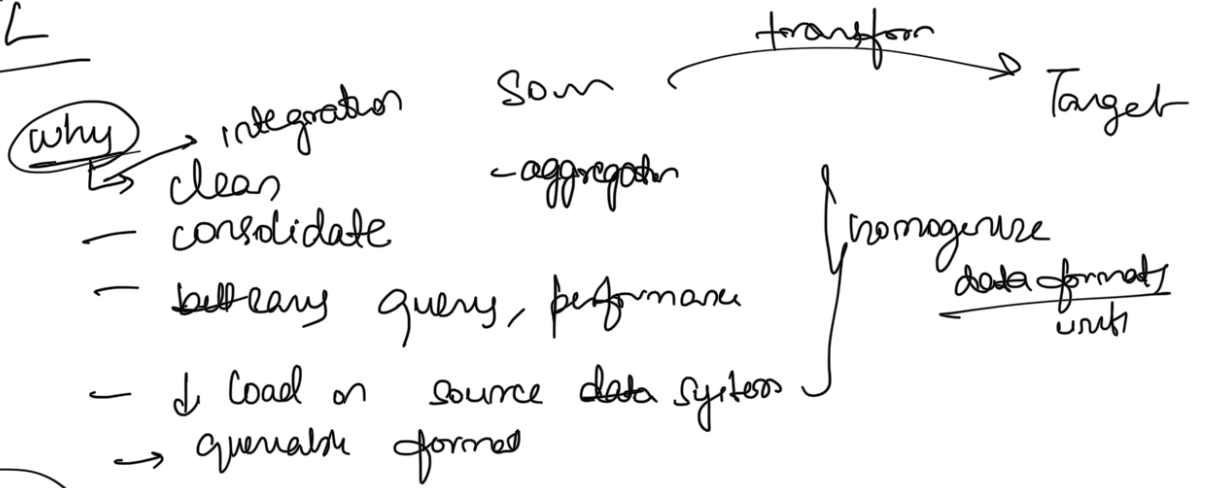


ETL Concepts

By Kusuma Mam

26th March 2022

ETL



Data

- collection of facts about obj/process
- Quant/Qual

OLTP

Online Transaction

ATM, POS (Point of Sale)

OLAP

Online Analytical processing

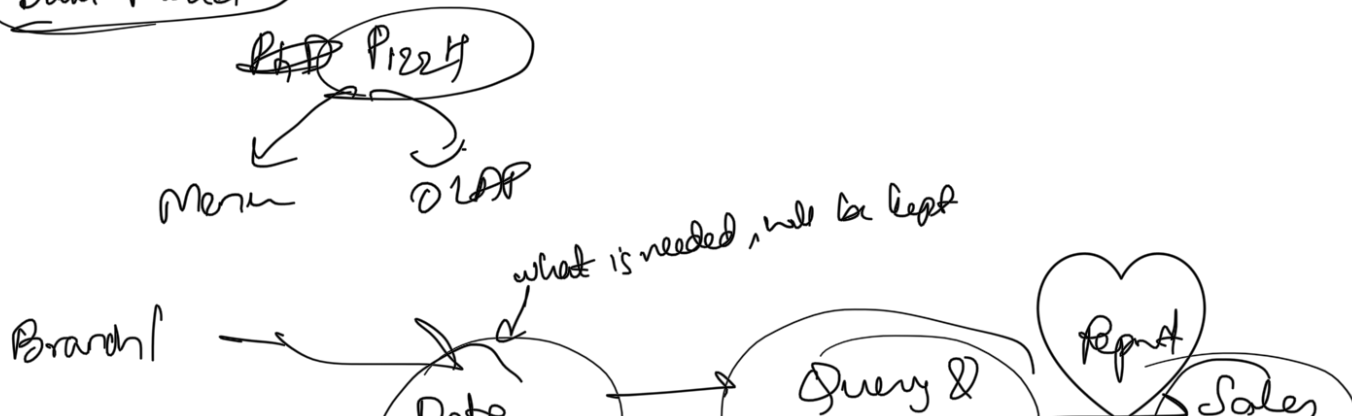
classification/computation

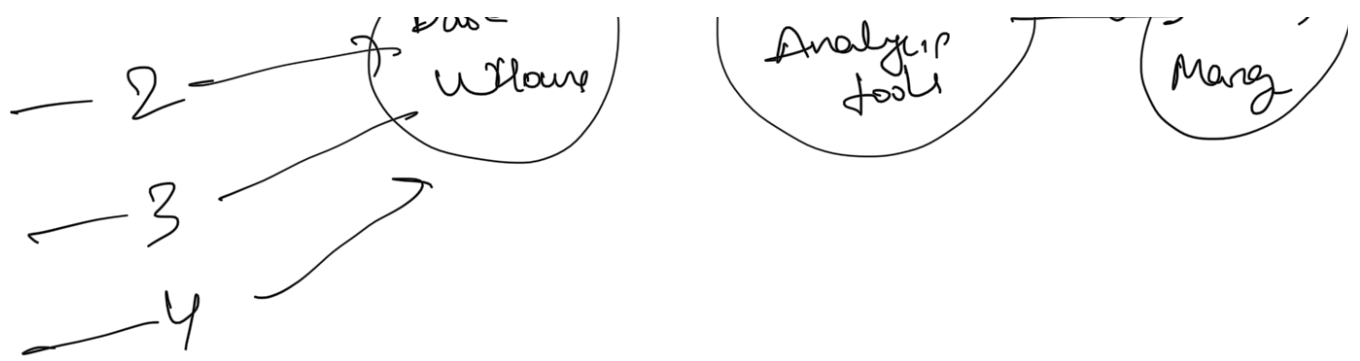
Ex: BI (Business Intelligence) Reports

Data Warehouse

→ large store of data accumulated from diff source

Data Model





ETL → copy data from source & des which represents differently from each other
 → based on consumers of data

Exa

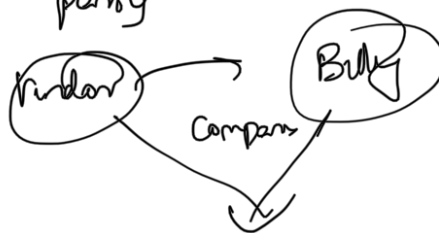
many data → Data warehouse

Migration

Unstructured data → store in BigData
Hadoop

Ex: Oracle to Hadoop
 ETL

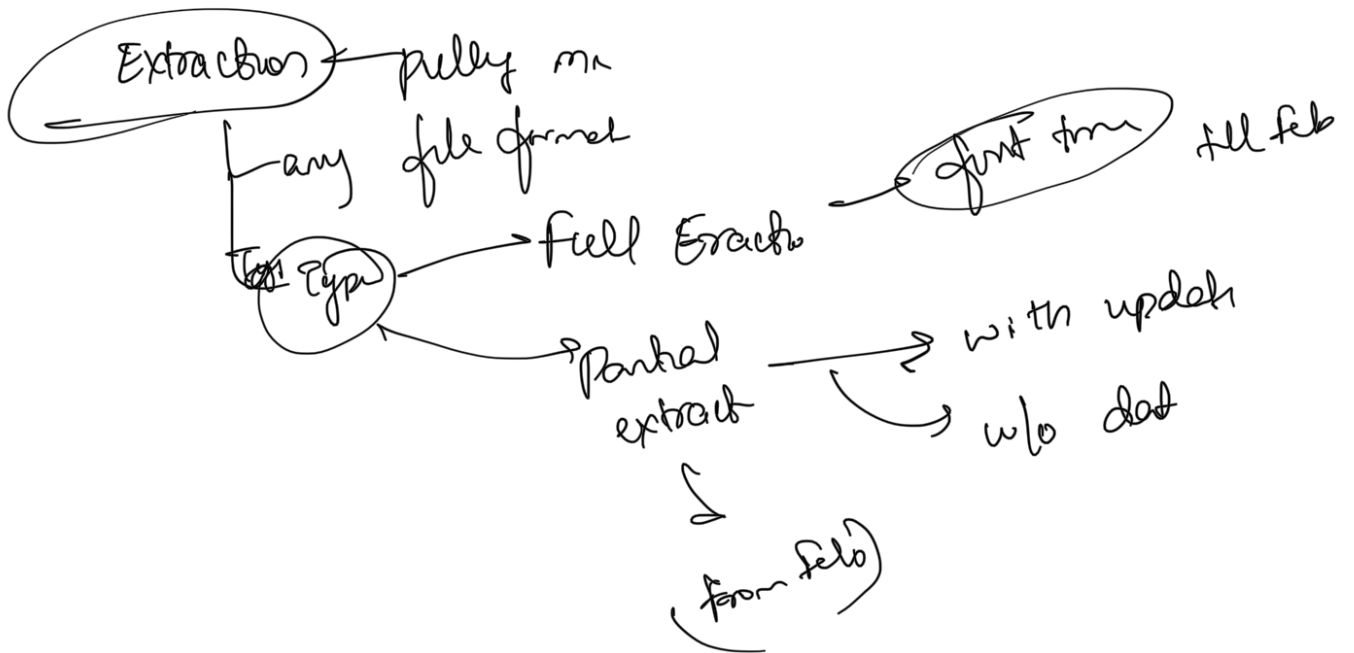
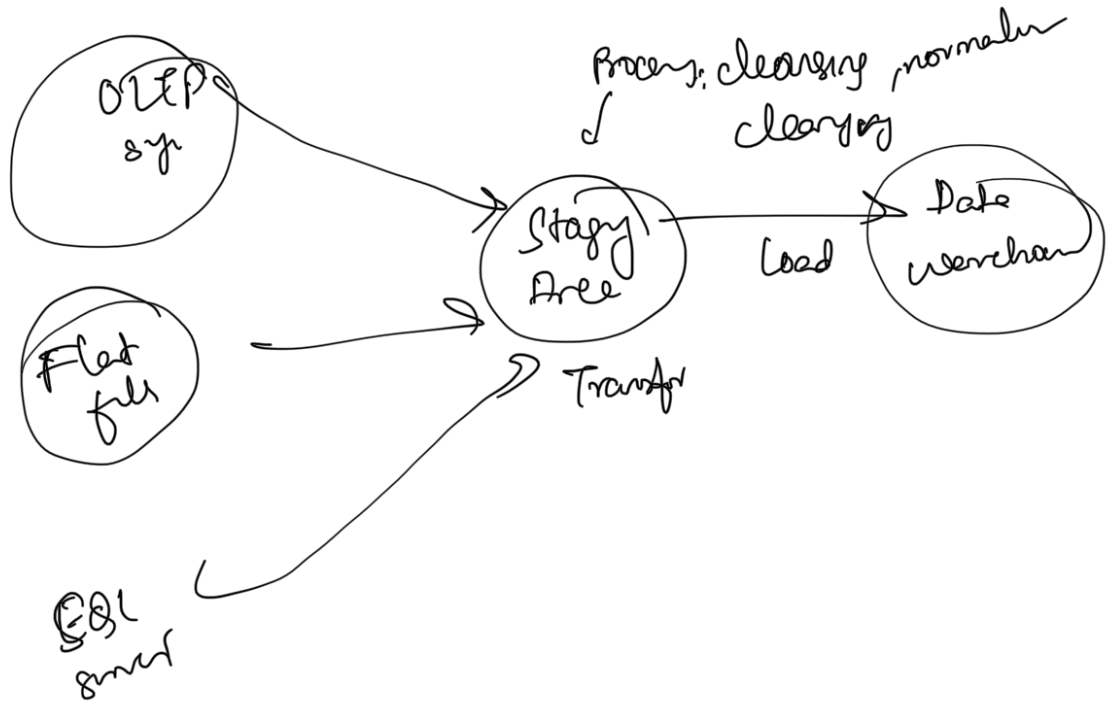
- mergers and acquisitions
- 3rd party



Data Integration

→ single view of Truth.

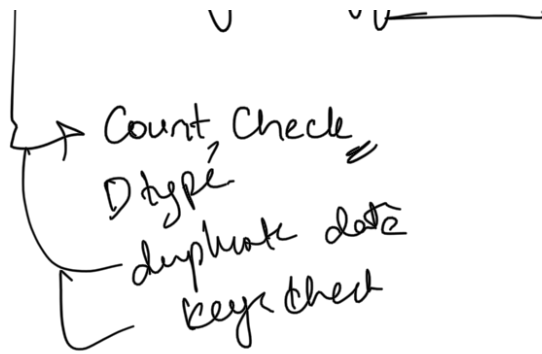
U
- Naming / Symbol. changes
- Abstraction



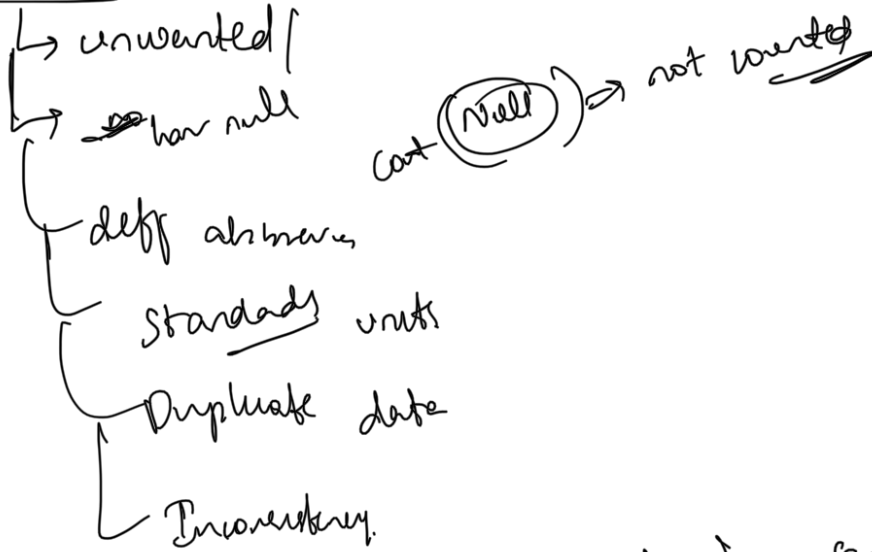
(Extract Validator.)

My First Bill

Mapping diff columns in diff table



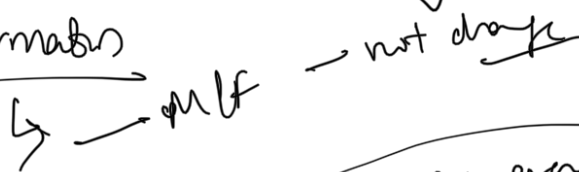
Dirty Data



Master Table entry should be seen as OLTP table

Data Transformation & Cleansing

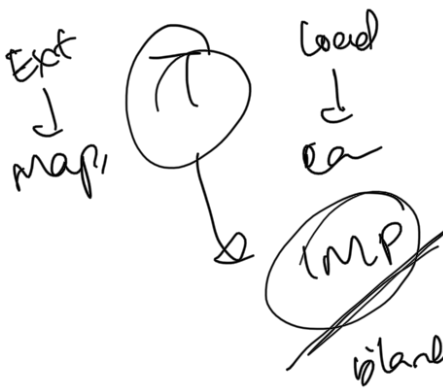
Transformations



~~Standard~~

currency changes everyday

Standardisation

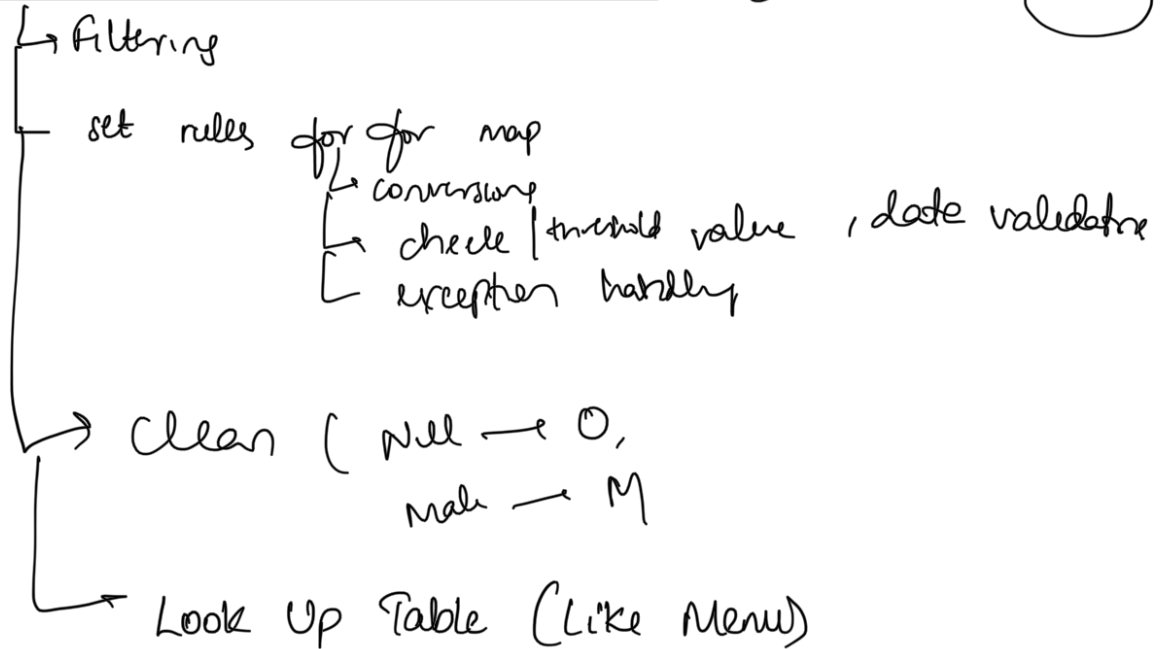


blank space to ~~AAA~~

concrete format name

D Transform & D Cleansing Scope

OLTP $\xrightarrow{\text{filter}}$ D. use



D Transform Validation

count key check

D Quality Check

- Null Value Check
- Referential Check
- Range Check

- List of value
Blank
Duplicate

- Data Gen

Before & After

~~Prefer~~

Exception
ch Table

~~outline~~

Blank
↓
space

VS Null

— 0

✓ I -
cont

Loading

↳ Load fails, Recover mode

Types

- Initial
- Incremental
- Full Refresh

Load

Techniques

- Batch Load
- Seq. Load
- Parallel
- Incremental

Loading Verification

↳ Key field

ETL tasks

- Understand Date to be used for report
- Review Date Mo (Target column needs)
- ...

Best Practices for ETL tools

ETL - Extract, Transform, Load

- Never try to clear all data
- Clear something
- Determine cost of clearing data
- Speed Up Query processing

Example: Delivery
holiday: done is

ETL tools

Enterprise class

Oracle Data Integrator
Informatica

Open Source

Informatica ETL

Apache Camel

Simple

Apache Nifi

AKA

How to choose ETL

Volume

Flexible data acquisition

Transform

ETL tools

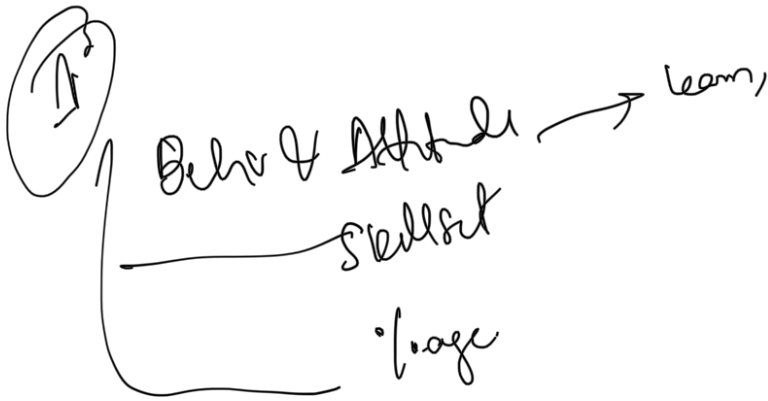
ETL tools

manage data in your system
source & derived

Summary

Core study

workload



Team players

imp than self-centred

Appraisal: Team Mindset

Team

Team completion vs

Whole prog
completion

! - "

after done,

help others

~~start~~ Company
Salary
Learning oppo }

Big Org. ↓ learn, 7 long roles