# Module-2
# The Data Ecosystem

In this module, you will learn about the different types of data structures, file formats, sources of data, and the languages data professionals use in their day-to-day tasks. You will gain an understanding of various types of data repositories such as Databases, Data Warehouses, Data Marts, Data Lakes, and Data Pipelines. In addition, you will learn about the Extract, Transform, and Load (ETL) Process, which is used to extract, transform, and load data into data repositories. You will gain a basic understanding of Big Data and Big Data processing tools such as Hadoop, Hadoop Distributed File System (HDFS), Hive, and Spark.

**Learning Objectives**

- Describe and differentiate between relational and non-relational database management systems.
- Classify data structures, file formats, and sources of data by their different types.
- Explain the features and use of the different languages used by data professionals.
- Describe how Data Warehouses, Data Marts, Data Lakes, and Data Pipelines work.
- Explain how the Extract, Transform, and Load process works to make raw data ready for analysis.
- Explain what Big Data is.
- Summarize the features and use of some of the Big Data processing tools.


# Before you start

## Practically Speaking: Where Data Actually Lives

- **Welcome & Transition to Module 2**

  - Recap: Completion of **Module 1 (Foundations)**

    - Participants now move on to a **crucial next step**: understanding the **Data Ecosystem**.

  - Focus of Module 2:

    - Core concepts like:

      - **Data structures**

      - **Repositories**

      - **Processing tools**

      - **Programming languages used in data**

    - Purpose:

      - Learn **how data exists**, **where to find it**, and **how to work with it effectively**


- **Why the Data Ecosystem Matters**

  - Not just theory:

- This is the **practical foundation** of all analysis and insights.

- **Every industry** relies on understanding data ecosystems:

  - Finance 💰

  - Healthcare 🏥

  - Pharma 🧪

  - Logistics 🚚

  - Retail 🛍️

- Analogy: **Navigating a city without a map**

  - Without understanding the ecosystem:

    - You might find useful things, but **not efficiently**.

  - Like mistaking a **bus stop** for a **subway station**—you're just lost.

---

- ◆ **Case Study: A Retail Company**

  - **Problem**: Tons of disconnected data:

    - Sales data

    - Website clicks

    - Social media mentions

    - Inventory levels

  - **Core Issue**: Not a lack of data, but a lack of understanding their **data ecosystem**

  - **Solution**:

    - Mapped out data sources and flow

    - Used integration tools to connect systems

  - **Outcome**:

    - Better customer insights

    - Improved marketing strategies

    - More efficient inventory management

- ◆ **Case Study: Logistics Company**

  - ● Needs to make **real-time decisions**:

    - ○ Factors include:

      - ■ Traffic

      - ■ Weather

      - ■ Road closures

    - ○ Streams of **high-velocity data**

  - ● **Requirement**:

    - ○ Big data tools that are built for **speed and agility**

  - ● **Outcome**:

    - ○ Improved efficiency

    - ○ Ability to react and adapt quickly

---

- ◆ **Exploring the Ecosystem: The Cave Analogy** 🗺️

  - ● Explorer 1:

    - ○ No map or light

    - ○ Finds things by **accident**

    - ○ **Lost in the dark**

  - ● Explorer 2:

    - ○ Has a **headlamp and detailed cave map**

    - ○ Can **strategically navigate**

  - ● Key Idea: Understanding the data ecosystem is like **having a map and flashlight**—it makes navigation purposeful, not random.

---

- ◆ **Developing the Right Mindset**

- 📌 **Practical Tip #1: Pause Before Diving In**

  - ● When encountering new data:

- ○ **Ask**: Where does it come from?

- ○ Determine structure:

  - ■ Structured (e.g., spreadsheets)

  - ■ Semi-structured (e.g., JSON)

  - ■ Unstructured (e.g., emails, documents)

📌 **Practical Tip #2: Trace Everyday Data Flows**

- ● Example: Ordering something online

  - ○ Track data from:

    - ■ Clicking "buy"

    - ■ Warehouse processing

    - ■ Delivery tracking

- ● Builds **intuitive understanding** of how data flows through systems

---

- ◆ **Understanding Repositories: Databases, Data Warehouses, Data Lakes**

- ● Action Step:

  - ○ For each data storage type, think of a **real-world example**:

    - ■ Local library → Simple database

    - ■ Streaming platform → Data lake for viewing behavior

- ● This helps ground **abstract concepts** in tangible examples

---

- ◆ **Final Reflection Question**

- ● As you learn about:

  - ○ **Data flow**

  - ○ **Storage**

  - ○ **Processing**

- ● Reflect on:

  - ○ Your own field or project

- ○ Identify **bottlenecks** or **opportunities for improvement**

- ○ Build **data awareness**

---

📌 **Summary**

This episode of *Practically Speaking* kicks off Module 2: **The Data Ecosystem**, emphasizing why understanding data's structure, storage, and movement is foundational to any data-driven task. Through vivid analogies (navigating cities and exploring caves) and real-world case studies (retail and logistics), it illustrates how lacking ecosystem knowledge can lead to confusion or inefficiency, while understanding it empowers smart, agile decisions. Learners are encouraged to pause before analyzing data, assess its type and origin, and begin mentally mapping how data flows—even in everyday scenarios.

🔑 **Key Takeaways:**

- The **data ecosystem** is the foundation of modern decision-making in every industry.

- Not knowing your data landscape = inefficiency.

- **Data structure and source awareness** saves time and improves analysis.

- **Mapping data flow** (even in simple activities) builds ecosystem intuition.

- Real-world examples help anchor complex repository types.

---

📊 **Table: What We Learnt in the Video**

| Topic/Section | What We Learnt |
| --- | --- |
| Module 2 Focus | Exploring how data is structured, stored, processed, and moved |
| Importance Across Industries | Data ecosystems are foundational in finance, health, retail, logistics, and more |
| City/Cave Analogy | Understanding the ecosystem = having tools/maps to navigate data confidently |
| Retail Case Study | Disconnected data hurts insight; integration and mapping unlock value |
| Logistics Case Study | Fast data requires specialized, agile tools and system understanding |
| Practical Mindset Tips | Pause and assess data structure & source; mentally map data flows in daily life |
| Repository Types | Link abstract terms like "data warehouse" or "lake" to real-world analogies |
| Reflection for Learners | Apply ecosystem thinking to personal projects to spot data flow or storage issues |

# The Data Ecosystem and Languages for Data Professionals

## Overview of the Data Analyst Ecosystem

🔧 **Introduction to the Data Analyst's Ecosystem**

- The **data analyst's ecosystem** consists of:

    - Infrastructure

    - Software

    - Tools

    - Frameworks

    - Processes

- Purpose:

    - To **gather, clean, analyze, mine**, and **visualize** data.

- This video provides a **quick overview**; future videos will cover topics in depth.

---

💾 **Types of Data**

**Based on structure, data can be classified as:**

**1. Structured Data**

- Rigid and organized format (rows and columns).

- Commonly stored in:

    - **Databases**

    - **Spreadsheets**

- Examples:

    - Tables in relational databases

    - Excel files

**2. Semi-Structured Data**

- Partially organized; mix of structured and unstructured components.

- Characteristics:

    - Has some consistent elements

    - Also includes irregular, free-form data

- Example:

    - **Emails**: Structured parts (sender, recipient), unstructured parts (email body)

## 3. Unstructured Data

- No defined structure

- Complex and qualitative in nature

- Cannot be reduced to rows and columns

- Examples:

    - **Photos**

    - **Videos**

    - **Text files**

    - **PDFs**

    - **Social media content**

---

## 💾 Data Sources & File Formats

- Data is collected in a **wide range of file formats**.

- Comes from a **variety of sources**:

    - **Relational databases**

    - **Non-relational databases**

    - **APIs**

    - **Web services**

    - **Data streams**

    - **Social platforms**

    - **Sensor devices**

---

## 📋 Data Repositories

- Data repositories include:

    - **Databases**

- - **Data warehouses**

  - **Data marts**

  - **Data lakes**

  - **Big data stores**

- Choice of repository depends on:

  - Data **type**

  - Data **format**

  - Data **source**

- Example:

  - **Big data** requires:

    - Big data warehouses

    - High-velocity storage systems

    - Real-time processing frameworks

---

## 💻 Languages in the Ecosystem

- **Types of languages** used by data analysts:

  1. **Query languages**

     - Example: SQL

     - Used to **query and manipulate** data

  2. **Programming languages**

     - Example: Python

     - Used to **build data applications**

  3. **Shell & scripting languages**

     - Used for **automating repetitive tasks**

---

## ⚙️ Tools, Frameworks, and Processes

- Covers all stages of data analytics:

- ○ **Gathering**

- ○ **Extracting**

- ○ **Transforming**

- ○ **Loading (ETL)**

- ○ **Cleaning**

- ○ **Mining**

- ○ **Analyzing**

- ○ **Visualizing**

- ● Examples of tools:

  - ○ **Spreadsheets**

  - ○ **Jupyter Notebooks**

  - ○ **IBM Cognos**

- ● Emphasis on:

  - ○ **Automation**

  - ○ **Streamlined processes**

---

📌 **Summary**

This video introduces the **data analyst's ecosystem**, covering the types of data (structured, semi-structured, unstructured), data sources, repositories, tools, languages, and frameworks involved in the data analysis lifecycle. It emphasizes how data characteristics influence the choice of storage, processing, and analysis tools. The ecosystem is diverse and includes everything from databases and programming languages to visualization platforms and automation tools.

**Key Takeaways:**

- ● Data can be structured, semi-structured, or unstructured.

- ● A wide variety of data sources and formats exist.

- ● Proper storage and tools depend on the nature of data.

- ● Analysts use query, programming, and scripting languages.

- ● Tools span the full data workflow: ETL, cleaning, analysis, visualization.

---

📊 **Table: What We Learnt in the Video**

| Topic/Section | What We Learnt |
|---|---|
| Types of Data | Structured (organized), semi-structured (mixed), and unstructured (free-form data) |
| Data Sources | APIs, sensors, social platforms, databases, and more |
| Data Repositories | Databases, data warehouses, marts, lakes, big data stores |
| Language Types | SQL for queries, Python for applications, scripting for automation |
| Tools and Frameworks | Used for ETL, cleaning, mining, analysis, visualization |
| Examples of Tools | Spreadsheets, Jupyter Notebooks, IBM Cognos |

## Types of Data

📌 **What is Data?**

- **Definition**: Data is unorganized information that is processed to make it meaningful.

- **Composition**:

    - Facts

    - Observations

    - Perceptions

    - Numbers

    - Characters

    - Symbols

    - Images

- **Purpose**: Can be interpreted to derive meaningful insights.

📁 **Classification of Data by Structure**

Data can be categorized into three types based on structure:

1. **Structured Data**

2. **Semi-Structured Data**

3. **Unstructured Data**

---

1️⃣ **Structured Data**

- **Definition**: Data that has a well-defined structure or adheres to a specified data model.

- **Storage Format**:

    - Well-defined schemas

    - Relational databases (SQL)

    - Tabular form: rows and columns

- **Characteristics**:

    - Objective facts and numbers

    - Can be easily stored, exported, and organized

    - Suited to standard data analysis tools

- **Examples of Sources**:

    - SQL Databases

    - OLTP (Online Transaction Processing) Systems

    - Spreadsheets (e.g., Excel, Google Sheets)

    - Online forms

    - GPS sensors

    - RFID tags

    - Web and network server logs

- **Tools for Analysis**:

    - Standard data analysis methods and tools

    - SQL-based querying systems

---

### 2 Semi-Structured Data

- **Definition**: Data with some organizational properties but no fixed or rigid schema.

- **Storage Format**:

    - Cannot be stored in simple tabular form

    - Organized using metadata and tags

    - Hierarchical structure

- **Examples of Sources**:

    - E-mails

- XML and JSON documents

- Binary executables

- TCP/IP packets

- Zipped files

- Integrated data from multiple sources

- **Technologies**:

  - XML and JSON:

    - Allow defining tags and attributes

    - Enable hierarchical organization

    - Widely used for data exchange and storage

---

## 3 Unstructured Data

- **Definition**: Data with no identifiable structure, sequence, or format.

- **Storage Format**:

  - Cannot be stored in relational databases

  - No specific rules or structure

  - Stored in files or NoSQL databases

- **Characteristics**:

  - Heterogeneous sources

  - Valuable for business intelligence and analytics

- **Examples of Sources**:

  - Web pages

  - Social media feeds

  - Images (JPEG, PNG, GIF)

  - Video and audio files

  - Documents and PDFs

  - PowerPoint presentations

  - Media logs

- ○ Surveys

- **Storage and Analysis**:

    - ○ Manual analysis (e.g., Word docs)

    - ○ NoSQL databases (with specific tools for analysis)

---

📌 **Summary**

This video introduces the concept of data and explains its classification based on structure into three categories: structured, semi-structured, and unstructured data. Each type is defined, its characteristics outlined, and examples provided. Structured data is highly organized and easily analyzed using standard tools, semi-structured data includes metadata and hierarchy without a rigid schema, while unstructured data lacks any clear structure but holds value in analytics with appropriate tools.

**Key Takeaways**:

- Data can be raw facts, figures, and media needing processing for meaning.

- Structured data fits databases; semi-structured data uses tags; unstructured data lacks a format.

- Different sources and technologies support each data type.

---

📊 **Table: What We Learnt in the Video**

| Topic/Section | What We Learnt |
|---|---|
| Definition of Data | Unorganized information composed of facts, numbers, symbols, and media |
| Structured Data | Highly organized, tabular format, stored in SQL databases; easy to analyze |
| Semi-Structured Data | Has some structure with tags/metadata (e.g., XML, JSON); no fixed schema |
| Unstructured Data | No clear structure; includes images, videos, web content; analyzed with NoSQL tools |
| Storage Formats | Structured: SQL DBs; Semi-structured: XML/JSON; Unstructured: Files/NoSQL |
| Examples of Sources | Structured: Spreadsheets, GPS; Semi-structured: E-mails, zipped files; Unstructured: Media, docs |

## Understanding Different Types of File Formats

◆ **Introduction to Data File Formats**

- Data professionals work with various file formats.

- Understanding each file format's structure, advantages, and limitations helps in choosing the right format for specific data and performance needs.

📄 **Delimited Text File Formats**

- **Definition**: Text files with values separated by a specific delimiter.

- **Common Delimiters**:

  - Comma (,): Used in **CSV** (Comma-Separated Values)

  - Tab (\t): Used in **TSV** (Tab-Separated Values)

  - Others: Colon (:), Vertical Bar (|), Space ( )

- **CSV vs. TSV**:

| Format | Delimiter | Use Case |
|--------|-----------|----------|
| CSV | Comma | Most common, but problematic if data includes commas |
| TSV | Tab | Alternative when commas exist in data |

-
  **Structure**:

  - Each **row** = a record

  - **Columns**: Represented in the first row (header)

  - Each column can hold **different data types** (e.g., date, string, integer)

- **Advantages**:

  - Flexible field lengths

  - Standard and widely supported format

  - Easily processed by most applications

  - Useful for simple, readable schemas

📊 **Microsoft Excel Open XML Spreadsheet (XLSX)**

- **Definition**: An XML-based file format developed by Microsoft for spreadsheets

- **Key Features**:

    - File extension: .xlsx

    - Contains **workbooks** with multiple **worksheets**

    - Data is stored in **cells** at the intersection of rows and columns

    - Based on **Open File Format**: Interoperable across different platforms

    - Can use and save all Excel functions

- **Benefits**:

    - Secure (cannot store malicious code)

    - Supports complex formulas and features

    - Widely supported by analytical tools

---

## ✖️ Extensible Markup Language (XML)

- **Definition**: A markup language with a standard for encoding documents in a format that is both human- and machine-readable

- **Characteristics**:

    - Self-descriptive language

    - Does not use predefined tags (unlike HTML)

    - Platform and language independent

- **Use Cases**:

    - Data transmission over the internet

    - Facilitates easy data sharing across systems

---

## 📑 Portable Document Format (PDF)

- **Definition**: A file format developed by Adobe for document presentation

- **Key Features**:

    - Platform-independent; consistent appearance across devices

    - Commonly used in **legal** and **financial** documents

    - Can include **fillable forms**

- **Advantages**:

  - Preserves formatting

  - Universally viewable and printable

  - Cannot be easily altered

---

🔧 **JavaScript Object Notation (JSON)**

- **Definition**: A lightweight, text-based data interchange format

- **Key Features**:

  - Language-independent

  - Easy to read and write

  - Highly compatible with web technologies

- **Use Cases**:

  - Web APIs and services

  - Transferring structured data (even multimedia types like audio and video)

- **Benefits**:

  - Easy to parse

  - Efficient for data exchange

  - Supported across all programming languages and browsers

---

📝 **Summary**

In this video, we explored **six standard data file formats** that are essential for data professionals. Each format has unique structures, advantages, and limitations that make them suitable for specific use cases. Understanding these formats enables professionals to choose the best one based on data type, system compatibility, and performance needs.

**Key Takeaways**:

- Delimited files like CSV and TSV are simple and versatile.

- XLSX supports complex calculations and is secure.

- XML ensures platform-independent structured data sharing.

- PDF is ideal for consistent document presentation and legal use.

- JSON is the go-to format for modern web data exchange.

---

📑 **Table: What We Learnt in the Video**

| Topic/Section | What We Learnt |
|---|---|
| Delimited Text Files | Store data using delimiters; CSV (comma) and TSV (tab) are most common |
| CSV vs. TSV | CSVs struggle with commas in data; TSVs are better when tabs are infrequent |
| XLSX | XML-based Excel format with multi-sheet support, formulas, and high security |
| XML | Human- and machine-readable, tag-based format for platform-independent data |
| PDF | Consistent formatting across devices; commonly used for forms, legal documents |
| JSON | Web-friendly, language-independent format; widely used in APIs and for structured data |

# Sources of Data

🌐 **Overview of Modern Data Sources**

- **Dynamic and Diverse Sources**:

    - Data sources today are more varied and constantly changing.

    - Crucial for analysis and decision-making in organizations.

---

🏢 **Internal Organizational Data**

- **Enterprise Applications**:

    - Support daily business activities, HR, workflows, and customer transactions.

- **Relational Databases**:

    - Examples: **SQL Server, Oracle, MySQL, IBM DB2**

    - Structured data storage using tables, rows, and columns.

    - Often used with internal applications and data warehouses.

- **Use Cases**:

    - Retail transactions → Regional sales analysis

    - CRM systems → Sales projections

🌐 **External Data Sources**

- **Public Datasets**:
    - E.g., demographic and economic data from government agencies.

- **Commercial Datasets**:
    - Examples: Point-of-Sale data, Financial data, Weather data
    - Used for: Strategy, demand prediction, distribution, marketing decisions

---

📄 **Flat Files, Spreadsheets, and XML Datasets**

**1. Flat Files**

- **Format**: Plain text with rows as records and delimiters (e.g., CSV, TSV)
- **Structure**: Single-table mapping
- **Common Format**: CSV (Comma-Separated Values)

**2. Spreadsheets**

- **Tools**: Excel (.XLS, .XLSX), Google Sheets, Apple Numbers, LibreOffice
- **Structure**:
    - Tabular data with rows and columns
    - Multiple worksheets → Map to different tables
    - Supports formatting and formulas

**3. XML Files**

- **Markup Language**: Uses tags to identify data values
- **Supports**: Hierarchical and complex data structures
- **Examples**: Online surveys, bank statements

---

🔗 **APIs and Web Services**

- **Definition**: Interfaces allowing interaction between users/apps and datasets

- **Data Formats**: XML, JSON, HTML, plain text, media files

**Popular Use Cases:**

- **Social Media APIs**: Twitter, Facebook

    - Used for **sentiment analysis / opinion mining**

- **Stock Market APIs**:

    - Share prices, EPS, historical data for trading

- **Lookup & Validation APIs**:

    - E.g., Postal code to city/state correlation

- **Database APIs**:

    - Internal/external data extraction

---

˅ **Web Scraping**

- **Also Known As**: Screen scraping, web harvesting

- **Purpose**: Extract structured data from unstructured web sources

**Uses:**

- Product comparisons from eCommerce sites

- Generating sales leads

- Extracting forum content

- Machine learning dataset collection

**Tools:**

- **BeautifulSoup**, **Scrapy**, **Pandas**, **Selenium**

---

🔁 **Data Streams**

- **Sources**: IoT devices, instruments, apps, social media, GPS, programs

- **Characteristics**:

    - Real-time, timestamped, geo-tagged data

**Use Cases:**

- Financial tickers for trading

- Retail streams for demand forecasting

- Video feeds for threat detection

- Social feeds for sentiment tracking

- Sensor data for machine monitoring

- Web click data for UX improvement

- Flight events for rebooking systems

**Tools:**

- **Apache Kafka**, **Apache Spark Streaming**, **Apache Storm**

---

📰 **RSS Feeds (Really Simple Syndication)**

- **Usage**:

  - Deliver real-time updates from forums, news sites

- **Access**:

  - Via **Feed Readers** converting text feeds into readable updates

---

**Summary**

This video discusses the diverse and evolving landscape of data sources vital for data analytics. It covers internal databases like SQL Server and external sources like government datasets, commercial data, flat files, XML, APIs, web scraping, and live data streams. Each source has distinct formats and tools, aiding in collecting, processing, and analyzing data for strategic and operational insights.

**Key Takeaways:**

- Internal systems use structured databases for operational data.

- External sources enrich analysis with diverse, sometimes unstructured, datasets.

- APIs and web scraping enable dynamic data extraction.

- Data streams and RSS feeds support real-time analytics.

---

📊 **Table: What We Learnt in the Video**

| Topic/Section | What We Learnt |
|---|---|
| Internal Data Sources | Data from internal systems stored in relational databases used for daily operations |
| Relational Databases | SQL Server, Oracle, MySQL, DB2 – structured data in tables |
| Flat Files | CSV, TSV – plain text with delimiters, maps to single table |
| Spreadsheet Files | Excel, Google Sheets – support tabular data, multiple sheets, formatting |
| XML Datasets | Tag-based structure, supports hierarchy and complex structures |
| APIs & Web Services | Data access via interfaces (e.g., JSON, XML); social media, stock, validation APIs |
| Web Scraping | Extracts data from unstructured sources using tools like BeautifulSoup, Scrapy |
| Data Streams | Real-time data from IoT, GPS, social media; used for predictions, alerts, analysis |
| RSS Feeds | Syndicated content from websites via feed readers for real-time updates |

## Languages for Data Professionals

🌐 **Overview of Languages for Data Professionals**

Data professionals commonly use three categories of languages:

- **Query Languages** (e.g., SQL)

- **Programming Languages** (e.g., Python, R, Java)

- **Shell & Scripting Languages** (e.g., Unix/Linux Shell, PowerShell)

📌 **Why It Matters**:
Proficiency in at least one language from each category is essential for handling data-related tasks efficiently.

---

💾 **Query Languages: SQL (Structured Query Language)**

◆ **Purpose:**

- Accessing and manipulating data in databases, especially **relational databases**

◆ **Capabilities:**

- Insert, update, delete records

- Create databases, tables, views

- Write **stored procedures** (reusable sets of instructions)

◆ **Advantages:**

- **Platform-independent & portable**

- Works with many types of databases (with some vendor-specific differences)

- Simple syntax, similar to English

- Requires fewer lines of code

- Efficiently handles large volumes of data

- **Interpreter-based**: immediate execution, suitable for prototyping

- Widely supported by a large global community

---

### 🐍 Programming Languages: Python

◆ **Overview:**

- Open-source, general-purpose, high-level language

- Known for its **readability, simplicity**, and **ease of learning**

- Great for handling large-scale computations and data tasks

◆ **Key Features:**

- Supports multiple paradigms: **object-oriented, imperative, functional, procedural**

- Few lines of code to accomplish tasks

- Platform-independent (runs on Windows, Linux)

- Open-source with a vast community

- Ideal for beginners and professionals alike

◆ **Popular Libraries:**

| Task | Libraries |
|---|---|
| Data analysis | Pandas |
| Numerical computing | Numpy, Scipy |
| Data visualization | Matplotlib, Seaborn |
| Web scraping | BeautifulSoup, Scrapy |
| Image processing | OpenCV |

📊 **Programming Languages: R**

◆ **Overview:**

- Open-source language focused on **data analysis, visualization, statistics, and machine learning**

- Popular in academia and research environments

◆ **Advantages:**

- Platform-independent and extensible

- Compatible with other languages like Python

- Can handle both **structured and unstructured data**

- Provides excellent **data visualization** tools

◆ **Key Libraries & Tools:**

| Purpose | Library/Tool |
|---|---|
| Visualization | ggplot2, Plotly |
| Reporting | Dynamic reports with embedded data/scripts |
| Interactivity | Create web apps for dynamic data exploration |

☕ **Programming Languages: Java**

◆ **Overview:**

- Object-oriented, class-based, platform-independent language

- High-performance and used in **big data** frameworks

◆ **Uses in Data Analytics:**

- Data cleaning, import/export, visualization

- Many big data tools/frameworks are Java-based:

  ○ **Hadoop, Hive, Spark**

- Ideal for **speed-critical** applications

# 🖥️ Shell & Scripting Languages

### ◆ Unix/Linux Shell

- Written as a series of commands in a text file
- Useful for:
    - File manipulation
    - Program execution
    - System tasks: disk backups, log evaluation
    - Running batch jobs and installation scripts
- Quick and efficient for **repetitive tasks**

---

### ◆ PowerShell (Microsoft)

- Cross-platform automation & configuration tool
- Ideal for working with structured data: JSON, XML, CSV, REST APIs
- Consists of:
    - Command-line shell
    - Scripting language

### ◆ Key Features:

- **Object-based** pipeline: manipulate data through sorting, filtering, grouping, etc.
- Supports:
    - Data mining
    - GUI development
    - Chart creation, dashboarding
    - Generating interactive reports

---

### 📌 Summary

This video introduces the essential types of languages used by data professionals — **Query Languages, Programming Languages**, and **Shell/Scripting Languages**. It explores SQL for database interactions, Python and R for data analytics and visualization, Java for high-performance applications, and

Shell/PowerShell for automation. Mastering at least one from each category enhances a data professional's ability to manage, analyze, and automate data-related tasks effectively.

🔑 **Key Takeaways:**

- SQL is indispensable for relational database operations.

- Python is versatile and beginner-friendly with powerful libraries.

- R excels at statistical analysis and data visualization.

- Java is robust and crucial in the big data ecosystem.

- Shell scripts and PowerShell automate repetitive tasks efficiently.

---

📊 **Table: What We Learnt in the Video**

| Topic/Section | What We Learnt |
|---|---|
| Query Languages (SQL) | SQL is used for database querying, supports portability, easy syntax, and efficient data retrieval. |
| Programming Language: Python | Python is open-source, easy to learn, and ideal for analytics using libraries like Pandas and Numpy. |
| Programming Language: R | R is statistical, open-source, highly extensible, and strong in data visualization. |
| Programming Language: Java | Java is fast, used in big data (e.g., Hadoop), and supports performance-heavy tasks. |
| Shell Scripting (Unix/Linux) | Shell scripts automate file and system tasks efficiently. |
| PowerShell | PowerShell handles structured data and supports automation, dashboards, and reporting. |

## Summary and Highlights

In this lesson, you have learned the following information:

A data analyst ecosystem includes the infrastructure, software, tools, frameworks, and processes used to gather, clean, analyze, mine, and visualize data.

Based on how well-defined the structure of the data is, data can be categorized as:

- Structured Data, that is data which is well organized in formats that can be stored in databases.
- Semi-Structured Data, that is data which is partially organized and partially free form.
- Unstructured Data, that is data which can not be organized conventionally into rows and columns.

Data comes in a wide-ranging variety of file formats, such as delimited text files, spreadsheets, XML, PDF, and JSON, each with its own list of benefits and limitations of use.

Data is extracted from multiple data sources, ranging from relational and non-relational databases to APIs, web services, data streams, social platforms, and sensor devices.

Once the data is identified and gathered from different sources, it needs to be staged in a data repository so that it can be prepared for analysis. The type, format, and sources of data influence the type of data repository that can be used.

Data professionals need a host of languages that can help them extract, prepare, and analyze data. These can be classified as:

- Querying languages, such as SQL, used for accessing and manipulating data from databases.
- Programming languages such as Python, R, and Java, for developing applications and controlling application behavior.
- Shell and Scripting languages, such as Unix/Linux Shell, and PowerShell, for automating repetitive operational tasks.

# Understanding Data Repositories and Big Data Platforms

## Overview of Data Repositories

📦 **What is a Data Repository?**

- **Definition**: A general term for data that is:

    - Collected

    - Organized

    - Isolated

- **Purpose**:

    - Supports business operations

    - Enables reporting and data analysis

- **Structure**:

    - Can be small or large

    - May consist of one or multiple databases

    - Manages and stores datasets

---

🗄 **Types of Data Repositories**

### 1. Databases

- **Definition**: A collection of data/information designed for:

    - Input

    - Storage

    - Search and retrieval

    - Modification

- **Managed by**: Database Management System (DBMS)

    - A set of programs to manage the database

    - Supports **querying** to extract specific information

- **Example Use**: Finding inactive customers (≥6 months)

💬 **Important Concepts:**

- **Database vs DBMS**:

    - **Database**: Stores data

    - **DBMS**: Manages and operates on the data

    - **Often used interchangeably**, though technically different

- **Querying**:

    - Essential function to retrieve targeted data

🔍 **Types of Databases:**

| Type | Description | Key Features |
|------|-------------|--------------|
| **Relational (RDBMS)** | Organizes data in tables (rows & columns) with defined schemas | Uses **SQL**, suitable for structured data |
| **Non-Relational (NoSQL)** | Schema-less, flexible data format; ideal for unstructured and varied data | High speed, scalability; used for **big data** |

💬 **Additional Notes:**

- **SQL (Structured Query Language)** is standard for relational databases.

- **NoSQL (Not Only SQL)** emerged due to:

    - Rise of IoT, cloud computing, and social media

    - Need for handling high data volume, variety, and velocity

---

## 2. Data Warehouse

- **Definition**: Centralized repository that consolidates data from various sources

- **Main Function**:

    - Supports **analytics** and **business intelligence**

- **Uses ETL Process**:

    - **Extract**: Gather data from multiple sources

    - **Transform**: Clean and format data for use

    - **Load**: Insert data into the central repository

🧱 **Related Concepts:**

- **Data Marts**:

- ○ Subsets of data warehouses, often department-specific
- **Data Lakes**:
  - ○ Store raw, unstructured, and structured data
- **Technology Evolution**:
  - ○ Initially relational (RDBMS-based)
  - ○ Now includes **NoSQL** systems for handling diverse data types

---

## 3. Big Data Stores

- **Definition**: Data repositories built for extremely large datasets
- **Characteristics**:
  - ○ Distributed storage and computing systems
  - ○ Enable scaling and efficient processing of massive volumes of data
- **Use Cases**:
  - ○ Real-time analytics
  - ○ Processing social media, sensor data, etc.

---

🔍 **Summary**

This video introduces the concept of data repositories—structures that store, manage, and isolate data to support business intelligence and analytics. It explores **databases**, **data warehouses**, and **big data stores**, distinguishing between **relational** and **non-relational** databases. A **Database Management System (DBMS)** is emphasized for managing data, and the **ETL process** is explained as essential to data warehousing. Modern repositories accommodate not just structured but also massive and unstructured datasets, thanks to advancements in NoSQL and distributed systems.

📌 **Key Takeaways:**

- Data repositories enable structured data storage and analytics.
- Relational databases use SQL and are schema-dependent.
- NoSQL databases are flexible and built for modern data needs.
- ETL is central to data warehousing for combining diverse sources.
- Big Data Stores handle high-volume, high-velocity data at scale.

📊 **Table: What We Learnt in the Video**

| Topic/Section | What We Learnt |
|---|---|
| Data Repository | Central place to collect, store, and manage data for analysis and operations |
| Database | Organized collection of data; DBMS manages and allows querying |
| Relational Database (RDBMS) | Uses tables with schema; best for structured data; uses SQL |
| Non-relational Database | Schema-less, scalable, and flexible; ideal for big data and real-time use cases |
| DBMS | Software that creates, maintains, and queries databases |
| Querying | Mechanism to retrieve specific data using query languages like SQL |
| Data Warehouse | Merges data from various sources for analytics via ETL |
| ETL Process | Extract, Transform, Load – to clean and consolidate data |
| Data Mart & Data Lake | Specialized or raw data storage variations used in business intelligence |
| Big Data Stores | Handle extremely large datasets with distributed storage and processing |

## RDBMS

📘 **Introduction to Relational Databases**

- A **relational database** is a structured collection of data organized into tables.

- **Tables** consist of:

  - **Rows (records)**: Each represents an individual data entry.

  - **Columns (attributes)**: Each represents a data field (e.g., name, ID).

- Tables can be **linked** via common data fields (e.g., Customer ID).

- Example:

  - **Customer Table**:

    - Columns: Company ID, Company Name, Address, Primary Phone.

  - **Transaction Table**:

    - Columns: Transaction Date, Customer ID, Amount, Payment Method.

  - Link: Based on **Customer ID**.

## 🔗 Table Relationships

- Linking tables helps in:

    - Creating new result tables through **queries**.

    - Generating consolidated **reports** (e.g., customer statements).

    - Understanding **data relationships** for better insights and decisions.

- SQL (Structured Query Language) is used to interact with and **query** relational databases.

## 📊 Comparison with Spreadsheets

| Feature | Spreadsheets | Relational Databases |
|---|---|---|
| Structure | Flat files with rows and columns | Structured schema-based tables |
| Capacity | Limited rows and columns | Can handle massive datasets |
| Relationships | No built-in relationships | Tables can be related via keys |
| Redundancy | High possibility | Minimizes redundancy |
| Data Integrity | Limited controls | Enforces consistency and data types |
| Performance | Slower with large data | Optimized for storage, retrieval, and processing |

## 🔐 Data Consistency and Security

- Supports **data types and value restrictions**.

- Ensures:

    - **Data consistency**

    - **Data integrity**

    - **Security architecture**:

        - Controlled access

        - Enforced governance policies

## ☁️ Types of Relational Databases

- **Deployment types**:

- Open-source (self-supported or commercial support)

- Commercial closed-source

- Cloud-based (Database-as-a-Service)

- **Examples**:

  - On-premise: IBM DB2, Microsoft SQL Server, MySQL, Oracle, PostgreSQL

  - Cloud: Amazon RDS, Google Cloud SQL, IBM DB2 on Cloud, Oracle Cloud, SQL Azure

---

## ✅ Advantages of Relational Databases

- **Data relationships**: Ability to create meaningful data by joining tables.

- **Flexibility**: Modify structure (add columns/tables) while running.

- **Reduced redundancy**: Centralized data storage and reference through keys.

- **Backup & Disaster Recovery**:

  - Easy export/import while database is running.

  - Continuous mirroring in cloud reduces recovery time to seconds.

- **ACID Compliance**:

  - **Atomicity**: All parts of a transaction succeed or none do.

  - **Consistency**: Database remains in a valid state.

  - **Isolation**: Transactions occur independently.

  - **Durability**: Completed transactions persist despite failures.

---

## 📈 Use Cases

1. **Online Transaction Processing (OLTP)**:

   - High transaction rate

   - Supports many users

   - Frequent insert/update/delete

   - Fast query responses

2. **Data Warehouses / OLAP**:

- Historical data analysis

- Supports business intelligence

3. **IoT Solutions**:

- Requires lightweight, fast databases

- Collects data from edge devices

---

⚠️ **Limitations of Relational Databases**

- Poor support for:

  - **Semi-structured** and **unstructured** data (e.g., images, logs)

- **Migration issues**:

  - Requires identical schemas/data types between databases

- **Field length limits**:

  - Cannot store values longer than defined field length

- Not ideal for **extensive big data analytics**

---

**Summary**

This module explains the foundational concepts of **relational databases** (RDBMS). A relational database organizes data into tables linked by shared fields (e.g., Customer ID), allowing users to generate complex reports and insights through SQL queries. RDBMSs are powerful tools for structured data due to their scalability, flexibility, data integrity, and security. Common use cases include **OLTP systems**, **data warehouses**, and **IoT solutions**. Despite their strengths, they are less effective for handling unstructured or semi-structured data, and migrations require consistent schemas.

**Key Takeaways**:

- Relational databases organize data in structured, interrelated tables.

- SQL enables powerful querying and data manipulation.

- Benefits include scalability, redundancy reduction, and ACID compliance.

- Used across a variety of data-intensive domains.

- Limitations include lack of support for unstructured data and migration complexity.

---

📋 **Table: What We Learnt in the Video**

| Topic/Section | What We Learnt |
|---|---|
| Structure of RDBMS | Tables with rows (records) and columns (attributes), linked by keys |
| SQL | Primary language for querying and managing relational databases |
| Advantages | Flexibility, reduced redundancy, ACID compliance, backup & disaster recovery |
| Use Cases | OLTP systems, OLAP in data warehouses, IoT solutions |
| Cloud RDBMS | Offers scalability and minimal data loss through mirroring |
| Limitations | Not suitable for unstructured data, migration challenges, field length limits |
| Examples | IBM DB2, MySQL, PostgreSQL, Oracle, Amazon RDS, Google Cloud SQL |

## NoSQL

🔍 **Introduction to NoSQL**

- **Definition**: NoSQL stands for **"Not Only SQL"**, indicating it is not limited to SQL-based querying and relational structures.

- **Nature**: NoSQL is a **non-relational** database model.

- **Schema**: Uses **flexible schemas** for data storage and retrieval.

- **Rise in Popularity**:

    - Due to the **cloud era**, **big data**, and **high-volume web/mobile applications**.

    - Focused on **scalability**, **performance**, and **ease of use**.

- **Common Misunderstanding**: "No" doesn't mean *no SQL at all*, but *not only SQL*.

- **Data Types Supported**:

    - **Structured**

    - **Semi-structured**

    - **Unstructured**

- **Schema-Free Design**: Enables **schema-less**, flexible, and free-form storage formats.

---

📦 **Types of NoSQL Databases**

1. 🔑 **Key-Value Stores**

- **Data Format**: Stored as **key-value pairs**

- **Key**: Unique identifier for data

- **Value**: Can range from simple types to complex formats like JSON

- **Best Use Cases**:

    - **User sessions**

    - **User preferences**

    - **Real-time recommendations**

    - **Targeted advertising**

    - **In-memory caching**

- **Limitations**:

    - Poor fit for:

        - Complex queries

        - Relationships between data

        - Multiple unique keys

- **Popular Examples**:

    - Redis

    - Memcached

    - DynamoDB

---

2. 📄 **Document-Based Databases**

- **Structure**: Each record is stored as a **document**

- **Benefits**:

    - **Flexible indexing**

    - **Powerful ad hoc queries**

    - **Efficient analytics**

- **Use Cases**:

    - **eCommerce**

    - **Medical records**

    - **CRM platforms**

- ○ **Analytics tools**
- **Limitations**:
  - ○ Not ideal for **complex queries** or **multi-operation transactions**
- **Popular Examples**:
  - ○ MongoDB
  - ○ DocumentDB
  - ○ CouchDB
  - ○ Cloudant

---

3. 📊 **Column-Based Databases**

- **Data Organization**: Data stored in **columns** rather than rows
- **Column Family**: Grouping of related columns (e.g., name & profile info)
- **Advantages**:
  - ○ Fast data access due to **columnar storage**
- **Use Cases**:
  - ○ **Heavy write operations**
  - ○ **Time-series data**
  - ○ **Weather data**
  - ○ **IoT data**
- **Limitations**:
  - ○ Not ideal for:
    - ■ **Complex querying**
    - ■ Frequently changing query patterns
- **Popular Examples**:
  - ○ Cassandra
  - ○ HBase

---

4. 🌐 **Graph-Based Databases**

- **Structure**: Data modeled using **nodes** (entities) and **edges** (relationships)

- **Best For**:

    - **Connected data**

    - **Visualizing relationships**

- **Use Cases**:

    - **Social networks**

    - **Fraud detection**

    - **Real-time recommendations**

    - **Network diagrams**

    - **Access management**

- **Limitations**:

    - Not optimized for **high-volume analytics** or **transactional throughput**

- **Popular Examples**:

    - Neo4j

    - CosmosDB

---

💡 **Advantages of NoSQL**

- **Supports Diverse Data Types**: Structured, semi-structured, and unstructured

- **Distributed Architecture**:

    - Scales across **multiple data centers**

    - Leverages **cloud computing infrastructure**

- **Cost-Effective**: Designed for **low-cost commodity hardware**

- **Scalability**:

    - **Horizontal scaling** with node addition

- **Design Simplicity**: Easier to manage, more agile development

- **Flexibility & Iteration**: Quicker iterations and adaptability

## 🆚 Relational (RDBMS) vs. Non-Relational (NoSQL)

| Feature | RDBMS | NoSQL |
|---|---|---|
| **Schema** | Rigid, pre-defined | Flexible, schema-agnostic |
| **Data Type** | Structured only | Structured, semi-structured, unstructured |
| **Scalability** | Vertical scaling (costly) | Horizontal scaling (cost-effective) |
| **Hardware** | High-end, commercial | Commodity hardware |
| **Query Language** | SQL | Varies (some support SQL-like) |
| **ACID Compliance** | Fully supported | Varies (some may relax ACID for performance) |
| **Technology Maturity** | Well-established | Newer, evolving |
| **Cost** | High | Lower |

---

## 📃 Summary

NoSQL databases are **non-relational** database systems designed for **flexibility, scalability, and high-performance** in modern application environments. Unlike relational databases, they allow **schema-less storage** and support a wide variety of data types, making them ideal for use cases like **real-time web applications, IoT, and big data analytics**. There are four main types—**key-value, document-based, column-based, and graph-based**—each tailored for specific scenarios. While RDBMS is still robust and reliable, NoSQL offers a **cost-effective and scalable alternative** for dynamic and rapidly changing data needs.

**Key Takeaways**:

- NoSQL ≠ No SQL; it means **Not Only SQL**

- Ideal for handling **large-scale**, **diverse data**

- Four main types: **Key-Value**, **Document**, **Column**, and **Graph**

- Offers **flexibility, scalability, and speed**, with trade-offs in **complex transactions** and **querying capabilities**

---

## 📘 Table: What We Learnt in the Video

| Topic/Section | What We Learnt |
|---|---|
| Introduction to NoSQL | NoSQL is a flexible, scalable, schema-less database model |
| Key-Value Stores | Store data as key-value pairs; great for caching and sessions |
| Document-Based Databases | Store each record as a document; ideal for eCommerce, CRMs |

| Column-Based Databases | Use columns instead of rows; good for write-heavy operations |
|---|---|
| Graph-Based Databases | Use nodes and edges to model relationships; ideal for connected data |
| Advantages of NoSQL | Scalability, cloud-readiness, cost-efficiency, flexibility |
| RDBMS vs. NoSQL Comparison | Differences in schema, scalability, cost, ACID compliance, etc. |

## Data Marts, Data Lakes, ETL, and Data Pipelines

### 📦 Data Warehouses

- **Definition**: A centralized, structured repository for massive amounts of data from operational systems.

- **Purpose**:

  - Analysis-ready data storage

  - Single source of truth for current and historical data

  - Supports operational and performance analytics

- **Characteristics**:

  - Cleansed, structured, and conformed data

  - Multi-purpose usage (reporting, analysis)

  - High reliability and availability for business insights

---

### 🧩 Data Marts

- **Definition**: A focused subsection of a data warehouse designed for specific business units or functions.

- **Purpose**:

  - Business-specific reporting and analytics (e.g., sales or finance)

- **Advantages**:

  - Isolated performance and security

  - Tailored to specific needs of users or teams

---

### 🌊 Data Lakes

- **Definition**: A repository for storing large volumes of raw data (structured, semi-structured, unstructured).

- **Purpose**:

  - Suitable for predictive and advanced analytics

  - Used when the organization doesn't want to restrict data to predefined use cases

- **Characteristics**:

  - Stores data in native format

  - Uses metadata tagging and unique identifiers

  - Retains all source data (complete fidelity)

  - May act as a staging area for a data warehouse

---

## 🔄 ETL Process (Extract, Transform, Load)

### 🟠 Extract

- **Definition**: Collecting data from identified source systems

- **Types**:

  - **Batch Processing**:

    - Moves data in large chunks at scheduled times

    - Tools: Stitch, Blendo

  - **Stream Processing**:

    - Real-time data pull and transform in transit

    - Tools: Apache Samza, Apache Storm, Apache Kafka

### 🔵 Transform

- **Definition**: Conversion of raw data into usable formats

- **Processes**:

  - Consistent formatting (e.g., date, units)

  - Removing duplicates and unwanted data

  - Enriching data (e.g., splitting full name)

  - Establishing relationships, applying rules and validations

## 🟢 Load

- **Definition**: Transporting processed data to a data repository

- **Types**:

    - **Initial Load**: First-time full data population

    - **Incremental Load**: Periodic updates and changes

    - **Full Refresh**: Erasing and reloading data tables

- **Load Verification**:

    - Checks for missing/null values

    - Monitoring performance and failure recovery

---

## 🔁 Data Pipelines

- **Definition**: Broader term than ETL; encompasses entire data flow from source to destination

- **Types**:

    - Batch data pipelines

    - Streaming data pipelines

    - Hybrid (batch + streaming)

- **Purpose**:

    - Facilitates both real-time and batch processing

    - Supports long-running and interactive queries

- **Destinations**:

    - Data lakes (common)

    - Other applications or visualization tools

- **Popular Tools**:

    - Apache Beam

    - Google DataFlow

---

## 📌 Summary

This lesson explores how organizations handle and process vast amounts of data using systems like data warehouses, data marts, and data lakes. It delves into the ETL process—Extract, Transform, Load—as the foundation of preparing data for analysis and how it's implemented through batch or real-time methods. Finally, it introduces data pipelines as a broader concept encompassing ETL, designed to handle continuous and large-scale data flow to various endpoints like data lakes or analytical platforms.

**Key Takeaways**:

- Data warehouses store structured, analysis-ready data.

- Data marts are business-specific segments of warehouses.

- Data lakes hold raw, flexible-format data for advanced use cases.

- ETL transforms raw data into usable information via structured steps.

- Data pipelines manage the full journey of data and support both batch and real-time processing.

---

📊 **Table: What We Learnt in the Video**

| Topic/Section | What We Learnt |
|---|---|
| Data Warehouse | Centralized, structured storage; supports business analytics; source of truth |
| Data Mart | Subset of a data warehouse; business-specific reporting; isolated access |
| Data Lake | Stores raw data of all formats; flexible use; used for predictive analytics |
| ETL Process | Extracts, transforms, and loads data into a repository for analysis |
| Extract (ETL) | Batch or real-time data collection using tools like Stitch or Kafka |
| Transform (ETL) | Cleansing and formatting of data for consistency and business logic |
| Load (ETL) | Moves data into a repository; includes full, incremental, and refresh loading |
| Load Verification | Checks for data integrity and system performance during loading |
| Data Pipelines | Broader concept including ETL; supports batch and streaming workflows |
| Streaming Tools | Apache Samza, Apache Storm, Kafka for real-time data streaming |
| Data Pipeline Tools | Apache Beam and Google DataFlow for scalable pipeline management |

# Foundations of Big Data

## Introduction to Big Data

- In the digital world, every interaction leaves a trace.

- Internet-connected devices capture vast amounts of data on daily habits (travel, workouts, entertainment).

- This massive collection of data is known as **Big Data**.

## Definition of Big Data (Ernst and Young)

- Big Data: dynamic, large, and diverse volumes of data created by people, tools, and machines.

- Requires innovative, scalable technology to:

    - Collect, host, and process data analytically.

    - Generate real-time business insights related to consumers, risk, profit, performance, productivity, and shareholder value.

- No single universal definition, but common elements exist.

## The 5 V's of Big Data

1. **Velocity**

    - Speed of data accumulation.

    - Data is generated continuously and very fast (near/real-time streaming).

    - Technologies: local and cloud-based systems for rapid processing.

2. **Volume**

    - Scale and amount of data stored.

    - Growth driven by increasing data sources, high-resolution sensors, and scalable infrastructure.

3. **Variety**

    - Diversity of data types.

    - Structured data: neatly organized (e.g., relational databases).

    - Unstructured data: irregular formats like tweets, blogs, pictures, videos.

    - Data sources include machines, people, internal and external processes.

    - Drivers: mobile tech, social media, wearables, geo technologies, video, IoT devices.

4. **Veracity**

- Quality and accuracy of data.

- Attributes: consistency, completeness, integrity, and ambiguity.

- Challenges around data truthfulness and traceability.

- Importance of differentiating real vs false data.

5. **Value**

- Ability to extract meaningful benefits from data.

- Value can be monetary, medical, social, or personal satisfaction.

- The key motivation behind investing in big data analysis.

## Examples Illustrating the V's

- **Velocity:** Every minute, hours of video are uploaded on YouTube.

- **Volume:** The global digital population (~7 billion people) generates ~2.5 quintillion bytes daily (equivalent to 10 million Blu-ray DVDs).

- **Variety:** Data types include text, images, videos, health data from wearables, IoT data.

- **Veracity:** 80% of data is unstructured; it needs classification, analysis, and visualization to be useful.

## Challenges and Tools for Big Data Analysis

- Traditional data analysis tools are insufficient for massive datasets.

- Alternative tools leverage distributed computing:

  - **Apache Spark**

  - **Hadoop** and its ecosystem

- These tools help extract, load, analyze, and process big data efficiently.

- Big data analytics enable organizations to:

  - Gain better customer insights.

  - Enhance service offerings.

  - Create new knowledge.

## Conclusion

- Everyday digital interactions (smartwatches, smartphones, workouts) contribute to big data.

- Data travels globally, undergoes big data analysis, and eventually delivers insights back to users.

**Summary**

Big Data refers to the immense, rapidly growing, and varied volumes of data generated by digital interactions in today's connected world. The five key characteristics—velocity, volume, variety, veracity, and value—help define big data's nature and challenges. To handle this data effectively, new scalable technologies like Apache Spark and Hadoop are used, allowing organizations to derive valuable insights for better decision-making, customer engagement, and innovation. Understanding these core concepts is crucial, as the data we produce daily contributes to a complex global ecosystem that can deliver personalized benefits and drive business success.

**Key Takeaways:**

- Big Data is massive, fast, diverse, and requires new tech for real-time processing.

- The 5 V's (velocity, volume, variety, veracity, value) are fundamental to understanding big data.

- Big data analysis uses advanced distributed computing tools.

- Data quality and meaningful value extraction are critical challenges.

- Everyday digital device use feeds into this global data system.

**Table: What We Learnt in the Video**

| Topic/Section | What We Learnt |
|---|---|
| Introduction to Big Data | Data from daily digital interactions accumulates massively, creating Big Data. |
| Definition of Big Data | Big Data involves large, dynamic, diverse data needing scalable technologies for real-time insights. |
| Velocity | Data is generated continuously and rapidly; real-time processing is essential. |
| Volume | Huge scale of data generated daily (2.5 quintillion bytes), driven by widespread device usage. |
| Variety | Data comes in many forms (structured/unstructured) from multiple sources like social media, IoT. |
| Veracity | Data quality issues exist; accuracy and reliability are important for meaningful insights. |
| Value | The purpose of Big Data is to extract value, not just profit but also social, medical, and personal. |
| Big Data Tools | Traditional tools fail; Apache Spark and Hadoop enable processing of large, distributed data sets. |

| | |
|---|---|
| Practical Examples | YouTube uploads, wearable data, and billions of devices illustrate Big Data in everyday life. |
| Conclusion | Daily data from devices starts a global journey through analysis, returning useful insights to users. |

# Big Data Processing Tools

**Introduction to Big Data Processing Technologies**

- Big Data technologies handle large datasets: structured, semi-structured, and unstructured data.

- Purpose: Derive value from big data.

- Technologies discussed previously: NoSQL databases, Data Lakes.

- Current focus: Three open-source technologies — **Apache Hadoop**, **Apache Hive**, **Apache Spark**.

---

**Apache Hadoop**

- A Java-based, open-source framework.

- Provides distributed storage and processing of big data.

- Runs on clusters of computers (nodes = single computers; clusters = collection of nodes).

- Scalability: from single node to many nodes, each with local storage and computation.

- Supports various data formats, including emerging types:

  - Streaming audio/video

  - Social media sentiment

  - Clickstream data

  - Structured, semi-structured, unstructured data

- Enables:

  - Real-time, self-service access for stakeholders.

  - Cost optimization by consolidating data and moving "cold" data to Hadoop systems.

**Hadoop Distributed File System (HDFS)**

- Core storage system in Hadoop.

- Runs on multiple commodity hardware connected via network.

- Features:

  - Partitioning files across nodes for scalability.

  - Splitting large files into blocks, enabling parallel processing.

  - Replicates blocks on multiple nodes (default 3 replicas) for fault tolerance.

  - Data locality: moves computation close to data storage nodes to reduce network congestion and improve throughput.

- Example:

  - A phonebook file split by last name (A on Server 1, B on Server 2, etc.)

  - Blocks stored across cluster, with replicas on multiple servers.

- Benefits:

  - Fast recovery from hardware failures.

  - Supports streaming data with high throughput.

  - Accommodates large datasets (scales to hundreds of nodes).

  - Portable across hardware and operating systems.

---

**Apache Hive**

- Open-source data warehouse software built on Hadoop.

- Used for reading, writing, and managing large datasets stored in HDFS or systems like Apache HBase.

- Optimized for:

  - Long sequential scans.

  - Data warehousing tasks such as ETL (Extract, Transform, Load), reporting, and data analysis.

- Provides SQL-like query interface.

- Limitations:

  - High latency queries due to Hadoop's batch processing nature.

  - Not suitable for transaction processing (high write operations).

  - Less ideal for applications needing very fast response times.

**Apache Spark**

- General-purpose data processing engine.

- Designed for large-scale data processing and analytics.

- Key use cases:

    - Interactive analytics

    - Stream processing

    - Machine learning

    - Data integration

    - ETL

- Performance:

    - Uses in-memory processing for speed.

    - Spills to disk only when memory is insufficient.

- Supports multiple programming languages:

    - Java, Scala, Python, R, SQL.

- Deployment:

    - Can run standalone or on top of Hadoop and other clusters.

- Data Access:

    - Can read data from HDFS, Hive, and various other sources.

- Strength:

    - Real-time complex analytics.

    - Fast processing of streaming data.

**Summary**

This video explained three key open-source big data technologies: **Apache Hadoop**, **Apache Hive**, and **Apache Spark**. Hadoop offers a scalable distributed storage and processing framework that handles various data formats efficiently using its HDFS component. Hive provides a SQL-like interface for querying and analyzing large datasets stored on Hadoop, optimized for batch data warehousing but with high latency. Spark complements Hadoop and Hive by enabling fast, real-time analytics and complex processing through its in-memory computation capabilities and multi-language support, making it ideal for streaming

data and interactive tasks.

**Key Takeaways:**

- Hadoop enables scalable, fault-tolerant distributed data storage and processing.

- Hive simplifies big data querying using SQL but is suited for batch analytics with high latency.

- Spark accelerates big data processing using in-memory analytics and supports real-time, complex workloads.

---

**Table: What We Learnt in the Video**

| Topic/Section | What We Learnt |
|---|---|
| Big Data Processing Tech | Technologies manage large datasets (structured, semi-structured, unstructured) for value extraction. |
| Apache Hadoop | Distributed storage and processing framework using HDFS; scalable, fault-tolerant, supports multiple data types. |
| HDFS | Splits large files across nodes; replicates blocks for fault tolerance; enhances parallel processing and data locality. |
| Apache Hive | Data warehouse on Hadoop for SQL-based querying; good for ETL, reporting; high latency; not suitable for transactions. |
| Apache Spark | Fast in-memory distributed analytics engine; supports real-time processing, machine learning, and streaming; multi-language support. |

# Summary and Highlights

In this lesson, you have learned the following information:

A Data Repository is a general term that refers to data that has been collected, organized, and isolated so that it can be used for reporting, analytics, and also for archival purposes.

The different types of Data Repositories include:

- Databases, which can be relational or non-relational, each following a set of organizational principles, the types of data they can store, and the tools that can be used to query, organize, and retrieve data.
- Data Warehouses, that consolidate incoming data into one comprehensive storehouse.
- Data Marts, that are essentially sub-sections of a data warehouse, built to isolate data for a particular business function or use case.
- Data Lakes, that serve as storage repositories for large amounts of structured, semi-structured, and unstructured data in their native format.
- Big Data Stores, that provide distributed computational and storage infrastructure to store, scale, and process very large data sets.

ETL, or Extract Transform and Load, Process is an automated process that converts raw data into analysis-ready data by:

- Extracting data from source locations.
- Transforming raw data by cleaning, enriching, standardizing, and validating it.
- Loading the processed data into a destination system or data repository.

Data Pipeline, sometimes used interchangeably with ETL, encompasses the entire journey of moving data from the source to a destination data lake or application, using the ETL process.

Big Data refers to the vast amounts of data that is being produced each moment of every day, by people, tools, and machines. The sheer velocity, volume, and variety of data challenge the tools and systems used for conventional data. These challenges led to the emergence of processing tools and platforms designed specifically for Big Data, such as Apache Hadoop, Apache Hive, and Apache Spark.