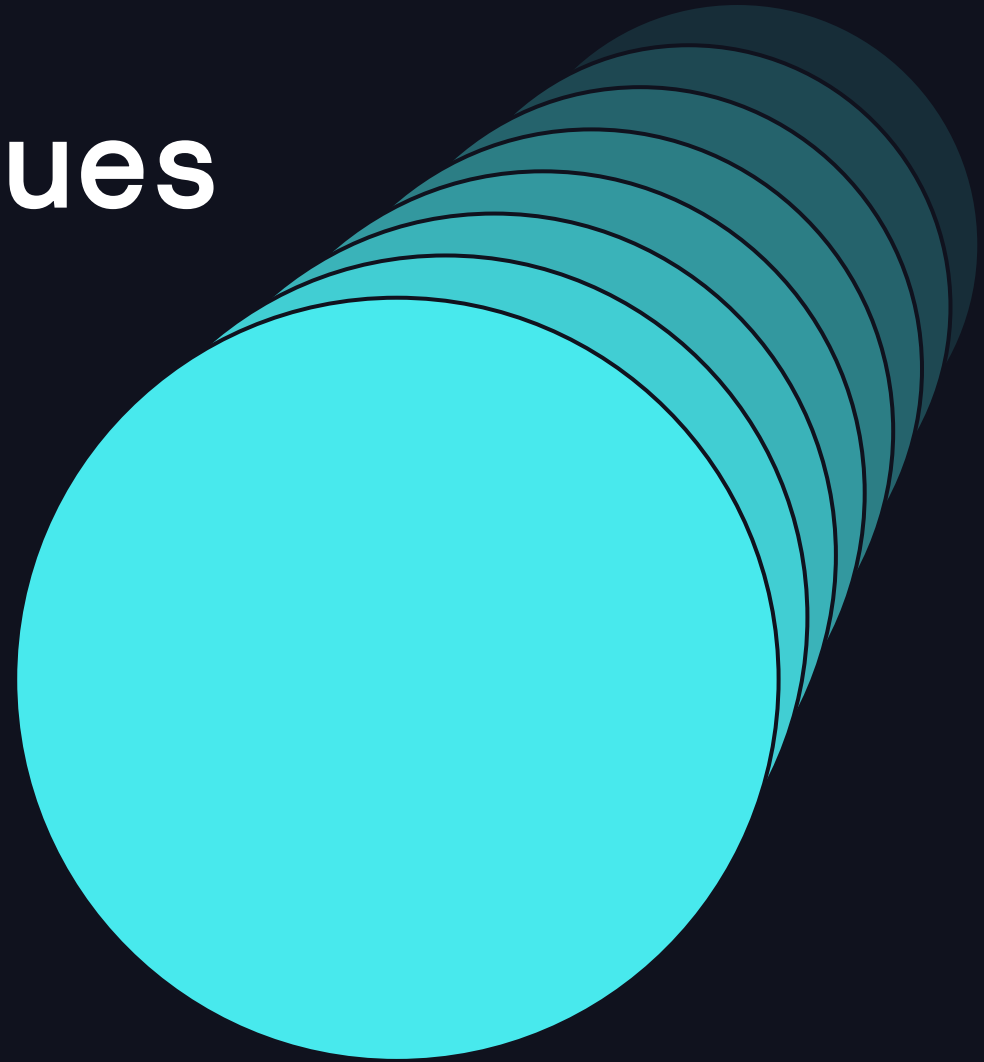# Practical techniques for applying data quality in the lakehouse

Liping Huang & Lara Rachidi
11 June 2024

# Meet the Speakers

**Liping Huang**

Senior Solutions Architect

https://www.linkedin.com/in/lipinght/

**Lara Rachidi**

Solutions Architect

https://www.linkedin.com/in/lara-rachidi/

# Agenda

❑ Six Dimensions of Data Quality

❑ Data Quality Management Lifecycle

❑ Crawl

❑ Walk

❑ Run

❑ Example Medallion Architecture

# Six dimensions model

## Dimensions of Data Quality

**Consistency**

No conflicts in data

**Accuracy**

No erros in data

**Validity**

Conform to set formats

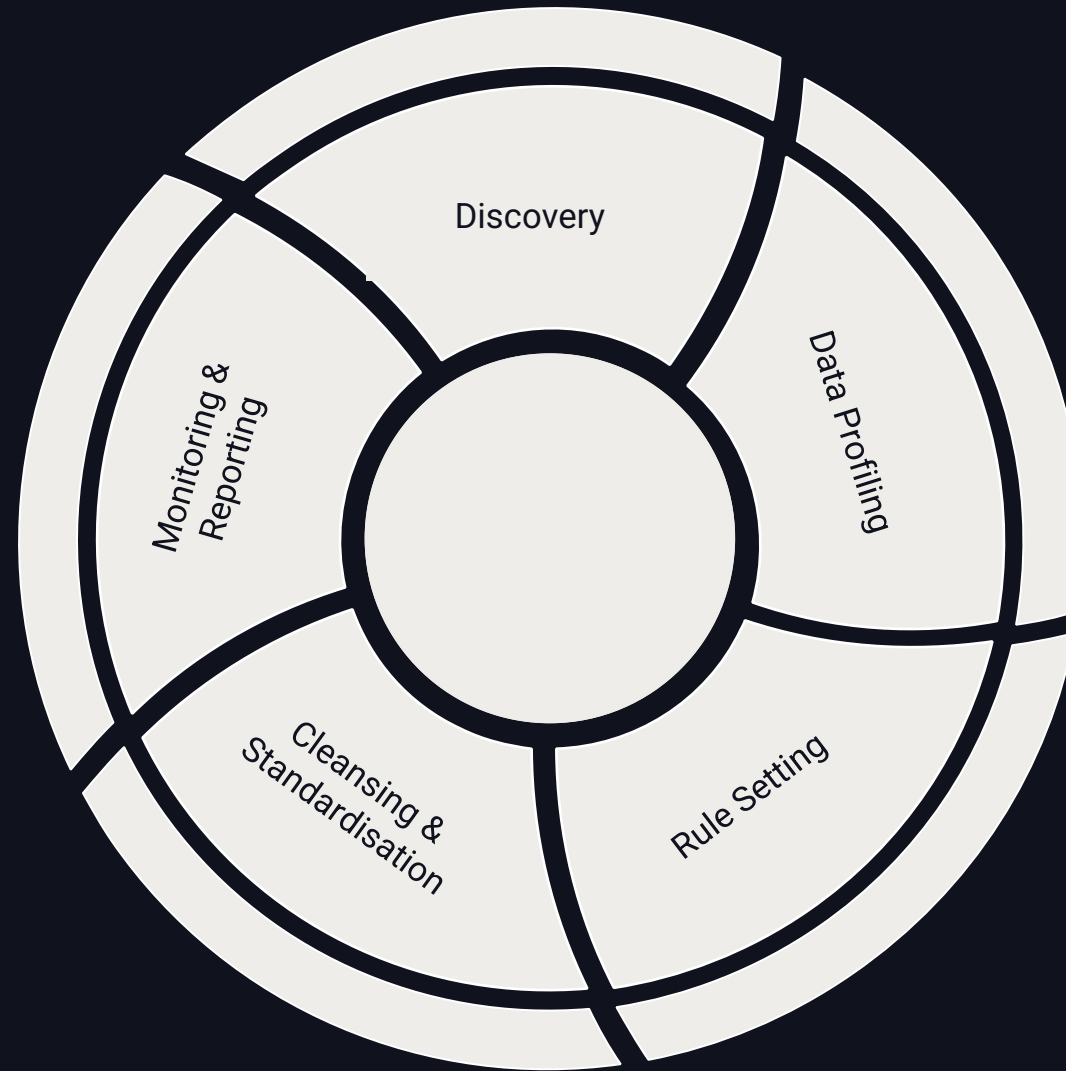**Completeness**

No missing data

**Uniqueness**

No duplicates in data
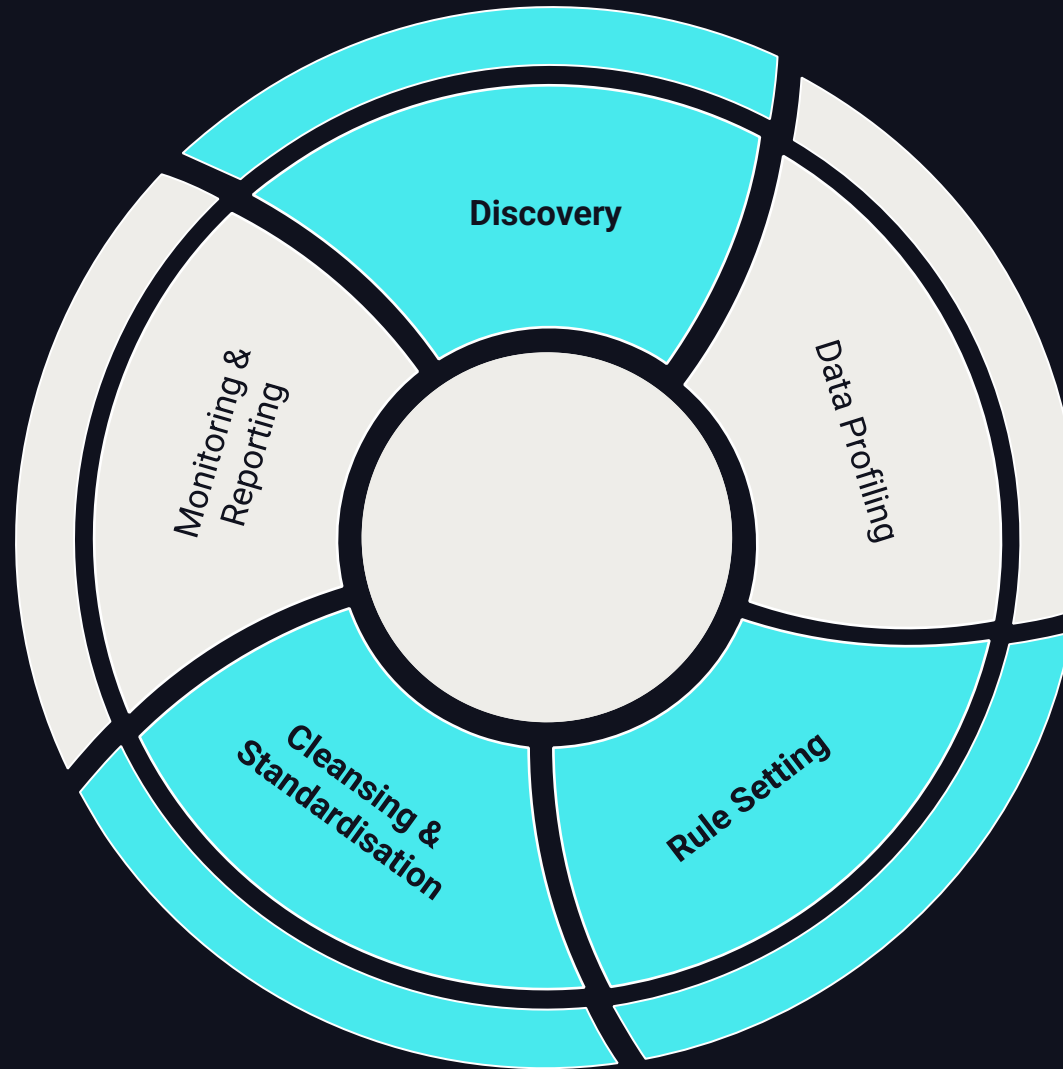
**Timeliness**

Up-to-date data

4

# Data Quality Management Lifecycle

―

# Data Quality Management Lifecycle

# Crawl

# Data Quality Management Lifecycle

# Discovery

# Discovery

## Data quality is a team sport

### Key Stakeholders

Business uers
Data scientists
Data engineers
Data analysts
Data stewards
Compliance officers
Executive stakeholders

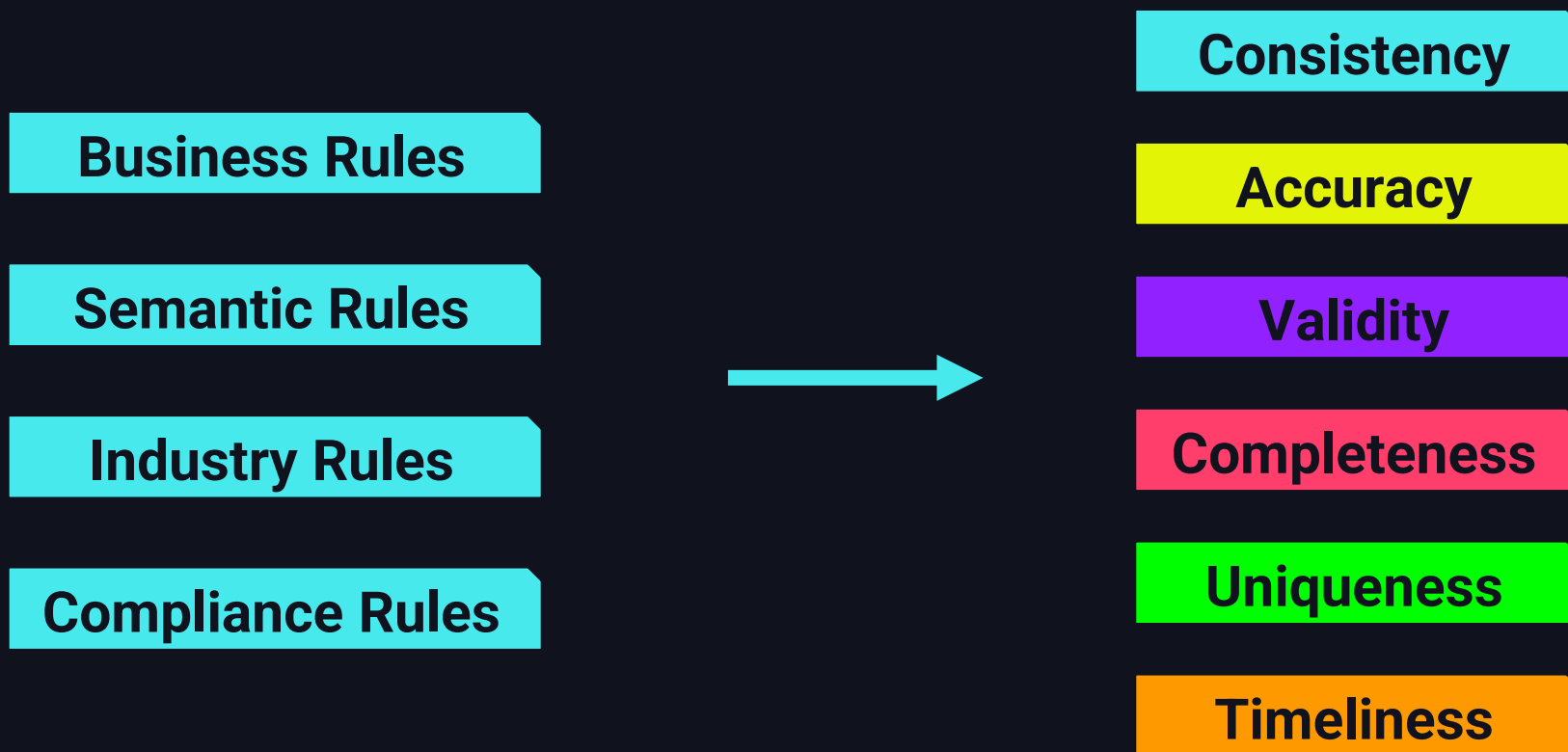### Requirement Gathering Techniques

Interviews
Surveys
Workshops
Obversations
Assessments

# Rule Setting, Cleansing & Standardization

# Data Quality Rule Setting

Different rules inform the requirements for six dimensions

Business Rules

Semantic Rules

Industry Rules

Compliance Rules

→

Consistency

Accuracy

Validity

Completeness

Uniqueness

Timeliness

# Detect inconsistencies

## Values in different datasets are not conflicting

Causes

Data sources
Transformation logic
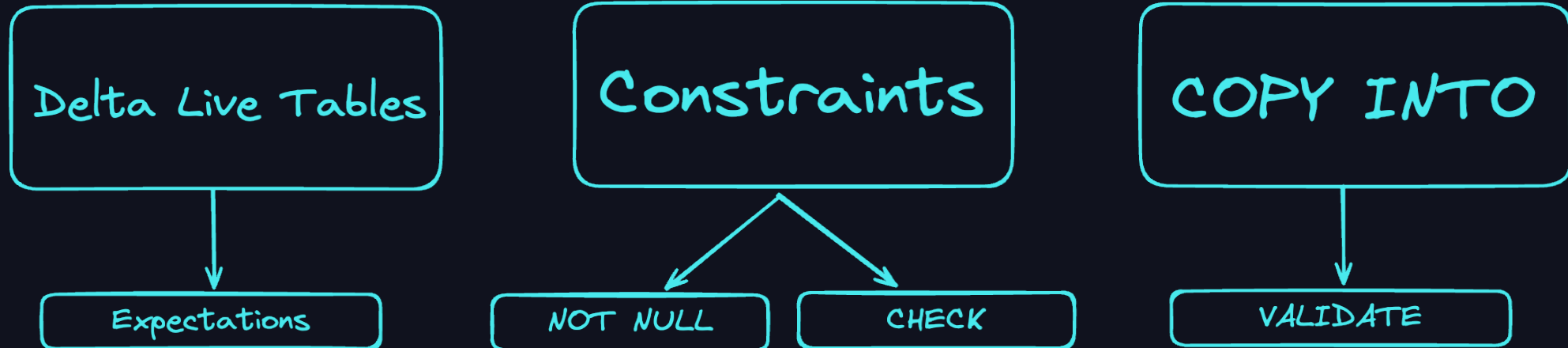Imputation techniques
Metric defitions

Techniques

MERGE
JOIN

# Detect inaccurate data
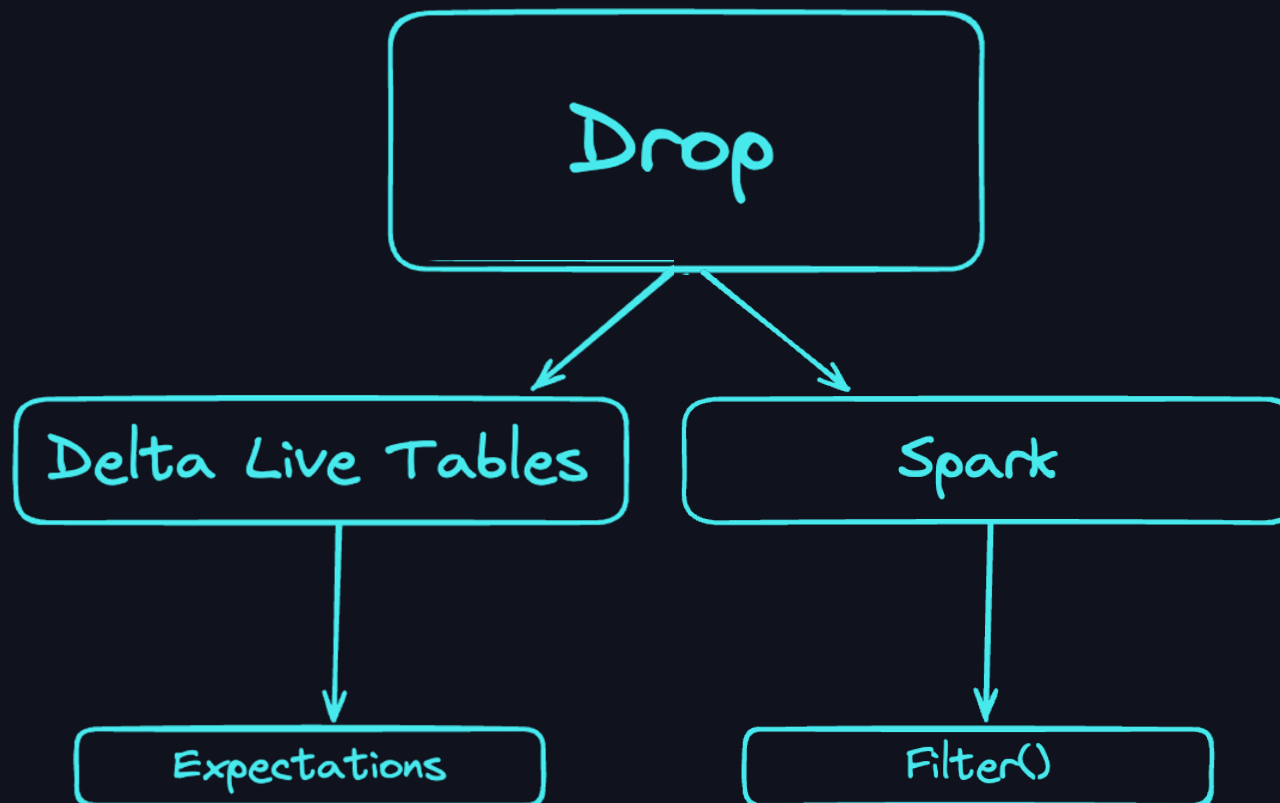
## Set rules and identify inaccurate data



Delta Live Tables → Expectations

Constraints → NOT NULL, CHECK

COPY INTO → VALIDATE

# Handle violations

Drop records failed against rule checks

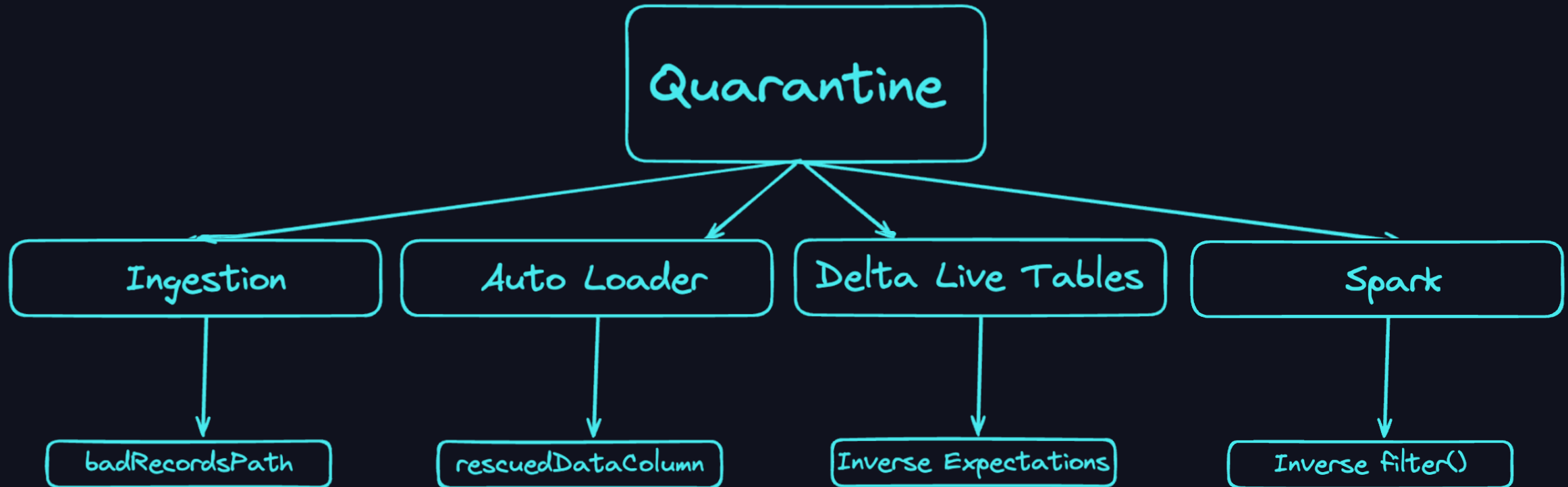# Handle violations (2)

Quarantine records failed against rule checks

# Schema Management

### Schema Enforcement

Delta
Auto Loader

### Schema Overwrite

overwriteSchema option

### Schema Evolution

mergeSchema option
AutoMerge
Auto Loader

### Schema Update

add new columns
add columns' comment
change column's ordering
replace columns
rename columns
drop columns
change columns' type

# Impute Missing Data

Impute missing records for data completeness

Min/Max

Fixed Value

K-nearest neighbors

Most frequent value

Next or previous value

Mean/Median/Moving Avg.

DATA+AI SUMMIT

# Drop Duplicates

## Remove duplicate records

Merge
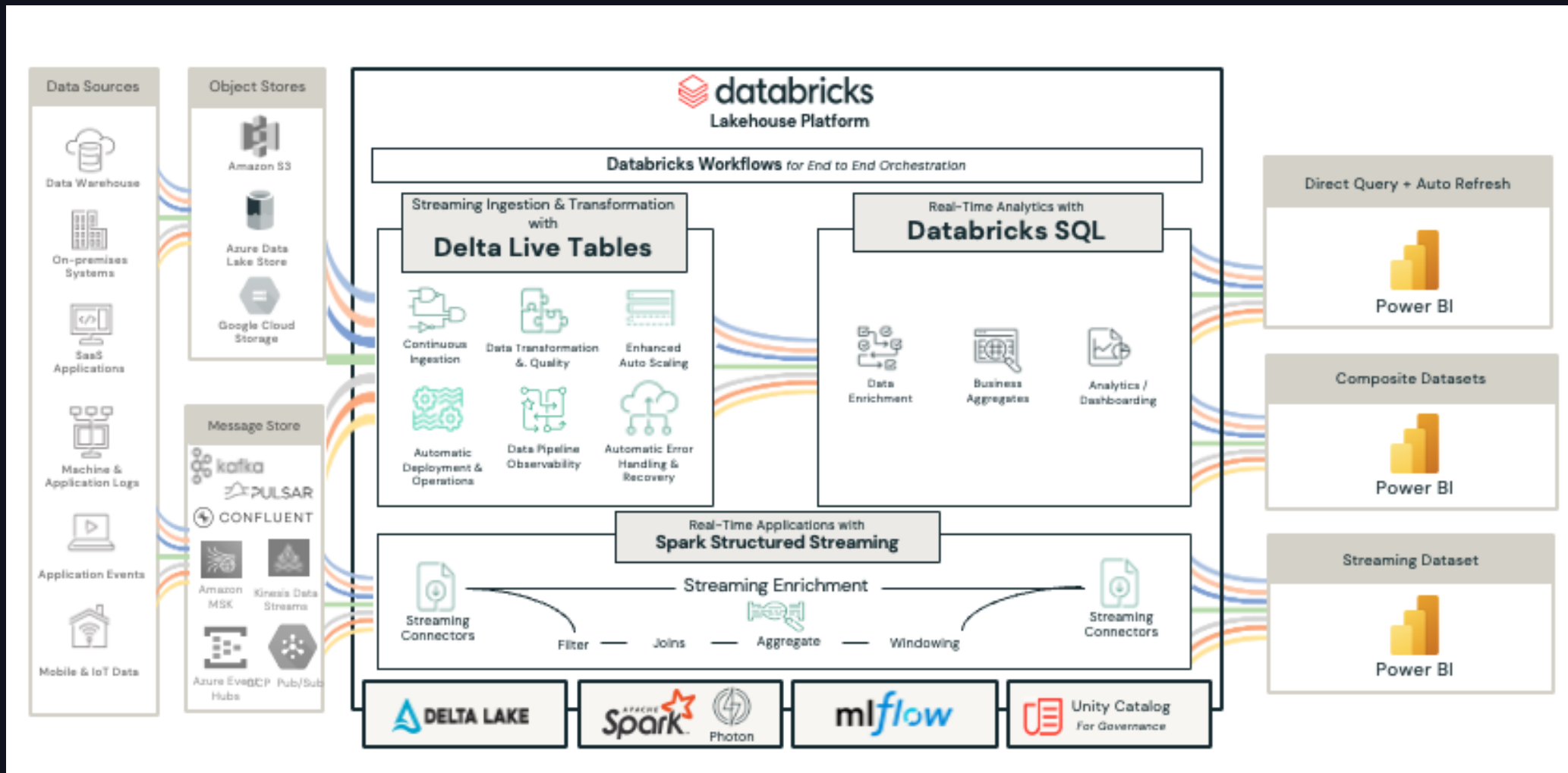
distinct()

dropDuplicates

Ranking Window

# Rollback and Vacuum

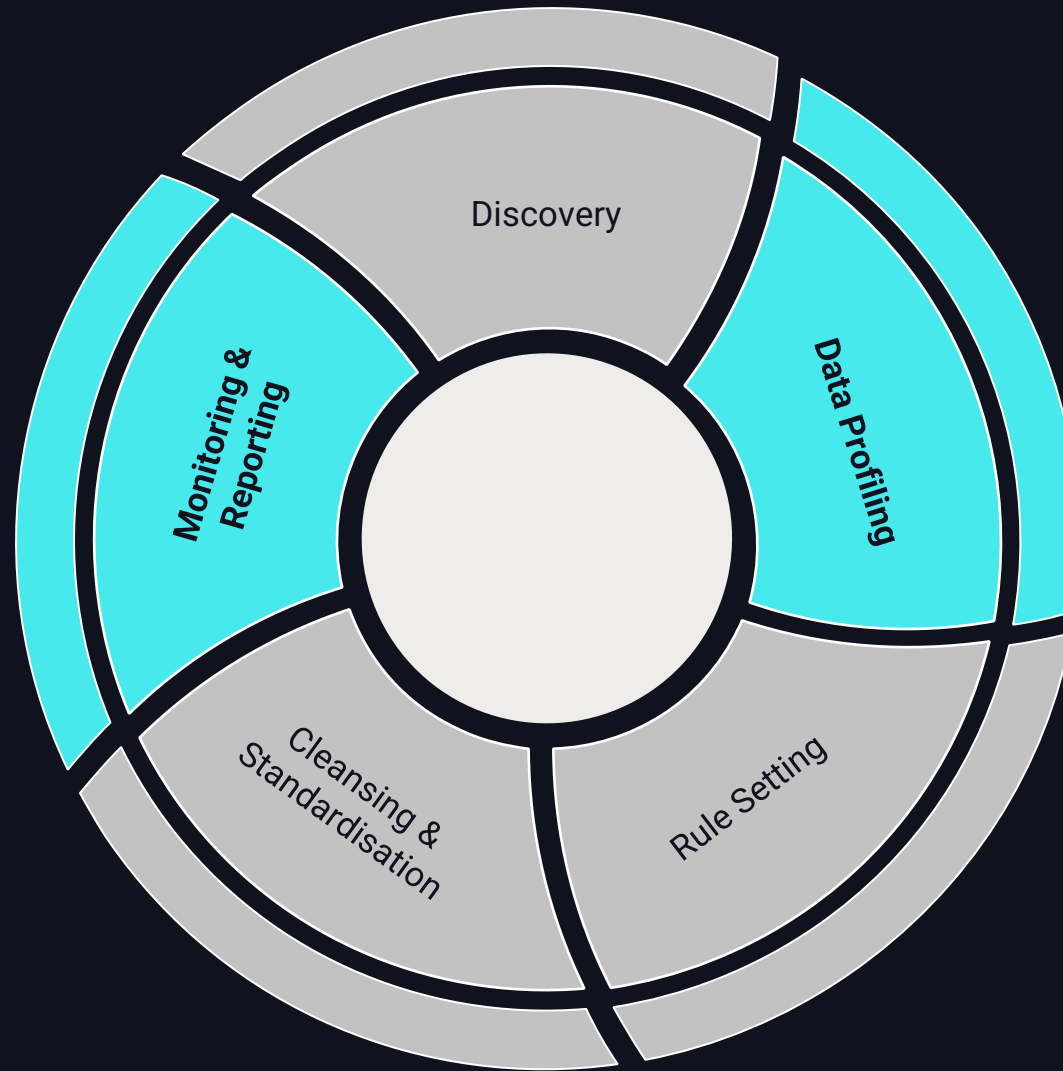Enable rollbacks but prevent access to invalid historical data

# Structured Streaming and DLT
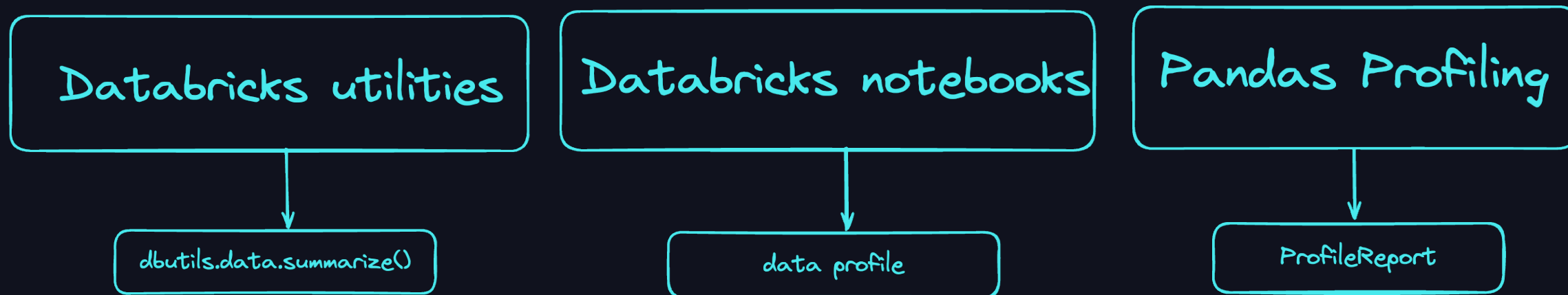
# Walk

# Data Quality Management Lifecycle

Discovery

Data Profiling

Rule Setting

Cleansing & Standardisation

Monitoring & Reporting

# Data Profiling

# Data Profiling

Identify trends and patterns in data and spot outliers and anomalies

Databricks utilities

Databricks notebooks

Pandas Profiling

dbutils.data.summarize()

data profile

ProfileReport

DATA+AI SUMMIT

# Metadata Management

# Data Augmentation

Add additional metadata to facilitate discovery

Current_timestamp()

ingestion

input_file_name()

ingestion

TBLPROPERTIES

Tables & Views

COMMENT

Tables & Columns

# AI-generated Descriptions and Comments

All Dimensions

DATA+AI SUMMIT

# Monitoring & Reporting

# Lakehouse Monitoring

Monitor data quality and model quality over time



❏ Time Series

❏ Model Inference

❏ Snapshot

❏ Optional: baseline table

# Lakehouse Monitoring

## Display automatically computed metrics in a fully customizable dashboard

# Dashboard and metrics tables

## Summary statistics and statistics related to data drift over time stored in UC

# Alerts with Lakehouse Monitoring

## Set up alerts for reporting

# Delta Live Tables Monitoring

**Data quality metrics in DLT are captured in event logs and can be reported using DBSQL**

# Root Cause Analysis leveraging Data Lineage

**All Dimensions**

**Roadmap**

# Run

DATA᛭AI SUMMIT

# Data Quality Management Lifecycle

# DataOps, MLOps, LLMOps

# Ensuring high quality with CI/CD

## Manual processes are error -prone, impacting data quality

CI: Continuous Integration

PROD

Ongoing monitoring and maintenance to ensure data quality
Data issues feed back into the DataOps cycle

Handover to Ops team and schedule deployment

STAGING

Mirrors prod environment for integration testing

CD: Continuous Deployment

Peer review (PR) + internal testing to ensure quality

DEV    DEV    DEV

Sandbox for experimenting in a controlled environment
Run unit tests

Tools for DataOps:

❑ Repos

❑ Databricks Asset Bundles

❑ Terraform

❑ Lakehouse Monitoring

# MLOps to ensure quality ML models

## Progress your code, data, models and pipelines through quality thresholds

**PROD**

Deploy the model in Prod
Real time inference with Model Serving
Monitor the model and retrain if required
Roll back to previous version if required

↑ Promote the model to the Model Registry

**STAGING**  Run tests against the prod datasets and make API level tests

↑ Peer review + internal testing to ensure quality

**DEV**  **DEV**  **DEV**  Data exploration, model training, track with MLflow

Tools:

- ☐ Repos
- ☐ Databricks Asset Bundles
- ☐ Terraform
- ☐ Lakehouse Monitoring
- ☑ MLflow

# LLMOps to ensure quality GenAI apps

## Progress your code, data, models and pipelines through quality thresholds



**PROD**

Deployment of LLM pipelines packaged using MLflow
Real time inference with Model Serving or AI Gateway
Track queries and responses using System Tables
Integrate human feedback (offline & online)
Enforce guardrails, e.g. against toxic responses

**STAGING**

Enforce guardrails, e.g. against toxic responses
Integrate human feedback (offline & online)
Offline evaluation using LLM as a judge

**DEV   DEV   DEV**

Store prompt and pipelines in a chain
MLflow: prompt iteration, pipeline packaging, tracking queries and responses
Enforce guardrails, e.g. filter toxic responses
Offline eval by integrating human feedback or using LLM as a judge
Low quality? Consider RAG or fine tuning instead of prompt engineering

Tools:

❏  MLflow

❏  LlamaIndex, LangChain

❏  Hugging Face

❏  DSPy

# Putting it all together

# Medallion Architecture

Detecting inconsistencies and inaccuracies

# Medallion Architecture

Handling Violations

DATA+AI SUMMIT

# Medallion Architecture

Schema management: enforcement, evolution, overwrite, update

# Medallion Architecture

Impute missing data

DATA+AI SUMMIT

# Medallion Architecture

Removing duplicate records

# Medallion Architecture

Enabling rollbacks but preventing access to invalid historical data

©2024 Databricks Inc. — All rights reserved

# Medallion Architecture

## Profiling data



Data Sources

**Bronze**

Databricks utilities

Databricks notebooks

Pandas profiling

**Silver**

Databricks utilities

Databricks notebooks

Pandas profiling

**Gold**

Databricks utilities

Databricks notebooks

Pandas profiling

# Medallion Architecture

Augmenting data

# Medallion Architecture

Monitoring and alerting

Data Sources

Bronze

Delta Live Tables

Lakehouse Monitoring

Silver

Delta Live Tables

Lakehouse Monitoring

Gold

Delta Live Tables

Lakehouse Monitoring

# Medallion Architecture

Conducting root cause analysis

# Medallion Architecture

Implementing DataOps, MLOps, LLMOps

©2024 Databricks Inc. — All rights reserved

# Key Takeaways

❑ Data quality is the foundation of everything

❑ Get the basics right: <u>Data Quality Management With Databricks</u>

❑ Start with requirements, work towards automation, iterate!

# Data Quality at Data+AI Summit

- **Sponsored by: IBM | Building Better Data Quality with IBM Data Fabric and Databricks**: Tuesday, Jun 11, 3:40 PM - 4:00 PM PDT
- **Sponsored by: Datafold | Shifting Data Quality to the Left: Automating Data Testing on Databricks**: Tuesday, Jun 11, 4:00 PM - 4:40 PM PDT
- **LLM Evaluation: Auditing Fine-Tuned LLMs for Guaranteed Output Quality**: Wednesday, Jun 12, 1:40 PM - 2:20 PM PDT
- **Building High-Quality and Trusted Data Products with Databricks**: Wednesday, Jun 12, 12:30 PM - 1:10 PM PDT
- **Lakehouse Monitoring GA: Profiling, Diagnosing, and Enforcing Data Quality with Intelligence**: Thursday, Jun 13, 11:20 AM - 12:00 PM PDT