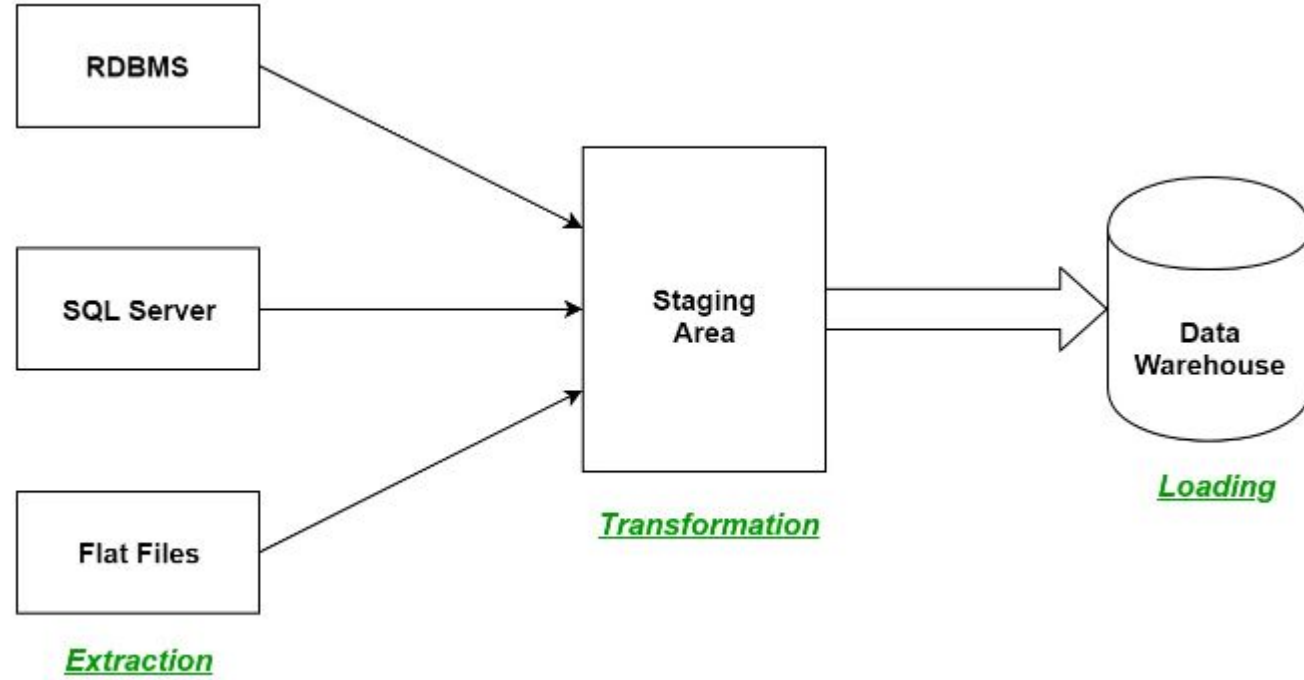


Data Staging (ETL) Design and Development

- ETL is a process in Data Warehousing and it stands for Extract, Transform and Load. It is a process in which an ETL tool extracts the data from various data source systems, transforms it in the staging area, and then finally, loads it into the Data Warehouse system.



Extraction:

- The first step of the ETL process is extraction. In this step, data from various source systems is extracted which can be in various formats like relational databases, No SQL, XML, and flat files into the staging area. It is important to extract the data from various source systems and store it into the staging area first and not directly into the data warehouse because the extracted data is in various formats and can be corrupted also. Hence loading it directly into the data warehouse may damage it and rollback will be much more difficult. Therefore, this is one of the most important steps of ETL process.

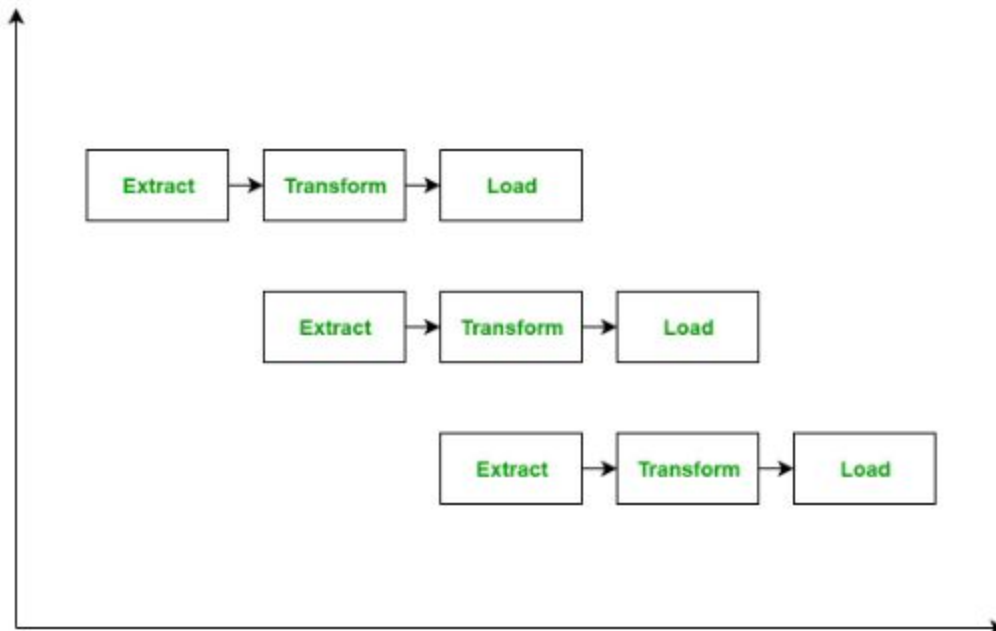
Transformation:

- The second step of the ETL process is transformation. In this step, a set of rules or functions are applied on the extracted data to convert it into a single standard format. It may involve following processes/tasks: **Filtering** – loading only certain attributes into the data warehouse.
- **Cleaning** – filling up the NULL values with some default values, mapping U.S.A, United States, and America into USA, etc.
- **Joining** – joining multiple attributes into one.
- **Splitting** – splitting a single attribute into multiple attributes.
- **Sorting** – sorting tuples on the basis of some attribute (generally key-attribute).

Loading:

- The **third and final step of the ETL process is loading**. In this step, the transformed data is finally loaded into the data warehouse. Sometimes the data is updated by loading into the data warehouse very frequently and sometimes it is done after longer but regular intervals. The rate and period of loading solely depends on the requirements and varies from system to system.

- ETL process can also use the pipelining concept i.e. as soon as some data is extracted, it can be transformed and during that period some new data can be extracted. And while the transformed data is being loaded into the data warehouse, the already extracted data can be transformed. The block diagram of the pipelining of ETL p



ETL Tools:

- Most commonly used ETL tools are
 - **Hevo**, Sybase, Oracle Warehouse builder, CloverETL, and MarkLogic.

Data Warehouses:

- Most commonly used Data Warehouses are
 - **Snowflake**, Redshift, BigQuery, and Firebolt.