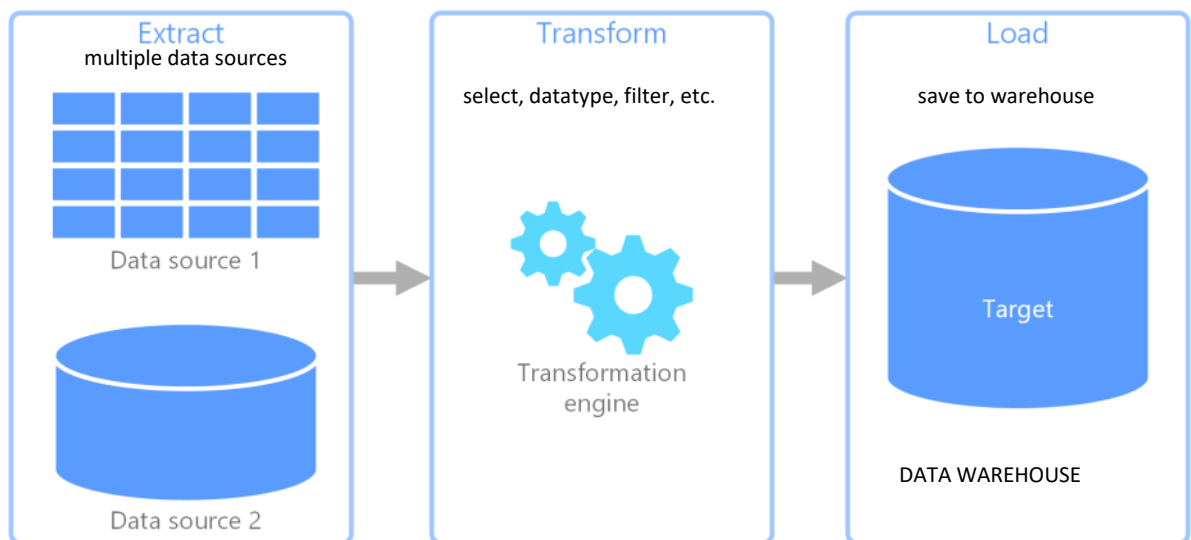
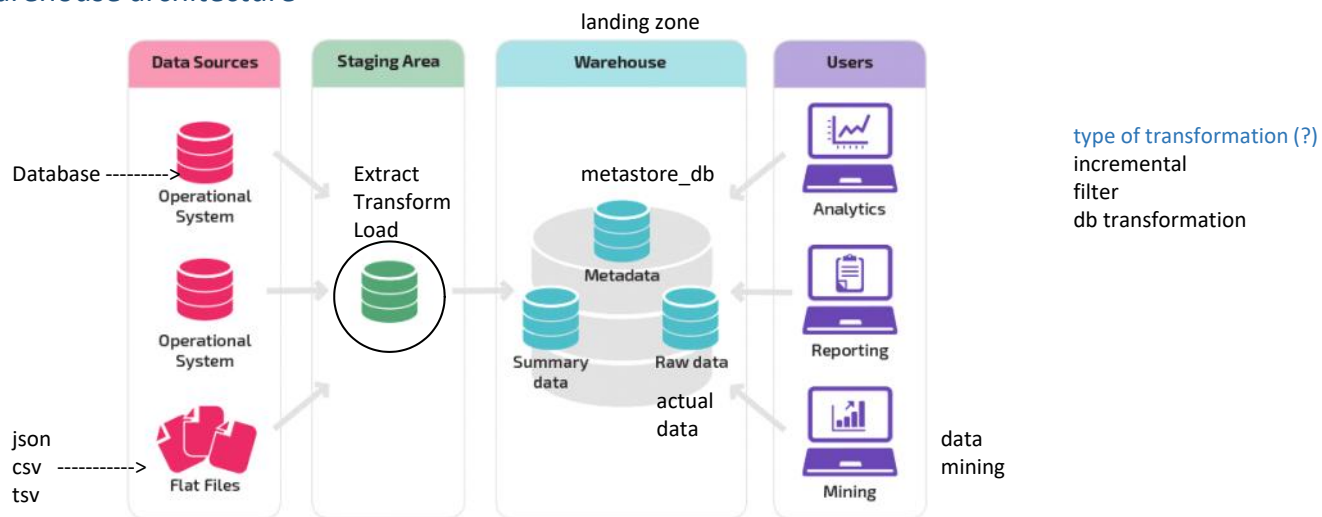


ETL

06 September 2024 15:53

Data warehouse architecture



Steps for ETL

1. determine all target data needed in data warehouse
2. determine all data sources (internal and external)
3. prepare data mapping for target data elements from sources
4. establish comprehensive data extraction rules
5. determine data transformation and cleansing rules
6. plan for aggregate tables
7. organize data staging area and test tools
8. write procedures for all data loads
9. ETL for dimension tables
10. ETL for fact tables

ETL key factors

- complexity of data extraction and transformation
- data loading functions

Data extraction techniques

Data in operational systems

current value	a single variable getting updated again and again over time	data about a person's residence getting updated over time when he moves
periodic status	a variable which (kind of) represents a timeline to show the changes	the status of a property being put for auction [put] --> [put, value changed] ----> [put, value changed sold]

Types of data extraction

- 'as is' i.e. static data extraction
 - o take capture of data at a given point in time
 - o mainly used for initial load
 - o the data capture would include each status/event at each point in time
- Data of revision
 - o aka incremental data capture
 - o look for and extract periodic new data
 - o 3 options for immediate data extraction capture through
 - transaction logs
 - dB triggers
 - source applications
 - o deferred data extraction capture through
 - date and time stamp
 - by comparing files

Data transformation

- data extracted must be useful
- we have to *enrich* and *improve* data quality before sending it
- tasks
 - o conversion
 - o summarization
 - o enrichment
 - o format revisions
 - o decoding of fields
 - o calculated and derived values
 - o splitting of single field
 - o merging of information
 - o character set conversion
 - o conversion of units of measurements
 - o data/time conversion
 - o deduplication
 - o key restructuring
- Data integration and consolidation
 - o entity identification problem

primary key

- uniquely identifies a row

foreign key

- common in two tables

surrogate key

- artificial key generated by db system
- ex. random UID generated by system

Data loading

- types
 - o initial load
 - o incremental load
 - o full refresh
- during loads, data has to be offline
- need to find a window of time when loads may be scheduled without affecting data warehouse users

- divide up whole load process into smaller chunks and populate few files at a time
 - o run smaller loads in parallel
 - o keep some parts of the warehouse up
- having staging area and data warehouse db on same server will save efforts
- how to apply data?
 - o writing special load programs
 - o load utilities that come with dbms
- modes for applying data
 - o load
 - if target table already exists, its removed and new data is added
 - else it just adds to the empty table
 - o append
 - add the new data to the existing table

ETL tools

- data transformation engines
- data capture through replication
- code generator

SQL server integration server

components

- control flow
- data flow
- event handler

pivot table

powerful tool to calculate,

Fuzzy lookup transformation

it's an SSIS component that can be used to clean/standardize or correct data in input data source with a reference dataset

this transformation uses fuzzy matching to return one or more matches in the reference table

Lookup cache modes

In SQL Server Integration Services (SSIS), the Lookup transformation component has three cache modes that determine how data is cached during package execution:

1. **Full Cache:** This mode loads the entire lookup table into memory before processing any data. It is the default mode and is efficient for small to moderately sized lookup tables. However, it can consume a lot of memory if the lookup table is large¹².
2. **Partial Cache:** In this mode, SSIS caches only the rows that are required during the execution of the package. If a row is not found in the cache, it is retrieved from the database and added to the cache. This mode balances memory usage and performance¹².
3. **No Cache:** This mode does not cache any data. Each lookup operation queries the database directly. It is useful when the lookup table is very large or when memory usage is a concern, but it can be slower due to the repeated database queries¹².

These modes help optimize the performance and resource usage of your ETL processes depending on the size and nature of your lookup data.