# 1) What is Data Engineering:

**the modern data Ecosystem:**
- the Core concept is that data is a valuable asset
- The Ecosystem's purpose is to turn raw data into "actionable insights" for decision making.
- this is achieved through a collaboration of several specialized roles.

**• Key players in data Ecosystem:**
- data Engineer: build the highway for data  *the Architect*
- data analyst: drives on the highway to report on traffic *the storyteller*
- data scientist: use the highway to predict where new cities will be built *predictor*
- Business Analyst: decides why we need the highway and where it should go for economic growth. *translator*

**• data Engineer tasks:**
- collecting source data
- processing data
- storing data
- making data Available to users securely.

# 2) the Data Ecosystem and languages for Data professionals:

**• Data types & Sources:**
- Structured: rigid Schema (DBtable, spread)
- Semi-structured: Flexib.. with tags (Json xml emails)
- unstructured: not defin schem (img, video, social media posts)
- * Source: rational/non rational DB, APIs, web Scraping, Datastreams (IoT, GPS)

**• Data repositories (storage):**
- OLTP (transactional): For day to day operation (e.g banking transaction) optimized for fast writes.
- OLAP (Analytical): For complex analysis and reporting. optimized for fast reads on large Datasets.
  * Data Warehouse: central repository for structured cleaned data from multiple Sources.
  * Data Mart: subsection of a data warehouse
  • Data lake: stores Vast amounts of raw data in its native format (struct - nonstruct - semistructured)
  entire process of moving data from source to destin
  • Data Process & pipelines:
  - ETL: data transformed before loaded (like Data warehouse)
  - ELT: data loaded first (datalake) by data

---

**• Essential Languages: SQL, python, R, shell scripting.**

**• comparison:**

| repository types | Description & purpose | Key characteristics | example |
|---|---|---|---|
| DB | collecting data for input Storage, search, modificat | relational RDBMS: table Format, SQL, ACID compliant. - non relational (No SQL) schemas less, scales easily | RDBMS: OLTP (mysQL PostgreSQL) NoSQL: Large, diverse fast changing data |
| D Warehouse | structured | - optimized for OLAP - Analysis ready | Business Analy (Business intell & reporting) googlebig Query, snowflake |
| Data (Mart) | | | |
| D Lake | massive amount of raw data (stru/unstur/semi) | - stores everything - highly scalable. | AWS S3 Azure Data-Lake. |
| Big data | handle high Volum, Velocity, Variety. | - handles high Volum Velocity, Variety | Hadoop HDFs |

**• No SQL DB types:**
- Key-Value store (Redis, Dynamo DB)
- Documented based (like JSON) (MongoDB, Couch DB)
- column-based (Cassandra HBase.
- Graph-based uses Nodes (data) and edges (Neo4j, Cosmos DB).

RDMS (SQL) when have struct data and need ACID compliance.
-notonly
NOSQL you have high Volume, Variety, or need high scalability and Flexibility (N+ACID)

**• Data Gathering & Import methods:**
- SQL: extract data from relational databases (Select, Filter, sort, Group.)
- APIs: how applic programm acces data from an endpoint (like web service or data marketplace).
- Web Scraping: process Automat to download data from websites (text, contact info, images, .. from unstructured web pages.
- Data streams: Handling continuous real time data flows from sources (IoT, sensor, social media...).

---

**• Data exchange:** using Formel platforms to securely share data.
- where import data?!
* structured: from relation DB
* unst/semi: Geos into NoSQL or Data lakes

**• Data Wrangling (Data cleani)** process to make raw Data usable it involves Exploration, transformation, Validation.
Key tasks: pandas-numpy
Join: combin column from diff tabl
Union: .. rows from diff tables
Normalization: removing dupl data
Denormalization: combine data from multiple tables.
cleaning tasks:
handle missing Values, remove duplicate data, remove irrele data, fix syntax error, standard data (make text lowercase) convert datatypes.

**• Governance & Compliance** to ensure that Data secure, private, trustworthy & other laws & regulation