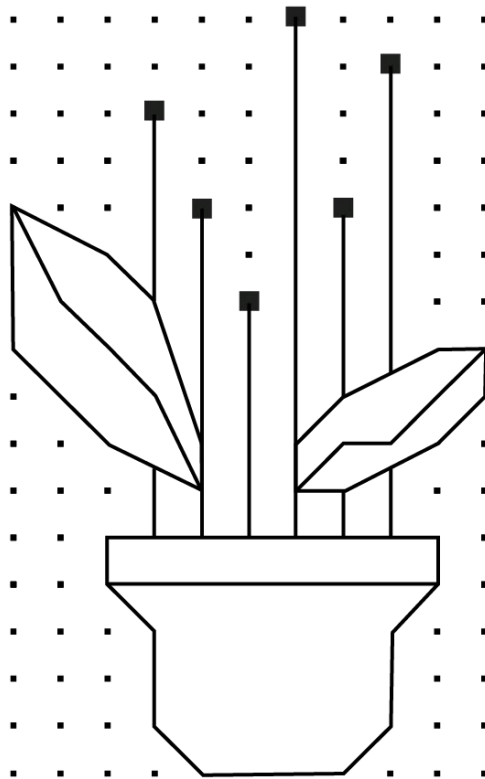# Invisible Interfaces

Considerations for Abstracting Complexities of a Real-time ML Platform

**Zhenzhong Xu**
Cofounder & CTO @ claypot.ai
July, 2023

The discovery of something invisible

# The Invisible Interface

Ubiquitous

Easy and responsive

Just works!

The endeavor to make things useful

*Real-time Decisions*

*that powers your business*

Fraud prevention     Personalization     Customer support     Dynamic pricing/discounting

Trending products     Risk Assessment     Account Take Over

Ads

ETA     Network analysis     Sentiment analysis     Object

detection

…

# The world is moving towards real-time

- [Instacart: The Journey to Real-Time Machine Learning](#) (2022)
  - Directly reduces **millions of fraud-related costs** annually.
- [LinkedIn's Real-time Anti-abuse](#) (2022)
  - LinkedIn moved from an offline pipeline (hours) to real-time pipeline (minutes), and saw **30% increase in bad actors caught online** and **21% improvement in fake account detection.**
- [How WhatsApp catches and fights abuse](#) (2022 | [slides](#))
  - **A few 100ms delay can increase the spam by 20-30%**.
- [How Pinterest Leverages Realtime User Actions in Recommendation to Boost Engagement](#) (2022)
  - According to Pinterest, this "*has been one of our most impactful innovations recently, increasing Home feed engagement by 11% while reducing Pinner hide volume by 10%*."
- [Airbnb: Real-time Personalization using Embeddings for Search Ranking](#) (2018)
  - Moving from offline scoring to online scoring grows bookings by **+5.1%**

# Real-time Decisions

| Exploration & Research | Model Architecture & Turning | Model Analysis & Selection | LLM Prompt Engineering |



# Data Fabric for Real-time AI

# Data Infrastructure

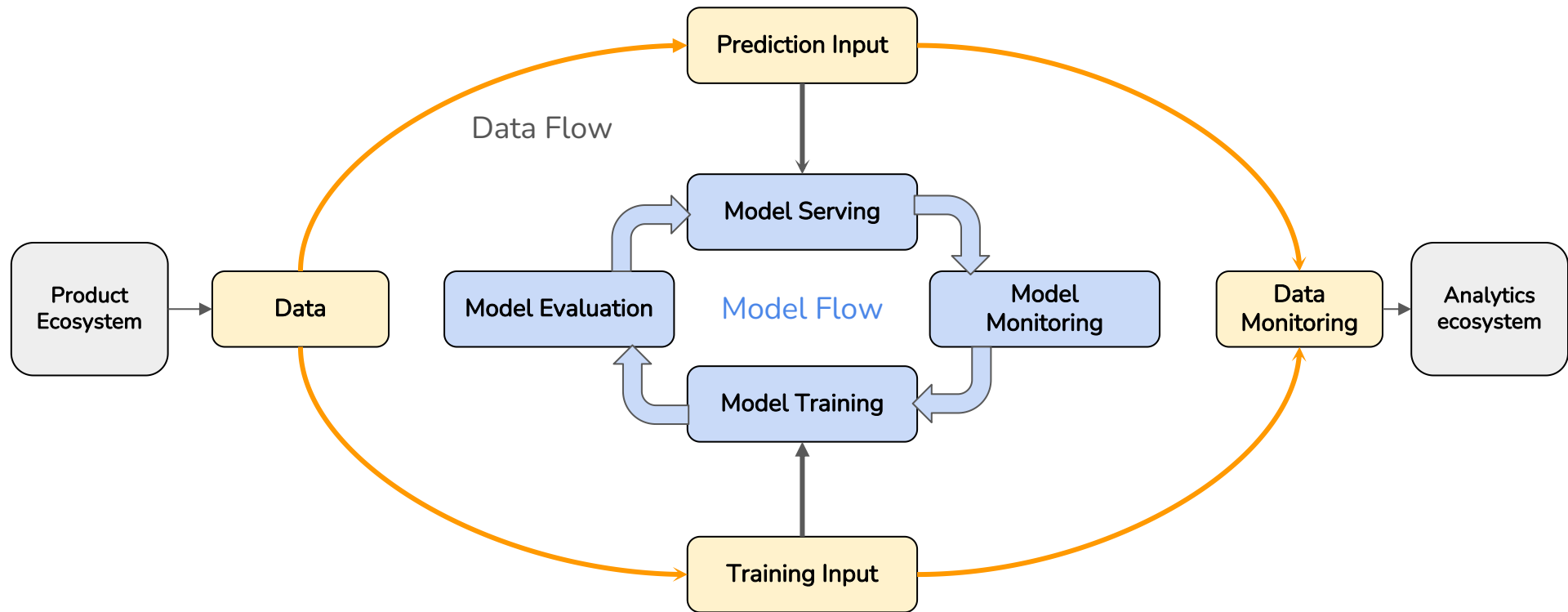| Data Sources | Ingestion & Transport | Storage | Query & Compute |

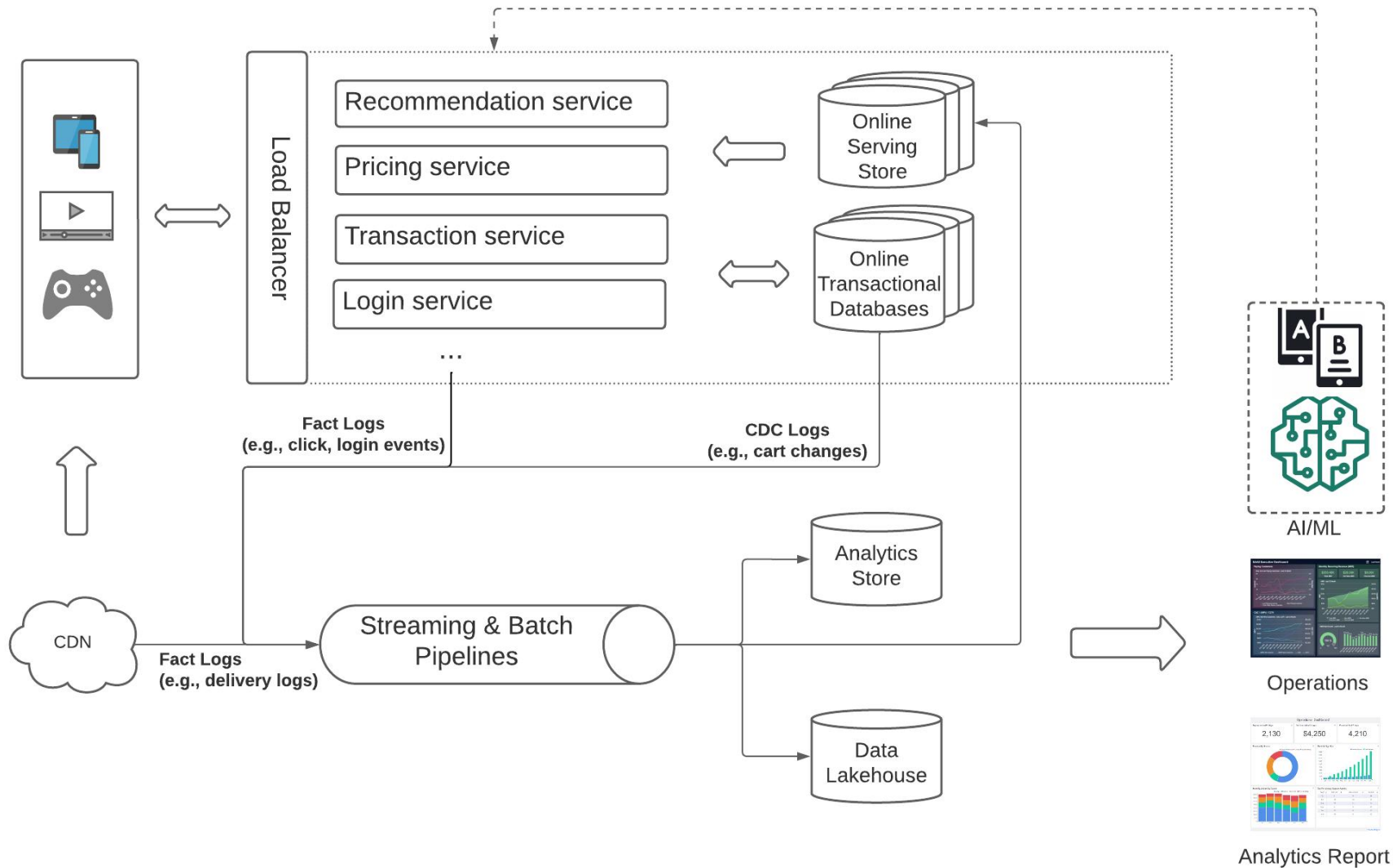| Workflow Orchestration | Analytics / Visualization | Multi-tenancy Isolation | Security & Governance |

Claypot

# The hard things towards real-time decisions

- Data silo and staleness
- Collaboration overhead
- Tech complexity

Recommendation service

Pricing service

Transaction service

Login service

Load Balancer

...

Online Serving Store

Online Transactional Databases

Fact Logs
(e.g., click, login events)

CDC Logs
(e.g., cart changes)

CDN

Fact Logs
(e.g., delivery logs)

Streaming & Batch Pipelines

Analytics Store

Data Lakehouse

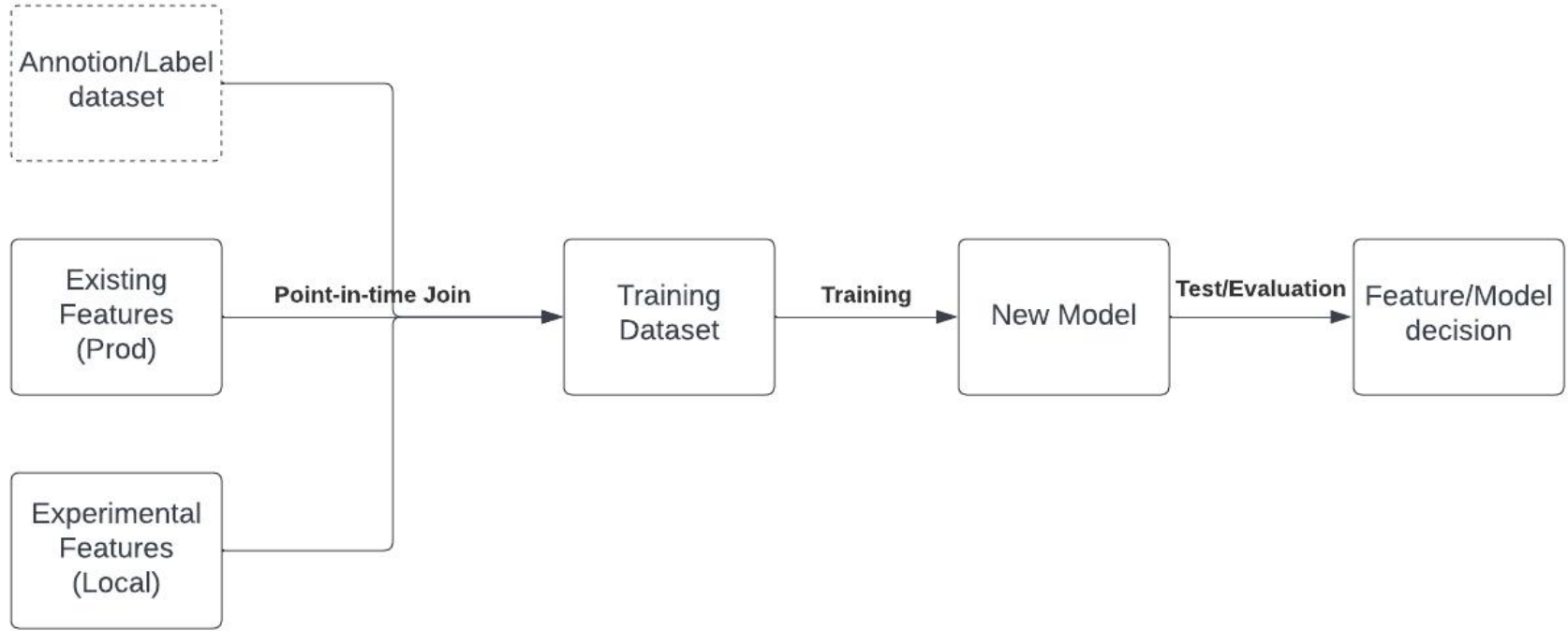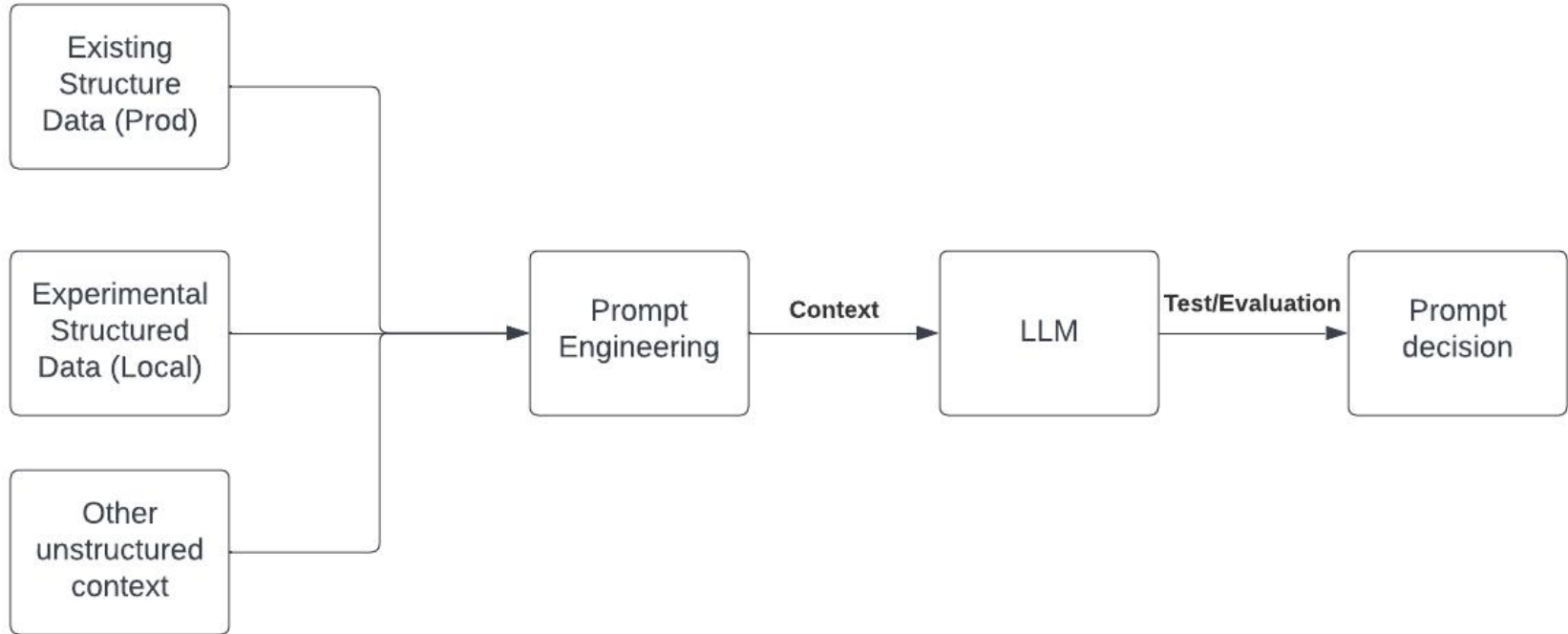AI/ML

Operations

Analytics Report

Claypot

## Challenge 1: From Experimentation to Production

- Slow prototyping
- Local vs. remote execution
- Divergent language & runtime

# Local Experimentation with Traditional Models
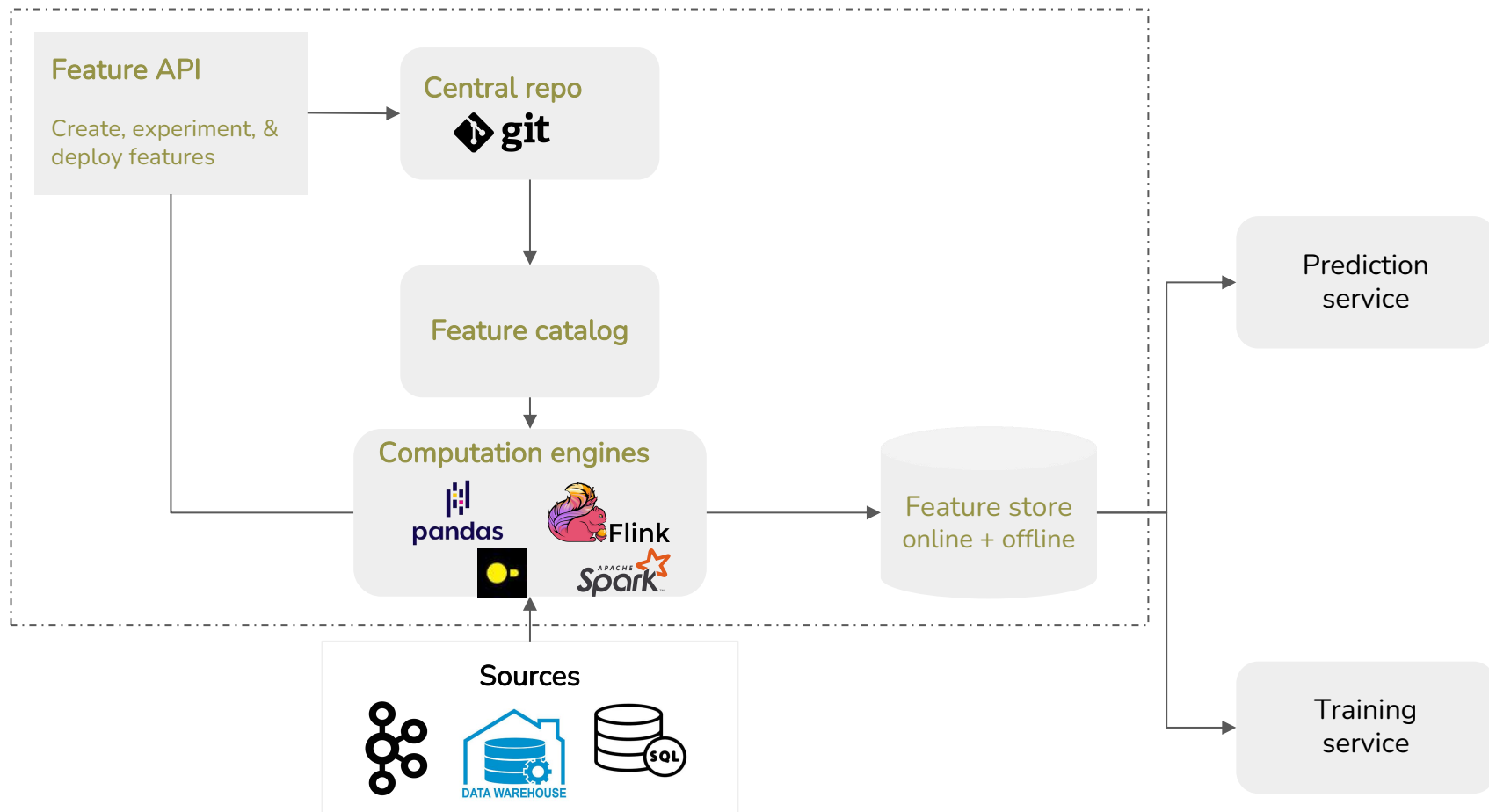
# Local Experimentation with LLMs

# Need an invisible interface to plug into compute ecosystems



Local/Single Machine

Remote/Distributed

# Declare features with familiar APIs

```python
@transformation
def average_transaction_amount_by_merchant(
    tx: Transactions,
    wspec: WindowSpec):

    return tx.groupby(["cc_num", "merchant"])["amt"].window(wspec).mean()
```
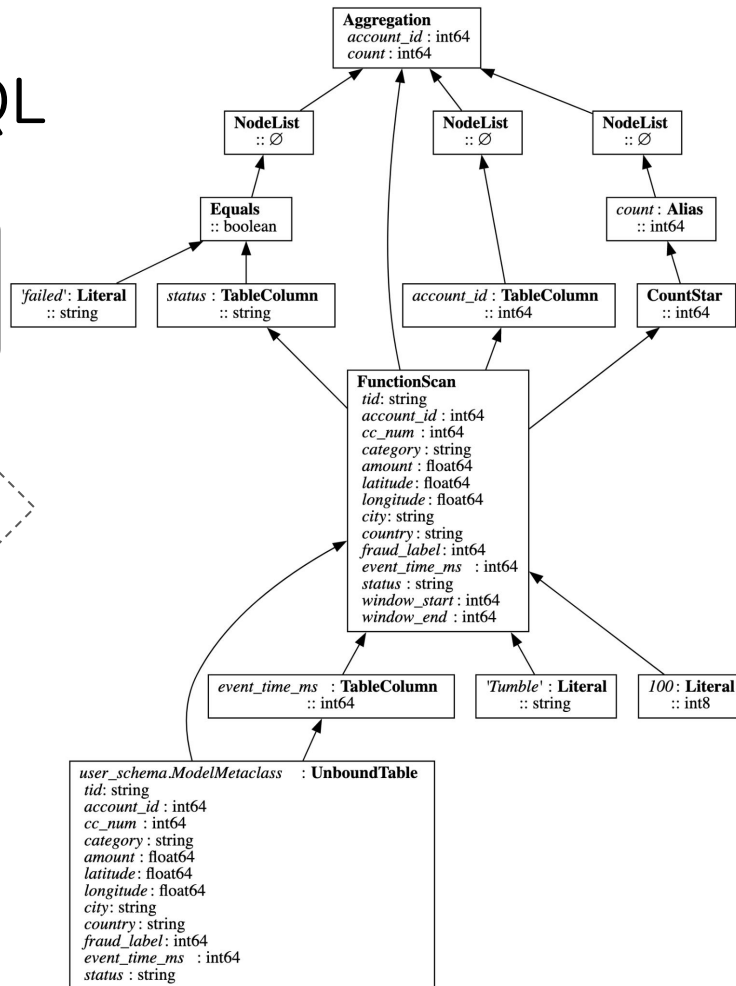
# Data Science Friendly: Python <> SQL

```python
@transformation
def transaction_count(tx: Transactions, wspec: WindowSpec):
    return tx[tx.status == "failed"].groupby("account_id").window(wspec).count()
```

**Relational Expression**

**Workload Compiler / Optimizer**

**Deployment**

**Aggregation**
*account_id* : int64
*count* : int64

**NodeList** :: ∅

**NodeList** :: ∅

**NodeList** :: ∅

**Equals** :: boolean

*count* : **Alias** :: int64

*'failed'* : **Literal** :: string

*status* : **TableColumn** :: string

*account_id* : **TableColumn** :: int64

**CountStar** :: int64

**FunctionScan**
*tid*: string
*account_id* : int64
*cc_num* : int64
*category* : string
*amount* : float64
*latitude*: float64
*longitude* : float64
*city*: string
*country* : string
*fraud_label*: int64
*event_time_ms*  : int64
*status* : string
*window_start* : int64
*window_end* : int64

*event_time_ms*  : **TableColumn** :: int64

*'Tumble'* : **Literal** :: string

*100* : **Literal** :: int8

*user_schema.ModelMetaclass*  : **UnboundTable**
*tid*: string
*account_id* : int64
*cc_num* : int64
*category* : string
*amount* : float64
*latitude*: float64
*longitude* : float64
*city*: string
*country* : string
*fraud_label*: int64
*event_time_ms*  : int64
*status* : string

# Same code can run on different computation engines

```
@transformation
def transaction_count(tx: Transactions, wspec: WindowSpec):
    return tx[tx.status == "failed"].groupby("account_id").window(wspec).count()
```

Intermediate Representation

**Relational Expression**

Compile into a relational expression (RE), which is SQL equivalent

**Workload Compiler/Optimizer**

Compile & optimize RE into the computation engine
(e.g., Panda, DuckDb, Flink, Spark) best suited for the job

**Deployment**
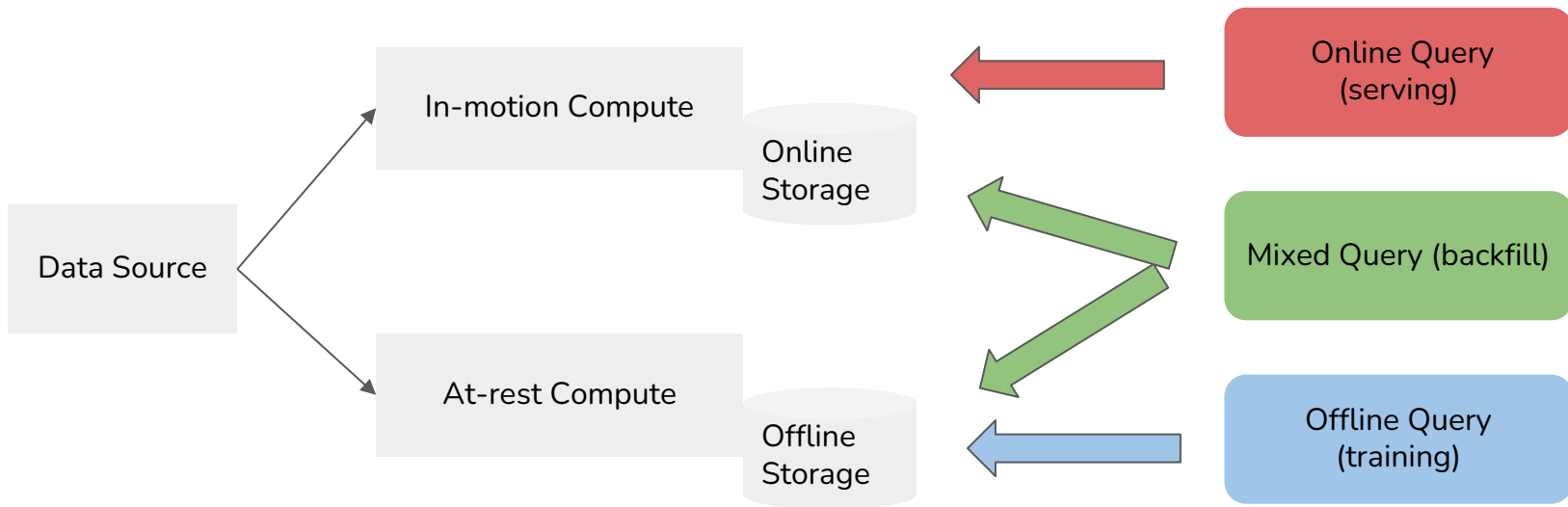
Spin up and manage computation jobs

Claypot

## Solution 1: Relational Expression based Compilation

- Unified yet familiar API
- Pluggable to many compute engines
- Minimize human error
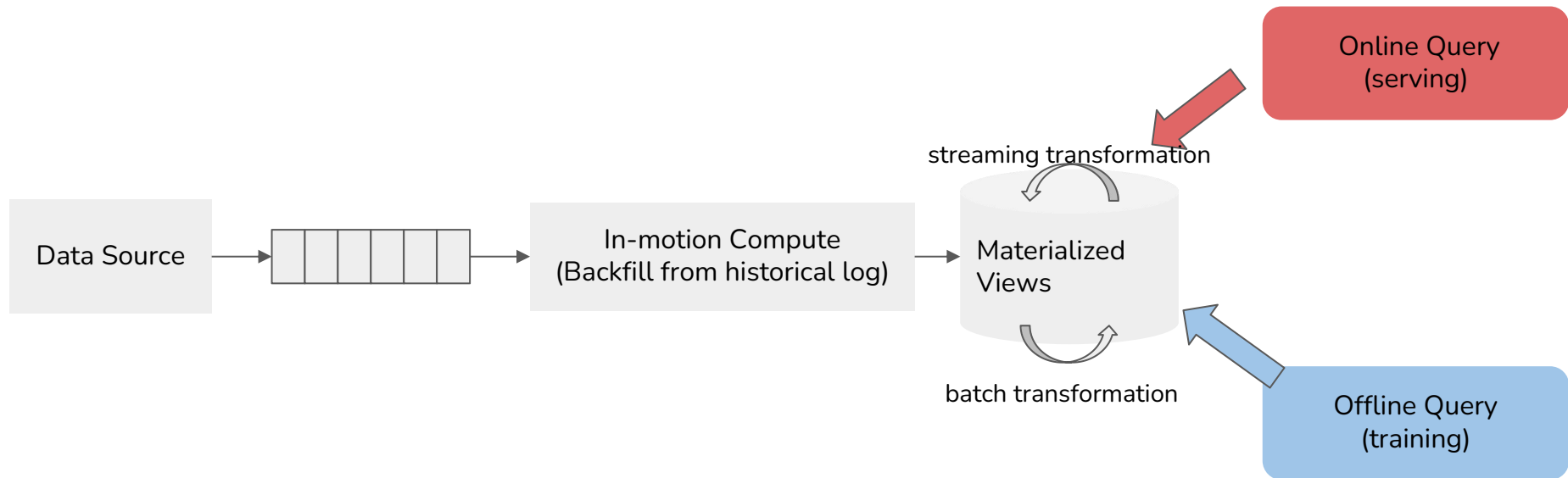- Prototype in minutes

Claypot

## Challenge 2: Streaming and Batch Divided

- Evolving architecture
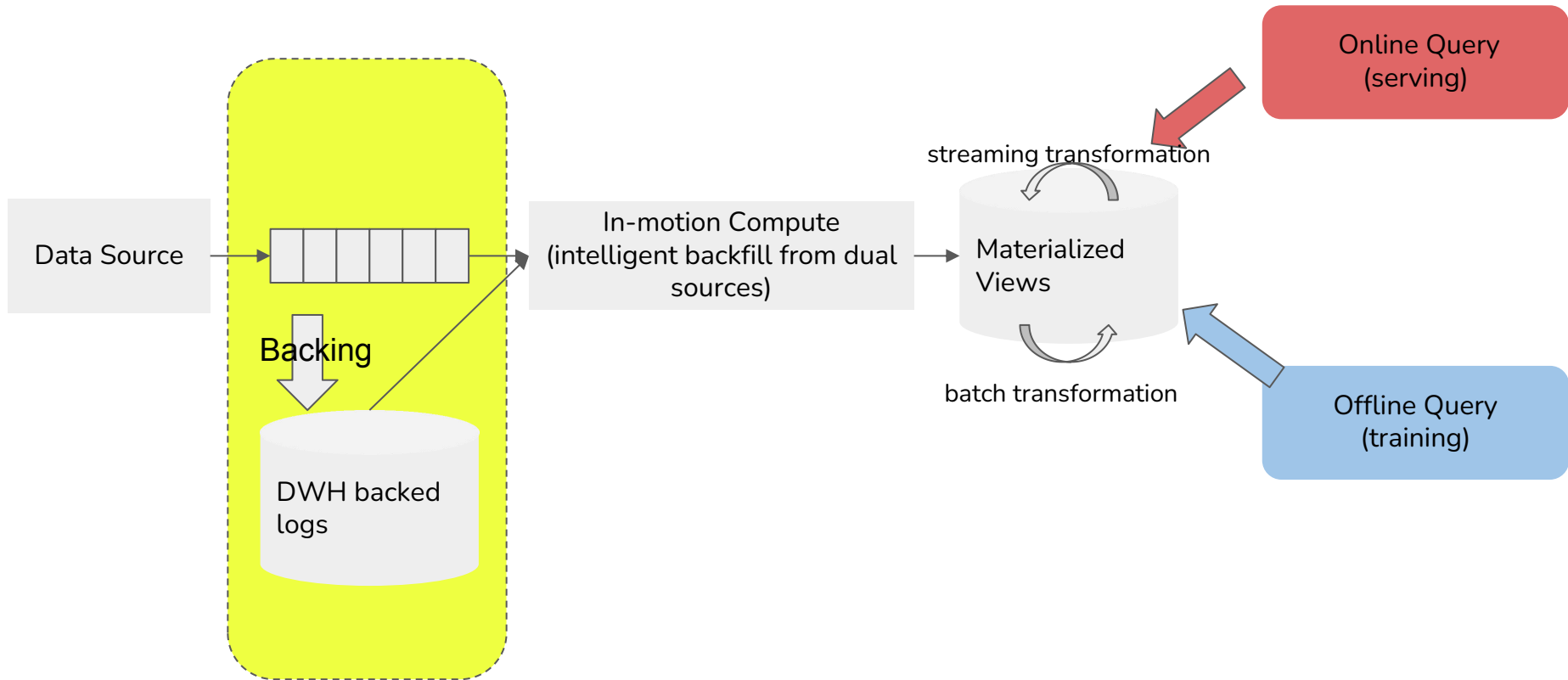- Difficult to backfill
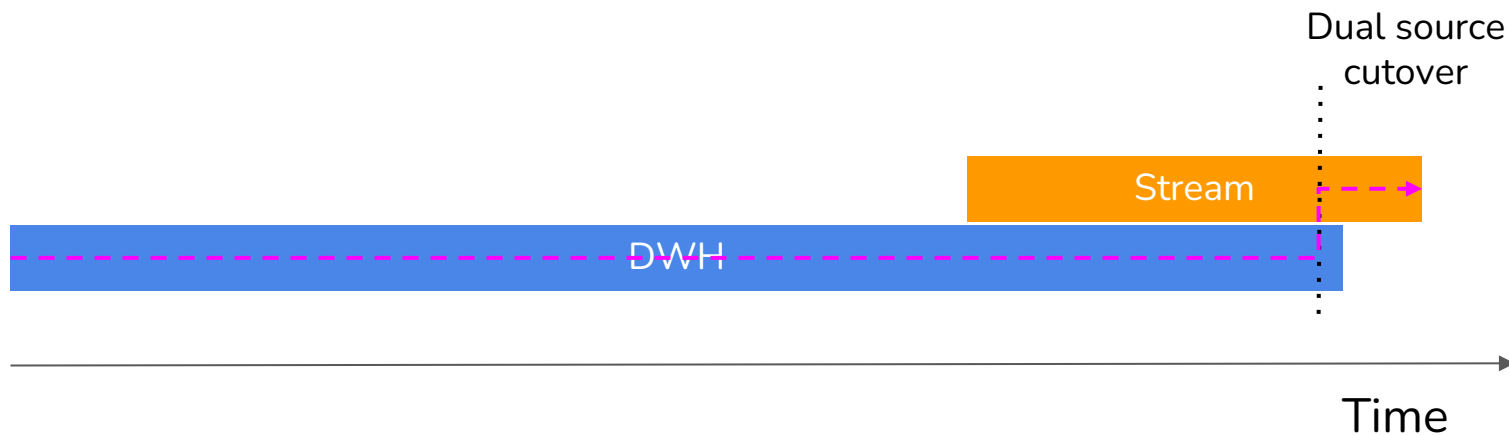- Train-predict inconsistencies

# Lambda Architecture

Data Source

In-motion Compute

Online Storage

At-rest Compute

Offline Storage

Online Query (serving)

Mixed Query (backfill)

Offline Query (training)

# Kappa (Streaming) Architecture

Data Source → [||||||] → In-motion Compute (Backfill from historical log) → Materialized Views

streaming transformation

batch transformation

Online Query (serving)

Offline Query (training)

# Unified Architecture



Data Source

Backing

DWH backed logs

In-motion Compute (intelligent backfill from dual sources)

Materialized Views

streaming transformation

batch transformation

Online Query (serving)

Offline Query (training)

# Batch and streaming source unified to simplify backfill

# Need an invisible interface to plug into storage ecosystems



Streaming Leaning

Batch Leaning

# Data Fabric for a Streaming Pipeline

**User Infra**

Schema Registry

Upstream Producer

**(1)**

Kafka Topic

**(2)**

**Streaming Pipeline**

Kafka Source Connector

**(3)**

Transformation

**(4)**

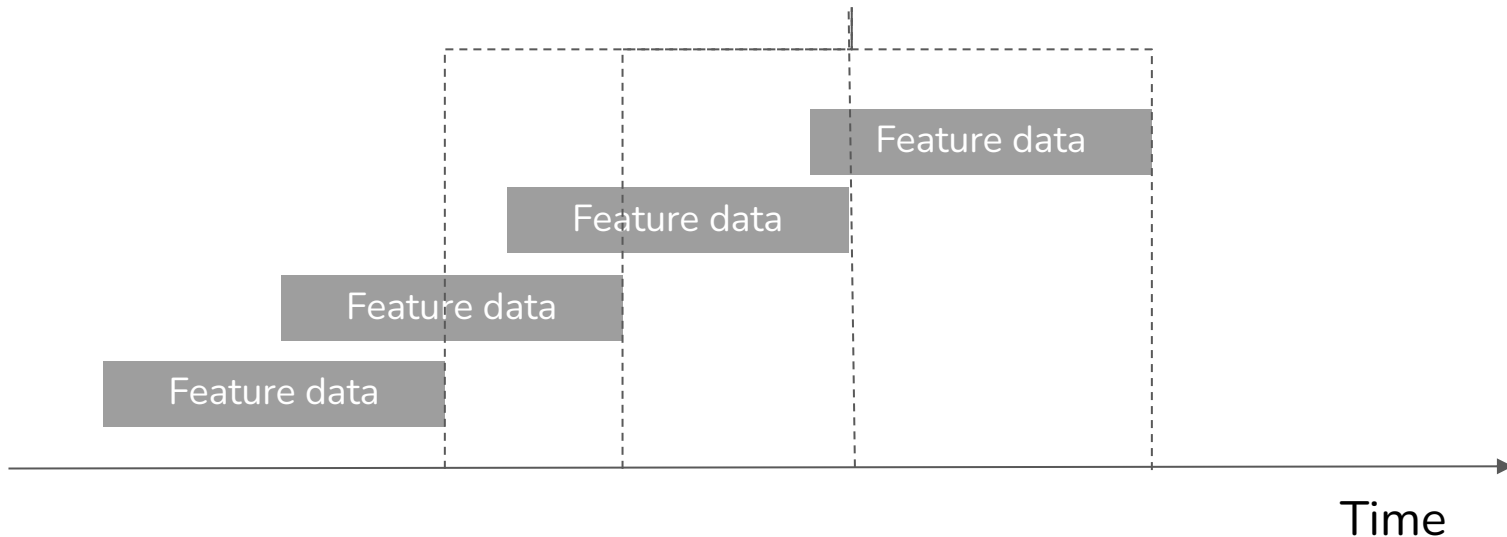Iceberg Sink Connector

**(5)**

**User Infra**

Catalog

Iceberg Table

# Data Fabric for a Unified Backfill Pipeline

# Training dataset backfill requires point-in-time correctness

# Point-in-time joins to generate training data

Given a spine (entity keys + timestamp + label), join features to generate training data

spine_df

| inference_ts | tid | cc_num | user_id | is_fraud |
|---|---|---|---|---|
| 21:30 | 0122 | 2 | 1 | 0 |
| 21:40 | 0298 | 4 | 1 | 0 |
| 21:55 | 7539 | 6 | 3 | 1 |

cc_num_tx_max_1h

| ts | cc_num | tx_max_1h |
|---|---|---|
| 9:20 | 2 | … |
| 10:24 | 2 | … |
| 20:00 | 4 | … |

user_unique_id_30d

| ts | user_id | unique_ip_30d |
|---|---|---|
| 6:00 | 1 | … |
| 6:00 | 3 | … |
| 6:00 | 5 | … |

```
train_df = pitc_join_features(
    spine_df,
    features=[
        "tx_max_1h",
        "user_unique_ip_30d",
    ],
)
```

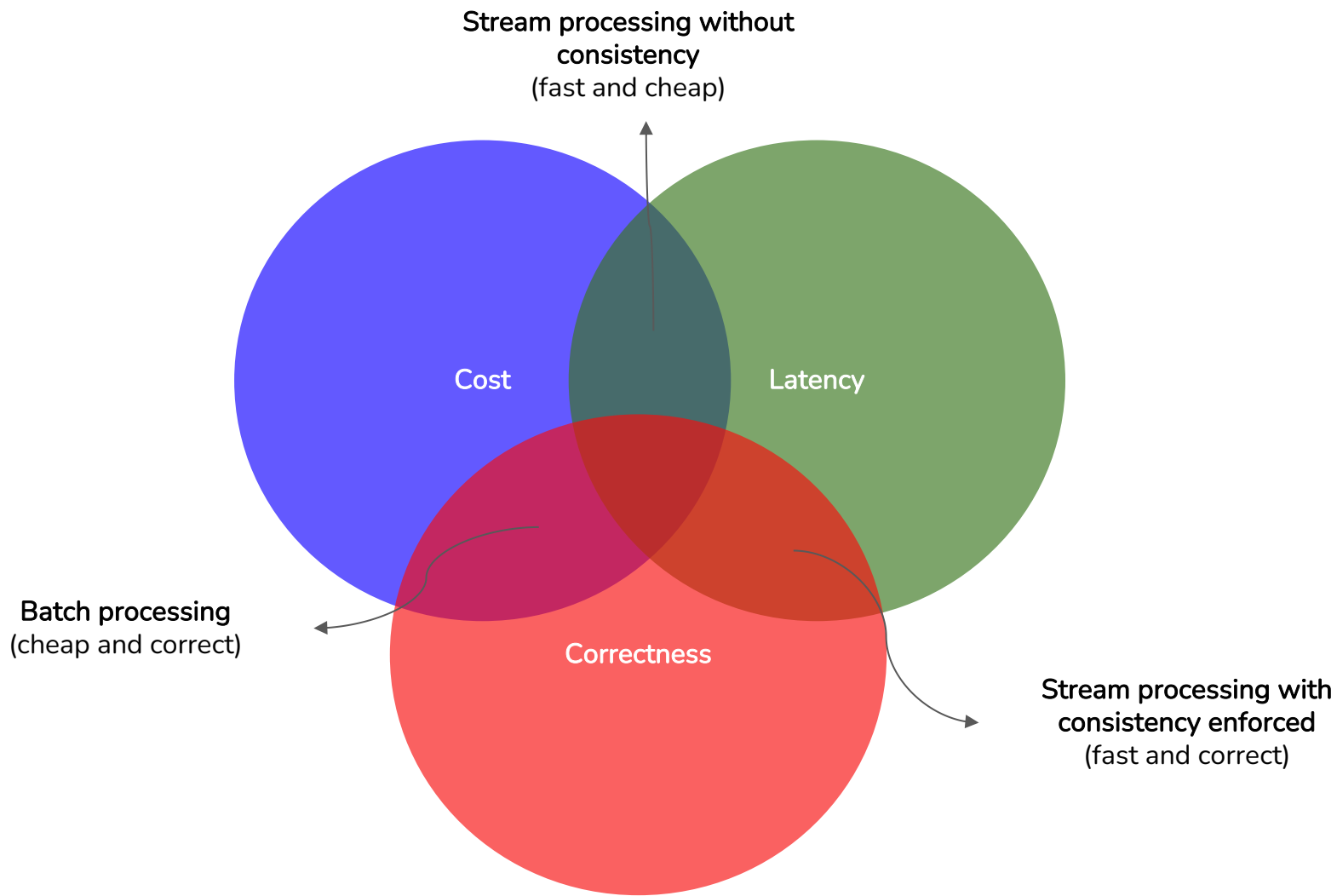| inference_ts | tid | cc_num | user_id | is_fraud | tx_max_1h | user_unique_ip_30d |
|---|---|---|---|---|---|---|
| 21:30 | 0122 | 2 | 1 | 1 | … | … |
| 21:40 | 0298 | 4 | 1 | 1 | … | … |
| 21:55 | 7539 | 6 | 3 | 3 | … | … |

Claypot

## Solution 2: Abstract streaming and batch data storage

- Unified streaming & batch source
- Unified online & offline feature stores
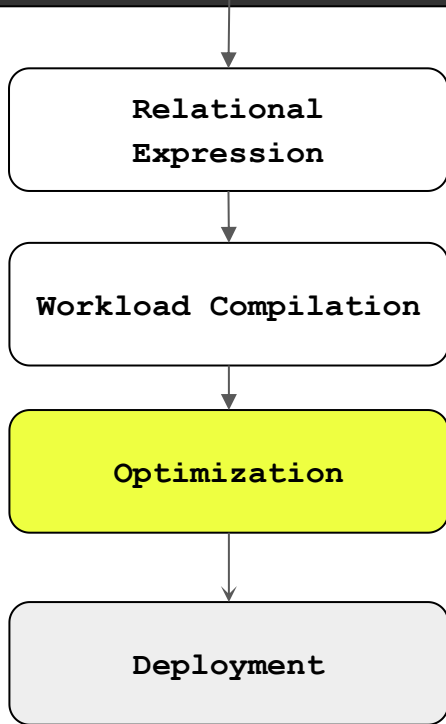- Pluggable to most storage technologies

Claypot

- Cost, latency, correctness surprises!
- Lack optimizations knobs

Stream processing without
consistency
(fast and cheap)

Cost

Latency

Batch processing
(cheap and correct)

Correctness

Stream processing with
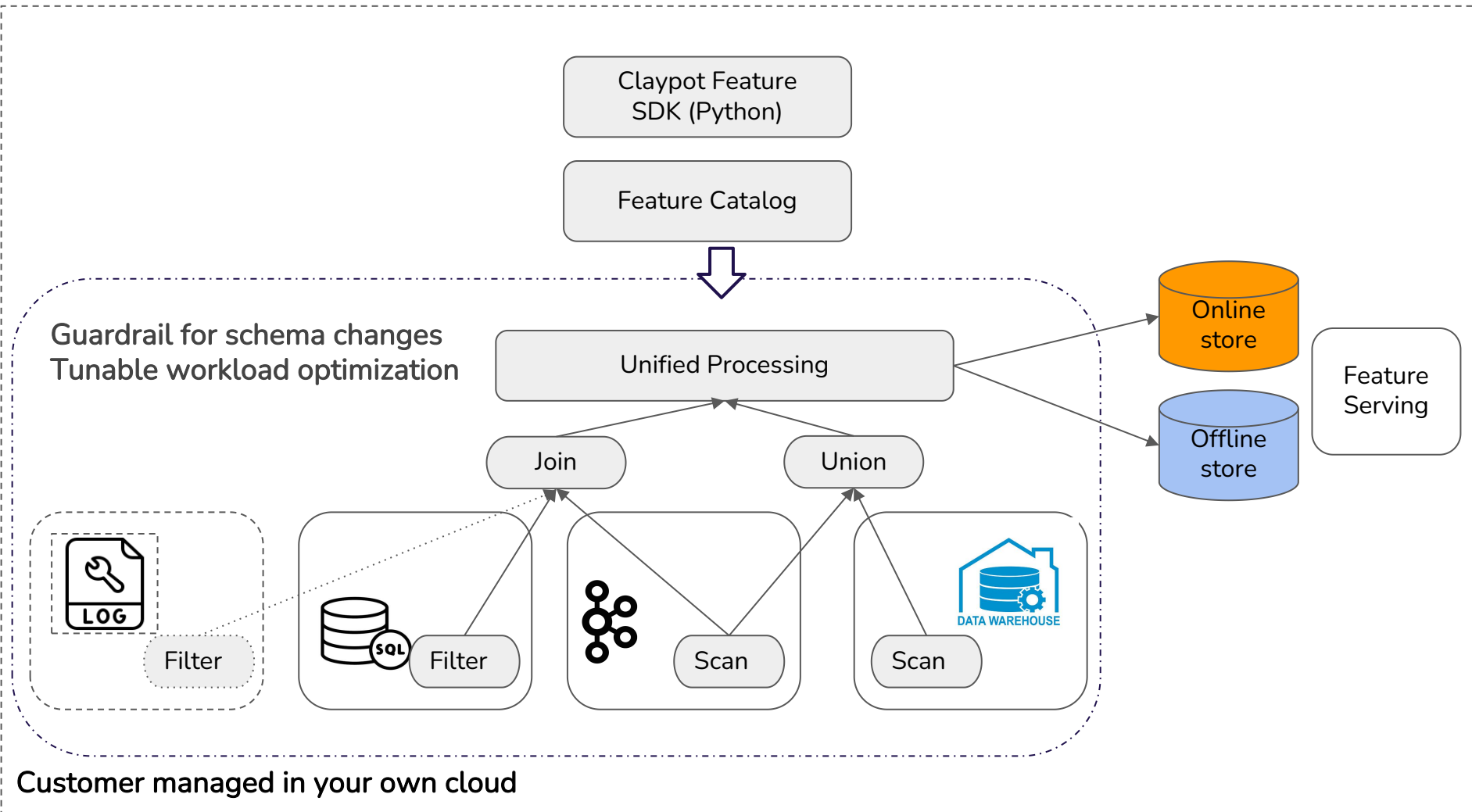consistency enforced
(fast and correct)

# Optimization

```
@transformation
def transaction_count(tx: Transactions, wspec: WindowSpec):
    return tx[tx.status == "failed"].groupby("account_id").window(wspec).count()
```

```
Relational
Expression
```

```
Workload Compilation
```

```
Optimization
```

```
Deployment
```

Various intelligent optimization can be done to make appropriate tradeoff across storage and compute systems.

Claypot

## Solution 3: Optimization knobs

- Abstract optimization complexity
- User controls with high level knobs
- Trust, no surprises!

# Claypot

# Make invisible interface possible!

- Ubiquitous
- Easy and responsive
- Just works!

https://zhenzhongxu.com/
zhenzhong@claypot.ai

the invisible interface