# Modernizing Banking ETL and Analytics with Delta Live Tables

Jian Zhou, Navy Federal Credit Union
Ricardo Portilla, Databricks

**NAVY FEDERAL**
**Credit Union**®

# Introduction

**Ricardo Portilla**

- Industry Principal Architect, 6 years Databricks
- Previously led Market Surveillance at FINRA
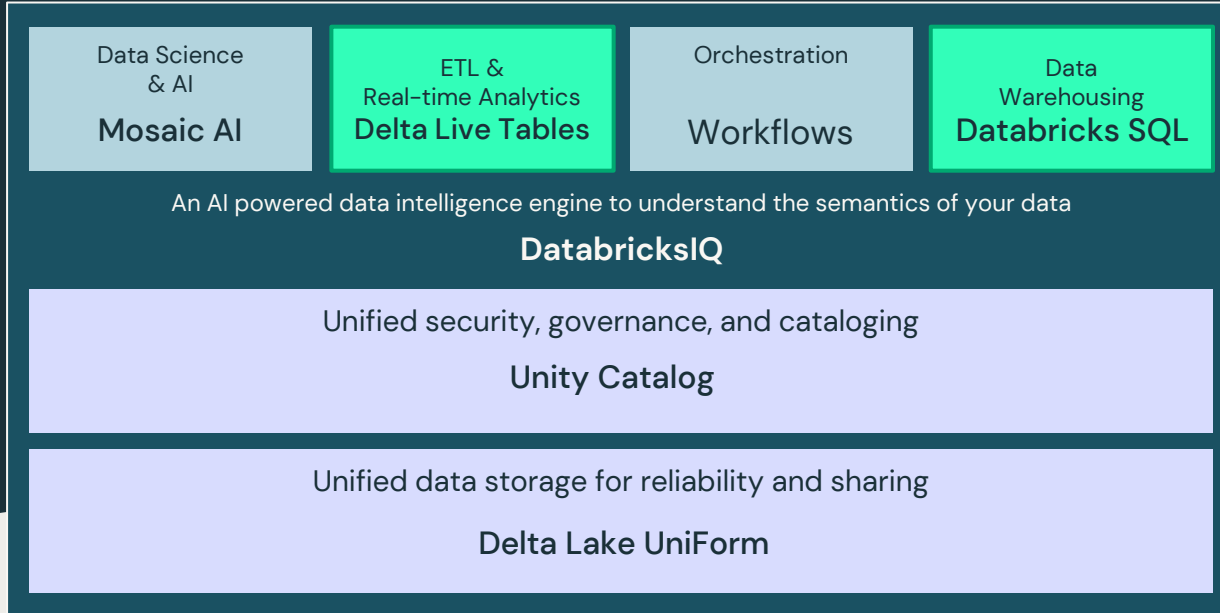- [Open source creator](#) and industry SME

2005 — amazon

2007 — NETFLIX

2010 — Uber / lyft

# Databricks Data Intelligence Platform

| Data Science & AI | ETL & Real–time Analytics | Orchestration | Data Warehousing |
|---|---|---|---|
| **Mosaic AI** | **Delta Live Tables** | **Workflows** | **Databricks SQL** |

An AI powered data intelligence engine to understand the semantics of your data

**DatabricksIQ**

Unified security, governance, and cataloging

**Unity Catalog**

Unified data storage for reliability and sharing
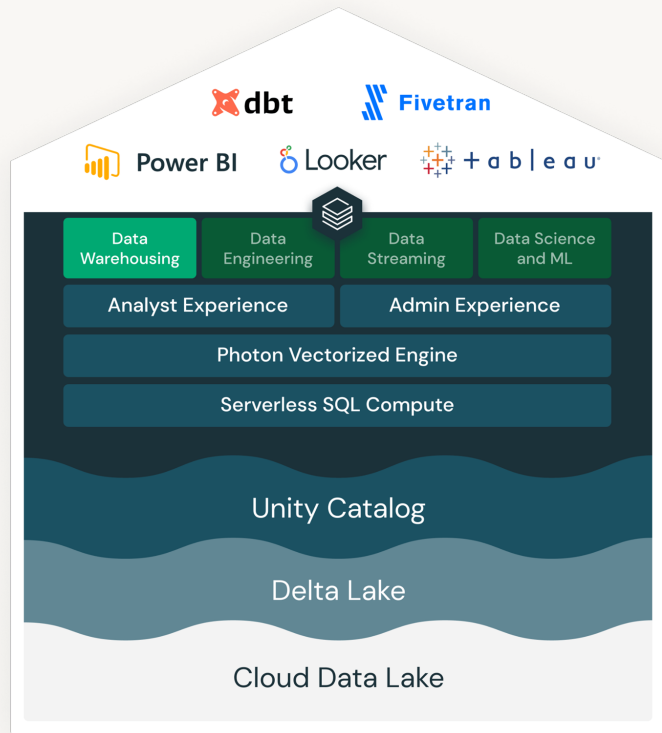
**Delta Lake UniForm**

**Open Data Lake**

All Raw Data
(Logs, Texts, Audio, Video, Images)

# The best data warehouse is a lakehouse

## Powered by Databricks SQL



Seamless Integration with the Ecosystem

Ease of Use

Real-world Performance

Centralized Governance

Open and Reliable Data Lake as the Foundation

# Introduction

**Jian (Miracle) Zhou**

- Sr. Manager, Digital Data Engineering

- Navy Federal Credit Union

- https://www.linkedin.com/in/miraclezhou/

DATA∣AI SUMMIT

June 2023

# MISSION:

# NEAR REAL-TIME INSIGHTS IN 6 WEEKS

# Our dataflow pattern



Serving & visualization

Data Source

Data Pipeline

ADLS Gen2

Databricks SQL

Power BI Dashboards

Ingestion & curation

Storage

Compute

Visualization

DATA AI SUMMIT

# The data source

## Online Banking Application Telemetry



Login / logout
Check Balances
Pay Bills
Make Transfer
Direct Deposit
Card Application

Application Insights

Since the initial rollout, Online v7 has logged over 8.5 billion events in Application Insights (as of April 2024)

# Turn on the telemetry firehose

DATA+AI SUMMIT

# NRT Pipeline



Serving & visualization

Application Insights → Azure Event Hubs (Kafka) → Near Real-time Data Pipeline → ADLS Gen2 (Storage) → Databricks SQL (Compute) → Power BI Dashboards (Visualization)

Ingestion & curation

DATA'AI SUMMIT

# The input

- Semi-structured JSON documents

- App events embedded in an array

  - 140+ event types

  - 100+ event attributes

- There are duplicates!

```
{
  "records": [
    {
      "time": "2024-04-08T14:00:03.3500000Z",
      "OperationId": "1f85fe13fa9cc460463b008bf94ca156",
      "SessionId": "40cf0031-a1db-426e-850b-1940a290b6a4",
      "Properties": {
        "activity": "login",
        "channel": "web",
        "ipAddress": "20.84.33.69"
      }
    },
    {
      "time": "2024-04-08T14:02:05.3110000Z",
      "OperationId": "88a4321463a92645bb5224238d5b7967",
      "SessionId": "40cf0031-a1db-426e-850b-1940a290b6a4",
      "Properties": {
        "activity": "ingestion",
        "channel": "web",
        "ipAddress": "20.84.33.69",
        "paymentOrderStatus":"SUCCESS"
      }
    }
  ]
}
```

App event 1

App event 2

A sample message from Event Hubs

DATA'AI SUMMIT

# The output



```
1  select *
2  from appevents_login
3  where channel is not null
4  order by event_date desc
5  limit 10;
```

dev data sample

Results

| # | t_time | | ls | gui | |
|---|---|---|---|---|---|
| | | 202...24 | | f39 | 3b... |
| | 2023-10-24 | | 6b125cdf | d06 | 8fe |
| 3 | 2023-10-24 | 2023-10-24 04:21:35.330 | web | 25198494-626f-4d20-a02a-df4b70d41689 | a4e699a6-672c-4c5b-9fe6-7b68692d... |
| 4 | 2023-10-24 | 2023-10-24 12:23:14.715 | mobile | 8a791733-9125-47b7-8469-4611b9db6627 | d729140c-2813-4f44-8e97-b6fb2bd1... |
| 5 | 2023-10-24 | 2023-10-24 12:25:08.771 | web | 62b346c3-4763-435f-a38f-d8a916fe0739 | 1d69c2ec-52ea-4996-af65-c5a73c799 |
| 6 | 2023-10-24 | 2023-10-24 13:24:56.989 | web | f5322ccc-688f-4d71-82ab-48ca3d84d7f5 | efaa771c-ec38-4e9b-8d05-949867b9... |
| 7 | 2023-10-24 | 2023-10-24 14:01:44.380 | mobile | 577a90f4-181f-402f-9a05-102445258fd3 | d5123f7a-4201-4093-ae86-09ddf013... |
| 8 | 2023-10-24 | 2023-10-24 12:41:03.562 | web | 949bbd15-dacd-4191-80b9-f134190df2fd | c1f6d04f-dfad-4671-ac4d-2af914c27c... |

**Persist** → **Flatten** → **De-dupe** → **Tailor**

DATA AI SUMMIT

# The pipeline

DATA'AI SUMMIT

# The iceberg of data engineering



20%
Functional

80%
Non-functional

Transformation

Data Quality

Scalability

Reliability

Observability

Error Handling

Data Persistence

Optimization

Monitoring

State Management

DATA AI SUMMIT

# Simplicity

DATA·AI SUMMIT

# Simplicity



| PERSIST | FLATTEN | DE-DUPE | TAILOR |
|---------|---------|---------|--------|
| Streaming table | Streaming table | Streaming table | View / Streaming table |

**PERSIST** — Streaming table
appevents_raw
Running · 190h 46m 51s
● 16M  ● 0

**FLATTEN** — Streaming table
appevents_flattened
Running · 190h 46m 51s
● 348M  ● 0

**DE-DUPE** — Streaming table
appevents_cleansed
Running · 190h 46m 48s
● 321M  ● 0  ● 0

**TAILOR** — View
appevents_login_t...

Streaming table
appevents_login
Running · 190h 46m 47s
● 41M  ● 0  ● 0

```python
@dlt.table()
@dlt.expect_or_drop(
  "valid records",
  "...")
def appevents_raw():
    options = get_kafka_config_options(
      spark,
      dbutils.secrets)
    return (spark
            .readStream
            .format("kafka")
            .options(**options)
            .load()
            .transform(parse_raw_appevents)
    )
```

# Simplicity



```python
@dlt.table()
@dlt.expect_or_drop(
  "valid records",
  "...")
def appevents_raw():
    options = get_kafka_config_options(
      spark,
      dbutils.secrets)
    return (spark
            .readStream
            .format("kafka")
            .options(**options)
            .load()
            .transform(parse_raw_appevents)
    )
```

```python
@dlt.table()
@dlt.expect_or_drop(
  "valid records",
  "...")
def appevents_flattened():
    return flatten_appevents(
      dlt.readStream("appevents_raw"))
```

# Simplicity

| PERSIST | FLATTEN | DE-DUPE | TAILOR |
|---------|---------|---------|--------|

Streaming table

appevents_raw ↻
Running · 190h 46m 51s
● 16M  ● 0

Streaming table

appevents_flattened ↻
Running · 190h 46m 51s
● 348M  ● 0

Streaming table

appevents_cleansed ↻
Running · 190h 46m 48s
● 321M  ● 0  ● 0

View

appevents_login_t...

Streaming table

appevents_login ↻
Running · 190h 46m 47s
● 41M  ● 0  ● 0

```python
@dlt.table()
@dlt.expect_or_drop(
  "valid records",
  "...")
def appevents_raw():
    options = get_kafka_config_options(
      spark,
      dbutils.secrets)
    return (spark
            .readStream
            .format("kafka")
            .options(**options)
            .load()
            .transform(parse_raw_appevents)
    )
```
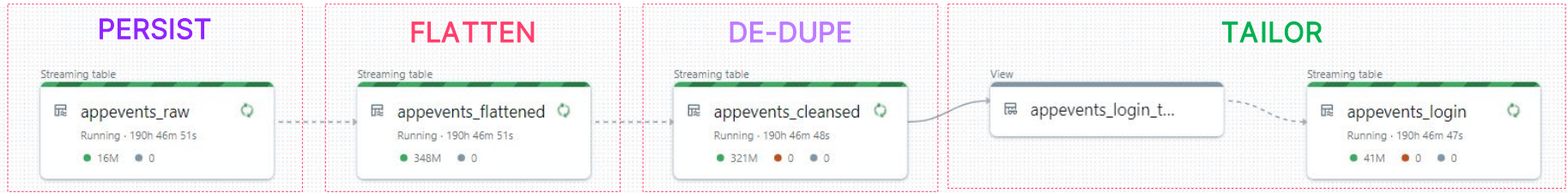
```python
@dlt.table()
@dlt.expect_or_drop(
  "valid records",
  "...")
def appevents_flattened():
    return flatten_appevents(
      dlt.readStream("appevents_raw"))
```

```python
dlt.create_streaming_table(
    "appevents_cleansed"
  )

dlt.apply_changes(
  target = "appevents_cleansed",
  source = "appevents_flattened",
  keys = ["event_date", ...],
  sequence_by = col("time"),
  except_column_list = [],
  stored_as_scd_type = 1
)
```

DATA·AI SUMMIT

# Simplicity



| PERSIST | FLATTEN | DE-DUPE | TAILOR |
|---|---|---|---|

**PERSIST** — Streaming table — appevents_raw — Running · 190h 46m 51s — 16M · 0

**FLATTEN** — Streaming table — appevents_flattened — Running · 190h 46m 51s — 348M · 0

**DE-DUPE** — Streaming table — appevents_cleansed — Running · 190h 46m 48s — 321M · 0 · 0

**TAILOR** — View — appevents_login_t... — Streaming table — appevents_login — Running · 190h 46m 47s — 41M · 0 · 0

```python
@dlt.table()
@dlt.expect_or_drop(
  "valid records",
  "...")
def appevents_raw():
    options = get_kafka_config_options(
        spark,
        dbutils.secrets)
    return (spark
            .readStream
            .format("kafka")
            .options(**options)
            .load()
            .transform(parse_raw_appevents)
    )
```

```python
@dlt.table()
@dlt.expect_or_drop(
  "valid records",
  "...")
def appevents_flattened():
    return flatten_appevents(
      dlt.readStream("appevents_raw"))
```

```python
dlt.create_streaming_table(
    "appevents_cleansed"
  )

dlt.apply_changes(
  target = "appevents_cleansed",
  source = "appevents_flattened",
  keys = ["event_date", ...],
  sequence_by = col("time"),
  except_column_list = [],
  stored_as_scd_type = 1
)
```

```python
@dlt.view
def appevents_login_temp():
    df = dlt.readStream(
      "appevents_cleansed")
    return filter_login_events(df)
```

# Simplicity



PERSIST | FLATTEN | DE-DUPE | TAILOR

Streaming table — appevents_raw — Running · 190h 46m 51s — 16M · 0

Streaming table — appevents_flattened — Running · 190h 46m 51s — 348M · 0

Streaming table — appevents_cleansed — Running · 190h 46m 48s — 321M · 0 · 0

View — appevents_login_t...

Streaming table — appevents_login — Running · 190h 46m 47s — 41M · 0 · 0

```python
@dlt.table()
@dlt.expect_or_drop(
  "valid records",
  "...")
def appevents_raw():
    options = get_kafka_config_options(
        spark,
        dbutils.secrets)
    return (spark
            .readStream
            .format("kafka")
            .options(**options)
            .load()
            .transform(parse_raw_appevents)
    )
```

```python
@dlt.table()
@dlt.expect_or_drop(
  "valid records",
  "...")
def appevents_flattened():
    return flatten_appevents(
      dlt.readStream("appevents_raw"))
```

```python
dlt.create_streaming_table(
    "appevents_cleansed"
)

dlt.apply_changes(
  target = "appevents_cleansed",
  source = "appevents_flattened",
  keys = ["event_date", ...],
  sequence_by = col("time"),
  except_column_list = [],
  stored_as_scd_type = 1
)
```

```python
@dlt.view
def appevents_login_temp():
    df = dlt.readStream(
        "appevents_cleansed")
    return filter_login_events(df)
```

```python
dlt.create_streaming_table(
    "appevents_login"
)

dlt.apply_changes(
  target = "appevents_login",
  source = "appevents_login_temp",
  keys = ["event_date", "event_id"],
  sequence_by = col("event_time"),
  except_column_list = [],
  stored_as_scd_type = 1
)
```

©2024 Databricks Inc. — All rights reserved

# Simplicity

| PERSIST | FLATTEN | DE-DUPE | TAILOR |
|---------|---------|---------|--------|

Streaming table
**appevents_raw**
Running · 190h 46m 51s
● 16M ● 0

Streaming table
**appevents_flattened**
Running · 190h 46m 51s
● 348M ● 0

Streaming table
**appevents_cleansed**
Running · 190h 46m 48s
● 321M ● 0 ● 0

View
**appevents_login_t...**

Streaming table
**appevents_login**
Running · 190h 46m 47s
● 41M ● 0 ● 0

```python
@dlt.table()
@dlt.expect_or_drop(
  "valid records",
  "...")
def appevents_raw():
    options = get_kafka_config_options(
      spark,
      dbutils.secrets)
    return (spark
            .readStream
            .format("kafka")
            .options(**options)
            .load()
            .transform(parse_raw_appevents)
    )
```

```python
@dlt.table()
@dlt.expect_or_drop(
  "valid records",
  "...")
def appevents_flattened():
    return flatten_appevents(
      dlt.readStream("appevents_raw"))
```

```python
dlt.create_streaming_table(
    "appevents_cleansed"
  )

dlt.apply_changes(
  target = "appevents_cleansed",
  source = "appevents_flattened",
  keys = ["event_date", ...],
  sequence_by = col("time"),
  except_column_list = [],
  stored_as_scd_type = 1
)
```

```python
@dlt.view
def appevents_login_temp():
    df = dlt.readStream(
      "appevents_cleansed")
    return filter_login_events(df)
```

```python
dlt.create_streaming_table(
    "appevents_login"
)

dlt.apply_changes(
  target = "appevents_login",
  source = "appevents_login_temp",
  keys = ["event_date", "event_id"],
  sequence_by = col("event_time"),
  except_column_list = [],
  stored_as_scd_type = 1
)
```

# Simplicity



| PERSIST | FLATTEN | DE-DUPE | TAILOR |
|---|---|---|---|

Streaming table

appevents_raw
Running · 190h 46m 51s
● 16M   ● 0

Streaming table

appevents_flattened
Running · 190h 46m 51s
● 348M   ● 0

Streaming table

appevents_cleansed
Running · 190h 46m 48s
● 321M   ● 0   ● 0

View

appevents_login_t...

Streaming table

appevents_login
Running · 190h 46m 47s
● 41M   ● 0   ● 0

```python
@dlt.table()
@dlt.expect_or_drop(
  "valid records",
  "...")
def appevents_raw():
    options = get_kafka_config_options(
      spark,
      dbutils.secrets)
    return (spark
            .readStream
            .format("kafka")
            .options(**options)
            .load()
            .transform(parse_raw_appevents)
    )
```

```python
@dlt.table()
@dlt.expect_or_drop(
  "valid records",
  "...")
def appevents_flattened():
    return flatten_appevents(
      dlt.readStream("appevents_raw"))
```

```python
dlt.create_streaming_table(
    "appevents_cleansed"
)

dlt.apply_changes(
  target = "appevents_cleansed",
  source = "appevents_flattened",
  keys = ["event_date", ...],
  sequence_by = col("time"),
  except_column_list = [],
  stored_as_scd_type = 1
)
```
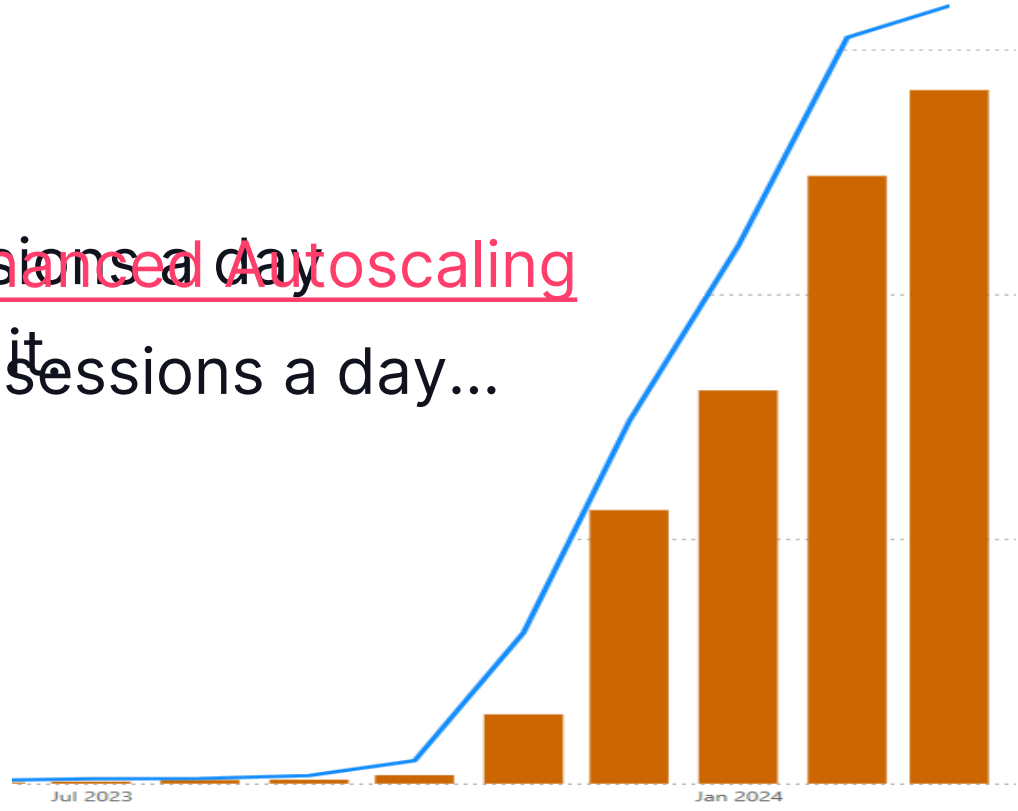
```python
@dlt.view
def appevents_login_temp():
    df = dlt.readStream(
      "appevents_cleansed")
    return filter_login_events(df)
```

```python
dlt.create_streaming_table(
    "appevents_login"
)

dlt.apply_changes(
  target = "appevents_login",
  source = "appevents_login_temp",
  keys = ["event_date", "event_id"],
  sequence_by = col("event_time"),
  except_column_list = [],
  stored_as_scd_type = 1
)
```

# Simplicity



| PERSIST | FLATTEN | DE-DUPE | TAILOR |
|---------|---------|---------|--------|

```
Streaming table
appevents_raw
Running · 190h 46m 51s
● 16M   ○ 0
```

```
Streaming table
appevents_flattened
Running · 190h 46m 51s
● 348M   ○ 0
```

```
Streaming table
appevents_cleansed
Running · 190h 46m 48s
● 321M   ● 0   ○ 0
```

```
View
appevents_login_t...
```

```
Streaming table
appevents_login
Running · 190h 46m 47s
● 41M   ● 0   ○ 0
```

```python
@dlt.table()
@dlt.expect_or_drop(
  "valid records",
  "...")
def appevents_raw():
    options = get_kafka_config_options(
      spark,
      dbutils.secrets)
    return (spark
          .readStream
          .format("kafka")
          .options(**options)
          .load()
          .transform(parse_raw_appevents)
    )
```

```python
@dlt.table()
@dlt.expect_or_drop(
  "valid records",
  "...")
def appevents_flattened():
    return flatten_appevents(
      dlt.readStream("appevents_raw"))
```

```python
dlt.create_streaming_table(
    "appevents_cleansed"
)

dlt.apply_changes(
  target = "appevents_cleansed",
  source = "appevents_flattened",
  keys = ["event_date", ...],
  sequence_by = col("time"),
  except_column_list = [],
  stored_as_scd_type = 1
)
```

```python
@dlt.view
def appevents_login_temp():
    df = dlt.readStream(
      "appevents_cleansed")
    return filter_login_events(df)
```

```python
dlt.create_streaming_table(
    "appevents_login"
)

dlt.apply_changes(
  target = "appevents_login",
  source = "appevents_login_temp",
  keys = ["event_date", "event_id"],
  sequence_by = col("event_time"),
  except_column_list = [],
  stored_as_scd_type = 1
)
```

3 queries. 2 CDC calls. ~100% for functional requirement.

# Scalability

We turned on Enhanced Autoscaling and forgot about it

From 1,000 sessions a day
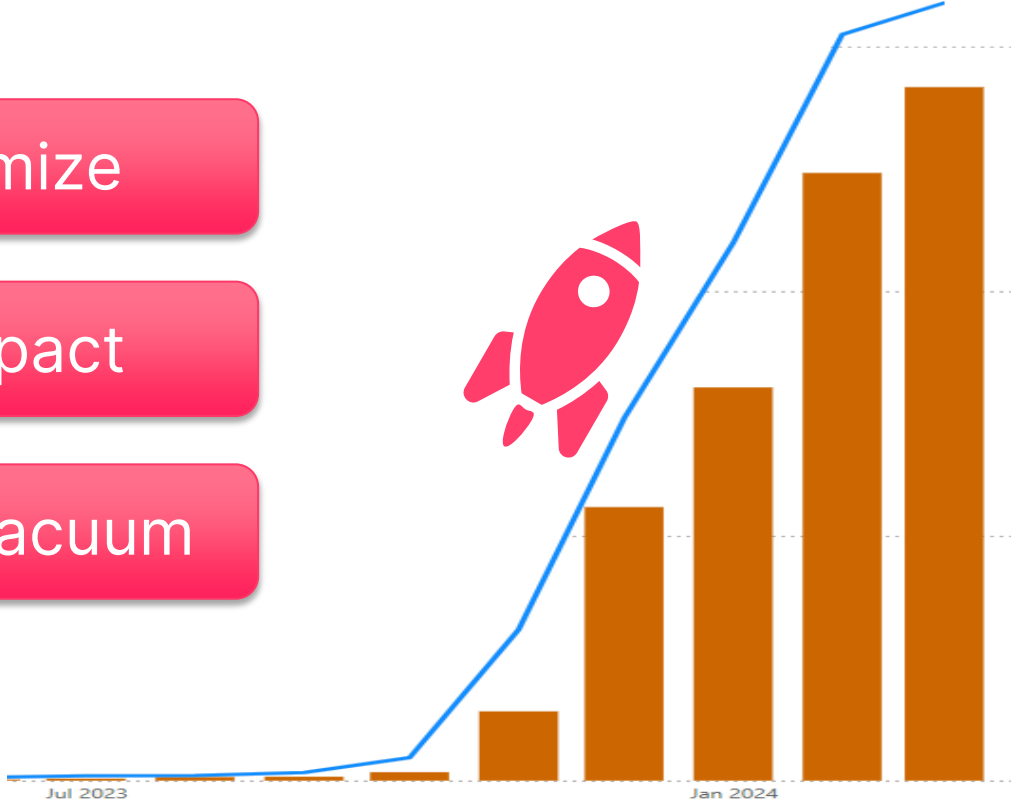to >1,000,000 of sessions a day...

# Reliability

**24*7**
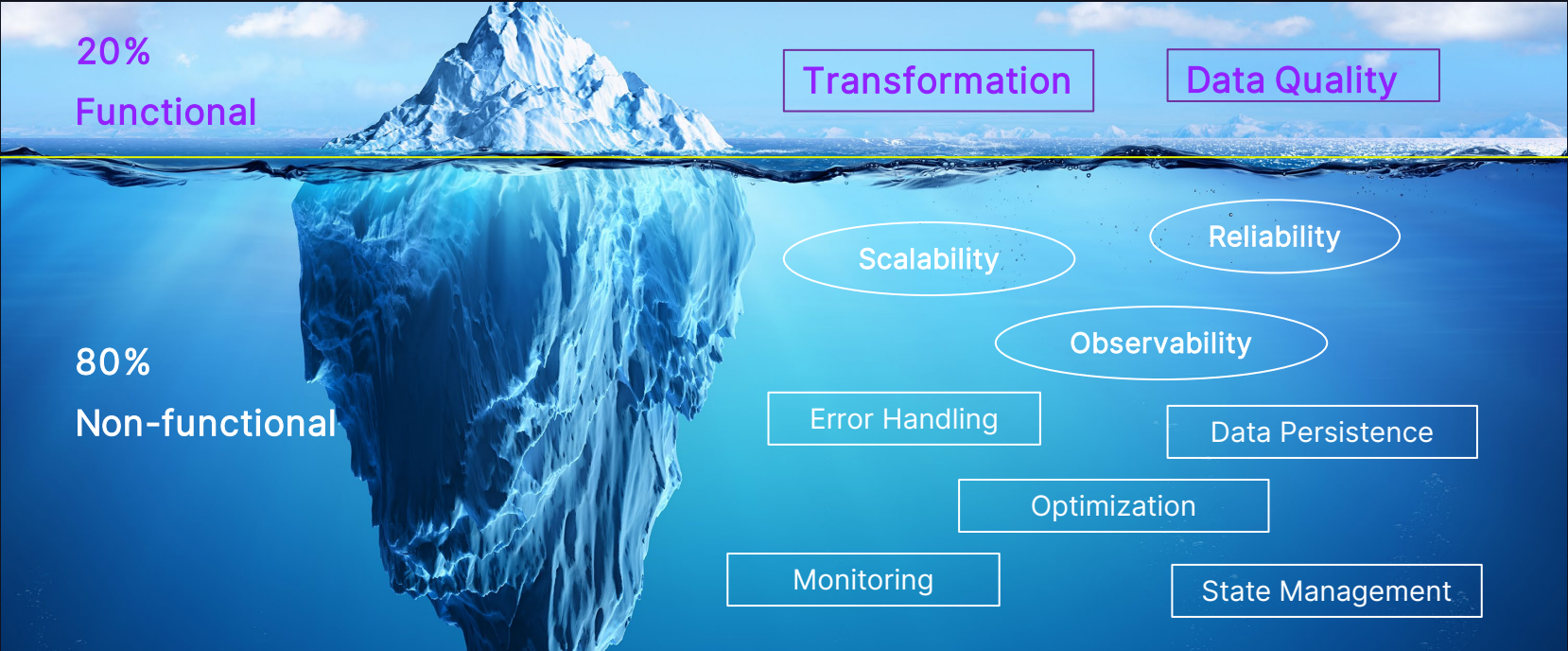continuous streaming for 9 month

**9** billion
events processed

**100**%
self-recovery from cloud failures

DATA AI SUMMIT

# Performance



Auto-Optimize

Auto-Compact

Scheduled Vacuum

Jul 2023

Jan 2024

# Where DLT shines



20% Functional

Transformation     Data Quality

80% Non-functional

Scalability

Reliability

Observability

Error Handling

Data Persistence

Optimization

Monitoring

State Management

- Time

Engineers

DLT

+ Scale

DATA AI SUMMIT

# Power BI + Databricks SQL



Databricks SQL

Power BI On-premises Data Gateway

Direct Query (real-time)

Power BI Semantic Model

DATA·AI SUMMIT

# v7 Online Wave 2 Member Migration

MRID Omni v7 Survey Dashboard

Last Refresh Date
10/27/2023

DIGITAL INFORMATION MANAGEMENT

**Micro-Wave**
- ☐ Wave 2a
- ☐ Wave 2b
- ☐ Wave 2c

9/21/2023    10/25/2023

### Sessions & Unique Members - Online Wave 2
● Wave 2 Sessions    ● Wave 2 Unique Members

## Omnichannel | Real Time Online v7 Activity

Refreshed:    10/27/2023 11:58:27 AM

DIGITAL INFORMATION MANAGEMENT

**Select Options**

Wave 3a Unique Members & Sessions

### Current Day

| Unique Members | Sessions |
|---|---|
| 476 | 608 |

Unique Members & Sessions in last 24 hours (Duration by the hour in EST time)

● Unique Members    ● Unique Sessions

*"For the first time, I can measure real time member activities across our high volume Digital channels, which allows for quick decision making on system availability and scaling in real time."*
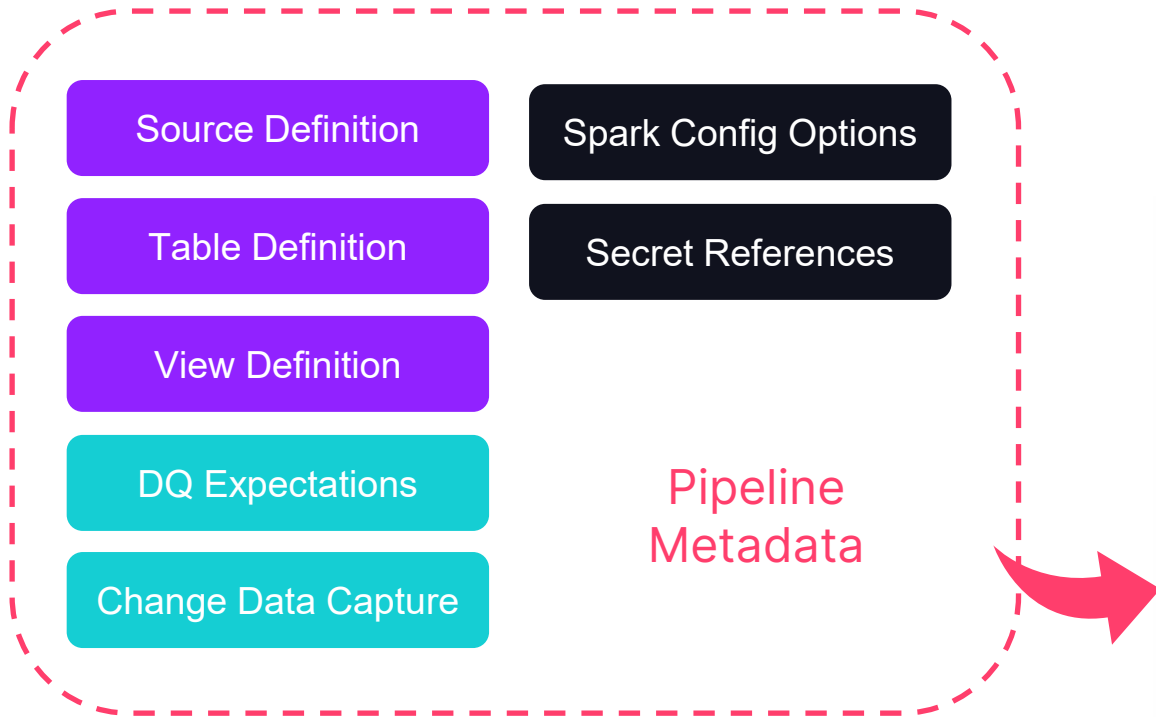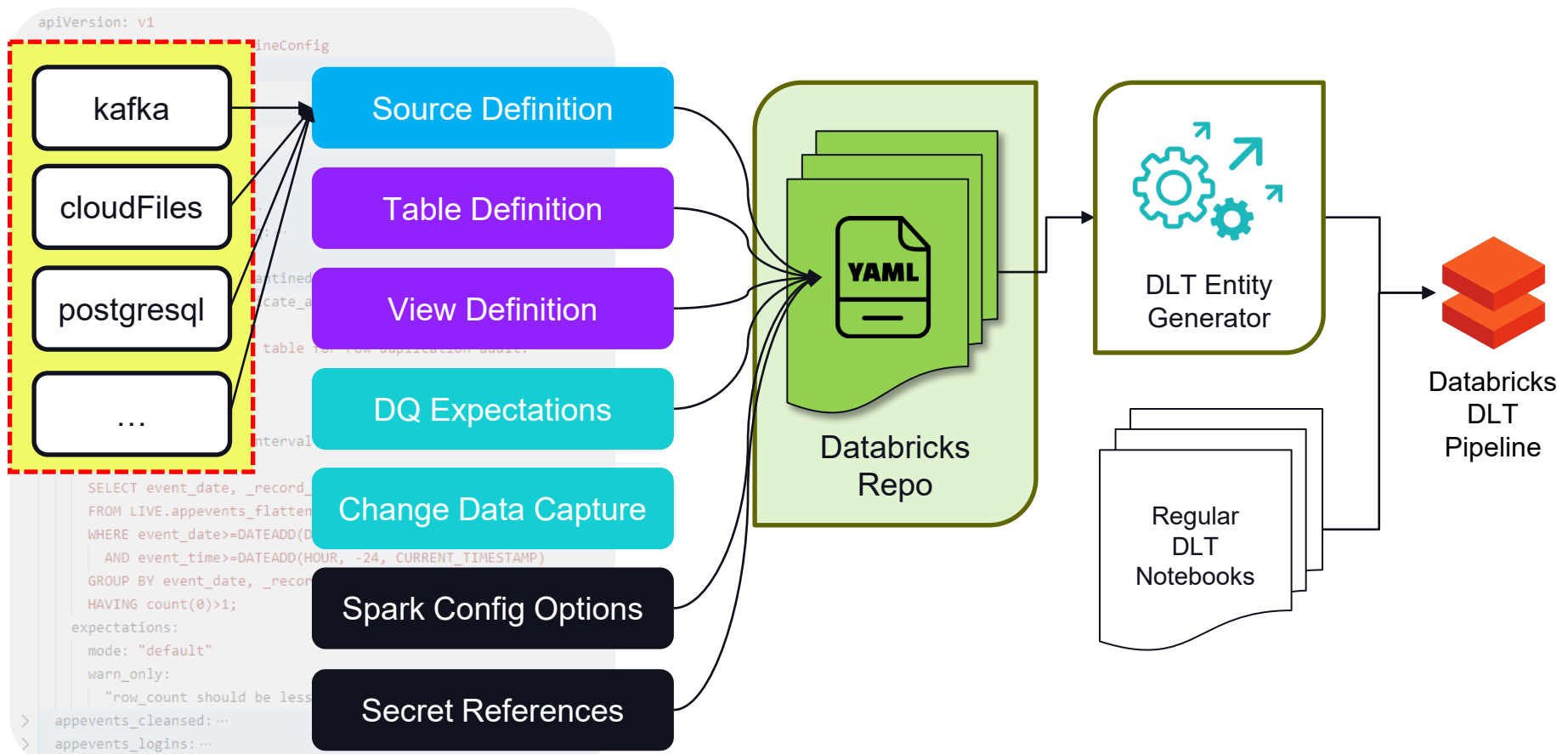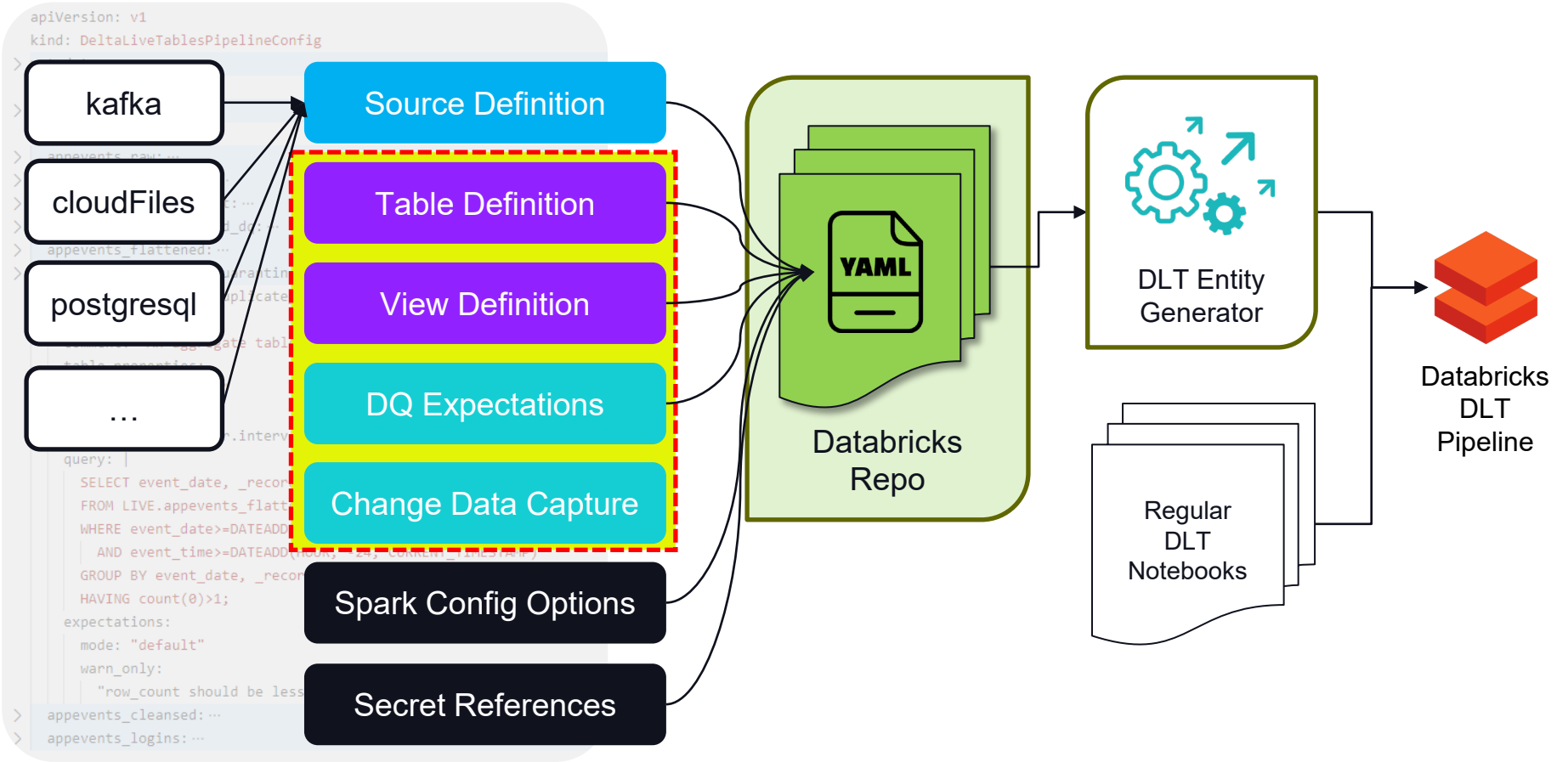
*- Gautam N., VP Digital Engineering*

DATA AI SUMMIT

# Speed to market

Online Banking v7
"Wave 0" Rollout

6/05/23

6/22/23

7/16/23

6/12/23

7/14/23

7/19/23

Proof-of-concept
(1 week)

Pipeline Dev, Test, & CI/CD
(3 weeks)

Pipeline
Go-live

Dashboard
Go-live

# METADATA-DRIVEN PIPELINE DEVELOPMENT

DLT API

Pydantic

Spark SQL API

YAML

DLT Entity Generator

a.k.a. "DEG"

Source Definition

Table Definition

View Definition

DQ Expectations

Change Data Capture

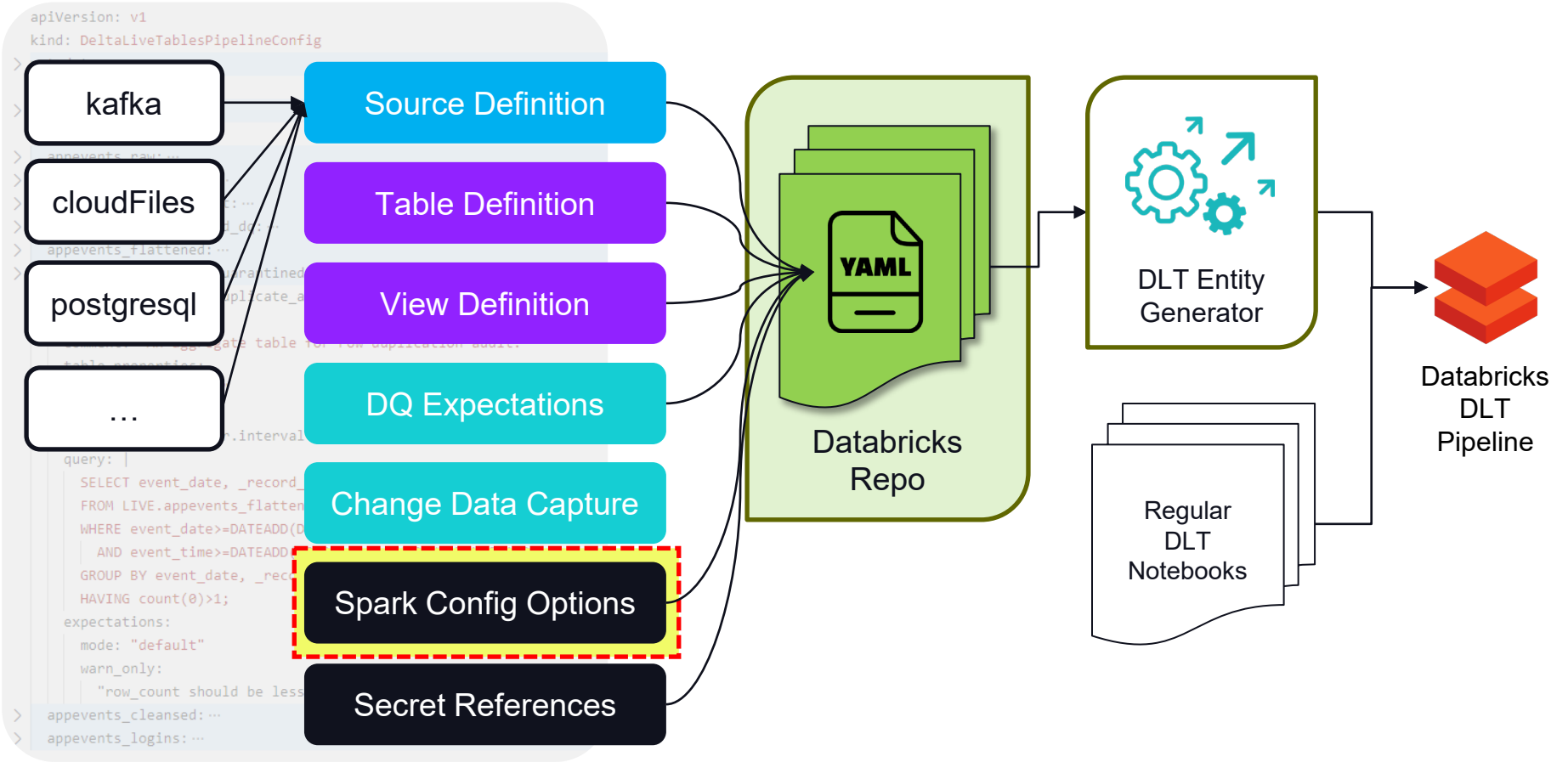Spark Config Options

Secret References

Pipeline Metadata

```yaml
apiVersion: v1
kind: DeltaLiveTablesPipelineConfig
> metadata: ⋯
  sources:
>   appevents_source: ⋯
  live_table_queries:
>   appevents_raw: ⋯
>   vw_appevents_raw_dq: ⋯
>   appevents_raw_dq_audit: ⋯
>   vw_appevents_flattened_dq: ⋯
>   appevents_flattened: ⋯
>   appevents_flattened_quarantined: ⋯
    appevents_flattened_duplicate_audit:
      type: "table"
      comment: "An aggregate table for row duplication audit."
      table_properties:
        "quality": "audit"
      spark_conf:
        "pipelines.trigger.interval" : "1 hour"
      query: |
        SELECT event_date, _record_hash_sha256, count(0) as row_count
        FROM LIVE.appevents_flattened
        WHERE event_date>=DATEADD(DAY, -1, CURRENT_DATE)
          AND event_time>=DATEADD(HOUR, -24, CURRENT_TIMESTAMP)
        GROUP BY event_date, _record_hash_sha256
        HAVING count(0)>1;
      expectations:
        mode: "default"
        warn_only:
          "row_count should be less than 2": "row_count<2"
>   appevents_cleansed: ⋯
>   appevents_logins: ⋯
```

DATA·AI SUMMIT

DATA'AI SUMMIT
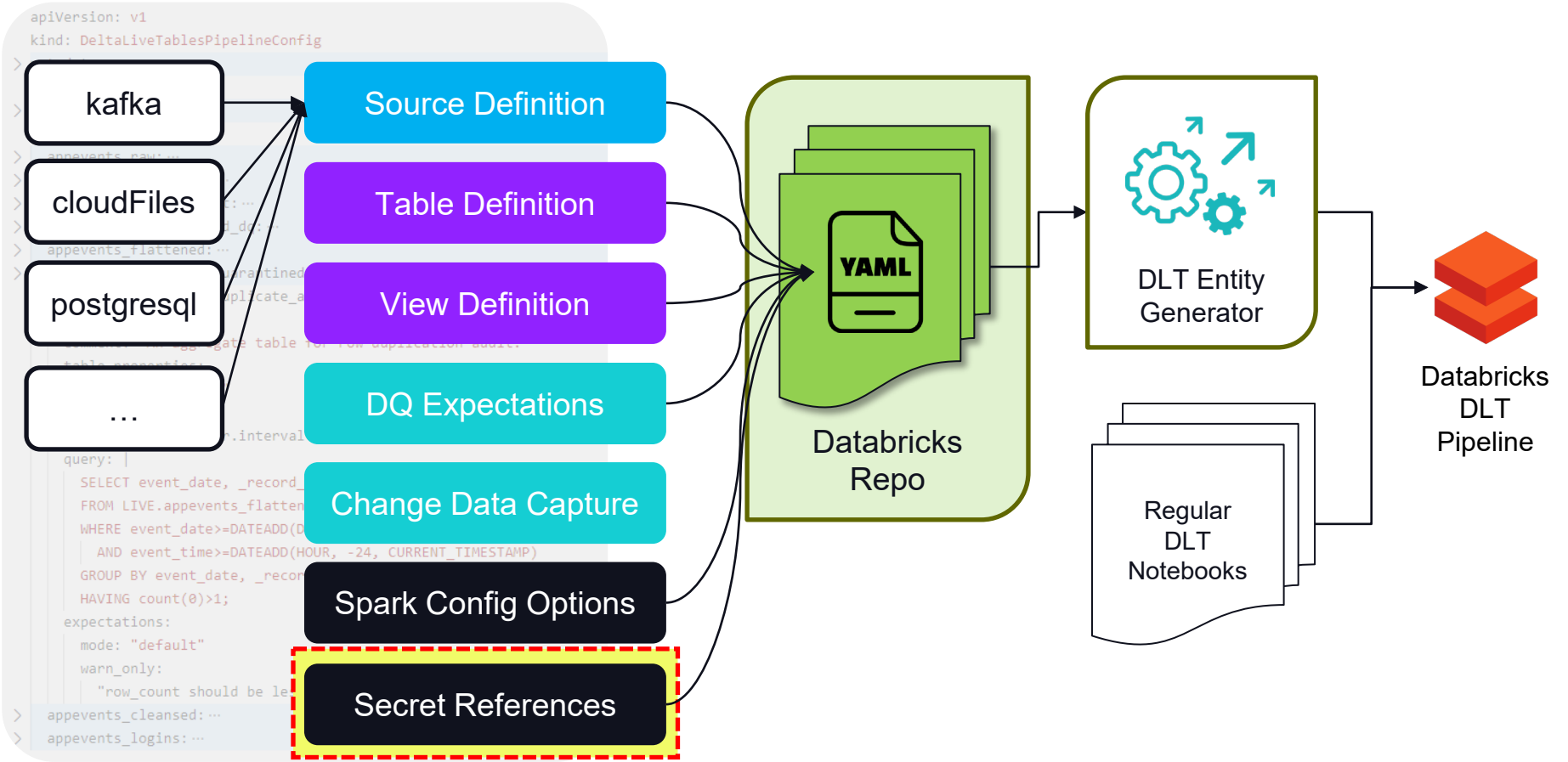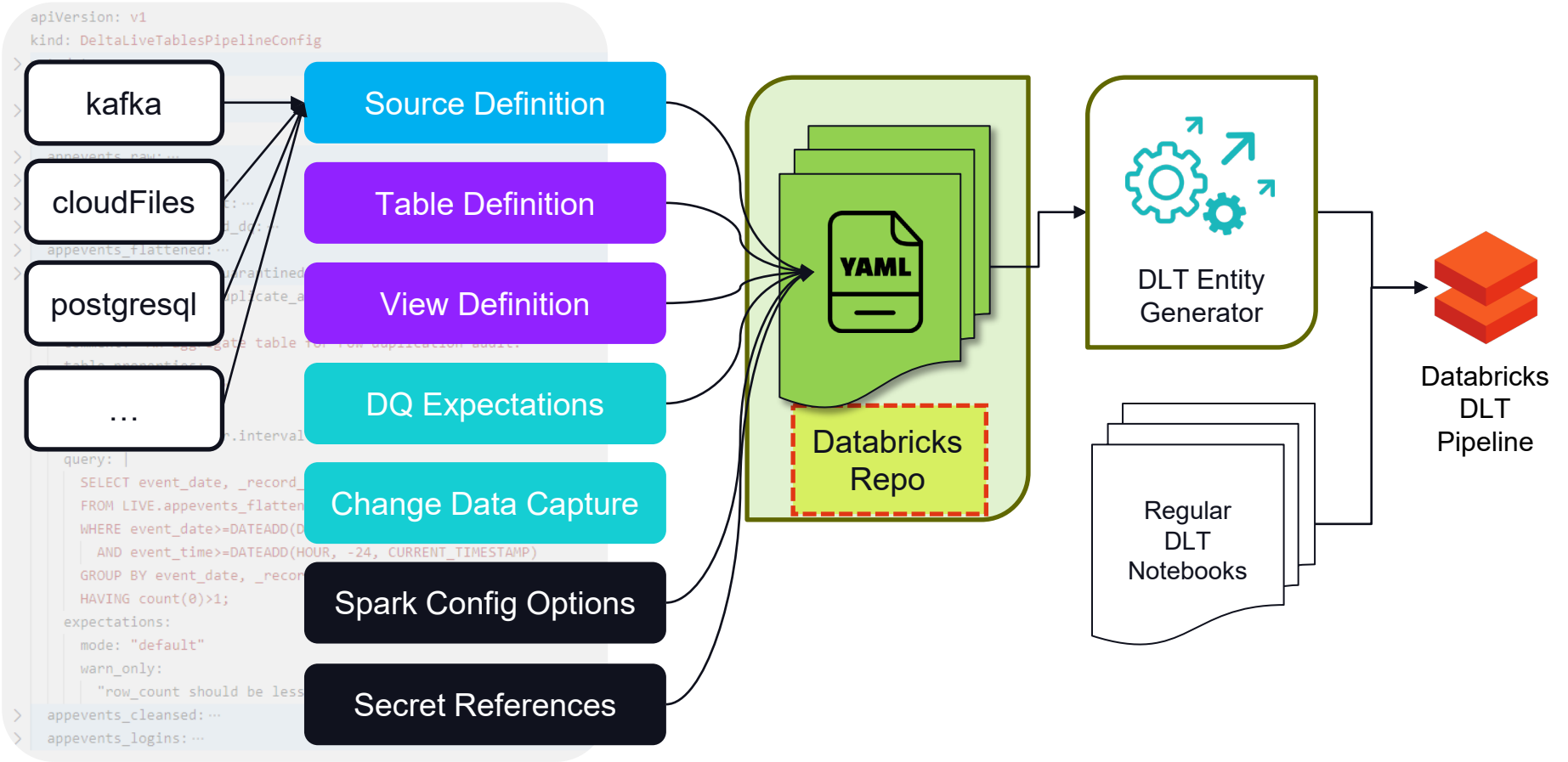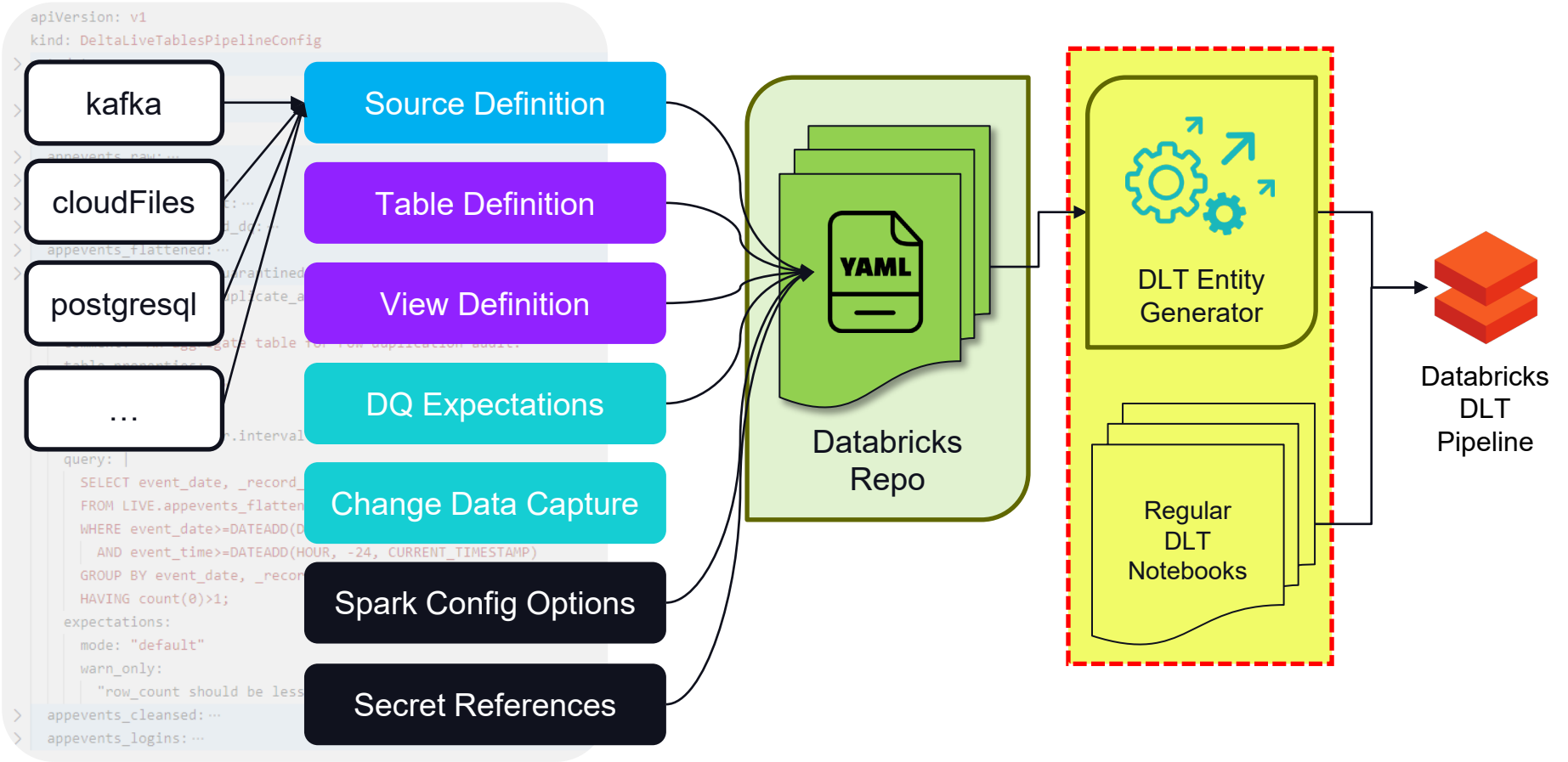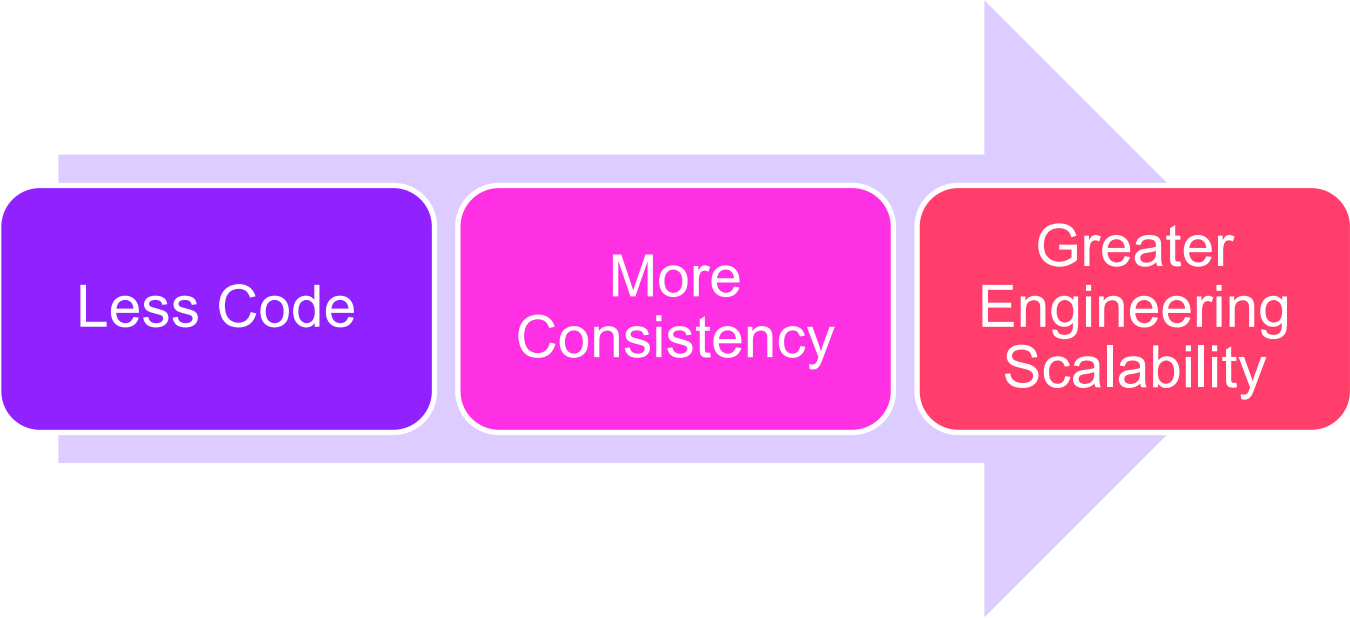
DATA'AI SUMMIT

DATA·AI SUMMIT

DATA·AI SUMMIT

# Metadata-driven solution for engineering scalability

**Less Code** → **More Consistency** → **Greater Engineering Scalability**

# SUMMARY

DLT isn't a silver bullet, but it worked for our use case wonderfully.

- **Simple** declarative programming model → speed to market
- Built-in **scalability**, **reliability**, and **optimization** → simplified operations
- **Programmability** → engineering creativity & scalability

# QUESTIONS?

# THANK YOU!

Connect on LinkedIn

Connect on LinkedIn