Data-Intensive Distributed Computing

CS 451/651 (Fall 2025)



Rubber, Meet Road: Orchestrators (v1.0)

Week 5: September 30

Khaled Ammar Rocket Innovation Studio



These slides are available at https://lintool.github.io/cs451-2025f/

Khaled Ammar

Director of Data Science – Applied Research Rocket Innovation Studio



This Week

Now: Orchestrator

DAGs, ML pipelines, Airflow, etc.

Next: Data Management in Production

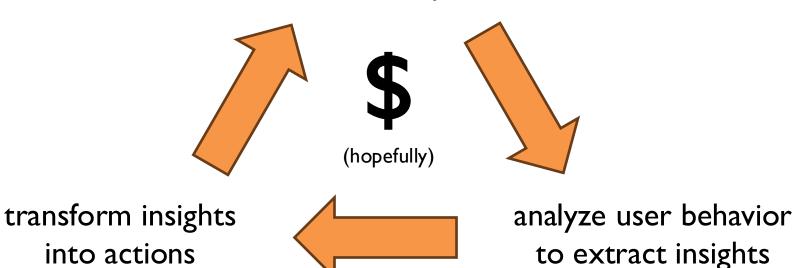
Data Governance, Metadata, Feature Stores

Context...

The Data Flywheel

(a virtuous cycle)

Build a useful product















Context...

What's this course about?

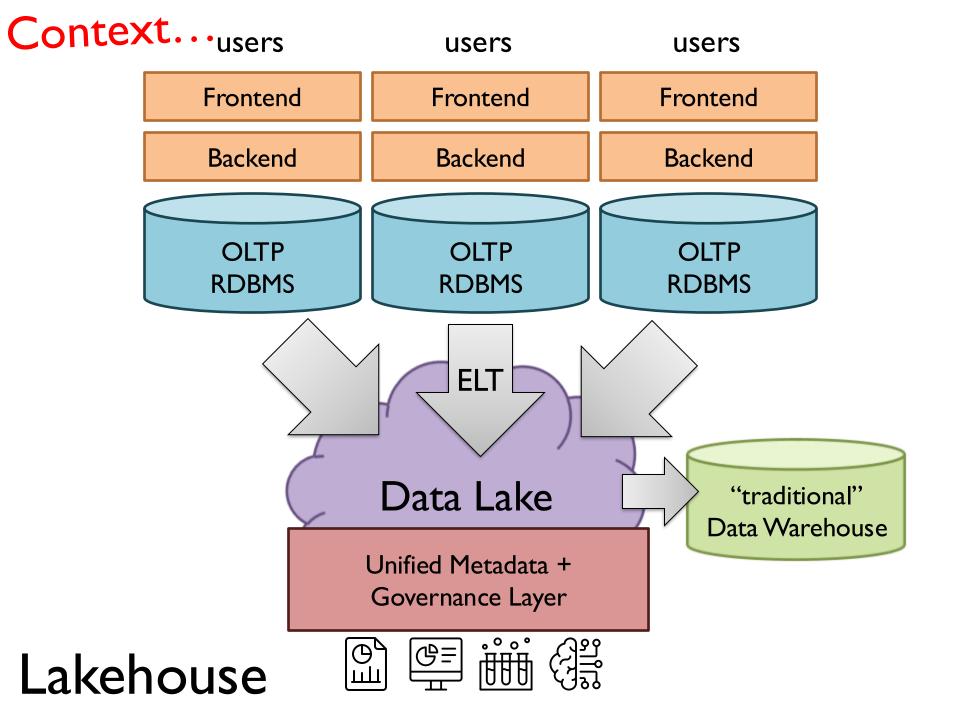
The infrastructure that supports the data flywheel.

data platforms + data engineering

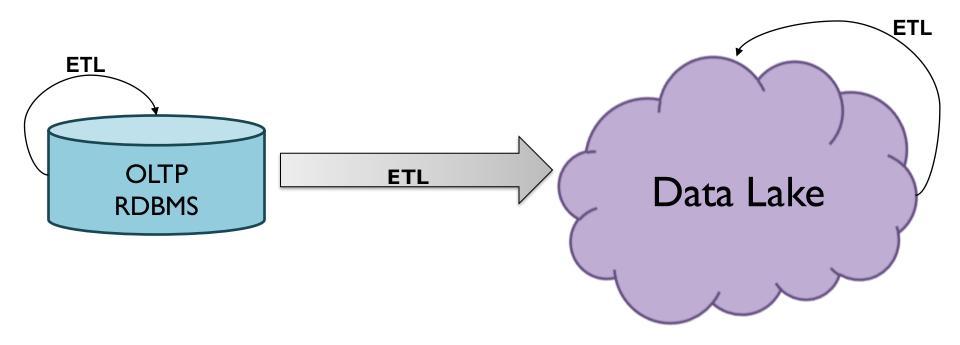
Context...

What problems do data platforms solve?

Ingesting, storing, manipulating, maintaining, serving... the data that supports the data flywheel.







Creating Spark jobs is not enough! How are you going to run it?

Creating Spark jobs is not enough!

How are you going to run it?



Script



Cron Jobs

Run multiple ETLs

Wake up in mid-night to run it

Manual retry

Multiple scripts are not easy to handle!

Run multiple ETLs

Scheduled to run automatically

No retry

Multiple jobs can run over each other

Workflow orchestrators

Declare a set of tasks, specify dependencies, retry policies, schedules, triggers, and monitor the whole pipeline.

Run multiple ETLs

Scheduled to run automatically – not only based on time but also using triggers!

Flexible retry policies

Clear dependencies between jobs

What is Workflow orchestration?

It is the process of managing the execution of independent tasks (e.g. Spark jobs) across multiple systems and time.

Common features:

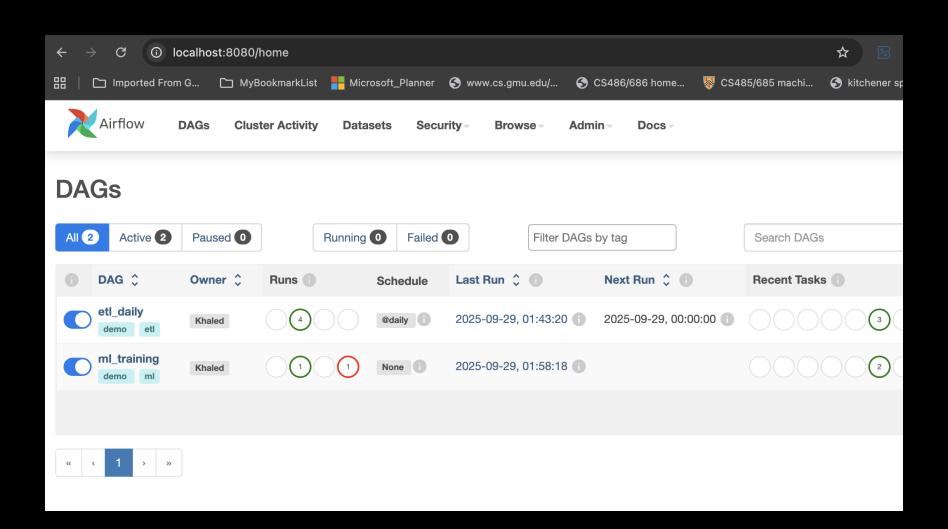
- **Dependency management:** specify that task B depends on task A's success.
- Scheduling: specify when a task can run and how often
- Retry policies: specify how many times to retry a failed job, any delay policies.
- Backfills / catchup: specify how to handle missed past intervals.
- Monitoring & Alerting: view the overall status of the system including failures and latency.

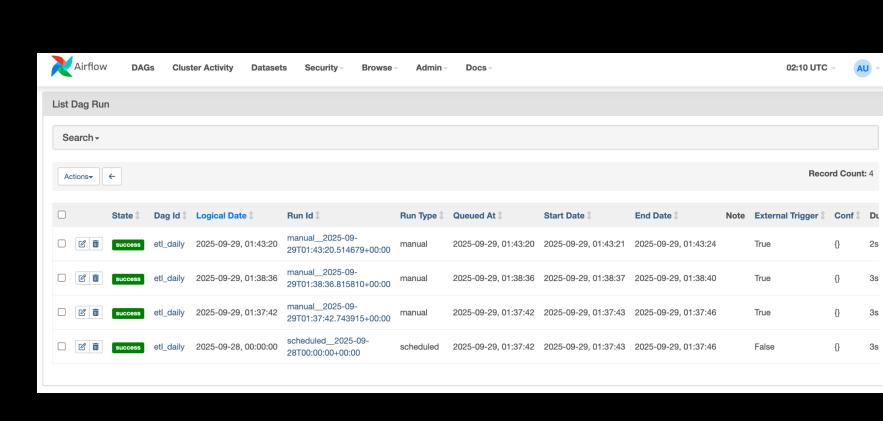
Airflow DAG Example

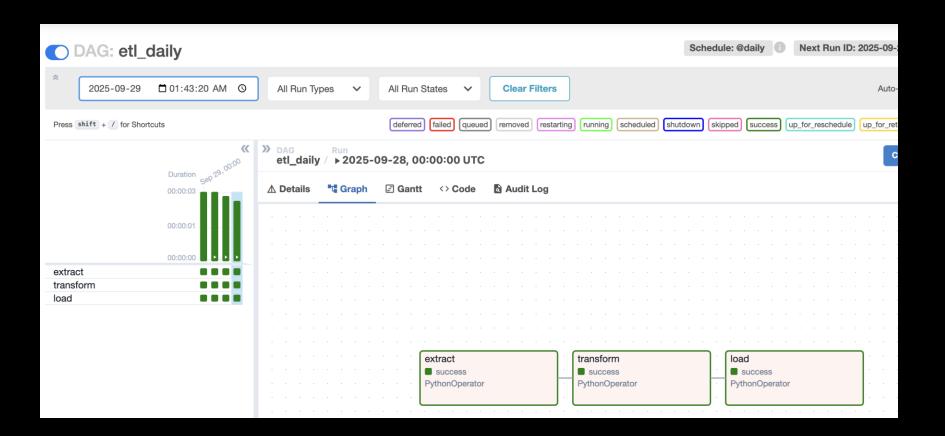
```
IZ
13
   default args = {
        "owner": "Khaled",
14
                                      Stateless DAG
        "depends_on_past": False,
15
        "retries": 2,
16
        "retry_delay": timedelta(seconds=10),
                                                  Retry policy
17
18
   }
19
20
   with DAG(
21
        dag id="etl daily",
22
        default args=default args,
                                              Define DAG details
23
        start_date=datetime(2025, 9, 1),
        schedule interval="@daily",
24
25
        catchup=False,
26
       tags=["demo","etl"],
   ) as dag:
27
28
29
        extract = PvthonOperator(
30
            task id="extract",
                                                  Define Extract step
31
            python callable=extract run,
            op_kwargs={"day": "2025-09-03"},
32
33
34
35
        transform = PythonOperator(
36
            task id="transform",
                                                  Define Transform step
37
            python callable=transform run,
            op kwargs={"day": "2025-09-03"},
38
39
40
41
        load = PythonOperator(
                                                  Define Load step
            task id="load",
42
43
            python_callable=load_run,
44
45
                                                 Order to run these steps
46
        extract >> transform >> load
```

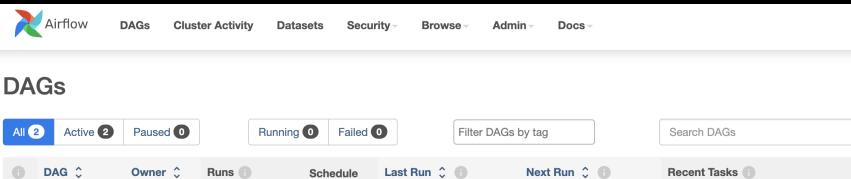
Airflow DAG Example

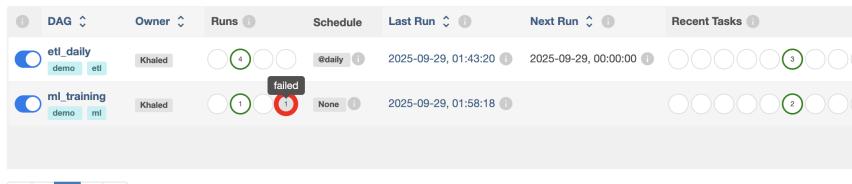
Live Demo

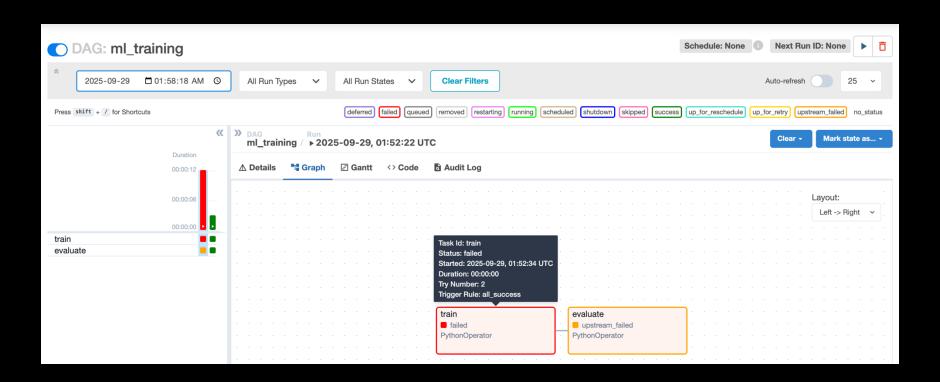




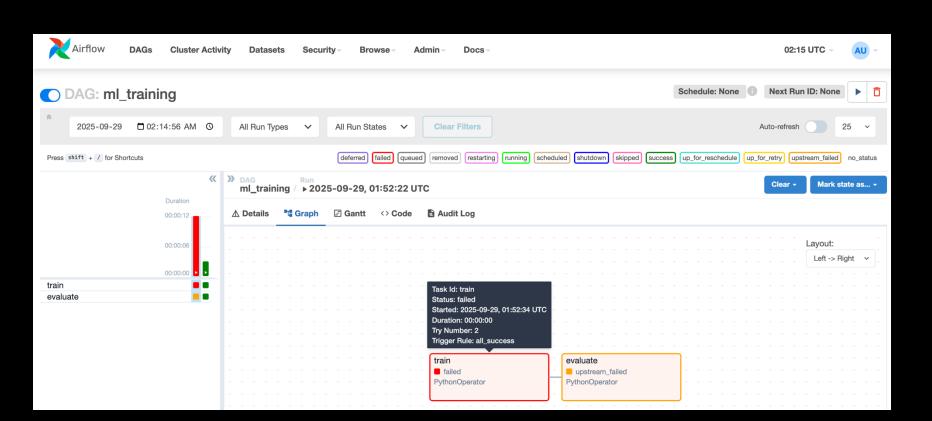


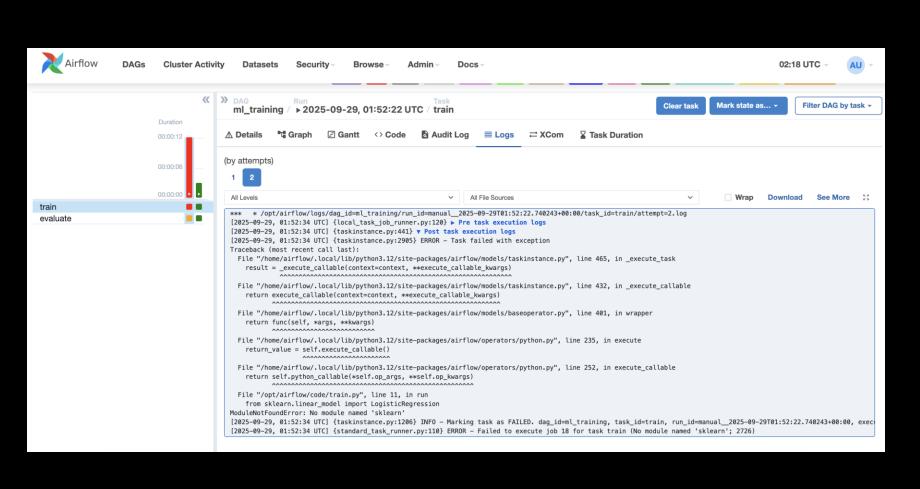




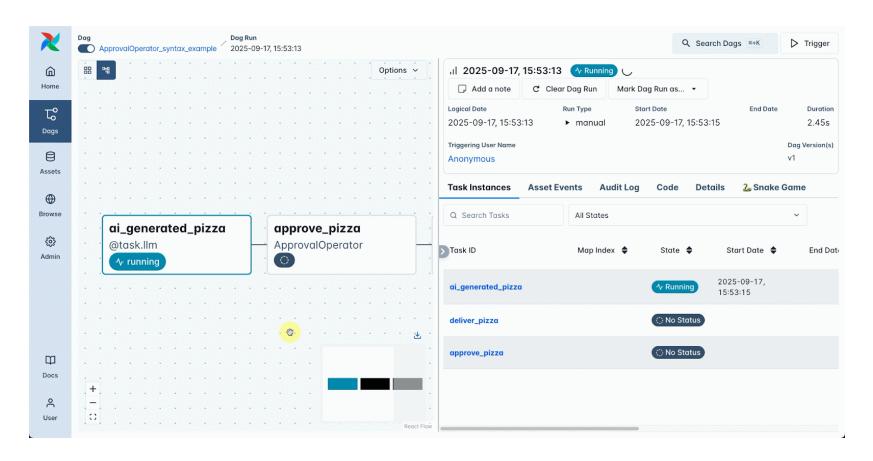


```
12
   default_args = {
13
        "owner": "Khaled",
14
        "depends_on_past": False,
                                                  Retry policy
15
        "retries": 1.
16
        "retry_delay": timedelta(seconds=10),
17
18
19
   with DAG(
20
        dag_id="ml_training",
21
        default_args=default_args,
22
        start_date=datetime(2025, 9, 1),
23
        schedule_interval=None,
24
        catchup=False,
25
        tags=["demo","ml"],
26
      as dag:
27
28
        train = PythonOperator(
29
            task_id="train",
30
            python callable=train run,
31
32
33
        evaluate = PythonOperator(
34
            task_id="evaluate",
35
            python_callable=evaluate_run,
36
            op_kwargs={"threshold": 0.6},
37
38
39
        train >> evaluate
40
```





Airflow Recent release: 3.1.0



https://airflow.apache.org/blog/airflow-3.1.0/

State of Airflow in 2025

Thousands of DAGs are running daily.

Airflow remains the scheduler and dependency manager, not the compute engine.

Orchestrate jobs across heterogenous systems.

DAGs often have data quality and validation steps.

External data availability often start DAGs in integration use cases.

Some teams use DAGs beyond data and ML, including sending emails, moving files, etc.

Limitation of Airflow



DAGs are defined when the DAG file is parsed, but complex workflows are data-dependent!

Example: One task discover 20 tables in an S3 path, then need to create 20 different tasks with different schema definition. Everything is doable with code, but it gets messy!



Airflow is fundamentally a batch scheduler; one task needs to finish for another to start!

Example: Executing a long running (streaming) task like Kafka consumer that continuously consume data and generate results cause the DAG to stuck!

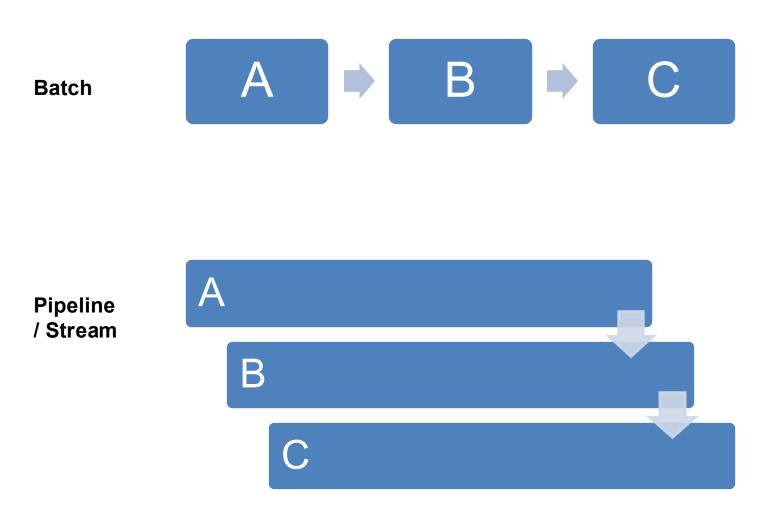
Modern Orchestrators

Tool	Strengths	Challenges
Airflow	Very mature, strong for batch ETL / Python orchestration	Static DAGs, no streaming support
Prefect	Supports dynamic workflow	Limited non-Python support, can support data flows but not optimized
dbt	Excellent for SQL and data lineage	SQL-only, and no streaming support
Kubeflow	Specialize in ML pipelines on K8; native TensorFlow/PyTorch integration	Requires K8 expertise, limited outside ML, heavy for small projects – no general support for streaming

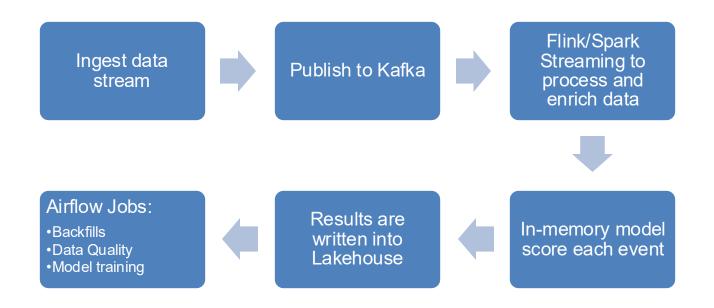
It is not uncommon to use multiple orchestrators in the same enterprise

Challenge with streaming workflows

A task can start processing the output before its predecessor finishes!



Streaming Orchestration



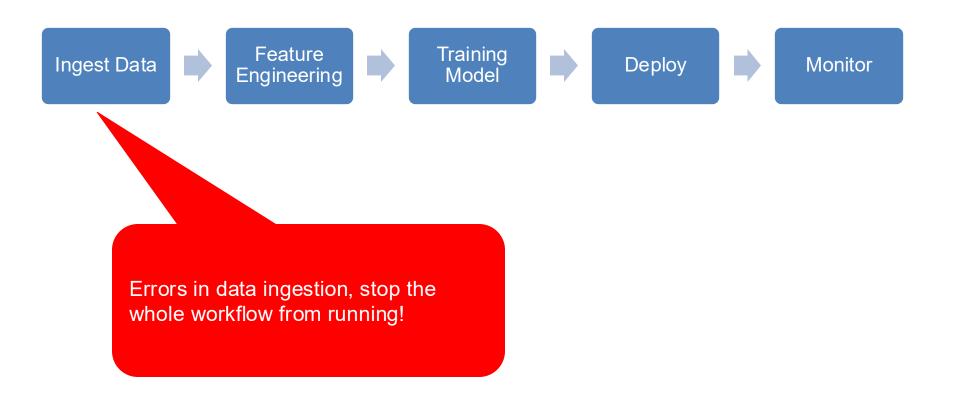
Ingestion: Kafka

Stream computing: Flink | Spark Streaming

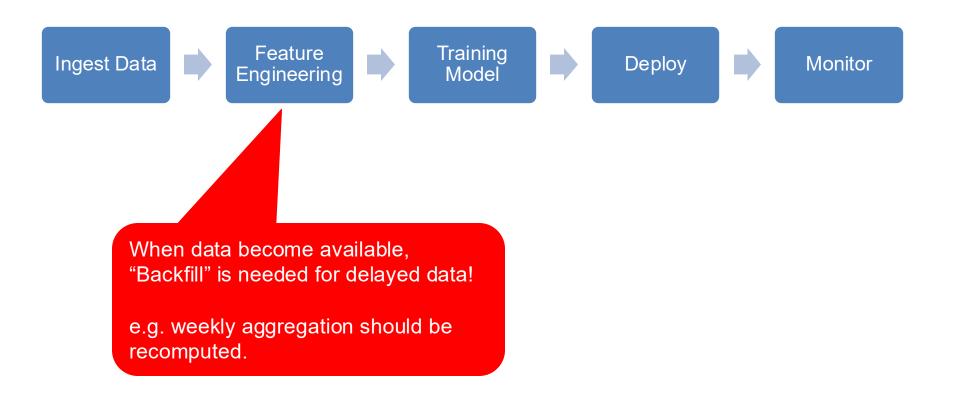
Model serving: Flink operator **Output:** Kafka | S3 | Lakehouse

Batch orchestration: Airflow | Prefect

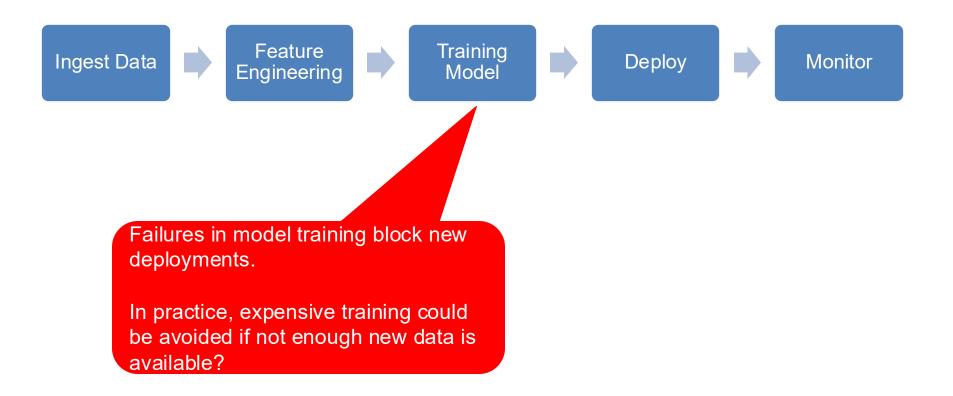
ML workflow



ML workflow

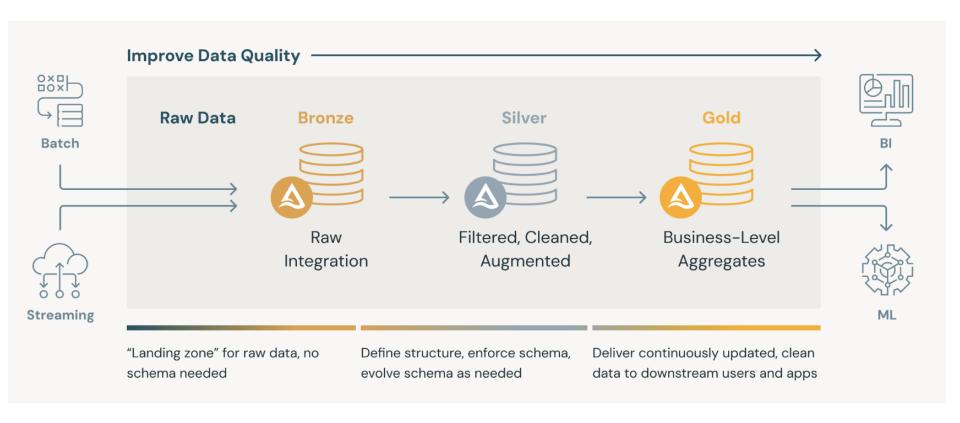


ML workflow



Data Science Setup in Enterprises:

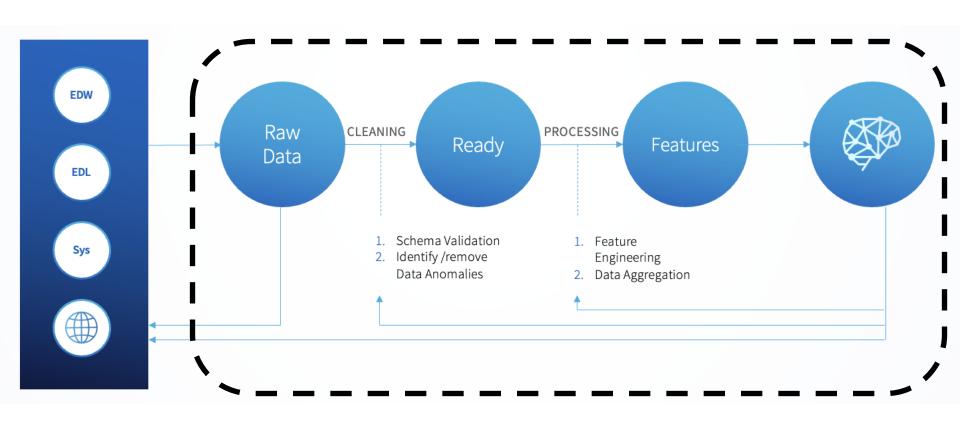
Medallion Architecture



Source: https://www.databricks.com/glossary/medallion-architecture

Data Science Setup in Enterprises:

Innovation Area



When not to use Orchestrator?

Is orchestrator enough for data pipeline?