

# ETL Process in Data Warehouse



Zain Shaukat

# Outline

- ETL
- Extraction
- Transformation
- Loading

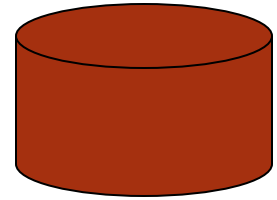
# ETL Overview

- Extraction Transformation Loading – ETL
- To get data out of the source and load it into the data warehouse – simply a process of copying data from one database to other
- Data is **extracted** from an OLTP database, **transformed** to match the data warehouse schema and **loaded** into the data warehouse database
- Many data warehouses also include **data from non-OLTP systems** such as text files, and spreadsheets; such data also requires extraction, transformation, and loading
- When defining ETL for a data warehouse, it is important to think of ETL as a **process, not a physical** implementation

# ETL Overview

- ETL is often a **complex combination of process and technology** that consumes a significant portion of the data warehouse development efforts and requires the skills of business analysts, database designers, and application developers
- It is not a one time event as **new data** is added to the Data Warehouse **periodically** – monthly, daily, hourly
- Because ETL is an integral, ongoing, and periodic part of a data warehouse

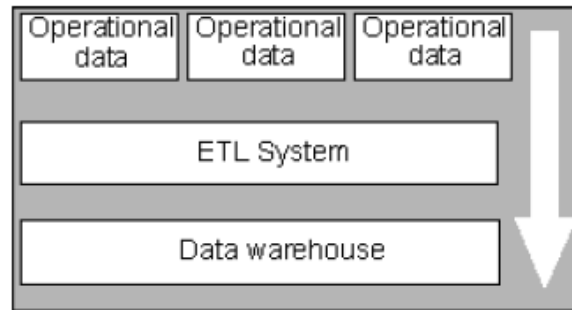
# ETL Staging Database



- ETL operations should be performed on a **relational database server separate** from the source databases and the data warehouse database
- Creates a logical and physical separation between the source systems and the data warehouse

# Extraction

# Extraction



- The **integration of all of the dissimilar systems** across the enterprise is the real challenge to getting the data warehouse to a state where it is usable
- Data is extracted from **heterogeneous** data sources
- Each data source has its distinct set of characteristics that need to be **managed and integrated into the ETL** system in order to effectively extract data.

# Extraction

- ETL process needs to effectively integrate systems that have different:
  - DBMS
  - Operating Systems
  - Hardware
  - Communication protocols
- Need to have a **logical data map** before the physical data can be transformed
- The logical data map **describes the relationship** between the extreme starting points and the extreme ending points of your ETL system usually presented in a table or spreadsheet



Target			Source			Transformation
Table Name	Column Name	Data Type	Table Name	Column Name	Data Type	

- The content of the logical data mapping document has been proven to be the critical element required to efficiently plan ETL processes
- The primary purpose of this document is to provide the ETL developer with a **clear-cut blueprint of exactly what is expected from the ETL process.**
- The transformation can contain anything from the absolute solution to nothing at all. Most often, the transformation can be expressed in SQL. The SQL may or may not be the complete statement

- The analysis of the source system is usually broken into two major phases:
  - The data discovery phase
  - The anomaly detection phase

# Extraction - Data Discovery Phase

- **Data Discovery Phase**

key benchmark for the success of the data warehouse is the cleanliness and cohesiveness of the data within it

- Once you understand what the target needs to look like, you need to identify and examine the data sources

# Data Discovery Phase

- It is up to the ETL team to drill down further into the data requirements to determine each and every source system, table, and attribute required to load the data warehouse
  - Collecting and Documenting Source Systems
  - Keeping track of source systems
  - Determining the System of Record - Point of originating of data
  - Definition of the system-of-record is important because in most enterprises data is stored redundantly across many different systems.
  - Enterprises do this to make nonintegrated systems share data. It is very common that the same piece of data is copied, moved, manipulated, transformed, altered, cleansed, or made corrupt throughout the enterprise, resulting in varying versions of the *same* data

# Data Content Analysis - Extraction

- Understanding the content of the data is crucial for determining the best approach for retrieval
  - **NULL values.** An unhandled NULL value can destroy any ETL process. NULL values pose the biggest risk when they are in foreign key columns. Joining two or more tables based on a column that contains NULL values **will cause data loss!** Remember, in a relational database NULL is not equal to NULL. That is why those joins fail. Check for NULL values in every foreign key in the source database. When NULL values are present, you must *outer* join the tables
  - **Dates in nondate fields.** Dates are very anomalous elements because they are the only logical elements that can come in various formats, literally containing different values and having the exact same meaning. Fortunately, most database systems support most of the various formats for display purposes but store them in a single standard format

# Determining Changed Data

- **Audit Columns** - Used by DB and updated by triggers
- **Audit columns** are appended to the end of each table to store the date and time a record was added or modified
- You must analyze and test each of the columns to ensure that it is a reliable source to indicate changed data. If you find any NULL values, you must find an alternative approach for detecting change – example **using outer joins**

# Determining Changed Data

## Process of Elimination

- Process of elimination preserves exactly one copy of each previous extraction in the staging area for future use.
- During the next run, the process takes the entire source table(s) into the staging area and makes a comparison against the retained data from the last process.
- Only differences (deltas) are sent to the data warehouse.
- Not the most efficient technique, but most reliable for capturing changed data

# Determining Changed Data

## Initial and Incremental Loads

- Create two tables: previous load and current load.
- The initial process bulk loads into the current load table. Since change detection is irrelevant during the initial load, the data continues on to be transformed and loaded into the ultimate target fact table.
- When the process is complete, it drops the previous load table, renames the current load table to previous load, and creates an empty current load table. Since none of these tasks involve database logging, they are very fast!
- The next time the load process is run, the current load table is populated.
- Select the current load table MINUS the previous load table. Transform and load the result set into the data warehouse.