

## Data Engineering Lifecycle

- Schema: Defines the hierarchical organization of data
- Schemaless: Application defines the schema as data is written
- Fixed schema: Enforced in the DB to which application writes must conform
- Ingestion:
  - Push: Source system writes data out to a target
  - Pull: Data is retrieved from a source system
- Featurization: Extract and enhance data features useful for training ML models
- BI:
  - Describe a business's past and current state
  - Data is stored in a clean but fairly raw form with minimal postprocessing business logic
- Operational Analytics:
  - Fine-grained details of operations, e.g., live view of inventory
  - Focused on present and doesn't concern historical trends
- Security: Principle of least privilege
- Data Management: Encompasses the set of best practices DE will use to accomplish the task of managing the data lifecycle technically and strategically
- Data Governance: Ensure quality, integrity, security, and usability of the data collected by an organization
- Metadata:
  - Business: Non-technical questions about who, what, where, and how and provides a DE with the right context and definitions to properly use the data
  - Technical: Describes the data created and used by systems across the DE lifecycle
  - Data Lineage: Tracks the origin and changes to data
  - Schema: Describes structure of data stored in a system
  - Operational: Used to determine whether a process succeeded or failed and the data involved in the process
- Data Accountability: Assigning an individual to govern a portion of data
- Data Quality:
  - Quality tests, ensuring data conformance to schema expectations, data completions, and precision
  - Accuracy: Is the data factually correct? Duplicated values? Are the numeric values accurate?
  - Completeness: Are the records complete? Do all required fields contain valid values?
- Master Data: Business entities (employees, customers etc.)
- Master Data Management: Practice of building consistent entity definitions known as golden records
- Data Modeling and Design: Process for converting data into usable form
- Data Lineage:
  - Provides a trail of data's evolution as it moves through various systems and workflows
  - Tracks both the systems that process the data and the upstream data it depends on
- Data Integration: Process of integrating data across tools and processes. Happens through general-purpose APIs rather than DB connections
- Orchestration: Coordinating many jobs to run as quickly and efficiently as possible on a scheduled cadence (DAGs)

## Designing Data Architecture

### Architecture Concepts

- Domain: Real-world subject area, can contain multiple services
- Service: Set of functionality whose goal is accomplish a task, each service has particular tasks
- Scalability: Increase capacity of a system to improve performance
- Elasticity: Ability of a scalable system to scale dynamically
- Availability: Percentage of time an IT service or component is in an operable state
- Reliability: System's probability of meeting defined standards in performing
- Horizontal scaling: Add more machines to satisfy load and resource requirements where a leader node distributes tasks to worker nodes → Increase availability and reliability
- Vertical scaling: Increase CPU, memory etc. → Cannot offer high availability and reliability
- Architecture tiers:
  - Single tier: Database and applications are tightly coupled, residing on a single server. Good for testing but not for production
  - Multi tier: Composed of multiple layers (data, applications, business logic etc.), that are bottom-up and hierarchical
  - Three tier: Consists of data, application logic, and presentation tiers
  - Monoliths: Consists of a single codebase running on a single machine that provides both the application logic and user interface
  - Microservices: Comprise separate, decentralized, and loosely coupled services where each service has a specific function and is decoupled from other services

### Types of Data Architecture

- Data Warehouse:
  - Central data hub used for reporting and analysis
  - Data is typically highly formatted and structures for analytics use cases
  - Organizational DW: Organizes data associated with certain business team structure and processes
  - ETL:
    - \* Extract: Pull data from source systems
    - \* Transform: Clean and standardize data, organize and impose business logic
    - \* Load: Push data into the DW target DB systems/data marts that serve the analytical needs
  - ELT:
    - \* Raw data gets moved directly from production into staging area in the DW
    - \* Transformations are handled directly in the DW
    - \* Intention: Take advantage of massive computational power of cloud DWs
    - \* Data is processed in batches, and the transformed output is written into tables, views etc.
- Cloud Data Warehouse: Data is housed in object storage, allowing virtually limitless storage
- Data Mart: More refined subset of a warehouse designed to serve analytics and reporting, focused on single sub-organizations, departments etc. (each organization has its own DM)
- Data Lake: Central location for (un-/semi-)structured data

- Data Lakehouse:
  - Incorporates the controls, data management, and data structures found in a DW while housing data in object storage and supporting variety of querying and transformation engines
  - ACID (atomicity, consistency, isolation, durability) transactions
- Data Mesh: Divide of data between operational and analytical data
- Lambda Architecture:
  - Systems operating independently of each other- batch, streaming, serving
  - Source system is immutable and append-only, sending data to two destinations for processing (stream and batch)
  - In-stream processing intends to serve data with the lowest possible latency
  - Batch layer processes data and transforms it into a system (e.g., DW)
  - Serving layer provides a combined view by aggregating query results from the two layers
- Architecture for IoT:
  - Distributed collection of devices with an internet connection
  - Devices: Collecting and transmitting data to a downstream destination
  - IoT gateway: Hub for connecting devices and securely routing devices to appropriate destination on the internet