

DS684
Cloud Computing
Week 08

Regarding Labs and Assignments

- Class participation means more than Zoom attendance. You must actively participate in the discussion and labs, and answer questions.
- Must hit Submit button, otherwise no grade
- If you need extension in time, must send written request (**email**). Otherwise no grade and no makeup. Requests sent over Zoom chat do not count.
- For any technical difficulty (installation, Azure access, etc), you must send written explanation (**email**) before the deadline. Otherwise no grade and no makeup.

Teaching Schedule

Week 7: Azure Synapse Analytics Part I: Data Warehouse

Week 8: Azure Synapse Analytics Part II: Data Engineering

Week 9: Visualization using Power BI

Week 10: Azure Machine Learning

Week 11: Final project presentation

Agenda

- ETL
 - Traditional Data Processing Flow
- ELT Data Processing Flow
 - Medallion Architecture
- Data Processing Services
 - Lab: Synapse Data Pipeline

General Tasks of Data Engineering

- Extracting (reading) from a source
- Transforming
 - Filtering
 - Calculation
 - Joining
 - Aggregation
 - etc.
- Loading (saving) into a target

Extract, Transform, Load (ETL)

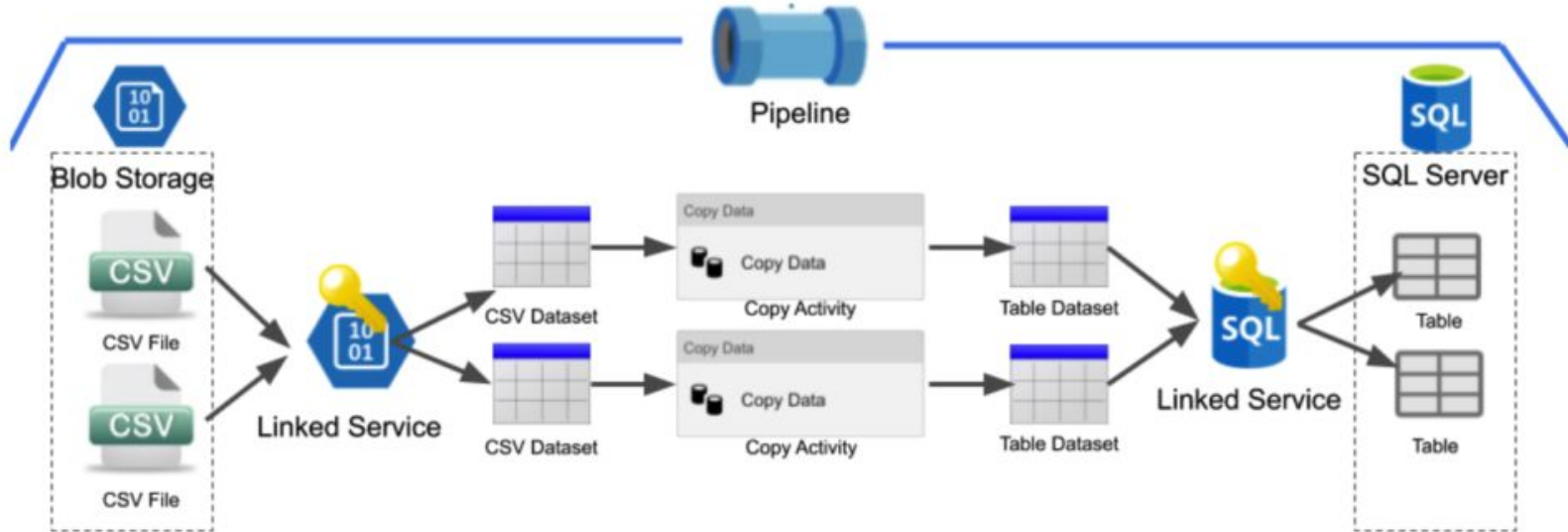
Different source systems will generate data that

- Comes in different format and frequency
 - Medicare history vs Medicaid history
- Is not joinable from different sources
 - Medical history vs medical device purchase history - different keys
- Is not clean
 - Rx usage history vs patient social network contents - lots of unrelated information

Extract, Transform, Load (ETL)

- Data from multiple sources (website, mobile, etc) are cleansed, consolidated, merged, and stored together
- Derived/aggregated values are calculated
- Loaded into more accessible schemas.

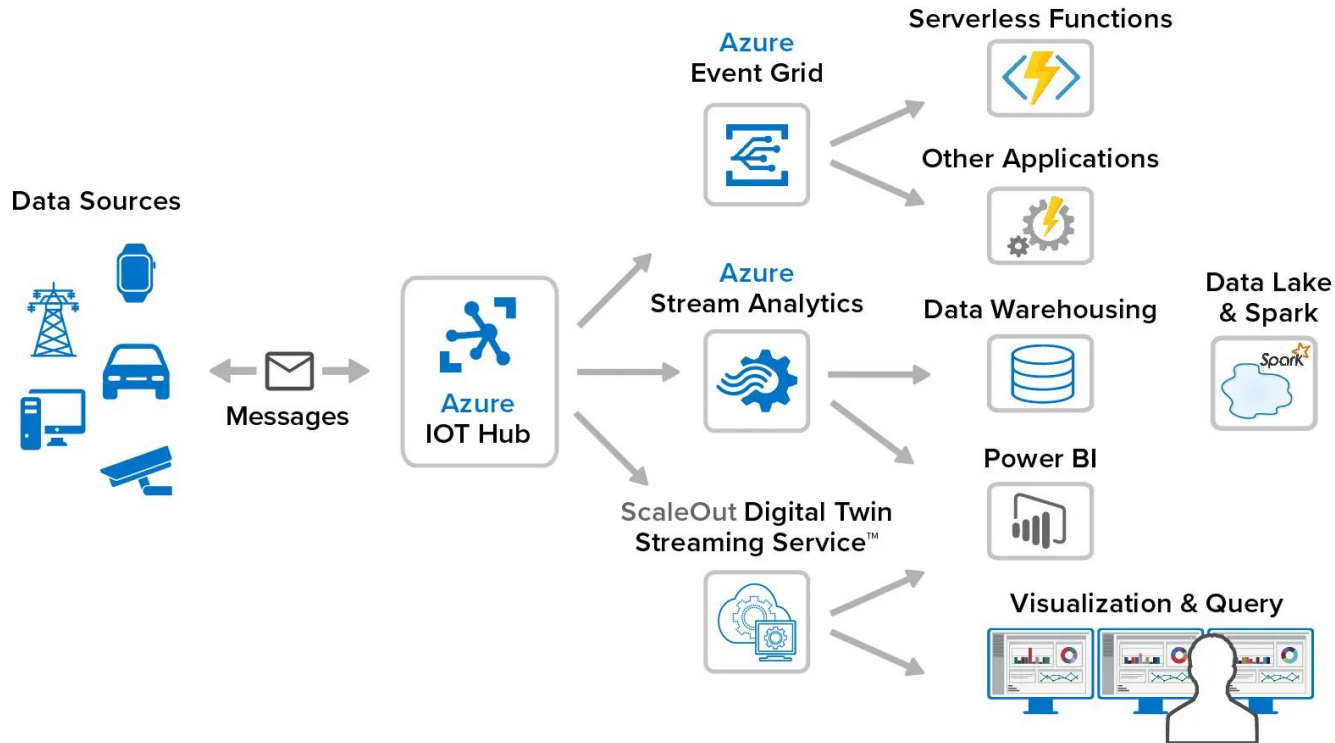
Extract, Transform, Load (ETL)



Batch vs Streaming

- Real world activities happen like a stream of events over time
 - Batch: Collect the events and process together
 - Streaming: Processing each event as it arrives
- Streaming is not a new concept/practice, but gets more attention as big data gains popularity
 - Velocity of data
 - Variety (source and format) of data
- Example: Internet of Things (IoT)

Azure IoT Streaming Example



Agenda

- ETL
 - Traditional Data Processing Flow
- ELT Data Processing Flow
 - Medallion Architecture
- Data Processing Services
 - Lab: Synapse Data Pipeline

Traditional Data Processing Flow

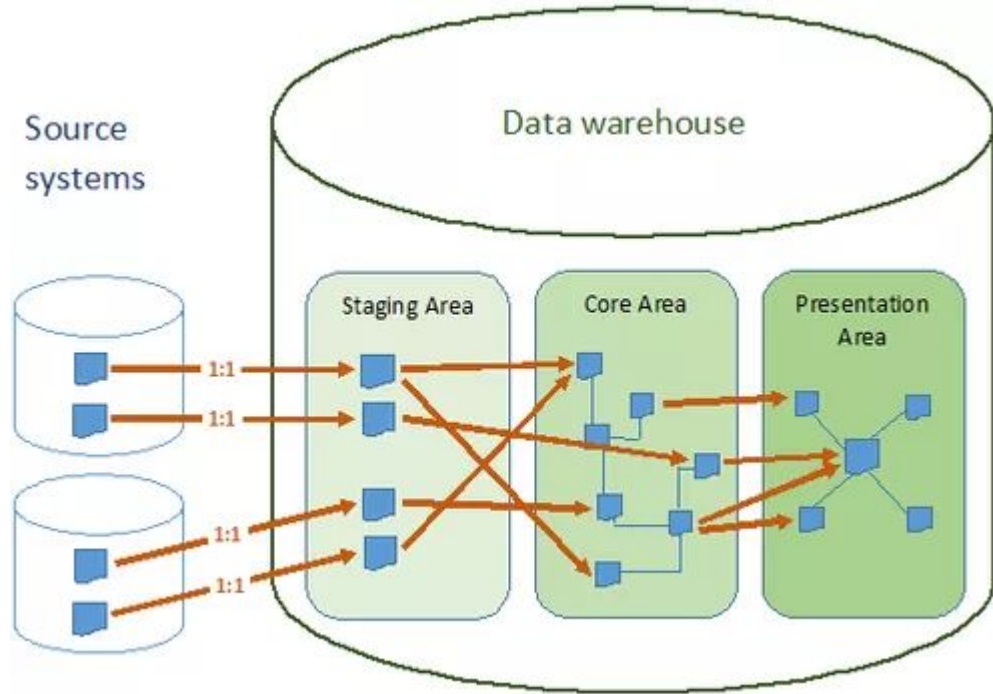
Source ->

Staging ->

Data Warehouse ->

Datamart

- Transformation happens between Staging area and Data Warehouse
- There might be multiple layers of staging



Traditional Data Architecture

ETL workflow will move data between different layers of staging tables, to data warehouse, and finally to datamarts

Look backwards, datamarts and data warehouses are usually designed first

Traditional Data Architecture

Datamart:

- Star schema
- Business oriented
- Organized around a particular business flow or activity

Traditional Data Architecture

Central storage

- Star/Snowflake schema Data Warehouse, or
- 3NF complaint ODS: Since datamart has been designed as star schema, data warehouse can be 3NF, or
- Logical data warehouse: a layer of views on top of ODS

Traditional Data Architecture

Staging:

1. Load data as is
2. Data manipulations
3. Intermediate results

There might be multiple layers of staging

ETL workflow will move data between different layers of staging tables, to data warehouse, and finally to datamarts

Agenda

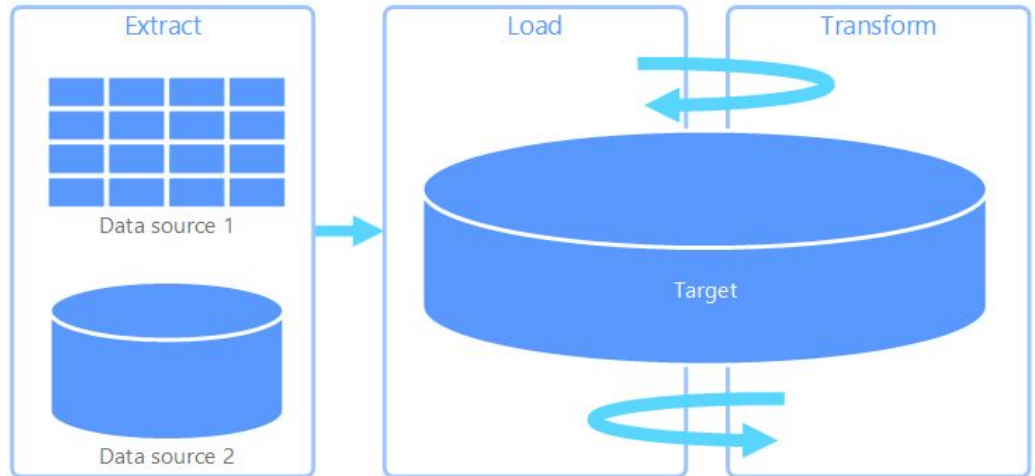
- ETL
 - Traditional Data Processing Flow
- **ELT Data Processing Flow**
 - Medallion Architecture
- Data Processing Services
 - Lab: Synapse Data Pipeline

ELT Data Processing Flow

Load into data lake first. Process when read (reducing processing needs)

Raw data lake based Lakehouse

Challenge: Not all datasets are equal. Some requires more attention.



Agenda

- ETL
 - Traditional Data Processing Flow
- ELT Data Processing Flow
 - Medallion Architecture
- Data Processing Services
 - Lab: Synapse Data Pipeline

Medallion Data Architecture

- Organize the data in a data lake

Bronze -> Silver -> Gold



Medallion Data Architecture

Bronze: raw data, all as is

Current practice is to collect data as detailed as possible

Medallion Data Architecture

Silver

- Cleansed: Handle missing, inconsistent, duplicated, and erroneous data
- Filtered: Remove unnecessary data
- Conformed: Make data type (date, string, integer) and data format (year, month, date e.g.) consistent
- Normalized/Denormalized: Convert between relational and non-relational schemas
- Feature engineering

Do you want to join (maintain foreign key) in silver stage?

- A matter of preference

Medallion Data Architecture

Gold: Ready for reporting

- Consumption-ready
- Project-specific
- Joined different datasets
- Denormalized into star schema
- Aggregated into statistics

Medallion Data Architecture

- Organize the data in a data lake

Bronze -> Silver -> Gold



Medallion Architecture vs Traditional Staging

- Medallion Architecture and Staging are not exclusive
- Staging approach is still an important part of data warehouse ETL design
- You will see a mixture of both in your future jobs

Agenda

- ETL
 - Traditional Data Processing Flow
- ELT Data Processing Flow
 - Medallion Architecture
- Data Processing Services
 - Lab: Synapse Data Pipeline

Data Processing Services

- General compute (VM, container, function)
- Synapse Data Pipeline (Azure Data Factory)

Data Processing Tools

- SQL
- Spark (python, scala, Java)
- Low-code/no-code ETL mapping

The first two are beyond the scope of this course. We will only introduce the third approach. However, keep in mind that all these three are important and are often required by employers.

Agenda

- ETL
 - Traditional Data Processing Flow
- ELT Data Processing Flow
 - Medallion Architecture
- Data Processing Services
 - Lab: Synapse Data Pipeline

Synapse Data Pipeline Demo and Lab

Bronze Stage

- Create sources
- Create sink

Silver Stage

- Create filter
- Create new derived column: Concatenate, Substring
- Create data type conversion

Gold Stage

- Create join

Designing the Final Project

How would you design your final project database?

Final Project

Review Assignment 07

- Table creation

Assignment 08

- Data processing
 - End result: a consolidated dataset