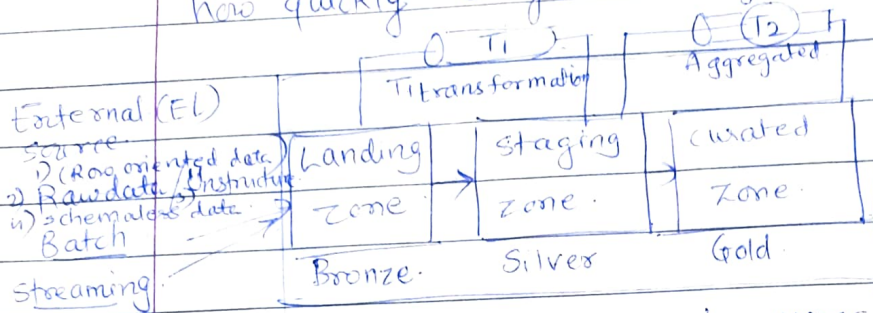# Data lake design

→ Objective — leverage data insights
how quickly u go to insights                    → horizontal split
                                                (partitioned / but
        ○ $T_1$              ○ $T_2$ ↑           2 -ordered
      $T_1$ transformation    Aggregated

| External (EL) | Landing | staging | curated | (Hive) |
|---|---|---|---|---|
| source: | Zone | Zone | Zone | Metastore |
| i) (Row oriented data | | | | → mapping bet[n] logical |
| 2) Raw data / Instructure | Bronze | Silver | Gold | name space to physical |
| ii) schema less data | | | | assets |
| Batch | | | | |
| Streaming | | | | |

5) Normalized          Datalake          on premise - HDFS
   (No redundancy)
6) compress data before ingestion.          cloud - Storage S3/ ADLS

<u>Landing Zone</u> — specific location on memory. → for particular
                    subject areas.
                    eg. 1) Amazon appparls landing zone.
                                      / subject
                        2) Retail business ₍ zones are landing zone.
$T_1$ transfor[n] ||     3) Saggrigated by time zone.

<u>Staging zone</u> — $T_1$ tran[n] ⟨ Scheduled
                                      triggered.

(more widely used) Triggered. — message queue. → publisher / Subscriber
    → queueing system.
    → because of variable injestion. triggered is used.
    → staging zone is called (source of all truth)
        all model u build, all dashboard u monitored, this
        is the zone. It has to be (versioned) (immutable)
    → Data is (denormalized) (join table and quickly get it)
        → have notion of downstream appl[n].
        → It is done in non linear fashion.
    → (columnar data). (parquette format)     OCR
    → horizontal split [(partitioned / bucketed)].
    → data here stays for longer duration.
    → during transformation, data cleaning is done.
        (missing value not in perspective of ML model)
    → keep marker to imputed data.
    → in healthcare missing data should not be imputed.
    → More imputation more noise.

→ compressed file can not be split form spark parallel processing
→ Splittable compression is used.

---

→ understand ingestion time, type of data, CPU uncompressing time and compressdata if needed.

[day end] → operationaldata ingestion→Trans — aggregate — available for next Busi Day.
              token/churning/coupons. prediction.

Benchmarking — find optimum Parameters.
Benchmarking — 1TB of data — how much time it takes.
2 TB — in same time →scale out → clouds are popular.
→ Access control list →. Landing zone access is subjective→file level read access
→ Staging zone → who has access to particular column, partition.
→ Landing zone access → Data engineers/staging zone — data analyst.

---

→ Data types conversion.

[star]
[snowflake] → schemabased datahere., [schema is enforced] here.
[scheema] → Two types of tables, facts and dimensions.
facts — Business operations are captured by facts.
Date, cid, tid$^x$, Pid$^{duct}$, sid$^{store}$, Amt. → T$^x$ data. [facts]

information about customer ⎫
              Product. ⎬ Reference data. → [dimensions]
              Store. ⎭

Curated Zone — pre computed statistics.
→ management wants to know KPIs. monthly. (aggregation).
→ freq of customer, regionwise
→ which store has— how many pitfalls.
→ moving average of revenue of last 4 weeks.
[Two core objective]→① Enable the lead time.
                    ② sources of features of ML model/statistical inference.
→ Some organization may push this data to datawareh/[feature store]
→ data here is in small size                          ↳ source of all truth ML
→ data is stored mostly row oriented fashion.

---

→ This is called referce architecture not standard. architecture
small organization may follow only two zones.
Large organization may have more than three zone.
→ metastores usually used in staging zone.