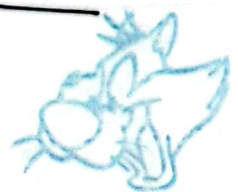# Week 2 : Data Ingestion

- Data Lake.
- Work flow orchestration.
- Airflow locally.
- Ingesting data to local postgres.
- Transfer service (AWS > GCP)

---

Data lake:
- Central repository that hold <u>big data</u> from many sources.
- Data may be structured or not.
- Main goal of DL to <u>ingest data as quickly</u> as possible <u>making it available</u> to others.
- DL should be <u>secure, scalable, Non-expensive</u>

| | Data Lake (DL) | Data warehouse (DW) |
|---|---|---|
| Processing | unstructured data, raw undergone minimal processing | structured data, refined data has cleaned, preprocessed |
| size | Large, petabytes(1024)TB Data transformed in need only data stored <u>indefinitely</u> ~'till | small compared, TB, alway processed before ingestion, data may <u>purged</u> abs periodically. |
| users | Data scientist, analysist | Business analysist |
| steps | ETL "Export Trasform load" extract ELT Transformation last layer Aws, GCP, Azure | ~~ETL~~ ETL |

. Data lakes come into existence because as companies
started to realize importance of data, they soon
founded they couldn't ingest data to DW, but
they didn't want to waste uncollected data when
their devs hadn't finished developing necessary relationship
for DW, so Data lake was born to collect any
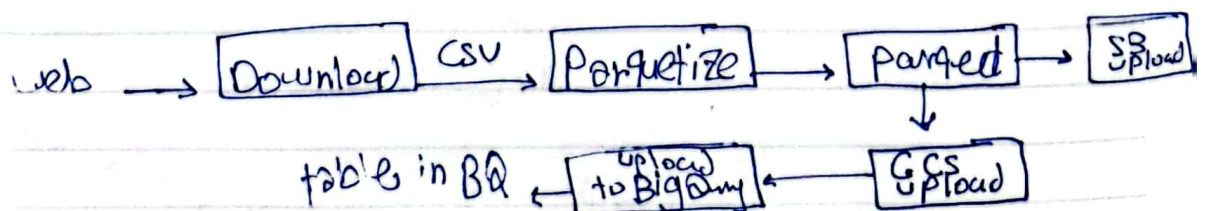potentially useful data till structure and organise
this data.

. Data swamp (مستنقع)
    ( Data lakes gone wrong)
    - Data lakes are only useful if data can be
      processed from it that happen when:
        - No versioning of data.
        - Incompatible schemes for same data.
        - No metadata associated.
        - Joins between different dataset not
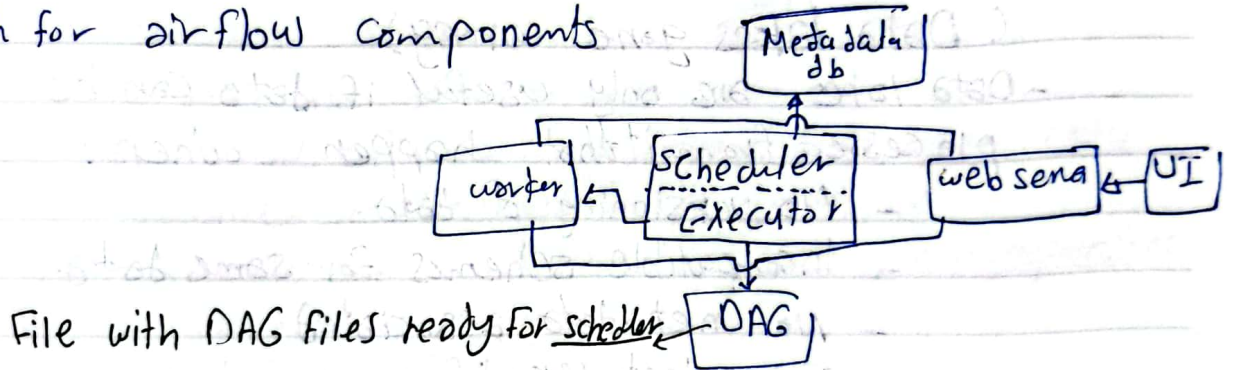          possible.

---

Orchestration (تنظيم) with Airflow:
For massive pipeline we call it workflow
in previous work we downloaded & process csv
file at same block, we should splet that
to different steps, separating help to avoid
problems like internet lose and repeat
download file.

web → Download → CSV → Parquetize → Parqued → S3 Upload
table in BQ ← Upload to BigQuery ← GCS Upload ←

- parquet (باركيت)
  Columnar storage datafile format, more efficient
  then CSV.

- The _previous_ graph are called Directed Acyclic Graph
  it miss for ops                                       (DAG)
  and flow is well defined.

---

Airflow Architecture:

Graph for airflow components

```
                                      ┌──────────┐
                                      │ Metadata │
                                      │   db     │
                                      └──────────┘
                                           ↑
        ┌─────────┐     ┌────────────┐     │     ┌──────────┐   ┌────┐
        │ worker  │ ←→  │ Scheduler  │ ────┘     │ web sena │ ← │ UI │
        └─────────┘     │ Executor   │           └──────────┘   └────┘
                        └────────────┘
                              │
                              ↓
                          ┌──────┐
File with DAG files ready for _scheduler_  │ DAG  │
                          └──────┘
```

- setting up docker ( Airflow )
  - full version            - lite version
  - we use docker compose for this task.

  ┌─────────────────────────┐
  │ all done on your pc )    │
  └─────────────────────────┘

# Creating DAG "Ingest Data to GCP with Airflow"

. DAG is created as <u>python script</u> importing series of libraries from airflow.

> with    DAG (dag-id="my-dag-name") as dag:
         op1 = Dummy operator (task-id= "task1")
         op2 . . . .
         op1 >> op2
    executing  op1 then op2 . . .


   . we will focous on <u>operator tasks</u> :
     - Bash operator : simple bash command.
     - Python operator : calls python method rather bash.


> download - dataset task = Bash operator (
    task-id =" download-dataset-task",
    bash_command = f "curl -sS $ dataset-url? > $path . . ..")
       instead of <u>wget</u>

> Format-to-parquet-task = Python Operator (
    task-id = "format-to-parquet-task",
    python-callable = format-to-parquet,
    op-kwargs = {
        "src-file": f " $path . . . }"


  - we can schedule intervals  to run
       code at a certain time.
          using   <u>Cron Job</u>
space q̇

          ┌─ * * * * * → day of week
   minute(0-59)  Hour  day  month  (0-6)   7(sundays) also
               (0-23)  month      (sun: sat)
               (1-31) (1-12)

0 6 2 * *
means **each month** **second day of this month**
at 6 am.

for example:
   Download this file each file and ingest it
   to dataset.
      • schedule_interval = "0 6 2 * * "

## Tricks

• URL = 'yellow-trip-data-{{ execution-data .strftime('\%Y-%m
                                                        -%m')}}.csv
      └──→ to download current data

with my-workflow as dag:
      download-task = Bash Operator (
         task-id = 'download'
         bash-command = (f) curl -sSL {URL-Template} >{output}
                        └ to pass      └ scilent
                          argument       show error

• Airflow offer series of pre-defined
   variable and **macros**.

• A DAG is idempotency( تسلسل ) , if end result
   is identical regardless multiples times.

• we have also **email operator** send emails.
   and even more , like GCP tools

2.3.3    ingest data to local postgres

- Using (Airflow) we can run certin code at certin time like
  - running python code to do something.
  - running bash.

- we are downloading NY- taxi data set using bash in airflow and want to repeat that every month using Cron.
- After that we would ingest this data to postgres & we will use python operator and pass python file to it.

Hint to write in multible line press ctrl +D

- we can creat Network between 2 docker compose

- URL TEMPLATE = URL-PREFIX + '/yellow-tripdata-
  {{ execution-data. strftime(\'%Y-%m\') }}. csv'

get current time when we execute code.

## 2.4.1 GCP Transfer Service

- GCP Enable us to trasfer files from different resources like ( AWS, Azure, url...)
  speed of transfer : 790 MB/S
  Charge : •012 $/GB

. we can do trasfer using
  - UI
  - Teraform
    "Just sheet which you edit variables"