Data-Intensive Distributed Computing

CS 451/651 (Fall 2025)



Rubber, Meet Road: Data Management (v1.0)

Week 5: October 2nd

Khaled Ammar Rocket Innovation Studio



This Week

Tuesday: Orchestrator

DAGs, ML pipelines, Airflow, etc.

Now: Data Management in Production

Data Governance, Metadata, Feature Stores

Is orchestrator enough for data pipeline?

Why Orchestration Alone is not enough?

- Discoverability
- Lineage
- Quality
- Governance

Tables: Transactions_cleaned_v2 & Clients_full_v4 Questions:

- I. Who owns this data?
- 2. What does each column mean?
- 3. How was the data prepared? cleaning process?
- 4. What is the refresh rate?
- 5. Can we legally use all columns?
- 6. ...

Who has chased a "mystery dataset" in a project before?

Data Catalog

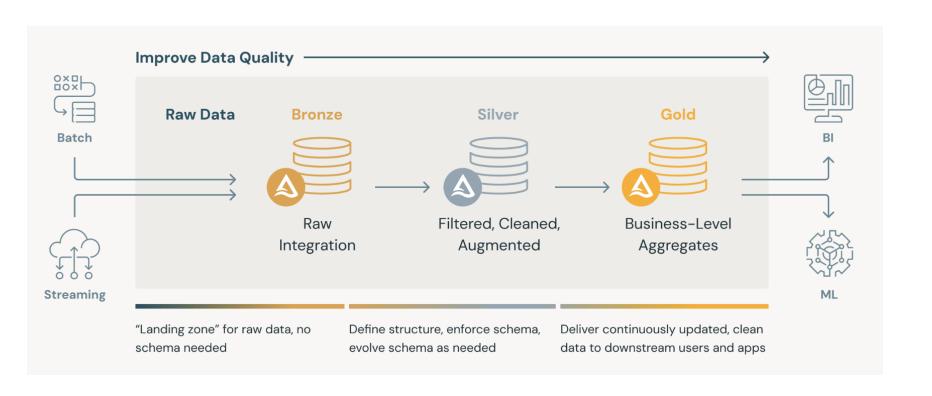
A data catalog is the system of record for your data assets.

It tracks:

- what exists (tables, views, files, ML models),
- who owns them,
- how they're structured (schema),
- how they're used (popularity),
- and where they came from (lineage).

Data Lineage

Record how data was constructed – useful for debugging and investigation.



Data Catalog Tools

Tool	Features	Limitations	Best fit
Unity Catalog	 Centralized governance for tables, files, ML models Fine-grained permissions Integrates with Delta Lake & MLflow 	Locked into Databricks ecosystem	Companies standardized on Databricks lakehouse
AWS Glue Data Catalog	Central metadata storeSchema inferenceLake Formation integration	weaker UX and less business metadata support.	AWS-native data lakes with Athena/Redshift
DataHub (LinkedIn OSS)	Active lineageRich metadata model (technical + business)Search & discovery UI	More engineering effort to run & maintain	Heavy engineering- focused companies.
Amundsen (Lyft OSS)	 Metadata search engine + lineage Simple interface (search-first) Integrations with warehouses + BI tools 	Less feature-rich than DataHub;	Quick-start OSS catalog for Bl/analyst teams





Search within a category using the pattern with wildcard support 'category: *searchTerm*', e.g. 'schema: *core*'. Current categories are 'column', 'database', 'schema', 'table', and 'tag'.

Browse Tags

tag1 1 tag2 1

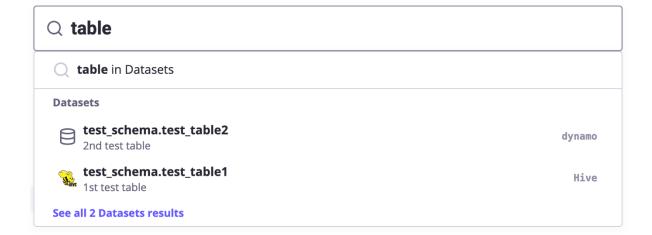
My Bookmarks

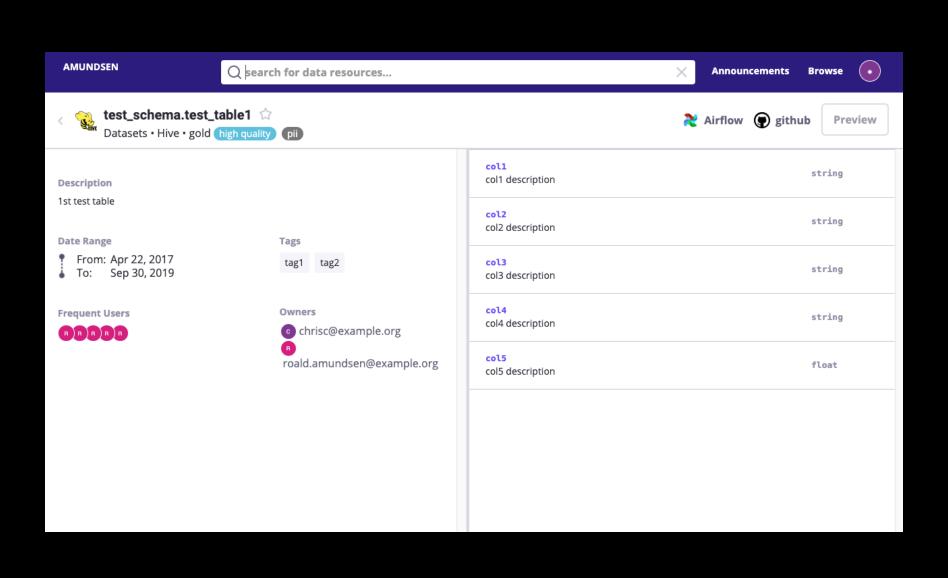


Popular Tables ①









col1

string

Description

This is an editable test description for the first column. This also supports **Markdown**.

Column Statistics Stats reflect data collected between May 22, 2015 and Jul 04, 2019.

distinct values	8	min
num nulls	500320	max
verified	230430	

min	aardvark
max	zebra

Unity Catalog Demo

Unity Catalog Search & Data Explorer

Databricks **Unity Catalog** is your security and governance layer across your entire <u>data stack</u>: <u>data pipelines</u>, tables, and ML models.

It provides:

- · Centralized metadata and user management
- Centralized data access controls
- Data lineage
- Data access auditing
- · Data search and discovery
- Secure data sharing with Delta Sharing

Unity Catalog comes with a **Data Explorer UI** and a **Search engine** to help you discover your assets. It increases data discoverability and **simplifies** data sharing between your teams.

Start

Data Catalog Rubber, meet road!

- Catalogs rot without ownership; make "owner" a required field.
- Generate lineage from ETL/code:
 - but validate critical edges; 100% automatic is rare.
- Don't flood users with noisy assets
 - Categorize data assets as bronze, silver, gold;
 - Most users need "gold" datasets.
- You may need to have multiple catalogs, not ideal but it happens.

Data Quality

Garbage In, Garbage Out



The cost of bad data represents an astonishing 15% to 25% of the revenue for most companies

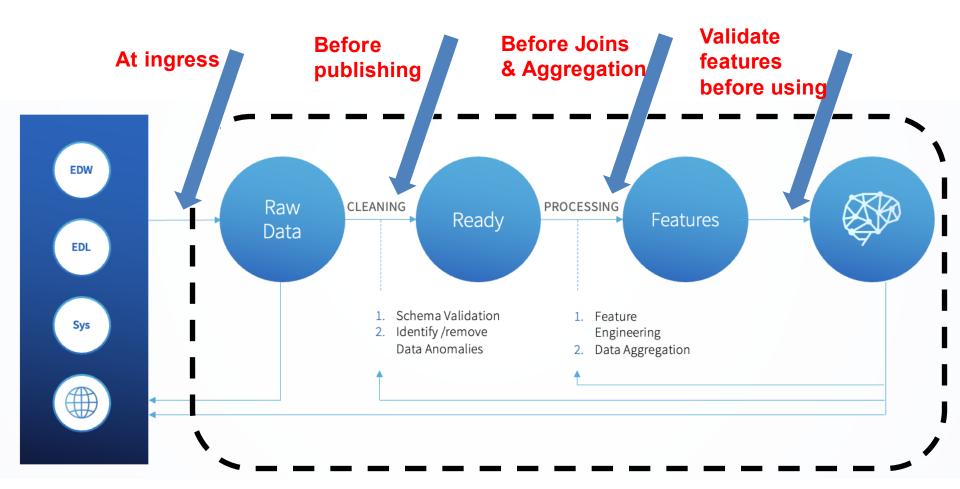
What should we check?

Short answer: everything!

Possible kind of checks:

- Schema checks: data types, not-null columns, unique keys, etc.
- Value checks: ranges, possible values, reference to other tables.
- Distribution checks: monitor drift as system progress, outliers, ratio of null values.
- Freshness: data update frequency.
- Business rules: Same level California employee makes more than Ohio.

When should we check?



Data Cleaning

Exercise:

For an events table like the assignment: (event_time, user_id, event_type, amount)

What expectations (aka quality checks) you should add?

Data Cleaning Rubber, meet road!

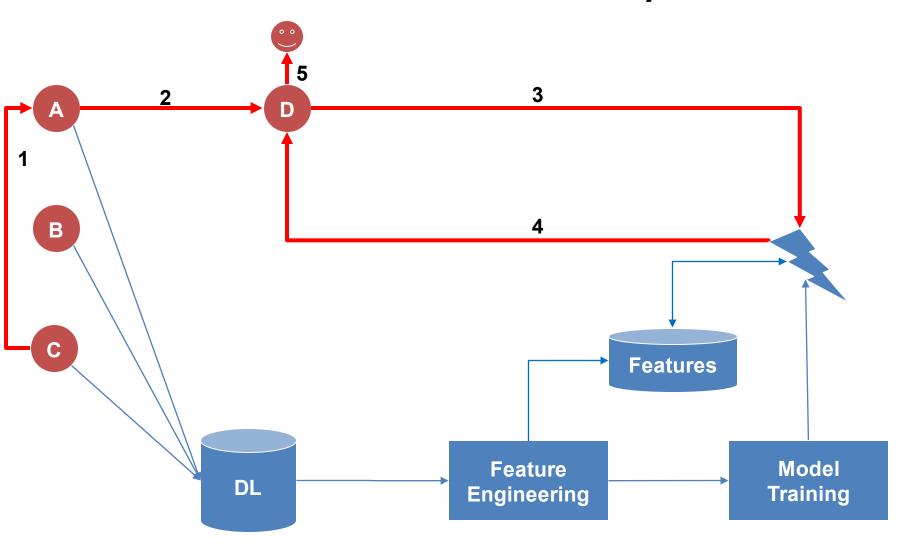
Same column may mean different things based on context.

Example:

For a company that manufacture toys, and sell them to retailers: What does revenue mean?

- Very strict rules cause many false alarms and often slow everything, or get ignored.
- Tests hidden in code are forgotten!
 - Centralize checks in configs
 - Surface checks in catalogs
 - Automate checks when possible

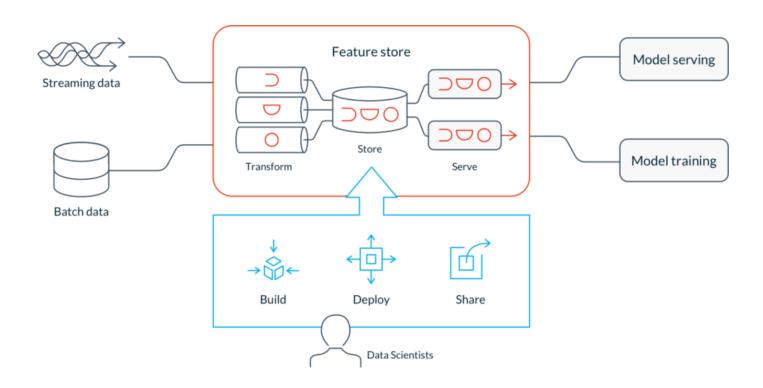
Feature Store, Why?



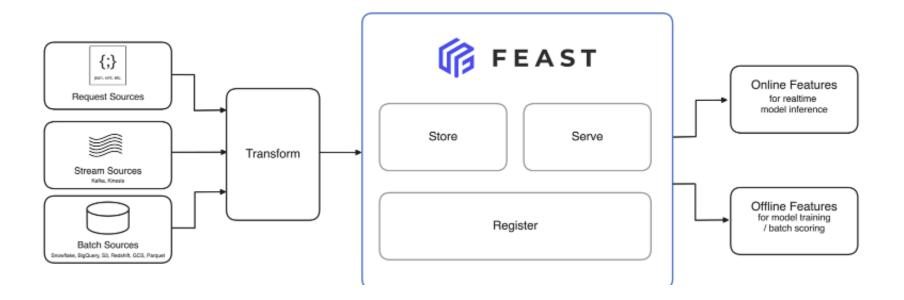
Feature Store

Features must be consistent between training and serving!

Feature Store Implementations (Tecton)



Feature Store Implementations (FEAST)



Feature Store Implementations (AWS SageMaker)



Feature Store Architecture

- Offline store: Parquet files, S3 buckets, etc.
- Online store: Redis/DynamoDB
- Feature Catalog & Lineage: obviously ©
- · Integration with model registry: part of lineage

Question: Why do we need an online store instead of using DW?

Feature Store Increase features quality

- Point-in-time correctness: prevent leakage
- Backfill when feature logic changes
- Cached features expires in online store

Feature Store Rubber, meet road

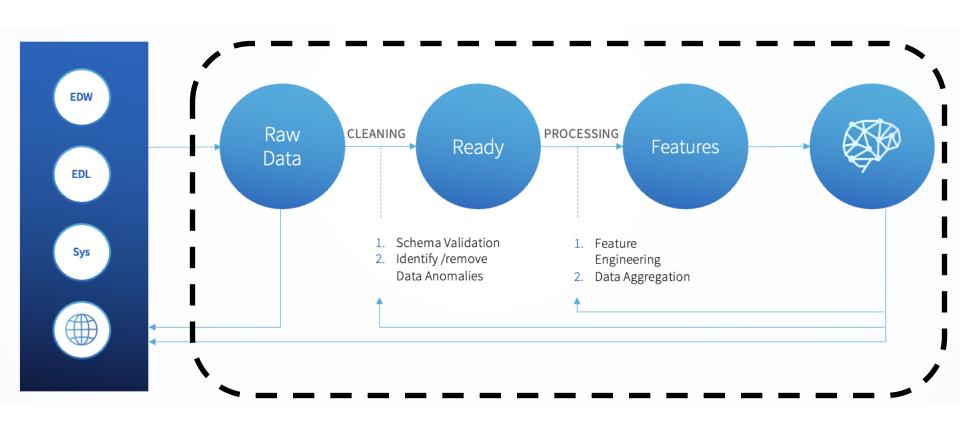
A big value for feature store, is cultural. It helps data teams share definitions and assumptions and start a dialogue about expectations.

Data Governance

- Legal vs Ethical
- Balance between "moving fast breaking things"; and legacy systems!
- You need governance to empower teams!
- PII data is sacred, in many industries.

Data Governance

safe innovation environment



Data Management

Rubber, meet road.

- Data management empower teams to build trustworthy data-enabled systems.
- Catalog & lineage: Important for data discovery & governance (owners)
- Quality: Bad data is the norm; stop it from being used in your system.
- Feature Store: A framework for consistent, low-latency ML features.
- Governance: empowers innovators to work in a safe environment.

Thank you!

