

## Week-2 Data Warehouse

Consolidate clean, accurate.

query = fast

⊙ GB to PB

☐ Serverless & no ops

⇒ ecosys. visualization & reporting tools  
" of ETL & data processing tools

Up to minute data

ML

Security & collab.

---

⊙ GB to PB

Bigquery

☐ Serverless & no ops

⇒ ecosys. visualization & reporting tools  
" of ETL & data processing tools

Up to minute data

ML

Security & collab.

X Data aging

X Storage management

- x fault recovery
- x Query engine optimization
- x Hardware
- x Updates

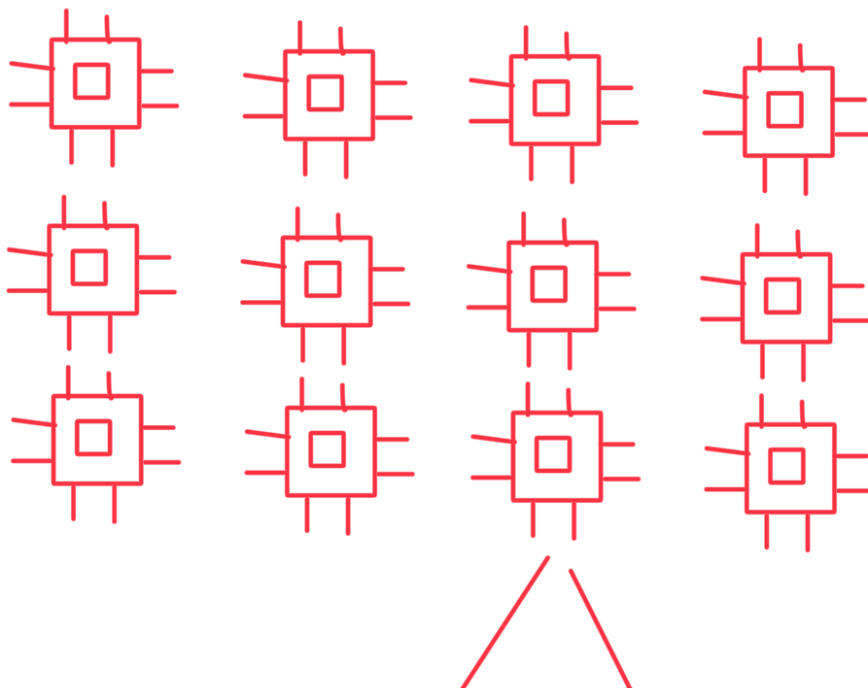
Big query = table-oriented (OLAP)

Storage engine      Analytic engine

Run length encoding  
Dictionary encoding

2000 slots

$\mu$ -service architecture



/ \  
CPU +  
RAM +  
Network

---

## Bigquery Services

Multiple dataset

billing → project  
Query → IAM

Region → multi zone

Google managed encryption key

predefined roles    row level security  
                                 ↕  
                                 policy

authorised view

Analytics UI

view → virtual table

materialized view

SQL

# pricing calculator

EL

ELT

ETL

CSV

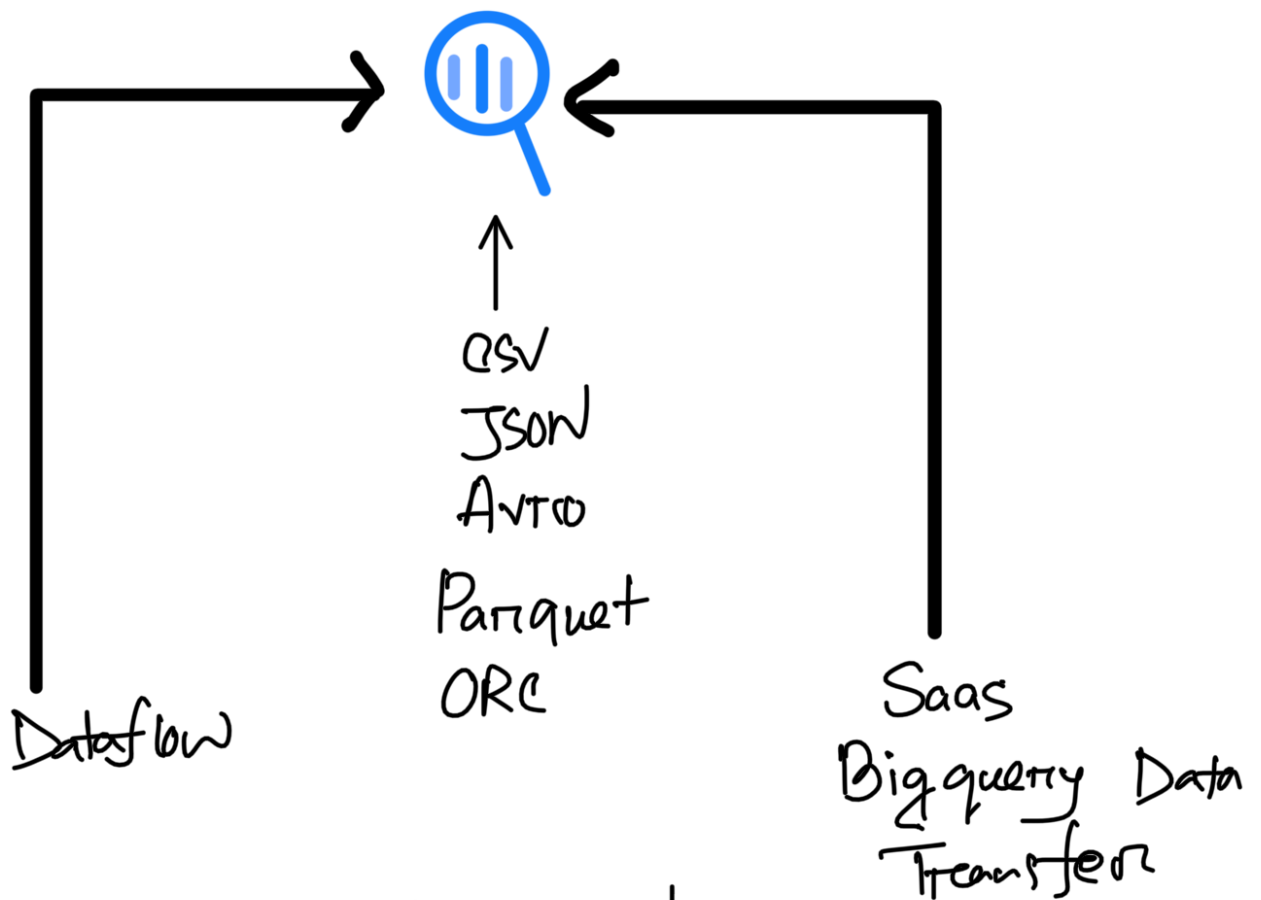
AVRO

Newline-delimited-JSON

DATASTORE-BACKUP

PARQUET

ORC

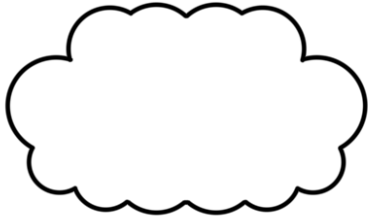


# Automate execution

point-in-time → 7 days

Transferring data →

backfilling



ideal

Support DML statement  
insert, update, delete, merge

DDL Statement

UDF → User defined function

---

### Practise

row\_count

Size\_of\_bytes

taxi-cabs

Bonus query

1.89

partitioning

Movie recommendation

Create function

---

### Schema design

Normalizing data

Denormalizing data (before loading)

Shuffle data

Nested & repeated → columns  
↓  
improve efficiency  
With relational data ←

---

Go-Jek 13+ PB

Normalizing data

Joins → costly

Data → repeated

Data → nested → repeated

ARRAYS → repeated

STRUCTS → pre joined table  
in a 1 to

Big query → column based

Type → Struct → Record

Arrays → regular field → Structs

Single table → STRUCTS (many)

---

nested & repeated

repeated values

Slot time compute

500

Multiple VCs  
unnest Command

18 GB

7.5 GB

Bytes shuffled

① → CROSS JOIN

inside dataset (gtc\_value)

---

① ✓ instead joins → advantage → nested & repeated fields  
denormalized

① ✓ keep dimension table (< 10 GB) normalized

① ✓ Denormalize table > 10 GB

---