtech) | App Link | v2vedtech.com

Method of data mining Classification

- **Decision Trees:** Splits data into branches to make decisions based on feature values.
- **K-Nearest Neighbors (KNN):** Classifies data by majority vote of its closest neighbors.
- **Random Forest:** Combines multiple decision trees to improve accuracy and reduce overfitting.
- **Support Vector Machine (SVM):** Finds the optimal boundary that separates different classes.

Decision trees Classifiers

A decision tree classifier is a flowchart-like structure used in data mining to classify data based on a series of decisions. It's one of the most intuitive and widely used classification methods.

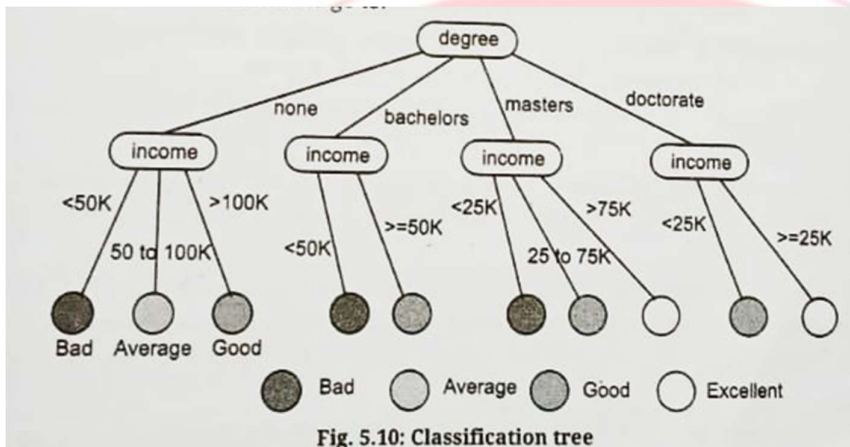
Structure of a Decision Tree:

- **Root Node:** The starting point, representing the first attribute to split on.
- **Internal Nodes:** Each node tests an attribute and branches based on its value.
- **Branches:** Represent outcomes of the attribute test.
- **Leaf Nodes:** Final decision or class label (e.g., "Yes", "No", "Average").

How It Works:

- The tree is built using a training dataset.

- At each node, the algorithm chooses the best attribute to split the data—often using metrics like information gain or Gini index.
- The goal is to create a tree that can accurately classify new data by following the path from root to leaf.



Regression

- Regression is a supervised learning technique used to predict continuous numeric values based on relationships between variables.
- It estimates how one variable affects another.
- Unlike classification (which predicts categories), regression predicts quantities.
- Example: Predicting house prices based on size, location, and number of rooms.

Types:

- **Linear Regression:** Assumes a straight-line relationship between input and output.
- **Nonlinear Regression:** Models more complex, curved relationships.

Applications:

- Forecasting trends (e.g., sales, stock prices).
- Identifying cause-and-effect relationships.
- Widely used in finance, economics, biology, and engineering.

Validating a classifier

- Validation ensures that a classifier performs accurately before it's used on real-world data.
- Prevents overfitting (too tailored to training data) and underfitting (too simplistic).
- Helps compare models and choose the best one.

Common Validation Methods:

- **Holdout Method:** Split data into training and testing sets.
- **Cross-Validation (k-fold):** Divide data into k parts, train on $k-1$, test on the remaining, repeat.
- **Bootstrap Method:** Train on random samples with replacement, test on the rest.
- These techniques help ensure your classifier is reliable and generalizes well to new data.

Association Rules (Data Mining)

- Association rules are used in data mining to uncover relationships between items in large transaction databases. They help answer questions like: *"If a customer buys item X, how likely are they to also buy item Y?"*

Rule Format:

- Expressed as: $X \rightarrow Y$
- Meaning: If itemset X appears in a transaction, itemset Y is likely to appear too.

Example:

In a supermarket:

- Transaction T1: {A, B, C}
- Transaction T5: {A, B, C, D}
- Rule: $A, B \rightarrow C$ This means customers who buy A and B are likely to also buy C.

Application:

- Market basket analysis
- Recommendation systems
- Inventory planning
- Customer behavior prediction

Clustering (Data Mining)

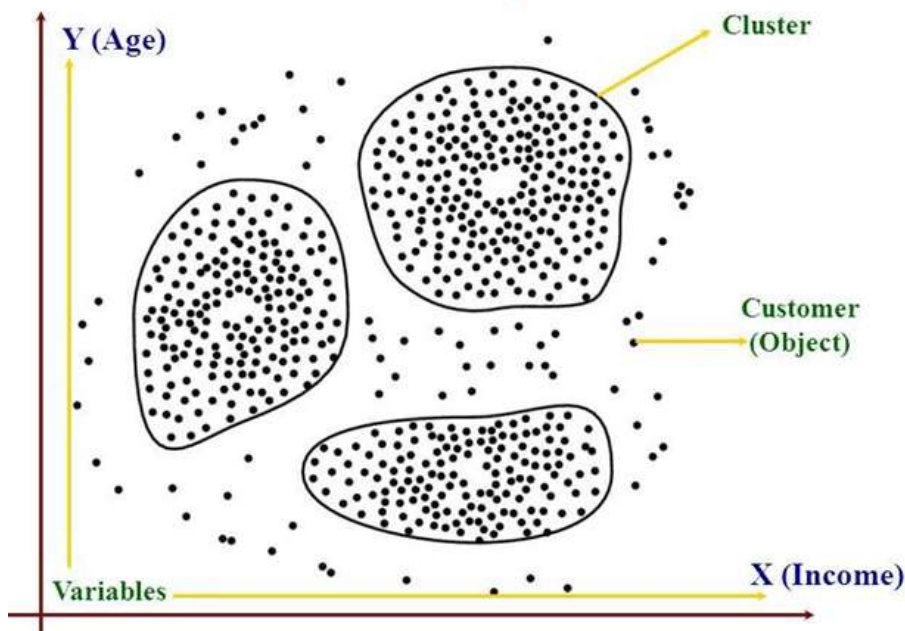
- Clustering is an unsupervised data mining technique that groups similar data objects into clusters, so that items in the same cluster are alike, and items in different clusters are distinct.
- No predefined labels—data is grouped based on similarity.
- Used in market segmentation, customer profiling, image processing, and pattern recognition.

- Helps discover hidden structures in data without prior knowledge.

Cluster analysis

- **Clustering analysis** is an unsupervised data mining technique used to group similar data objects into clusters, where:
 - Objects **within a cluster** are highly similar.
 - Objects **between clusters** are significantly different.
 - To discover hidden patterns or natural groupings in data without predefined labels.
 - Commonly used in market segmentation, customer profiling, image processing, and anomaly detection

Cluster Analysis



Clustering method in data mining:

1. **Density-Based Clustering Method:** Groups data based on regions of high density separated by regions of low density.
2. **Grid-Based Clustering Method:** Divides the data space into a grid structure and performs clustering on the grid cells.
3. **Model-Based Clustering Method:** Assumes a statistical model for each cluster and finds the best fit for the data.
4. **Constraint-Based Clustering Method:** Incorporates user-defined constraints to guide the clustering process.
5. **Partitioning Clustering Method:** Divides data into k clusters by optimizing a criterion like distance (e.g., K-means).
6. **Hierarchical Clustering Methods:** Builds a tree of clusters using either bottom-up (agglomerative) or top-down (divisive) approaches.

Other forms of Data Mining

1. **Association Rule Mining:** Discovers relationships between items frequently occurring together in transactions.
2. **Anomaly Detection:** Identifies rare or unusual data patterns that deviate from expected behavior.
3. **Sequential Pattern Mining:** Finds recurring sequences or ordered events in time-series or transactional data.
4. **Text Mining:** Extracts meaningful information and patterns from unstructured textual data.
5. **Neural Networks and Deep Learning:** Mimics brain-like structures to learn complex patterns and make predictions.
6. **Summarization:** Generates compact representations of data to highlight key insights and trends.

7. Data Visualization: Data Visualization is the graphical representation of information and data using charts, graphs, maps, and other visual formats. It transforms raw numbers into visual stories that are easier to interpret and act upon.

Benefits of Data Visualization

1. Improved understanding

Visuals simplify complex data, making patterns, trends, and outliers instantly recognizable. This helps users grasp information faster than reading through spreadsheets or reports

2. Faster Insights

By presenting data visually, decision-makers can quickly identify key metrics and anomalies, speeding up the analysis process and enabling real-time responses.

3. Better Communication

Charts and graphs make it easier to share findings with teams, stakeholders, or clients. They bridge the gap between technical data and non-technical audiences, fostering clearer discussions.

4. Enhanced decision Making

When data is visualized effectively, it supports evidence-based decisions. Leaders can make strategic choices with confidence, backed by clear, visualized data trends.

Introduction to Data Lake house

A Data Lake House is a modern data architecture that combines the flexibility of data lakes with the structured performance of data warehouses, offering a unified platform for storing, managing, and analyzing both structured and unstructured data.

- Handles all types of data—structured, semi-structured, and unstructured—in one place.
- Optimized for fast data processing and analytics.
- Reduces infrastructure costs by eliminating the need for separate systems.
- Supports various data formats and processing engines.
- Designed to manage massive volumes of data.

A data lake is a centralized storage system that holds vast amounts of raw data—structured, semi-structured, and unstructured—in its native format.

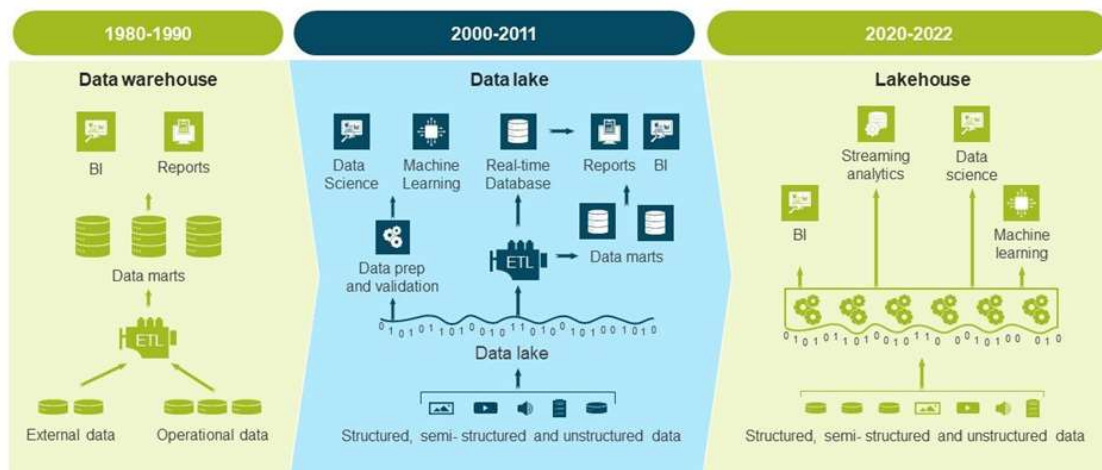
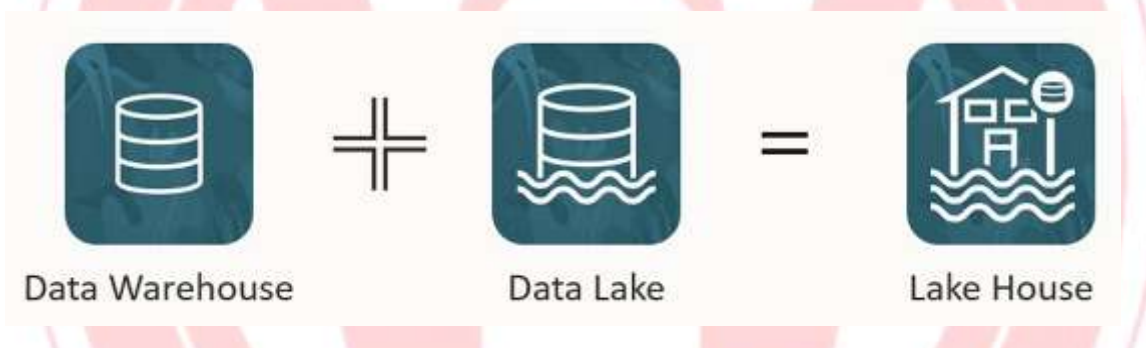
- You can store data as-is without needing to organize it first.
- Easily handles massive volumes of data, often using cloud-based platforms.
- Supports diverse data types like text, images, videos, logs, and sensor data.
- Enables various types of analysis—big data processing, machine learning, real-time analytics, and visualizations.

Comparison between Data Warehouse, Data Lake, Lakehouse

Fig 1.10 Concept of Lakehouse

Comparison between Data Warehouse, Data Lake and Data Lakehouse:

Parameters	Data Warehouse	Data Lake	Lakehouse
Data	Rational data from transactional systems, operational databases and business applications.	All data including structured, semi-structured, and unstructured.	Query every kind of data image, audio, video, and others.
Data Access	SQL only.	Open API, SQL, Python.	Open API, SQL, Python.
Data Format	Proprietary Format.	Open format.	Open Format.
Governance	Fine-grained Security.	Weak Governance and security.	Fine-grained Security.
Reliability	High with ACID Transactions.	Low Quality- Data Swamps.	High with ACID Transactions.
Performance	Fast query results using local language.	Faster query results, decoupling of computing and storage.	Faster and deeper insights without data movement.
Scalability	Scaling becomes expensive.	Low cost of scaling Regardless of the data-type.	Low cost of scaling regardless of the data-type.



Features of Data Lakehouse

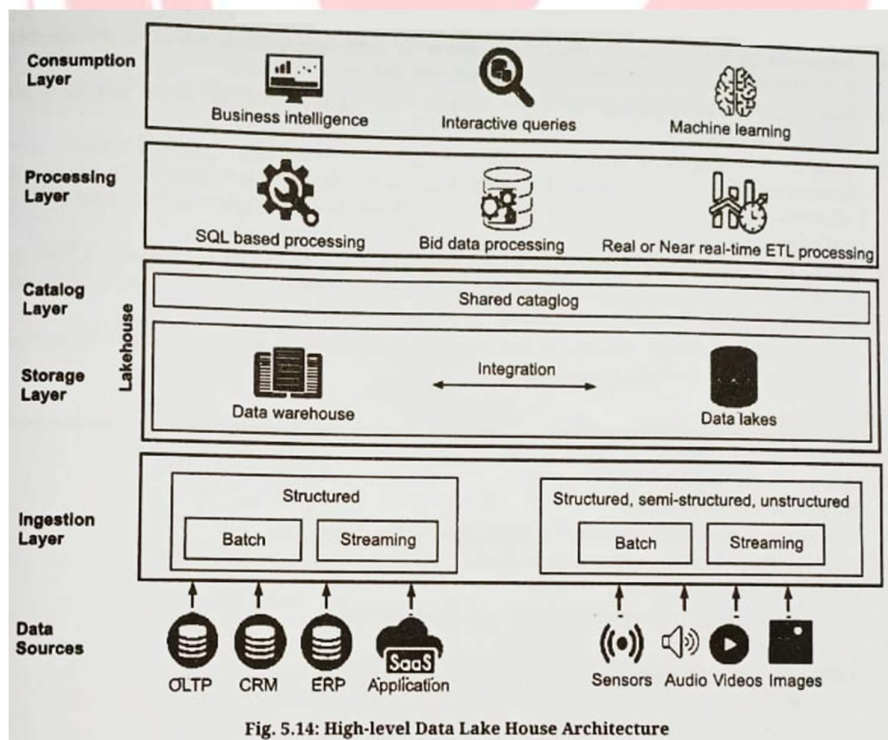
Mob No : [9326050669](tel:9326050669) / [9372072139](tel:9372072139) | Youtube : [@v2vedtechllp](https://www.youtube.com/@v2vedtechllp)

Insta : [v2vedtech](https://www.instagram.com/v2vedtech) | App Link | v2vedtech.com

1. **Transaction Support:** Ensures reliable and consistent data operations using ACID properties.
2. **Schema Enforcement:** Maintains data quality by enforcing structured schemas during storage and access.
3. **BI Support:** Seamlessly integrates with business intelligence tools for reporting and analytics.
4. **Openness:** Uses open formats and interfaces for broad compatibility and flexibility.
5. **Support for Diverse Data Types:** Handles structured, semi-structured, and unstructured data in one platform.
6. **Support for Diverse Workloads:** Efficiently manages batch processing, streaming, machine learning, and analytics.

Architecture of the Data Lakehouse

Its architecture is layered to support scalable, secure, and efficient data operations.



1. Ingestion Layer

Components: Data Connectors, Streaming Pipelines, ETL Tools

- This layer ingests data from various sources—databases, files, IoT devices, social media, etc.
- It supports batch and real-time streaming, ensuring timely and consistent data flow.
- ETL (Extract, Transform, Load) processes clean and prepare data before storage.

2. Storage Layer

Components: Object Storage (e.g., S3, ADLS), Table Formats (e.g., Delta Lake, Apache Iceberg)

- Stores raw and processed data in open formats, enabling flexibility and scalability.
- Supports structured, semi-structured, and unstructured data.
- Table formats enable ACID transactions, schema enforcement, and time travel.

3. Catalog Layer / Metadata & Governance Layer

Components: Data Catalogs, Security Policies, Lineage Tracking

- Manages metadata about datasets—schemas, ownership, access controls.
- Ensures data governance, including compliance, auditing, and lineage tracking.
- Tools like Hive Metastore or Unity Catalog help organize and secure data assets.

4. Processing Layer / Compute & API Layer

Components: Processing Engines (e.g., Spark, Presto), Query Interfaces (e.g., SQL, REST APIs)

- Executes data transformations, analytics, and machine learning workflows.
- Offers interactive and batch processing capabilities.
- APIs allow developers and analysts to query and manipulate data programmatically.

5. Consumption Layer

Components: BI Tools (e.g., Power BI, Tableau), Dashboards, Applications

- Delivers data insights to end-users through visualizations, reports, and apps.
- Supports real-time dashboards, predictive models, and embedded analytics.
- Enables cross-functional access for business, data science, and engineering teams.

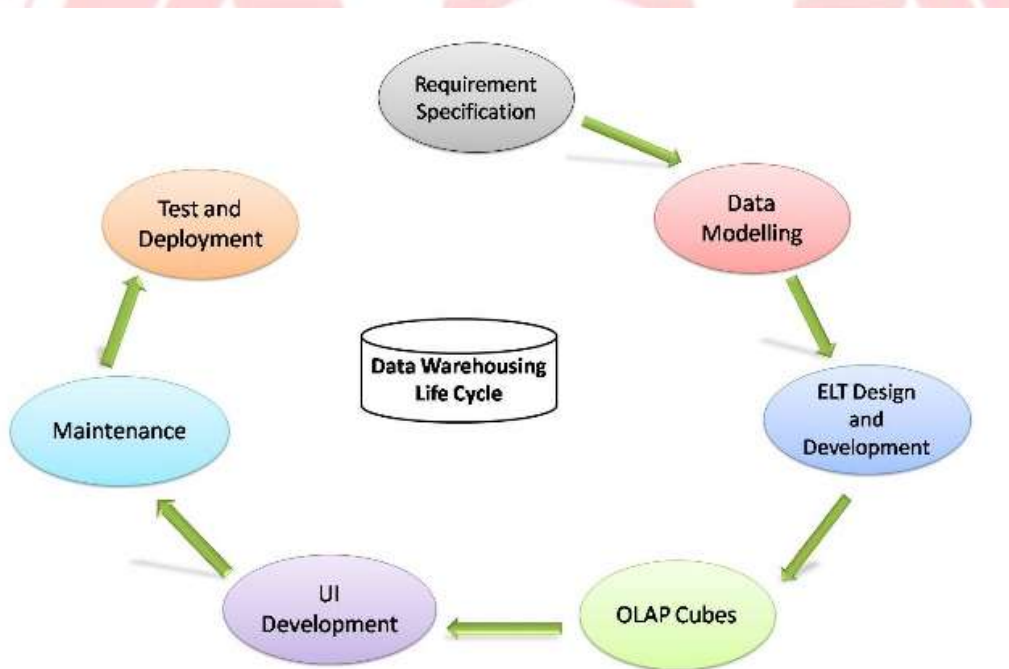
This layered architecture ensures that data is ingested, stored, governed, processed, and consumed efficiently—all within a unified platform.

Benefits of Data Lakehouse

1. **Unified Platform:** Combines storage and analytics in a single system for seamless data management.
2. **Supports Diverse Data Types:** Handles structured, semi-structured, and unstructured data effortlessly.

3. **Improved Data Governance:** Enforces policies, access controls, and lineage tracking for secure data use.
4. **Enables Advanced Analytics:** Powers machine learning, real-time insights, and predictive modeling.
5. **Simplified Architecture:** Reduces complexity by merging lake and warehouse capabilities.
6. **Better Data Quality:** Ensures consistency and accuracy through schema enforcement and metadata management.
7. **Lower Cost:** Cuts infrastructure and operational expenses by eliminating duplicate systems.

Data Warehousing life cycle



1. Requirement Gathering

- Firstly, Business analysts, the local technical lead, and the customer complete it.
- A business analyst creates a business requirement specification (BRS) document during this stage.

- Finally, after gathering the requirements, the data modeler begins to identify the dimensions, facts, and combinations depending on the needs.
- We can describe this as the data warehouse's general design.

2. Data Modelling

- This is the second step.
- This is the process of designing databases and visualizing data distribution.
- Star Schema, Snowflake Schema and Galaxy Schema are the three data models of data warehouse.
- It is the reasoning behind the storage of data in relation to other data.

3. ELT Design and development

- Thirdly a data lake can include the data that an ETL (Extract, transfer, load) tool has extracted from different source systems.
- ETL process extract the data from the lake. After that it transforms and load it into a data warehouse for reporting.
- A solid ETL procedure can be useful in creating a straightforward yet useful data warehouse that is beneficial at all organizational levels.
- ELT application design to fulfill specifications documents created during the analytical phase

4. OLAP Cubes

- This is the Fourth step.

- Also called hypercube or multidimensional cube.
- It is a data format that enables quick data analysis in accordance with the several aspects that characterize a business problem.
- From various data sources and file types, such as text files, excel sheets, multimedia files, etc., a data warehouse would extract information.
- The retrieved data is modified and cleaned before being placed into an OLAP server (or OLAP cube) for preliminary processing in preparation for additional analysis.

5. UI Development

- This is the Fifth step.
- A user interface is required for the communication between user and a computer system.
- A user interface's primary goal is to give users the ability to efficiently control the equipment or device they're using.
- There are several tools available on the market to aid in UI development. For those using BigQuery, excellent options include BI tools like Tableau or PowerBI.

6. Maintenance

- This is the Sixth step.
- In this step we can update or change the schema and the application domain or requirement of it.
- Systems for maintaining data warehouses must include tools for tracking changes to the schema, such as updates.

7. Test and Deployment

- This is the ultimate step.
- Finally, Business and organization evaluate data warehouses to determine implementation of necessary business problems.
- The warehouse testing requires examining vast amounts of data.
- Similarly, different types of data sources, including relational databases, flat files, operational data, etc., provide the data that must be compared.
- At the time of implementation, the majority of its capabilities are implemented. Additionally, both on-premises and cloud deployment options are available for the data warehouses.