

Structured → predefined model
fields / tables / columns / spreadsheet

Numerical
[interval] [ratio]

Categorical
[nominal] [ordinal]



Types

unstructured

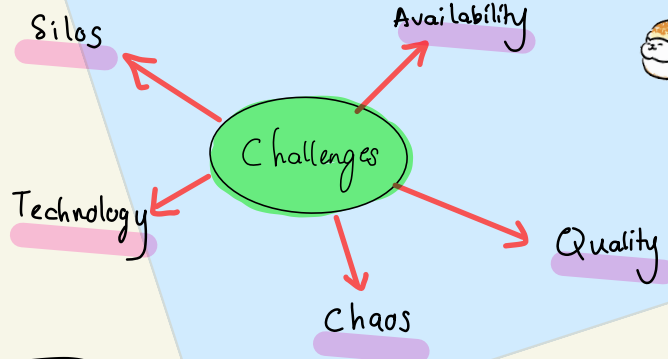
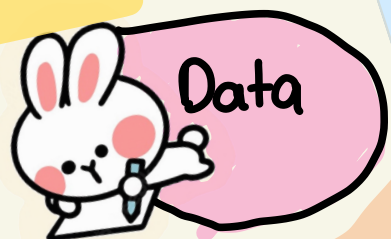
no field / attribute
binary files - audio / video / images

[Textual]
[Multimedia]
[XML / JSON]

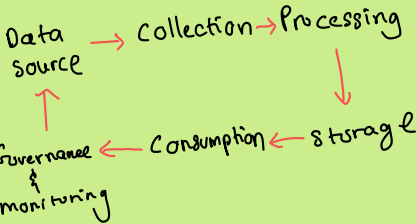
semi-structured

have markers / tags to separate elements

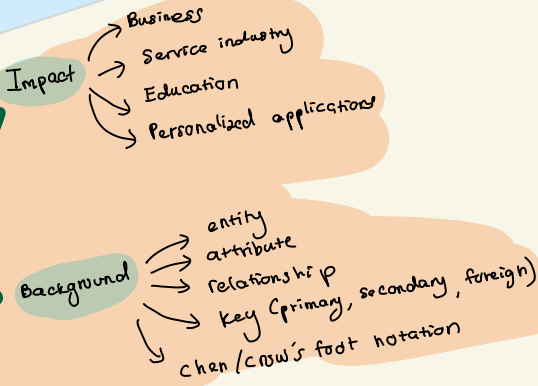
xml / HTML / email / webpage



Pipeline



Database



Extracting

- ① extract from sources
- ② Reading & comprehension
- ③ Copying & transfer

Issues

- source ident.
- method extrac.
- extraction freq.
- time window
- job sequence.
- exception handling

Transformation

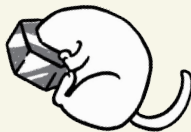
- ① selection
- ② Splitting / joining
- ③ Conversion
- ④ Summarization
- ⑤ Enrichment



Loading

- * Initial load
- * Incremental load
- * Full refresh

- Suitable for monolithic, legacy data sources



Data warehouse → enterprise DW
 ↓ → data mart → dependent
 metadata repo → independent

Source system → Component → OLAP

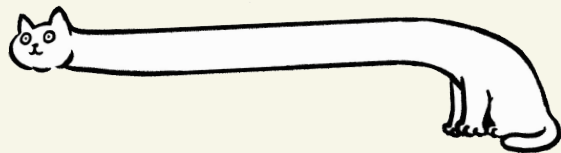
Component

Source system

End user

→ data-mining tools
 → reporting tools
 → statistical tools

Data staging → extraction
 → transform [clean, integ., aggreg.]
 → Loading



Data cube

2D 3D

OLAP

Operation

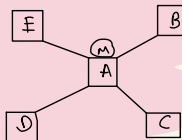
drill through
 → drill across
 → roll up (drill-up)
 → Drill down (roll down)
 ↓ slice & dice
 pivot (rotate)



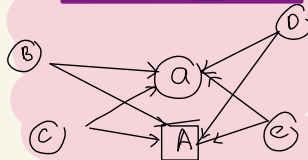
DW

Schema

Star Schema



Fact constellation schema



Snowflake Schema

