



Data Warehousing



Who am I?

Data Engineer @ TwentyToo AI

Ex Data Engineer @ GMC

Data Consultant @ Dynamu

Instructor @ Information Technology Institute (ITI) — Egypt

Instructor @ National Telecommunication Institute (NTI) — Egypt

Instructor @ CREATIVA — Egypt

Instructor @ DEPI — Egypt

Teaching Assistant @ Faculty of Artificial Intelligence, Cairo University

Coached students @ Faculty of Engineering, Helwan University

Lecturer @ Joseph Institute for Wireless Officers

 [LinkedIn.com/in/MohamedARoshdy](https://www.linkedin.com/in/MohamedARoshdy)



Agenda

- **Introduction to the Program**
- **Understanding Data-related Job Roles**
- **Data Journey**
- **Data Engineer: Bridging TEch and business.**
- **OLTP Vs OLAP.**
- **History of Data Warehousing > Bill Inmon Vs Kimball.**
- **Dependent Vs Independent data mart.**
- **DWH Architectures.**
- **ETL VS ELT**
- **Surrogate vs Natural Key**
- **Preparing the Environment for Labs**
- **Q&A**

Icebreaker / Warm-up

What do you think a Data Engineer actually does?



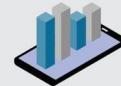
Data Engineer



Data Scientist



Data Analyst



Guess the Role

Overview of Data Roles

Data Analyst vs Data Scientist vs Data Engineer

	Data Analyst	Data Scientist	Data Engineer
Responsibilities	<ul style="list-style-type: none">• Accurate data• Visualize & report data	<ul style="list-style-type: none">• Source data• Analyze data• Run experiments	<ul style="list-style-type: none">• Build data pipelines and warehouse• Manage scalability of data products
Skills	<ul style="list-style-type: none">• Analytics• Communication & Visualization	<ul style="list-style-type: none">• Analytics• Model building• Math & Coding	<ul style="list-style-type: none">• Coding• Model implementation
Tools	<ul style="list-style-type: none">• SQL• Excel• Tableau	<ul style="list-style-type: none">• Python• Machine learning• Tableau & SQL	<ul style="list-style-type: none">• Java• C++ & Python• Hadoop & Spark

Be a Careful Data Engineer – It’s a Mindset



"Data Engineering is not just about tools and coding – it's a mindset."



State of Data Engineering 2022 map

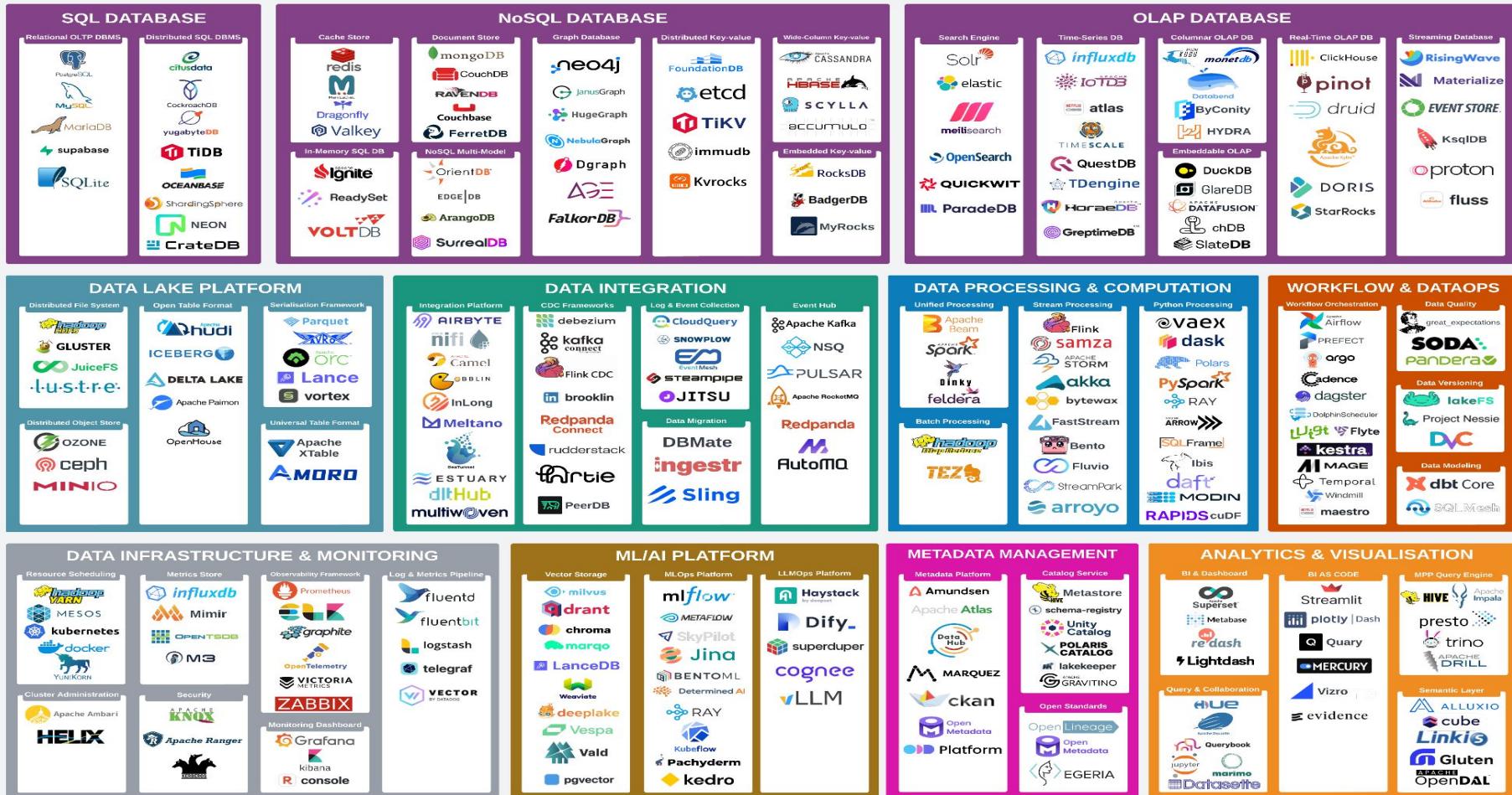


Presented by lakeFS

This image is a collage of logos for various data and machine learning platforms, organized into categories:

- Ingest SaaS**: Stitch, Airbyte, Fivetran, Segment, Rivery, Azure Data Factory, Dataddo, datacoral, MATILLION, SNOWPLOW.
- Object Storage**: amazon S3, Microsoft Azure Blob Storage, Google Cloud Storage, ORACLE CLOUD, IBM Cloud Object Storage, Alibaba Cloud, hadoop, MINIO, zadara, SEAGATE, filebase, kafka, PULSAR, beam, BENEATH, Flink, upsolver, StreamNative, CONFLUENT, tinybird, Stream Analytics.
- Metastore**: HIVE, Cloud Dataproc, Azure Purview, Ambari Glue, CLOUDERA, databricks Unity, Git for Data: lakeFS, Project Nessie, Open Table Formats: ONEHOUSE, hudi, CrowdStorage, wasabi, DigitalOcean, filebase, SwiftStack, ceph, LYVE, Azure Data Lake Storage Gen2, PURE STORAGE.
- Compute**: databricks, Spark, Cloudwick, RAY, Amazon EMR, Cloud Dataproc, Azure HDInsight, akka, trino, DASK, ASCEND.IDO, bodoai, CLOUDERA, Analytics Engine: snowflake, amazon ATHENA, amazon REDSHIFT, Google Big Query, Synapse, databricks, dremio, druid, pinot, star.tree, ClickHouse, Starburst, pentaho, ICEBERG, Tabular, FIREBOLT, ORC.
- Orchestration**: Airflow, Flyte, PREFECT, Wigt, dagster, ASTRONOMER, MONTE CARLO, lightup, Datafold, KENSU, great_expectations, Databand, HoloClean, Bigeye, unravel, WHYLABS, griffin, Metaplane, awslabs/deequ, redata, SODA, elementary, acceldata, timeseek.AI.
- MLOps End-to-End**: colab, OctoML, ABACUS.AI, METAFLOW, Driadat, Verta, cnvrg.io, mlflow, snorkel, SELDON, Hugging Face, DataRobot, W&B, HOPSWORKS, Kubeflow, Google Data Studio, CLEARML, Valohai, nephele.ai, DOMINO, data iku, Amazon SageMaker, FLOYDHUB, iguazio, hydrophere.io, ZenML, Michelangelo, Kedro, Qwak, comet, H2O.ai.
- Data Centric AI/ML**: DAGsHub, DVC, activeloop, Pachyderm, GRAVITI, LIGHTLY.
- Feature Stores**: RASGO, TECTON, FEAST, KASKADA, MOLECULA, HOPSWORKS, scribbleData, redis.
- ML Observability**: deepchecks, mona, Superwise, fiddler, Apres, arize, WHYLABS, Arthur, galileo, ROBUST INTELLIGENCE, GANTRY, truera, aporia.
- Notebooks**: Deepnote, HEX, dataform, dbt, jupyter, count, noteable, databricks, Querybook.
- Analytics Workflow**: Acryl Data, SELECT STAR.
- Discovery & Governance**: MARQUEZ, boom!, Qwak, Collibra, BigID, magda., Apache Atlas, Open Metadata, ckan, OKERA, Alation, MetaCat, Amundsen, data.world, IBM Watson, Cloud Dataproc, IMMUTA, DataHub, atlan, privacer, Platform, Acryl Data, SELECT STAR.

OPEN SOURCE DATA ENGINEERING LANDSCAPE 2025



Whom You'll Work With!!

Data Engineer: Bridging TEch and business.

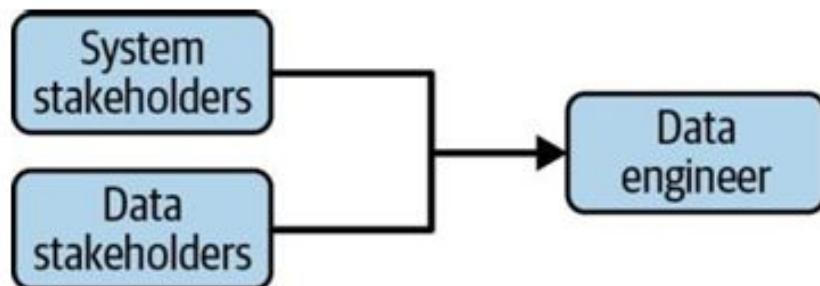


Figure 5-12. The data engineer's upstream stakeholders

A great Data Engineer !

Key Points:

A great Data Engineer understands the business as much as the technology.

Think beyond pipelines and storage – focus on delivering real value.

Bridge the gap between business needs and technical solutions.

Always ask: “How will this data solve a real problem?”

Data Journey

Files Era (Flat Files, CSV, Logs)

Data stored in files.

Simple but hard to integrate.

No standard structure → inconsistency problems.



From Files to Data Engineering – Supporting Analytics

Files Era (Flat Files, CSV, Logs)

Data stored in files.

Simple but hard to integrate.

No standard structure → inconsistency problems.



From Files to Data Engineering – Supporting Analytics

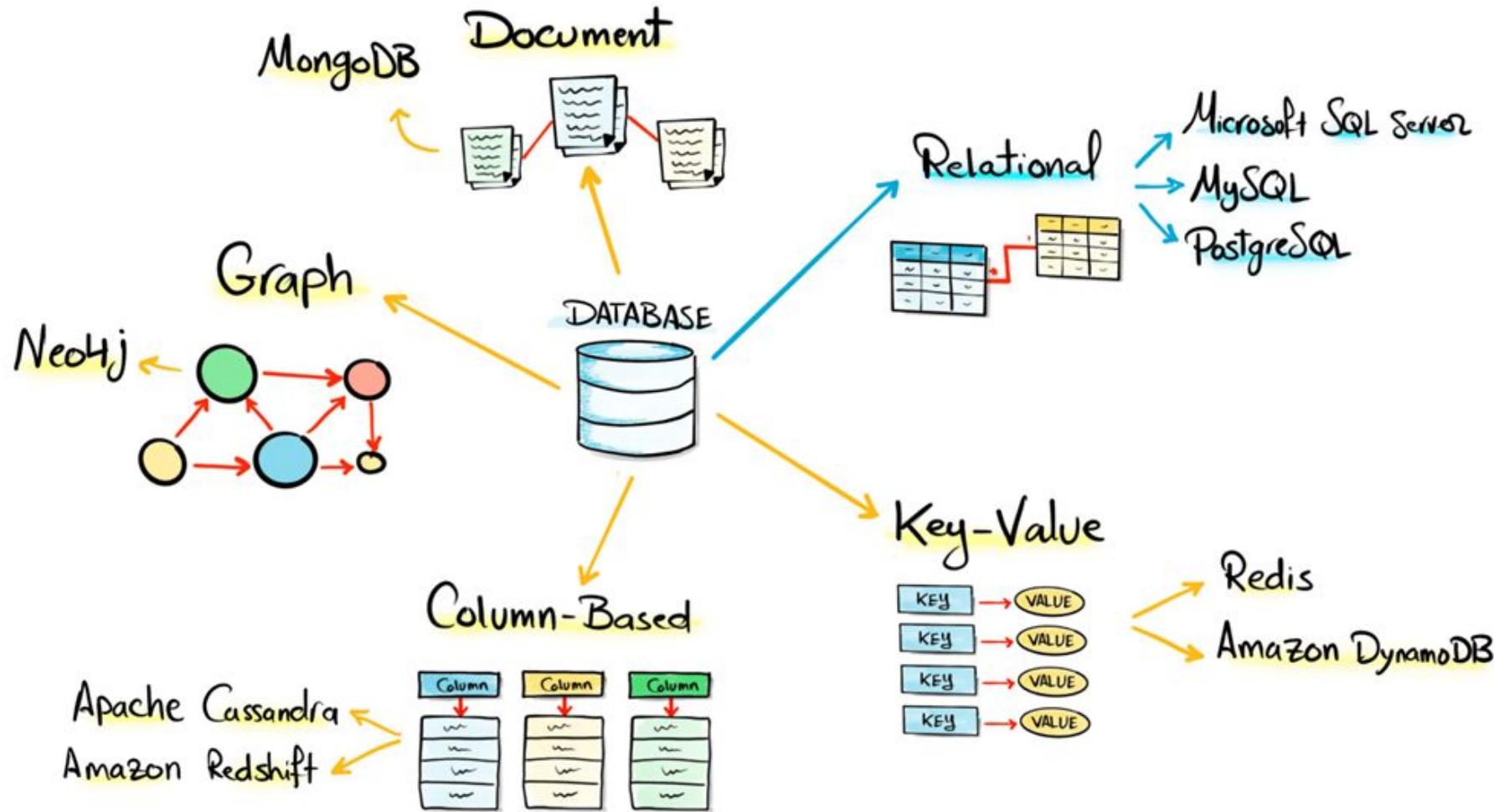
Databases Era (SQL Databases)

Solved integration & consistency issues.

Structured schema, easy querying.

Great for operational systems, but limited for large-scale analytics.





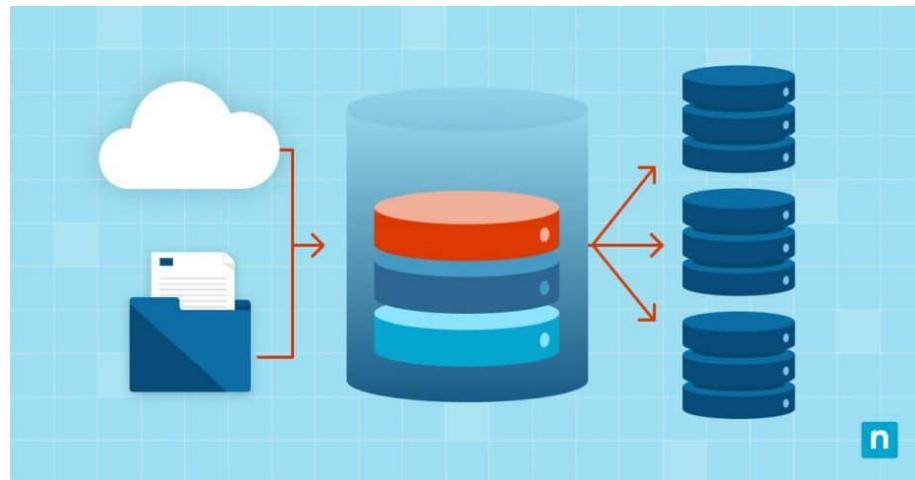
From Files to Data Engineering – Supporting Analytics

Data Warehouse Era

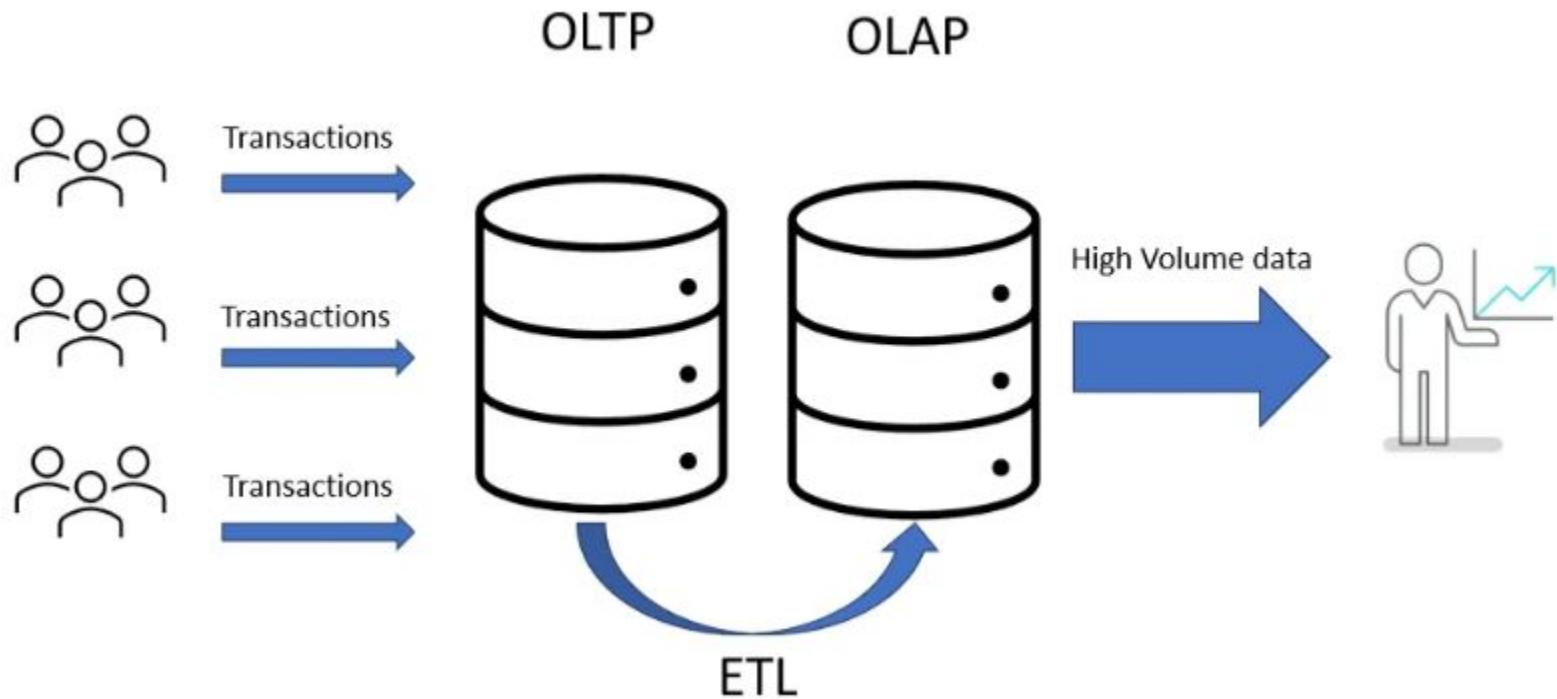
Designed for analytics, not operations.

Integrates data from multiple sources.

Optimized for large-scale queries & historical analysis.



OLTP Vs OLAP



History of Data Warehousing and “Bill Inmon”



INTERVIEW

BILL INMON

FATHER OF DATA
WAREHOUSE

History of Data Warehousing and “Bill Inmon”

1. Before Data Warehousing (Pre-1980s)

Data stored in separate operational (OLTP) systems.

Difficult to consolidate and analyze data across systems.

2. Birth of the Concept (1980s)

Bill Inmon introduced the idea of the Enterprise Data Warehouse (EDW).

Definition: "A subject-oriented, integrated, time-variant, and non-volatile collection of data to support decision-making."

History of Data Warehousing and “Bill Inmon”

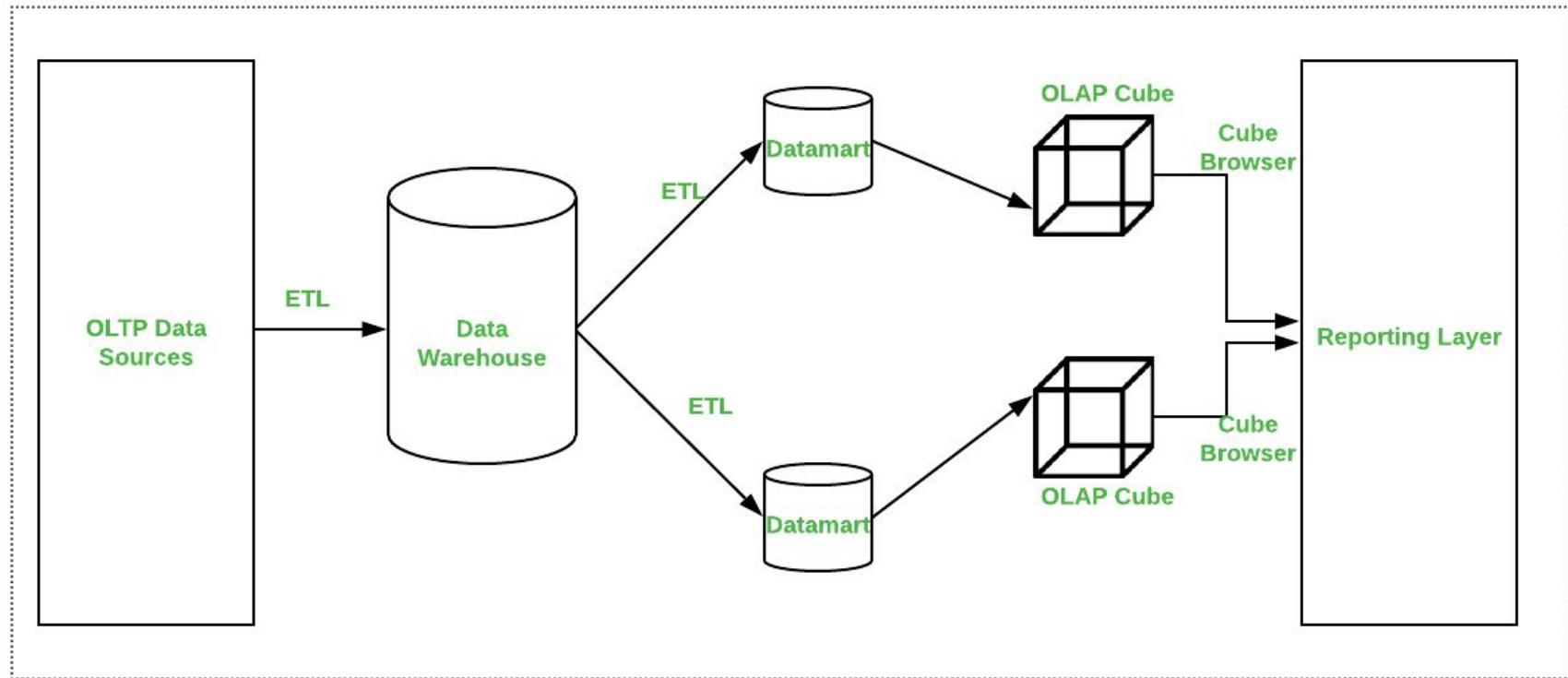
3. The Inmon Approach

Build EDW first (Normalized – often up to 3NF).

From EDW, derive Data Marts (Star or Snowflake) for analytics.

Pros: High data integration, single source of truth for the enterprise.

History of Data Warehousing and “Bill Inmon”

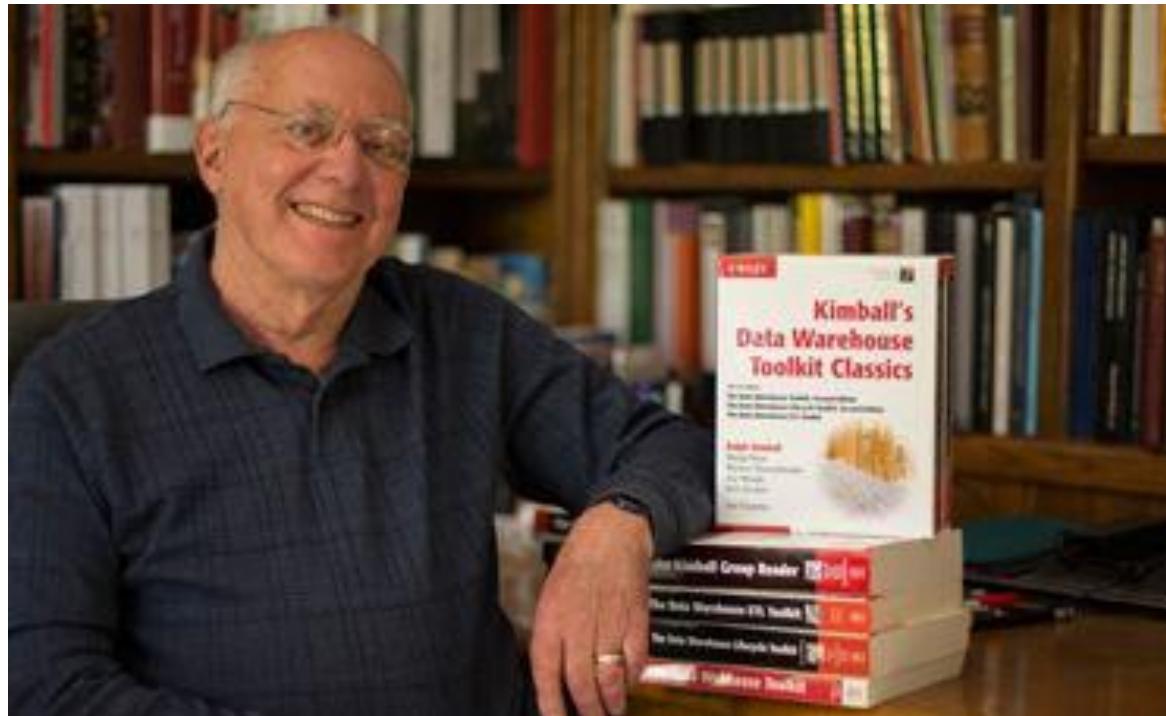


Inmon Model

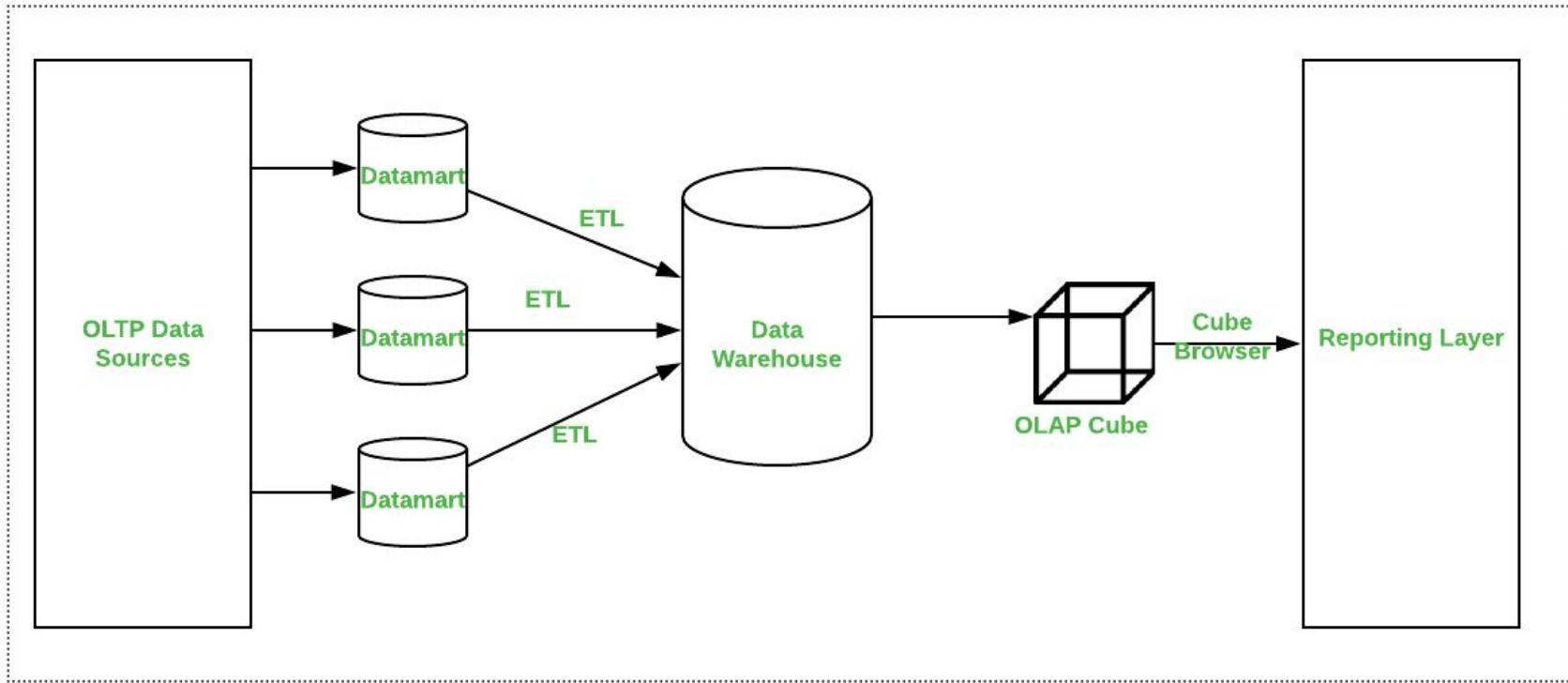
After Inmon – The Kimball Approach (1990s)

Ralph Kimball proposed building Data Marts first (Star Schema).

Faster analytics but less integrated at the start.

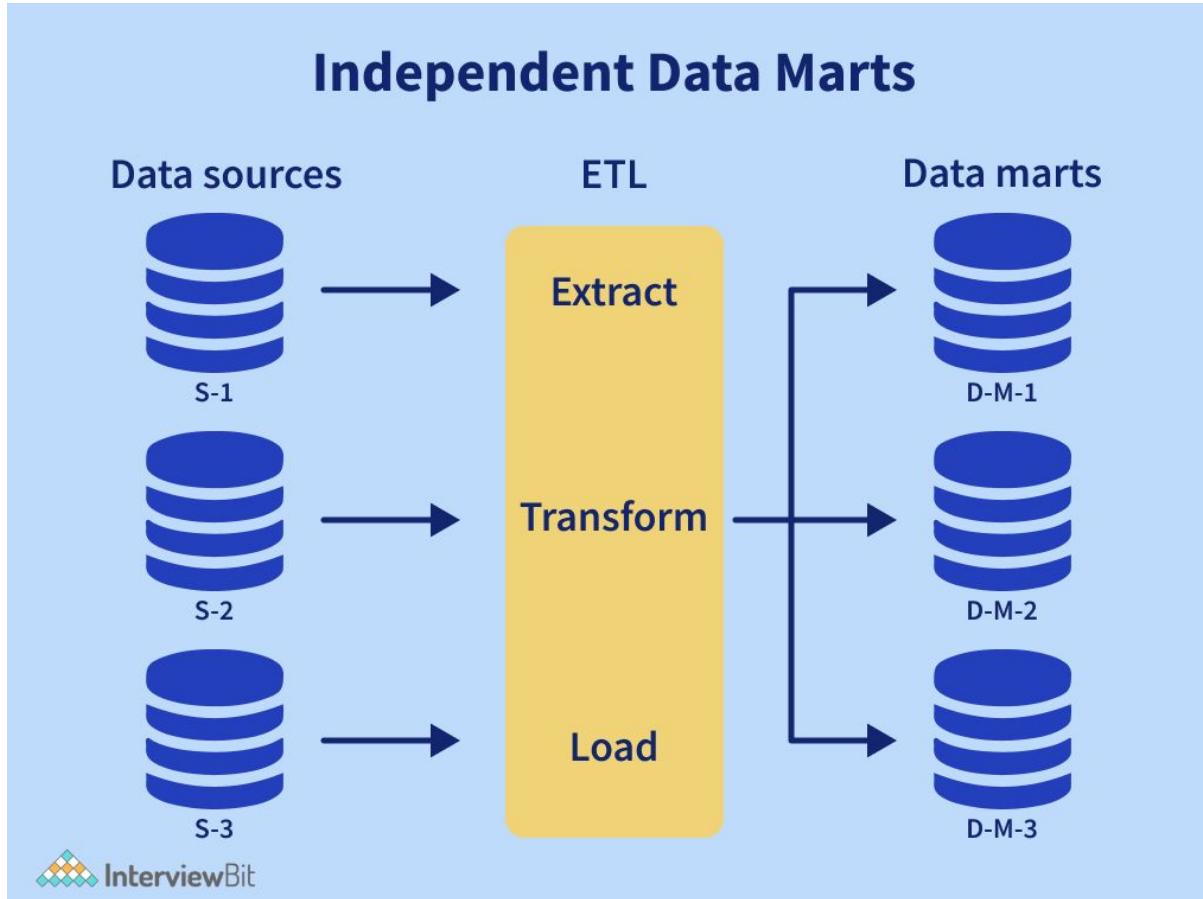


Kimball Approach

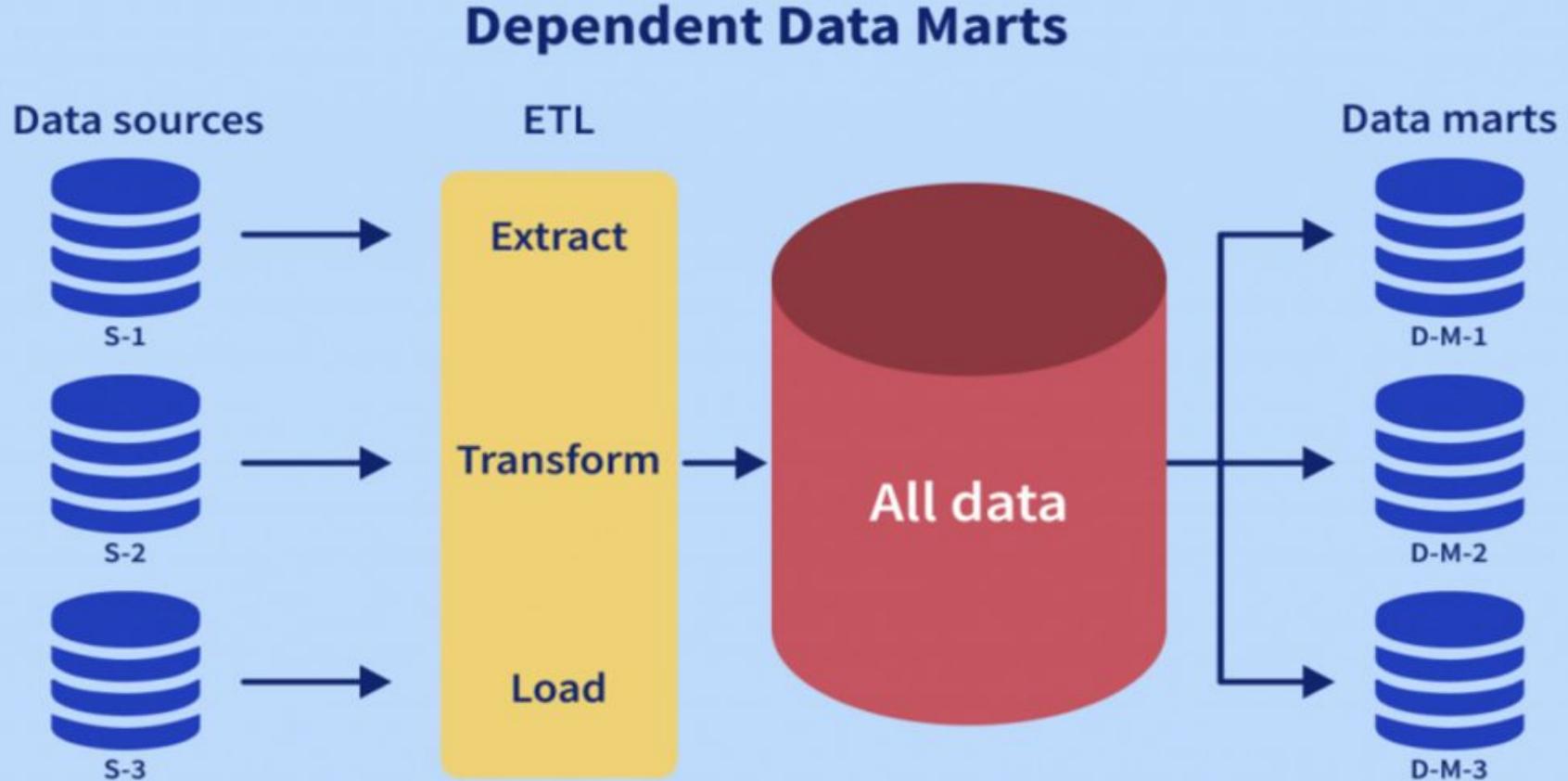


Kimball Model

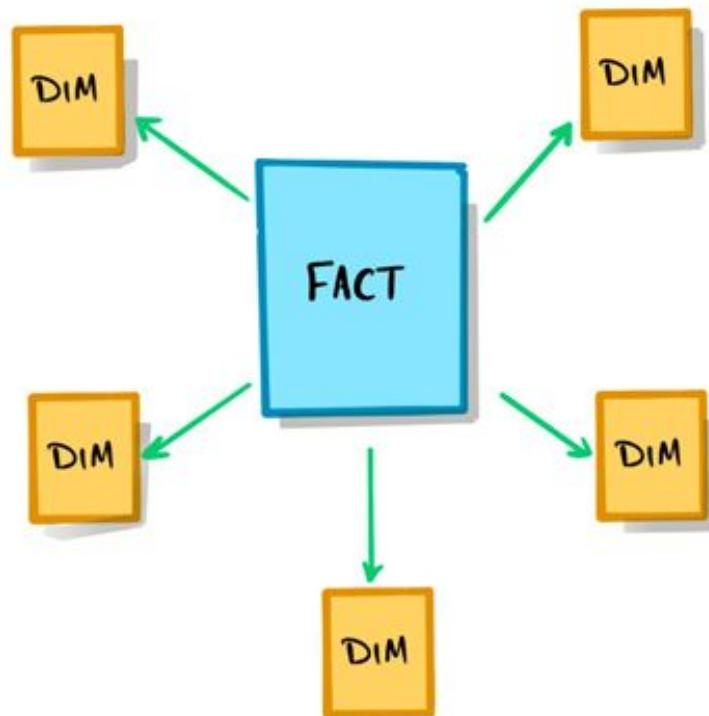
Dependent Vs Independent data mart “Which Approach?”



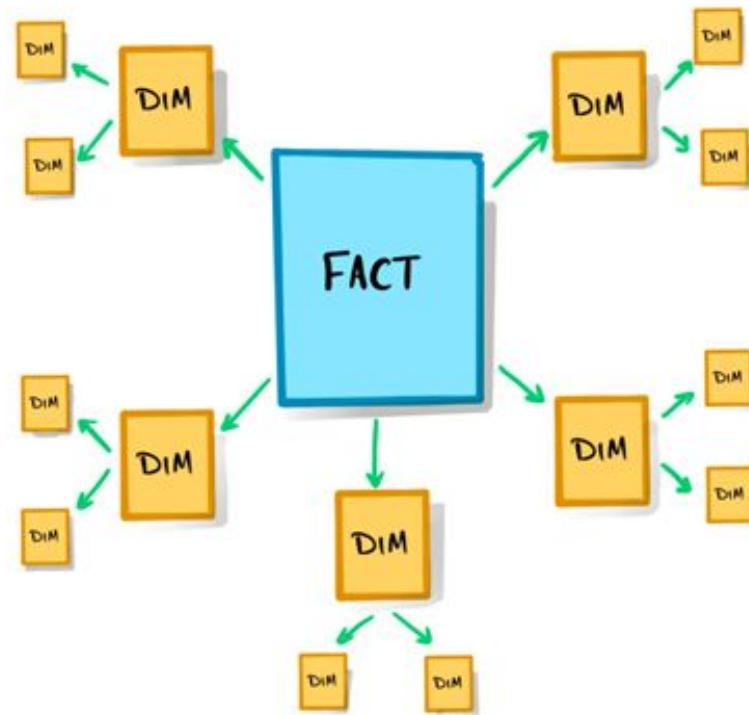
Dependent vs Independent data mart



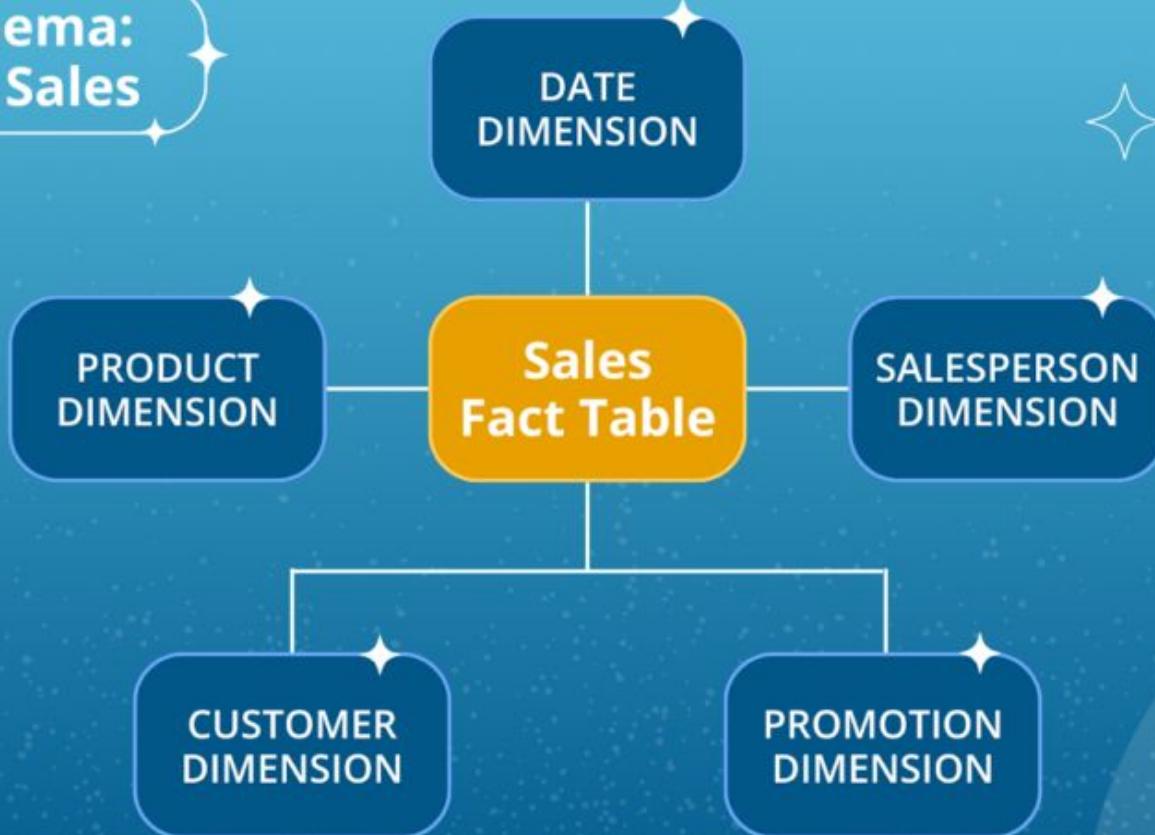
STAR SCHEMA



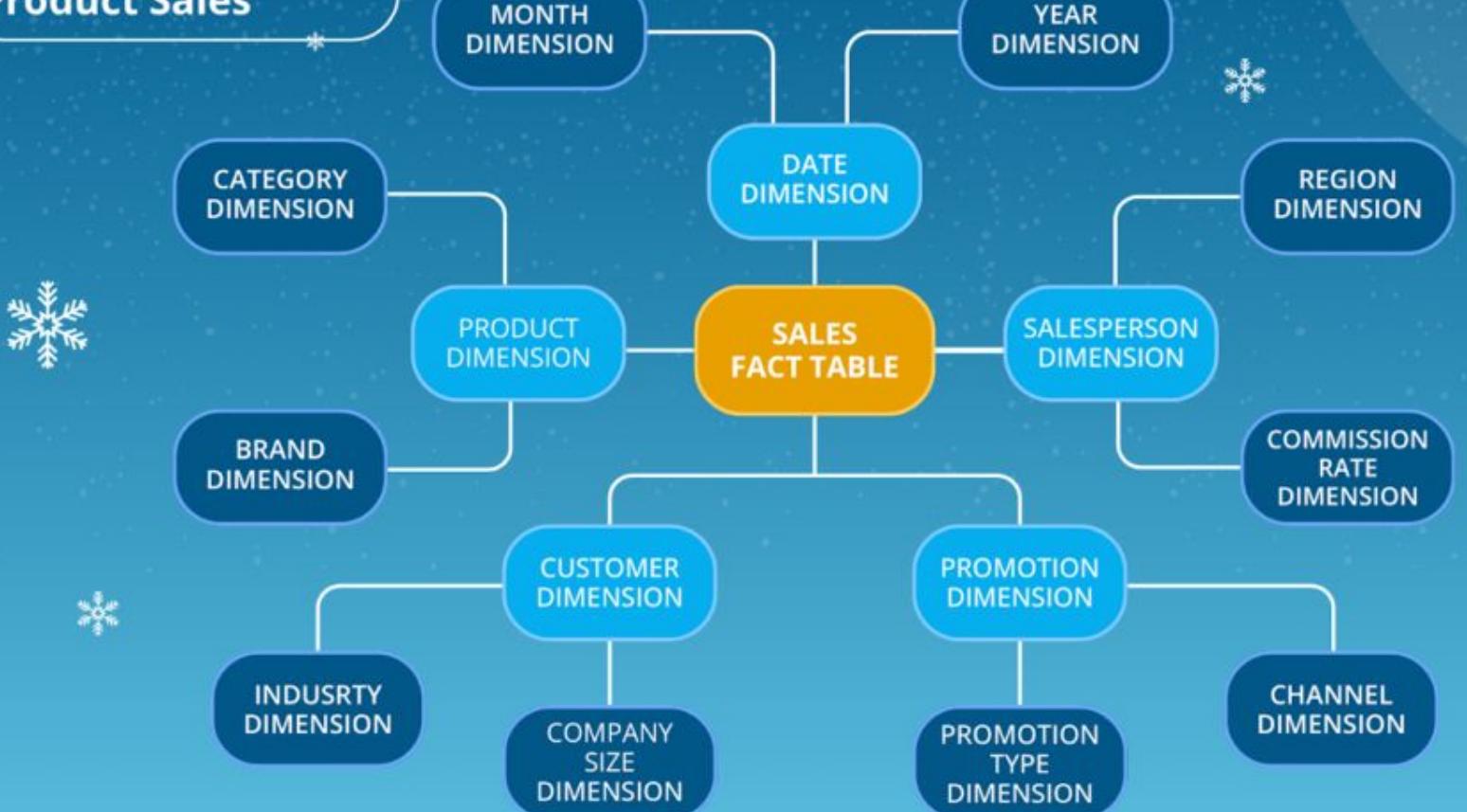
SNOWFLAKE SCHEMA



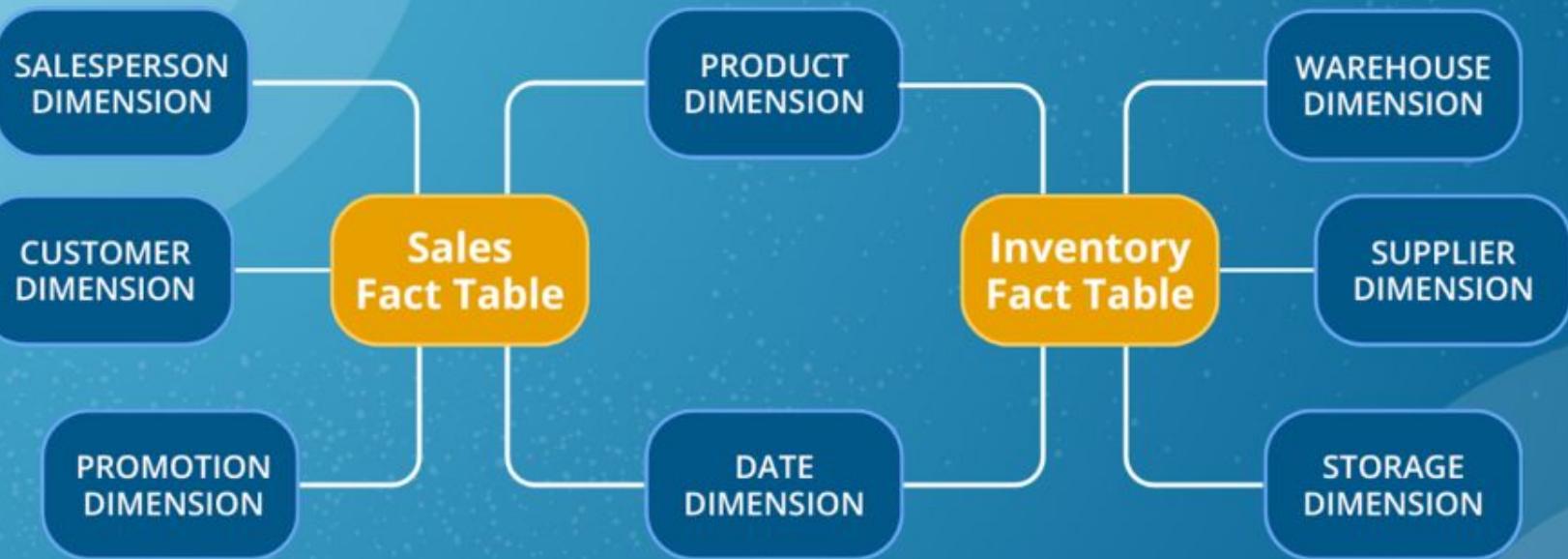
Star Schema: Product Sales



Snowflake Schema: Product Sales



Galaxy Schema



From Files to Data Engineering – Supporting Analytics

Modern Data Engineering

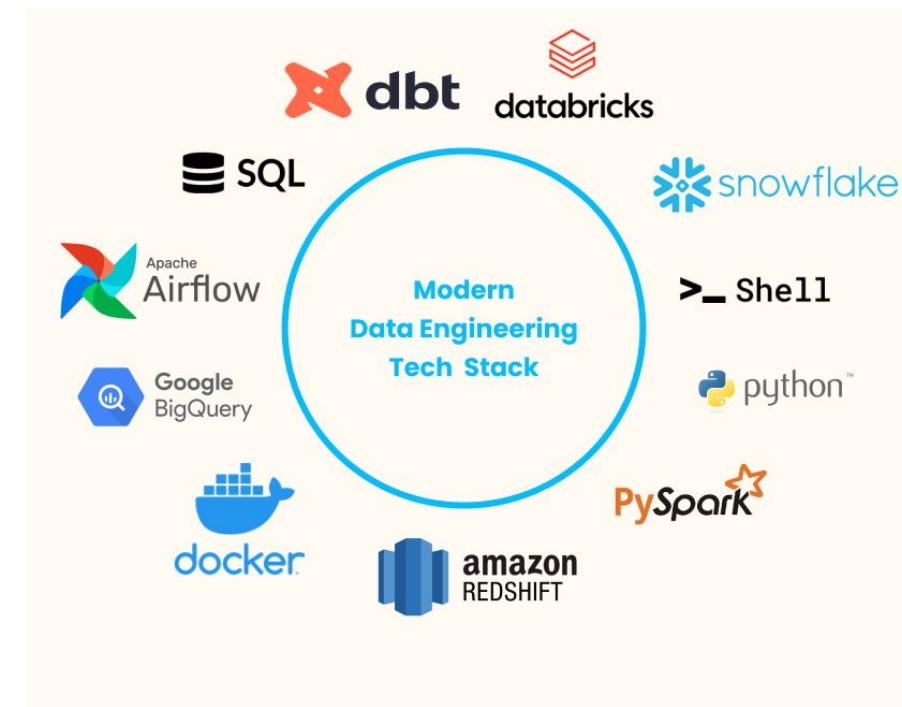
Combines databases, NoSQL, data lakes, and warehouses.

Builds pipelines for clean, integrated, analytics-ready data.

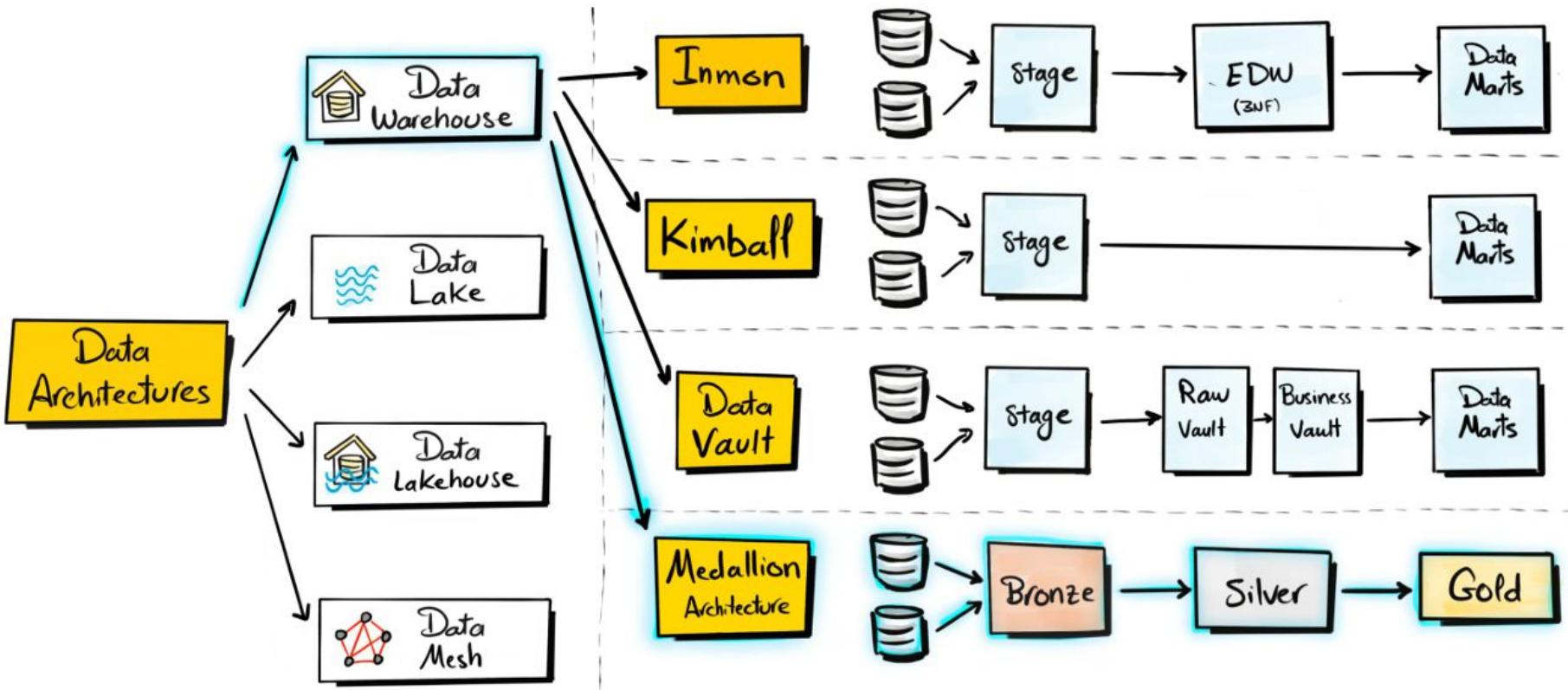
1980s: Bill Inmon → EDW (3NF)

1990s: Kimball → Star Schema

2000s+: Hybrid & Cloud DW



Data Architectures



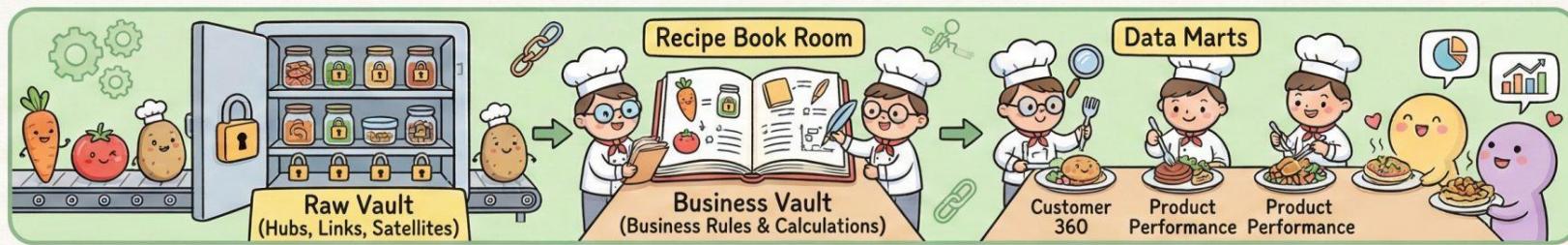
Inmon (Enterprise Data Warehouse)



Kimball (Dimensional Modeling)



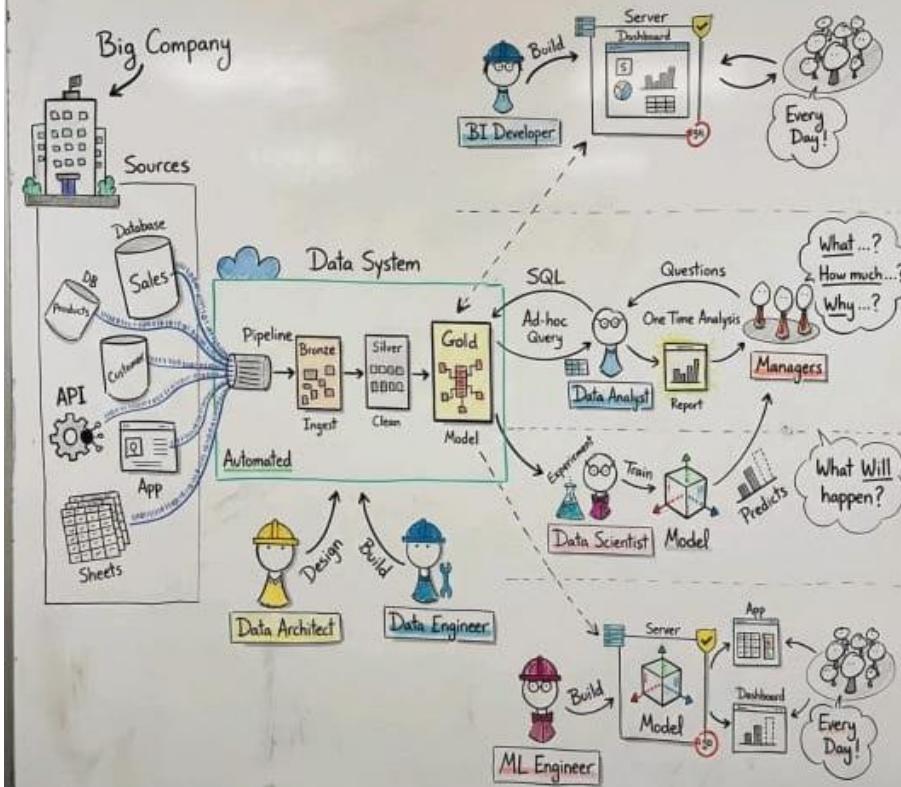
Data Vault (Agile & Auditable)



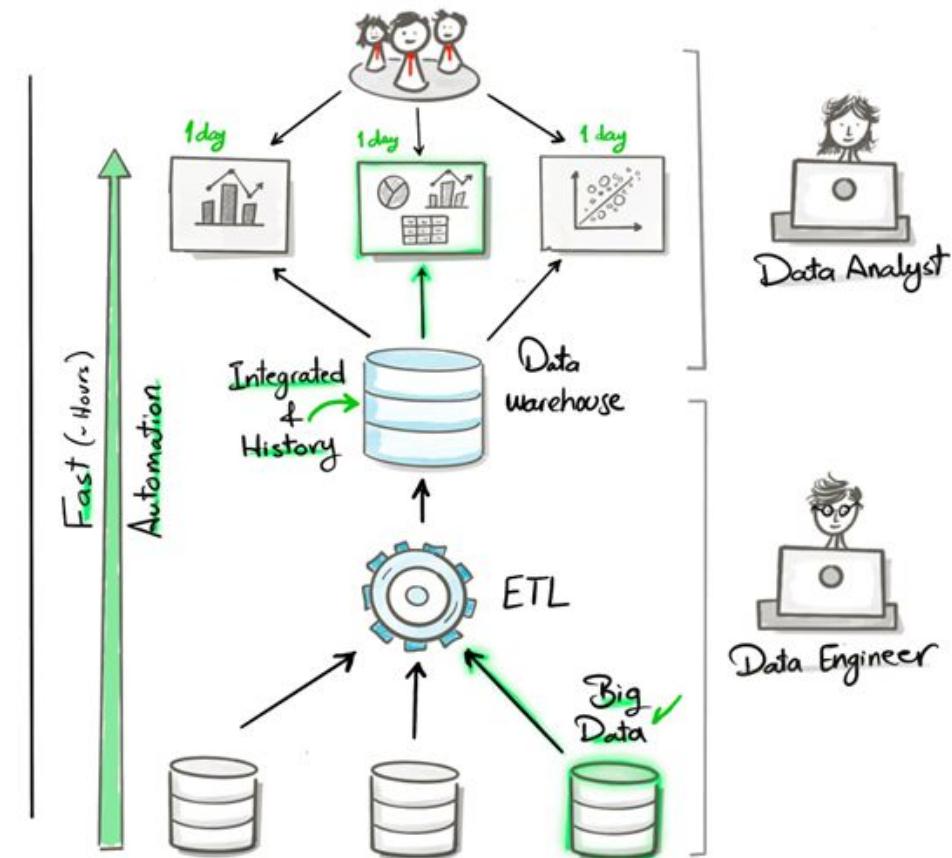
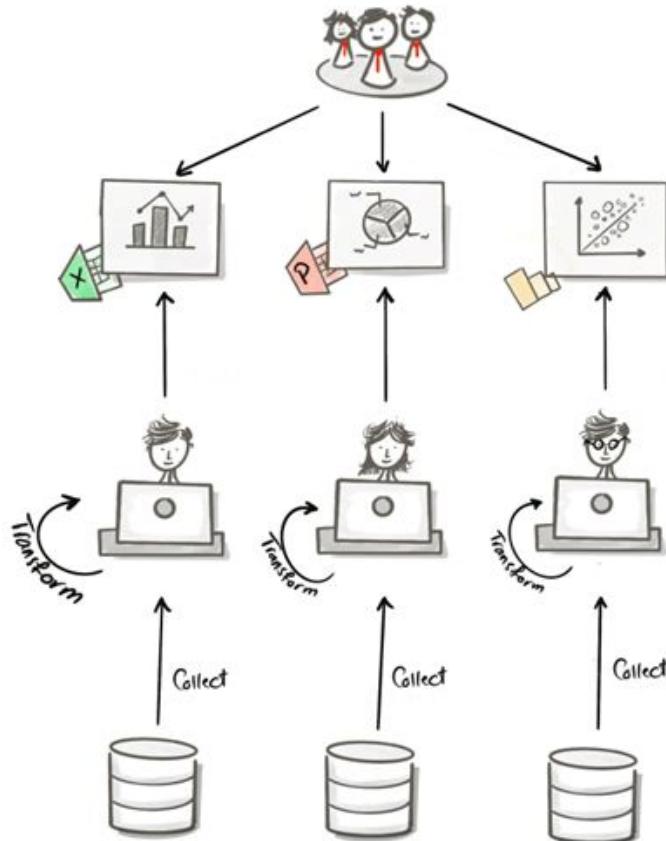
Medallion Architecture (Bronze, Silver, Gold)



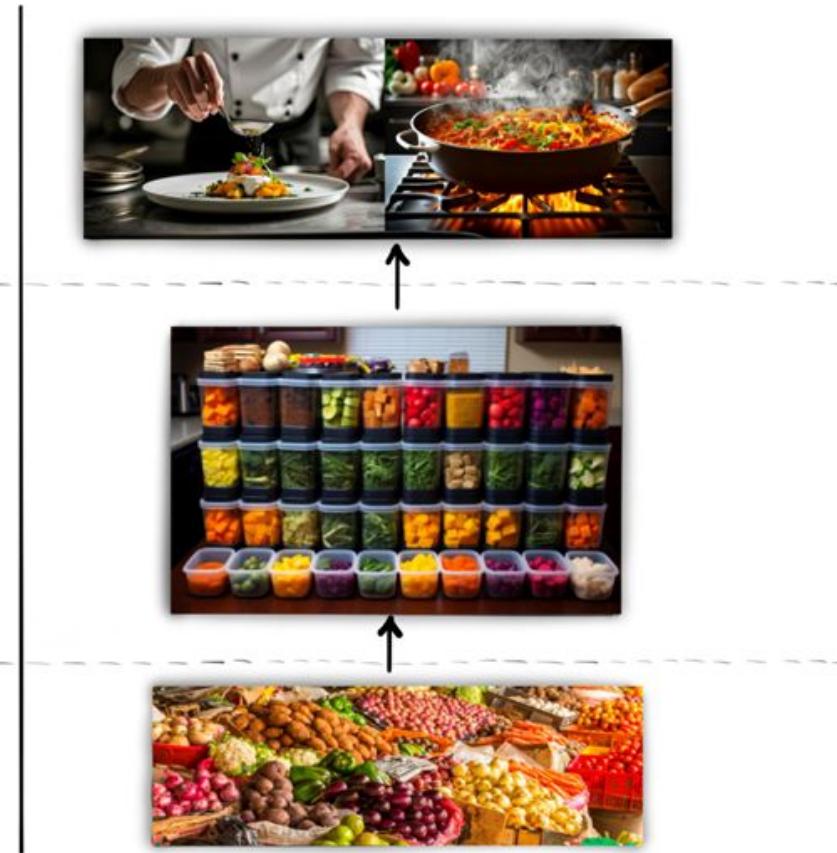
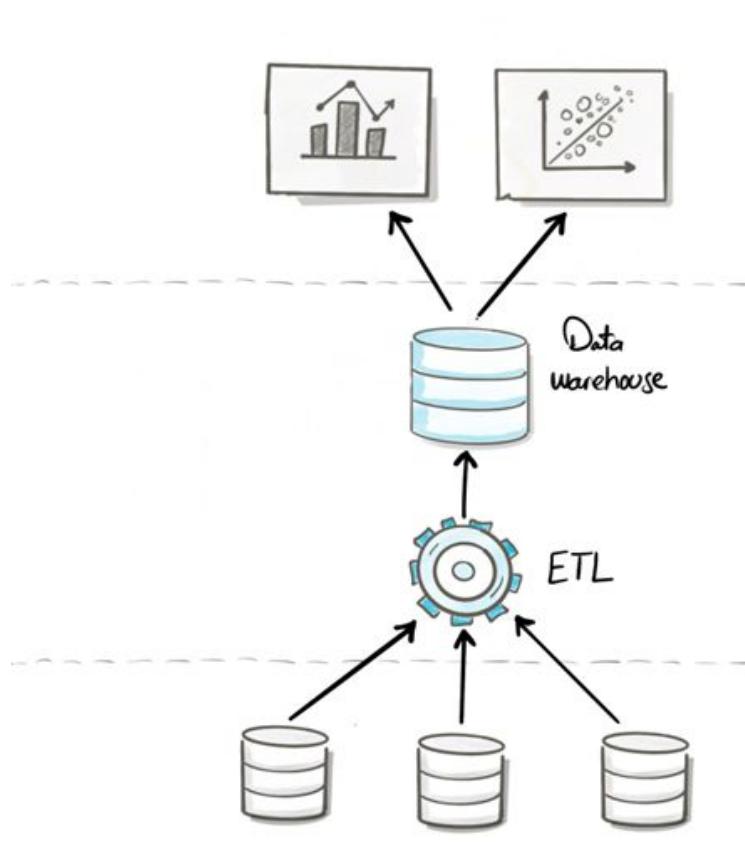
How modern data teams work



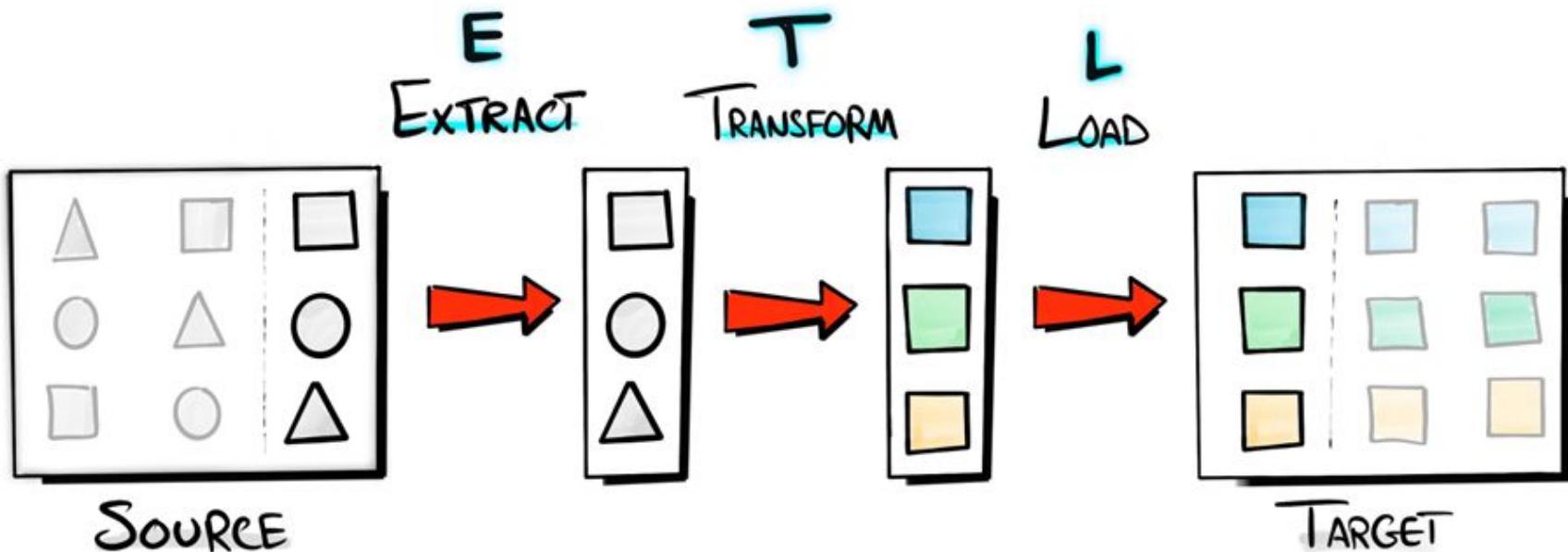
Why we need DWH !



Why we need DWH !



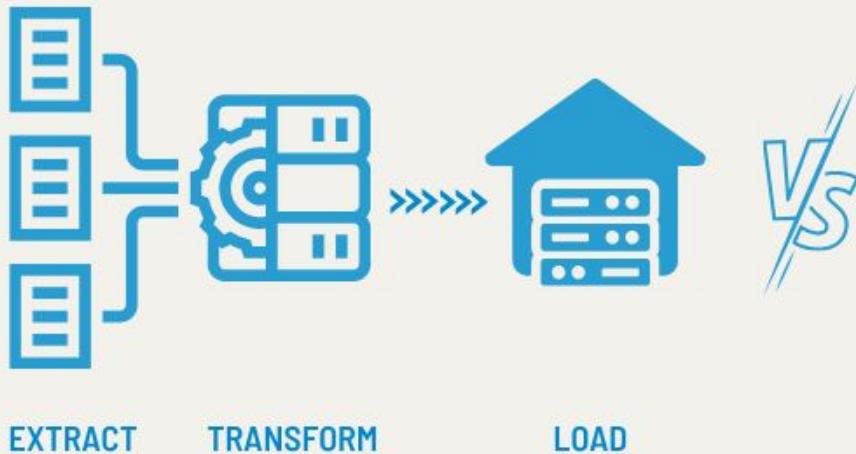
ETL Process



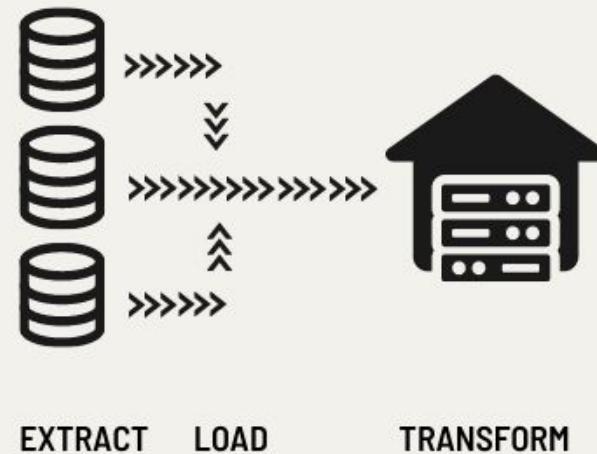
ETL VS ELT Guess??



ETL



ELT



Lambda Architecture

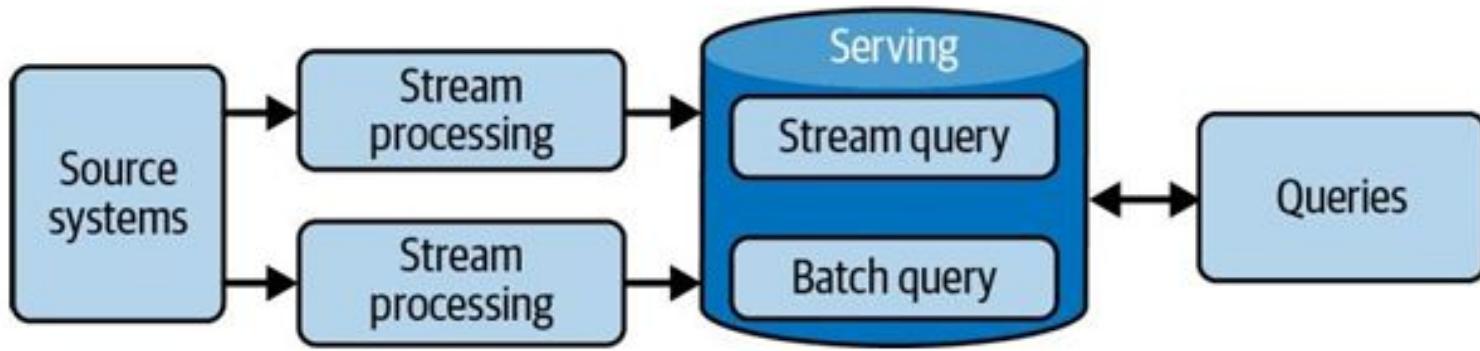


Figure 3-14. Lambda architecture

Kappa Architecture

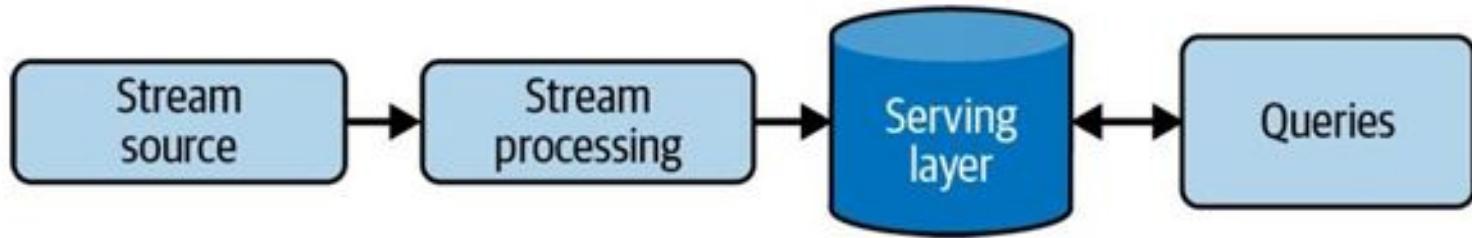
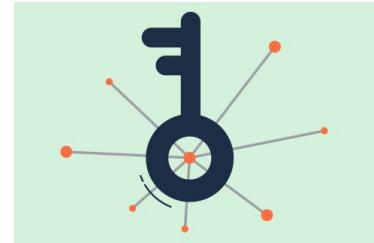


Figure 3-15. Kappa architecture

Surrogate vs Natural Key



It depend on several factors.

- The natural of the data.
- Database (DWH) platform.
- The group who uses this data.

Do we need to remove the natural key to use surrogate key ?

- NO, we will keep both in the table and treat the surrogate key as primary key.

Reasons to Learn SQL

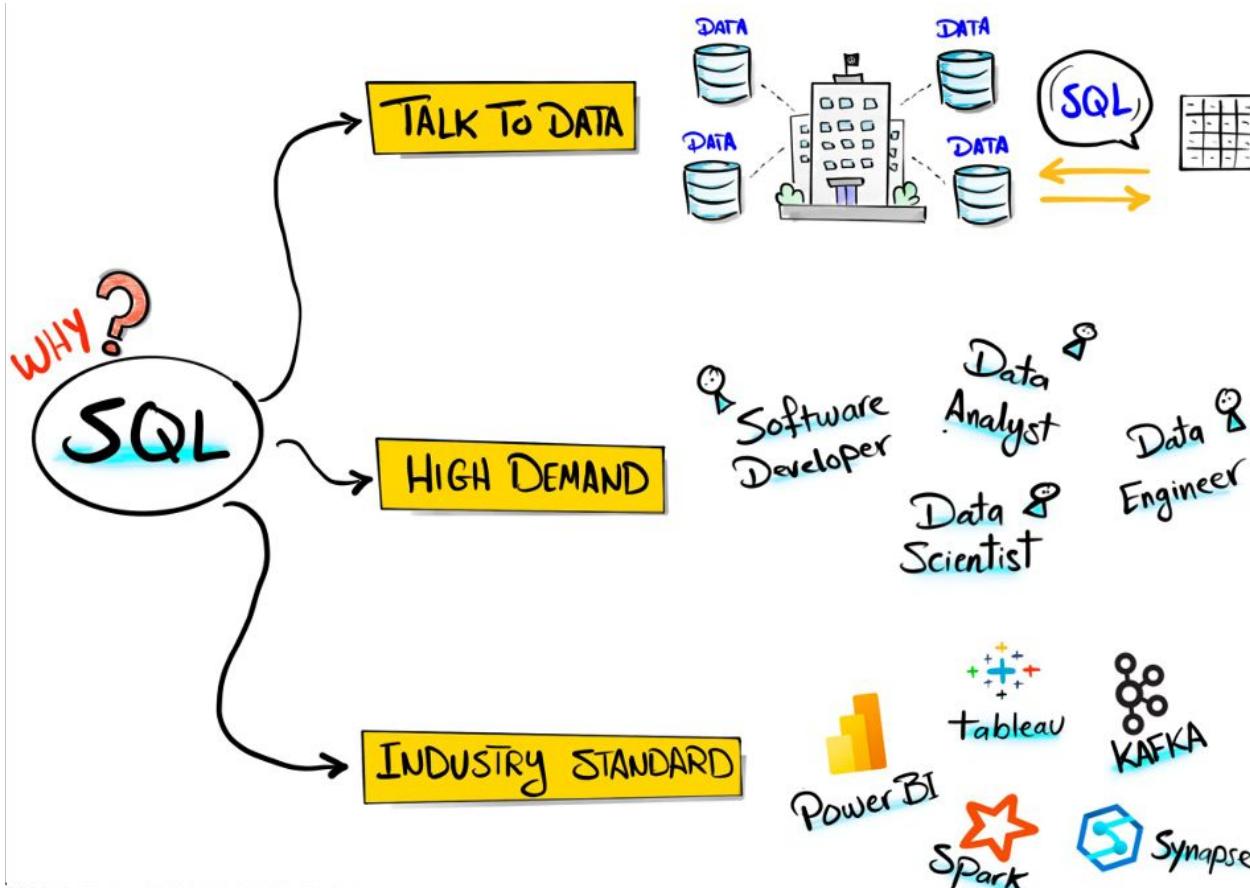




Data Engineer



Talk To Data!



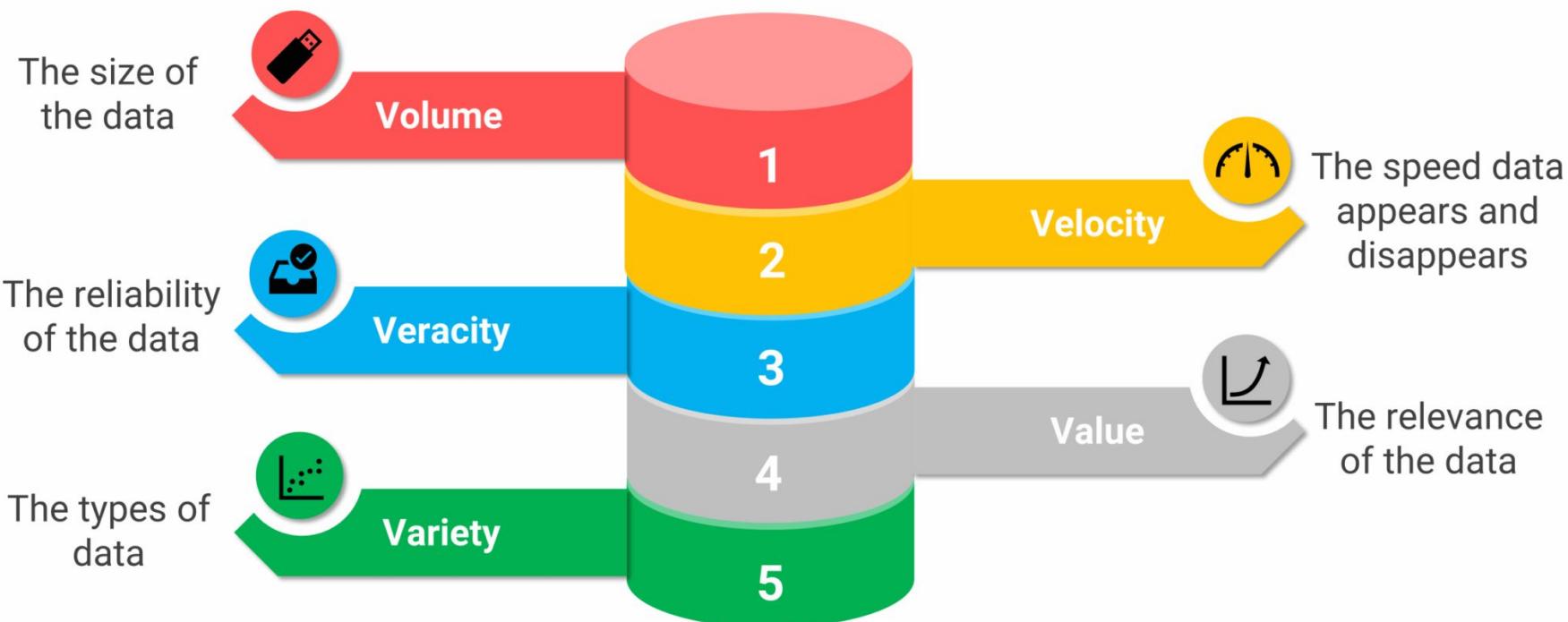
When Does Data Become “Big Data”

The 3 V's of big data

Big data is a collection of data from various sources, often characterized by what's become known as the 3 V's: *volume*, *velocity* and *variety*.

Volume	Velocity	Variety
<p>The amount of data from myriad sources.</p> 	<p>The speed at which big data is generated.</p> 	<p>The types of data: structured, semistructured, unstructured.</p> 

The 5 Vs of Big Data



SQL Server Overview

- **What is SQL Server?**
- **Editions (Developer, Express, Enterprise)**
- **Why we use it in this course**



Installation Steps Overview

- **Download SQL Server (Developer Edition)**
- **Download SQL Server Management Studio (SSMS)**
- **Custom Installation Configurations.**



Any
QUESTIONS?