

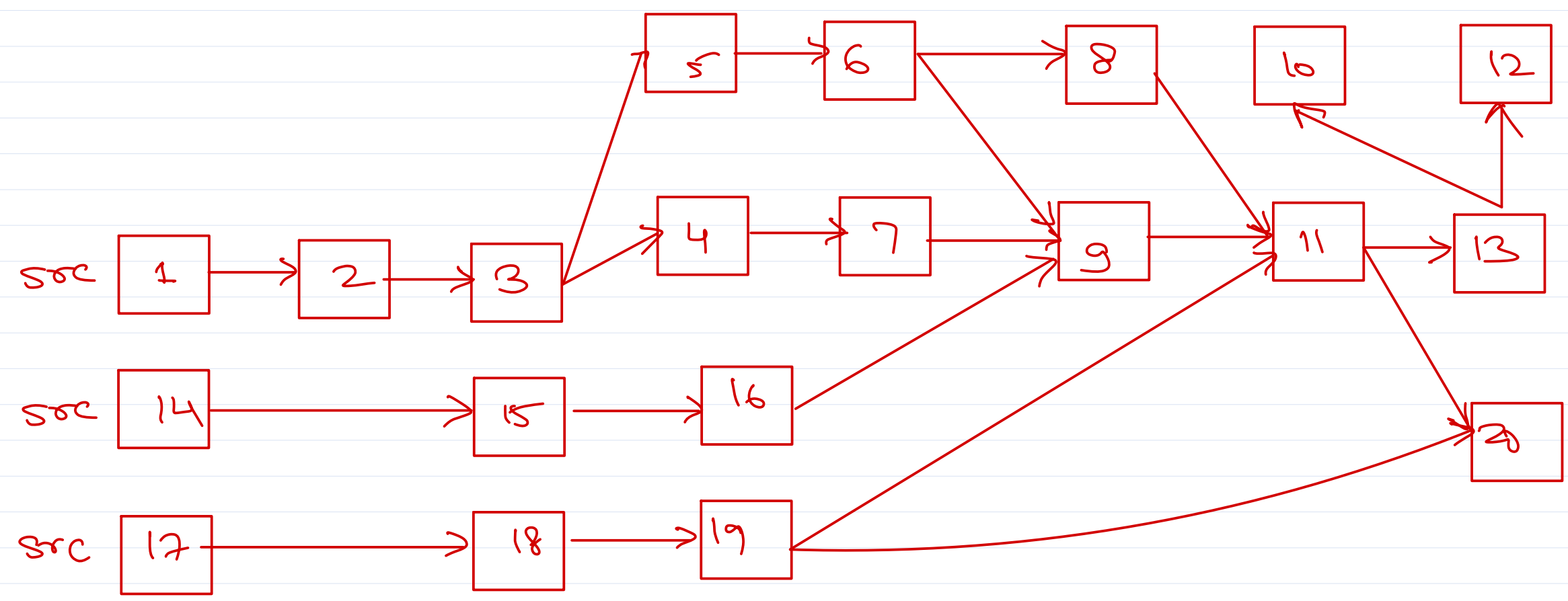


Big Data Technologies

Trainer: Mr. Nilesh Ghule.



Complex DAGs



DAG

task1 (download)

t1. >> t2 >> t3 >> t4

weather readings



Weather Station

Bash Operator / Python Operator (wget)

LOAD DATA
...
INTO ...
ncdc_staging

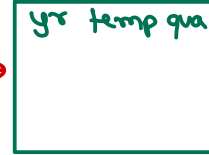
ncdc_staging



Hive table

INSERT ...
ncdc_orc ...
SELECT ...
ncdc_staging

ncdc_orc



Hive table

ALTER MATERIALIZED
VIEW mv ...
REBUILD;

mv_ncdc
yr avgtemp



Hive MV.

Power BI

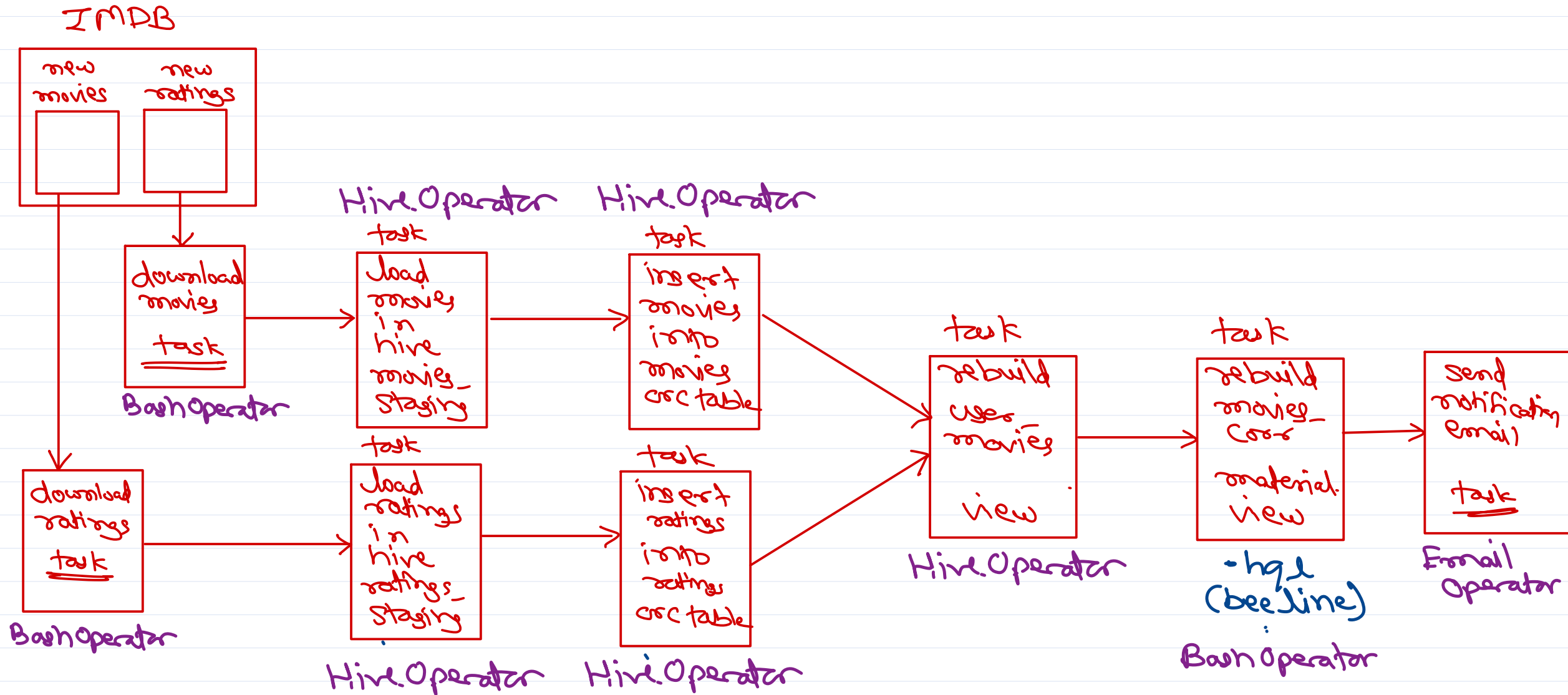
task2 (execute
hive q1 to
load data into
staging table)
Hive Operator

task3 (execute
hive q2 to
insert data
from staging
to main ncdc
table - orc).
Hive Operator

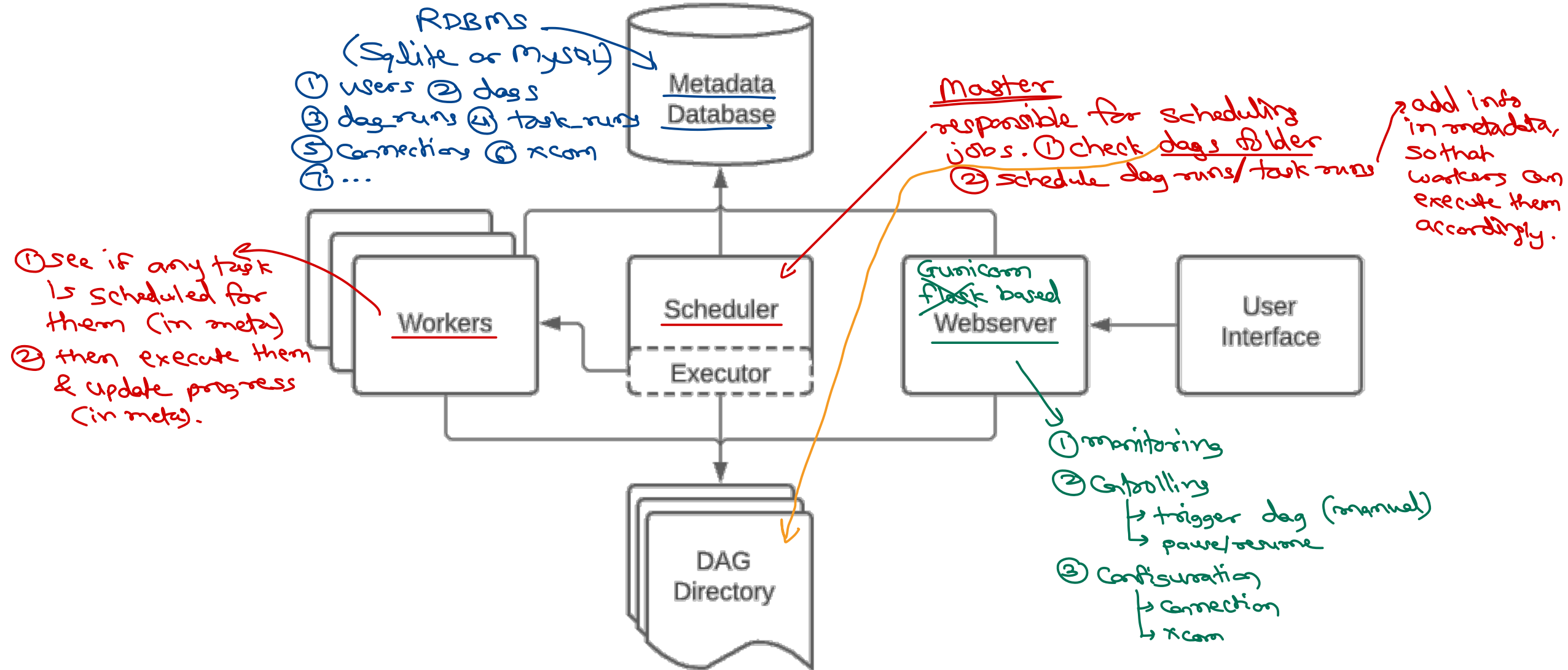
task4 (execute
hive q3 to
update/rebuild
materialized
view - yearly avg
temp)
Hive Operator



Movie Recommendation Workflow

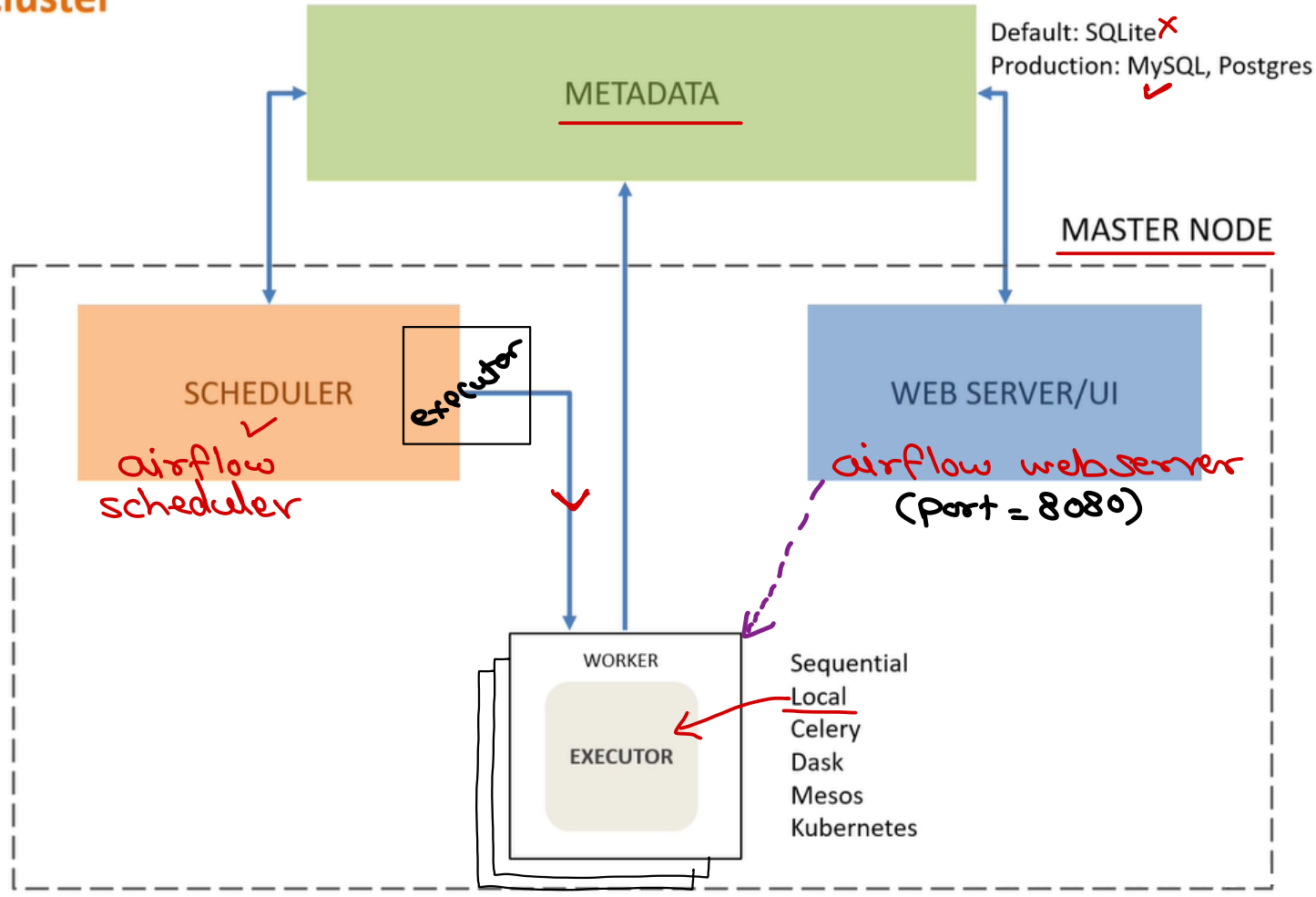


Airflow Architecture



Airflow – Single node installation

Single Node cluster

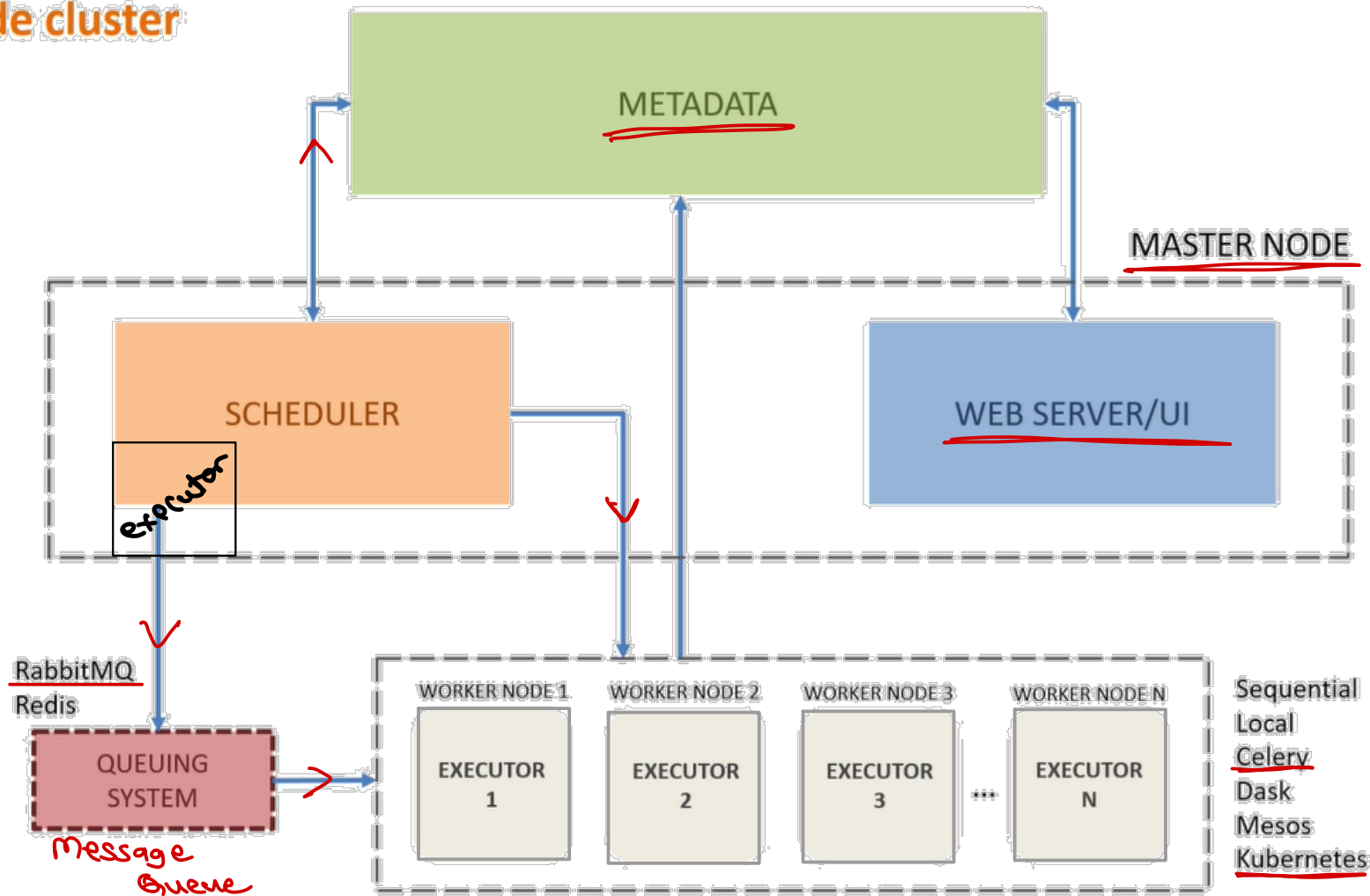


- Single worker node.
- Not scalable (max all resources available on the system).
- Quite fast for limited number of DAGs and data size.
- Use local executor.
- Direct communication between scheduler and ~~executor~~ workers



Airflow – Multi node installation

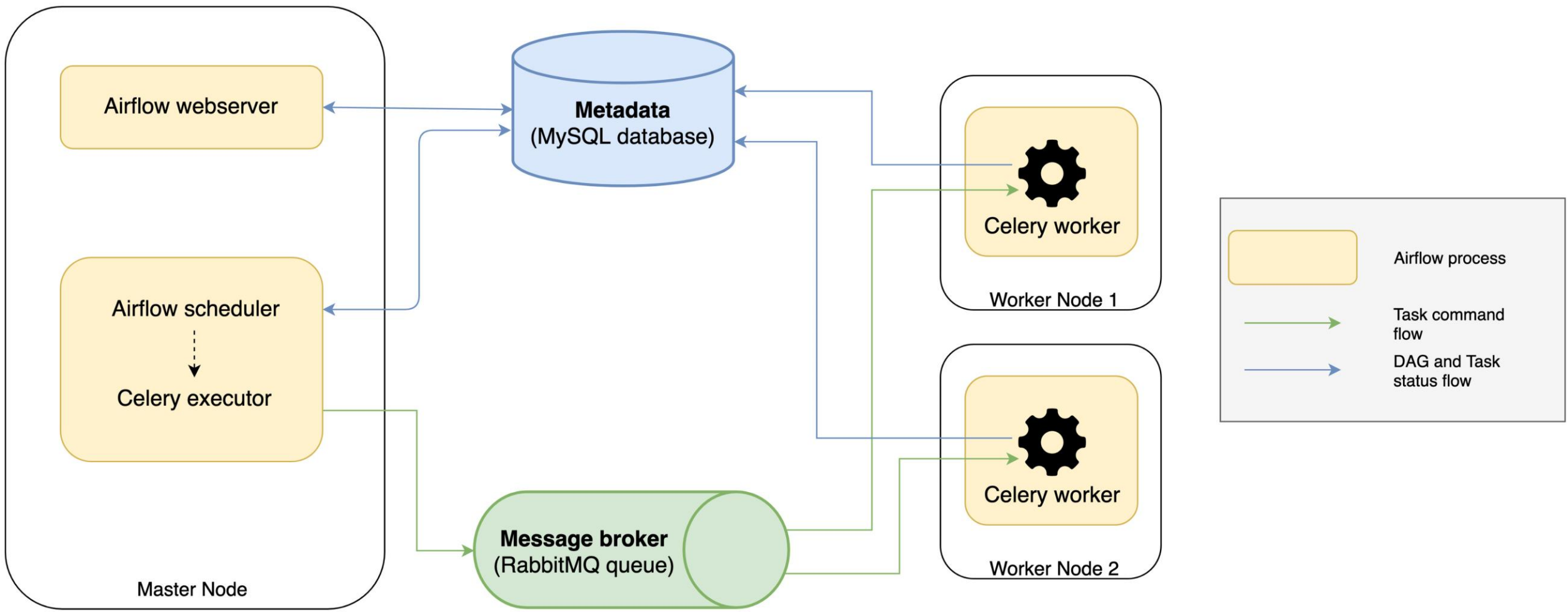
Multi Node cluster



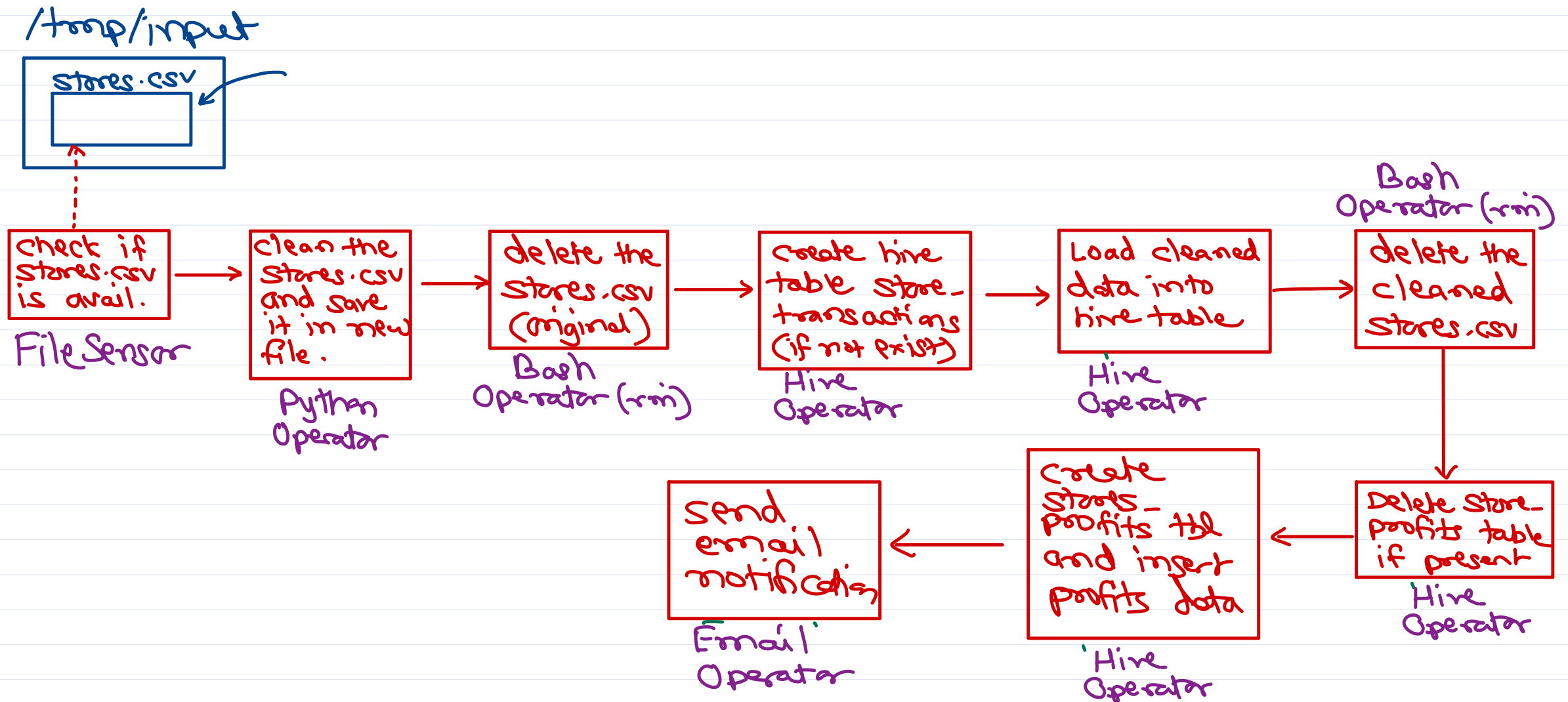
- Master node runs scheduler and web-server.
- Multiple worker nodes - running executors/worker
- Highly scalable for huge data and many nodes.
- Recommended executor is Celery.
- Scheduler and workers communicate with external queue system like Rabbit MQ / Redis.



Airflow – Multi node installation



Data ingestion and Reporting workflow





Thank you!

Nilesh Ghule <nilesh@sunbeaminfo.com>

