

Chronon, Airbnb's open source feature engineering framework

Pengyu Hou, Software Engineer, Airbnb



FEATURE STORE SUMMIT 2024

DATA FOR AI:
REAL-TIME, BATCH, AND LLMS

Organized by  **HOPSWORKS**



Agenda

- Why Chronon?
- What is Chronon?
- Point in time correct
- Life before Chronon
- Architecture
- Examples
- Takeaways + Call to action



Why Chronon?



FEATURE STORE SUMMIT 2024

DATA FOR AI:
REAL-TIME, BATCH, AND LLMS



Why Chronon?

What are some common challenges?

- Feature definition in multiple places
- Online offline inconsistency
- Slow backfills \Rightarrow slow iteration speed
- Repetitive glue code for the pipelines
- Duplicate work across multiple orgs



What is Chronon?



FEATURE STORE SUMMIT 2024

DATA FOR AI:
REAL-TIME, BATCH, AND LLMS

Chronon is an open-source platform enabling ML practitioners to efficiently develop, deploy, manage, and monitor ML data pipelines.



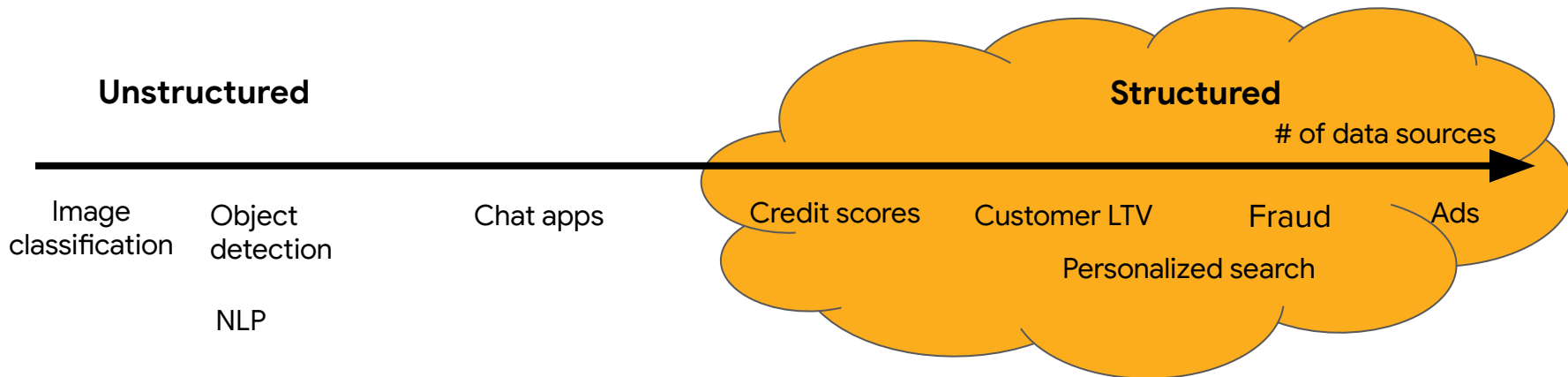
FEATURE STORE SUMMIT 2024

DATA FOR AI:

REAL-TIME, BATCH, AND LLMS



ML Applications



- Most of the data is available at once: full image
- Features are *automatically* extracted from *few* (often one) data stream:
 - words from a text
 - pixels from an image

- Data arrives steadily as user interacts with the platform
- Features extracted from *many* event streams:
 - logins
 - clicks
 - page views, etc
- **Iterative *manual* feature engineering**

Journey to open source



FEATURE STORE SUMMIT 2024

DATA FOR AI:
REAL-TIME, BATCH, AND LLMS



Journey to open source

- 2017-2022: iterated over four internal versions
- 2022: private beta partnered with Stripe
 - Fully adopted within Stripe
 - They are also in this summit
- 2024: announced open source
 - Stripe is the co-maintainer



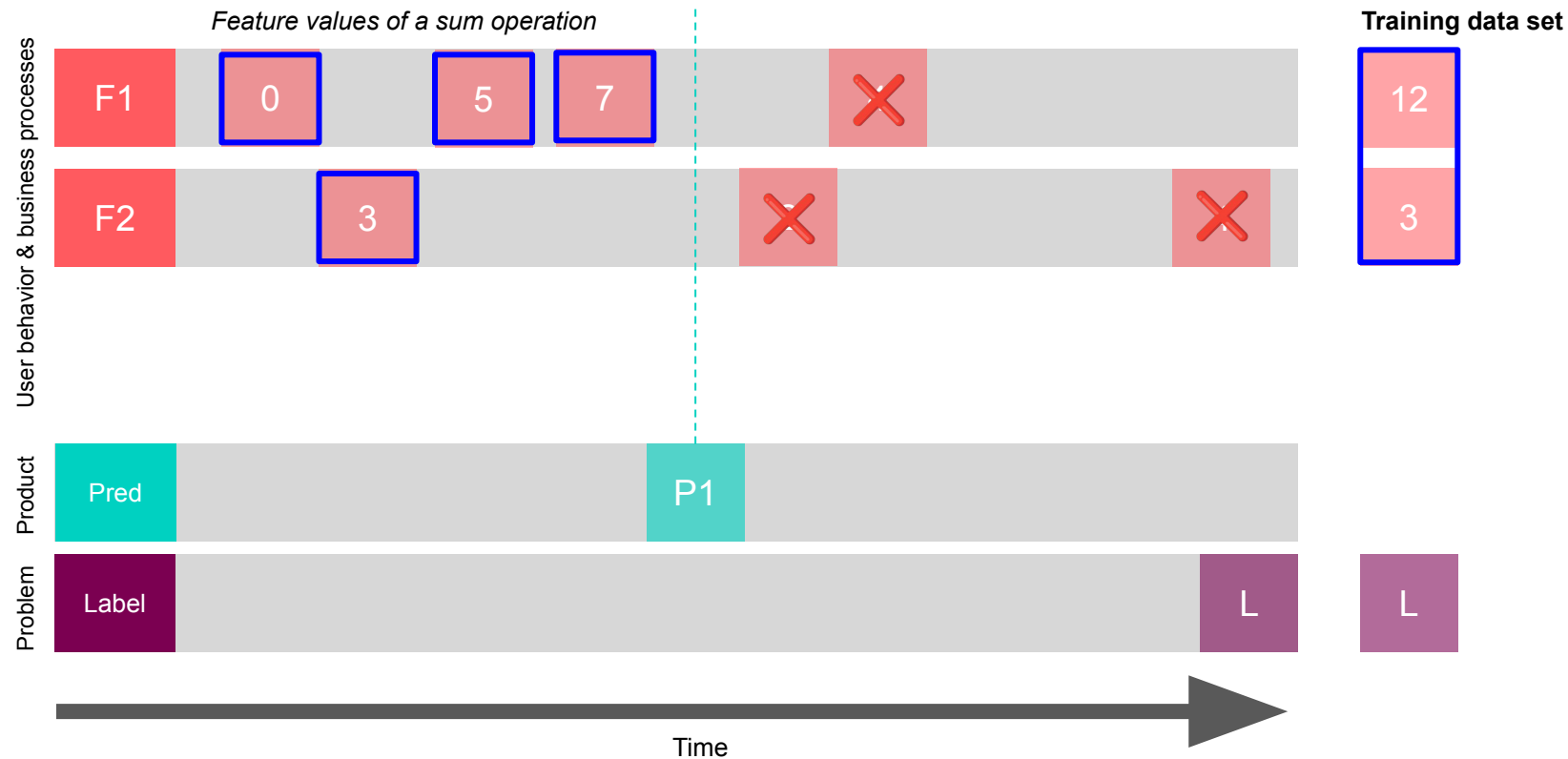
What is Point in Time Correct? (PITC)?



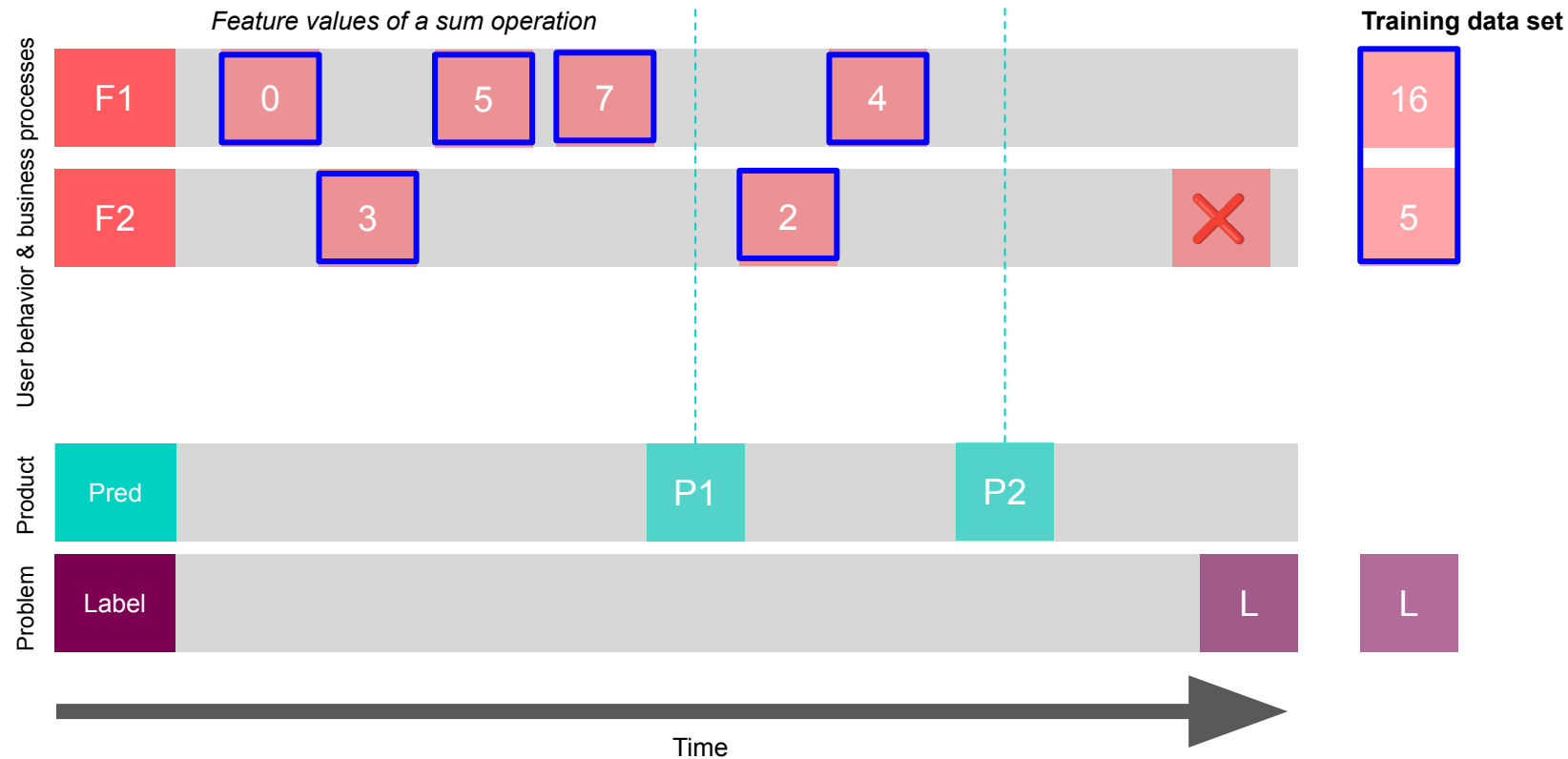
FEATURE STORE SUMMIT 2024

DATA FOR AI:
REAL-TIME, BATCH, AND LLMS

Point in time correct feature values



Point in time correct feature values





Point in time correct matters

What are some common symptom?

- Works well in local
- Performs poorly in prod





Point in time correct matters

What are some common symptom?

- Works well in local
- Performs poorly in prod

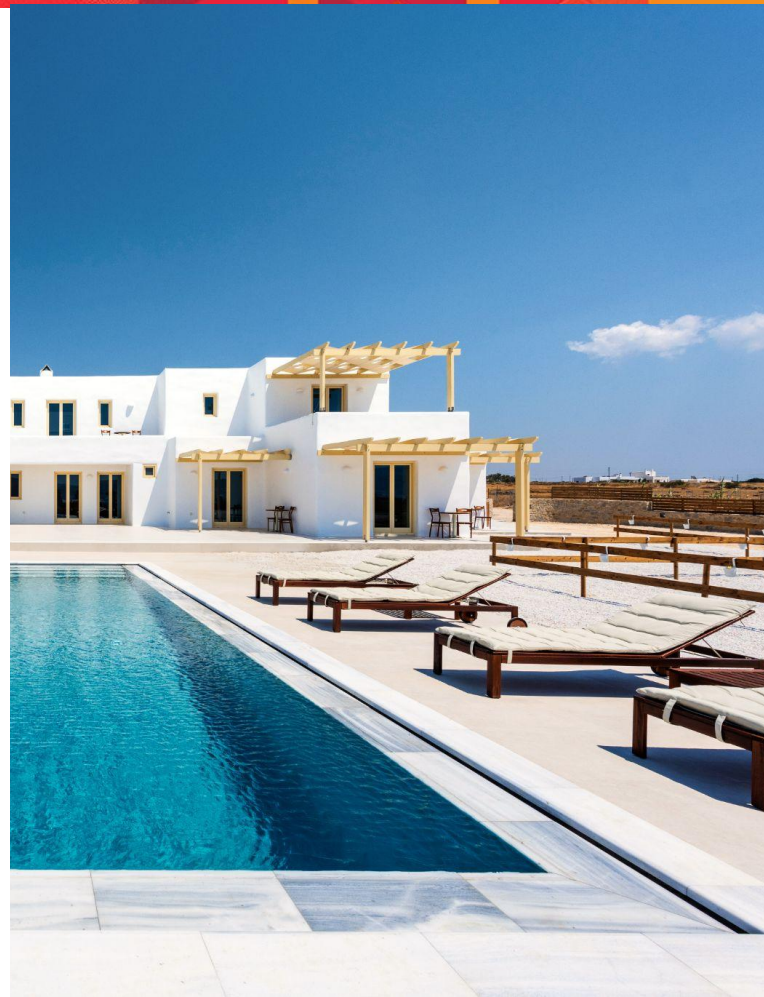
Data leakage example:

A house price prediction model can be based on:

- Rooms, location, size, age, etc

But not:

- Offers received
- Buyer's info



What was life like before Chronon?



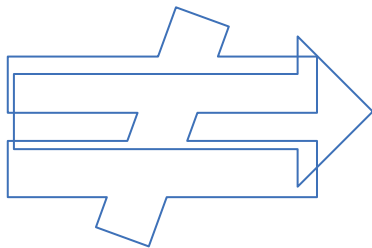
FEATURE STORE SUMMIT 2024

DATA FOR AI:
REAL-TIME, BATCH, AND LLMS



Replicate offline -> online

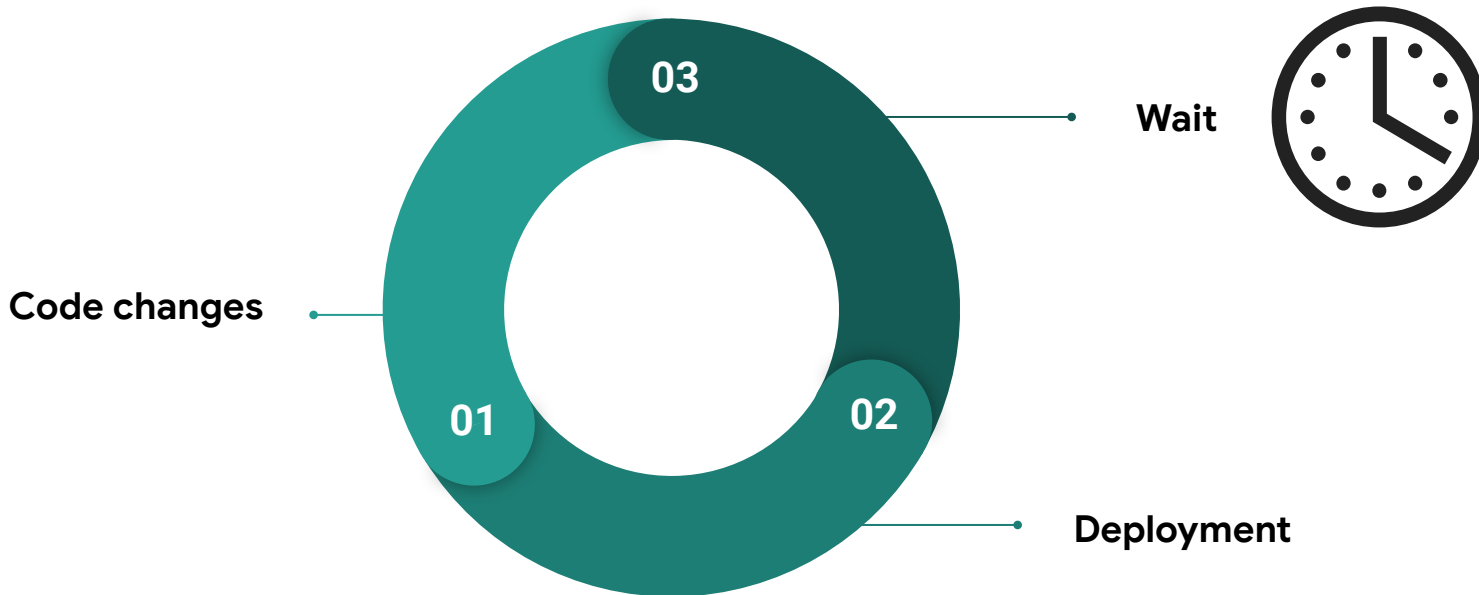
SQL



Program
Language of
the service



Log and wait



How can Chronon address consistency without waiting for months?



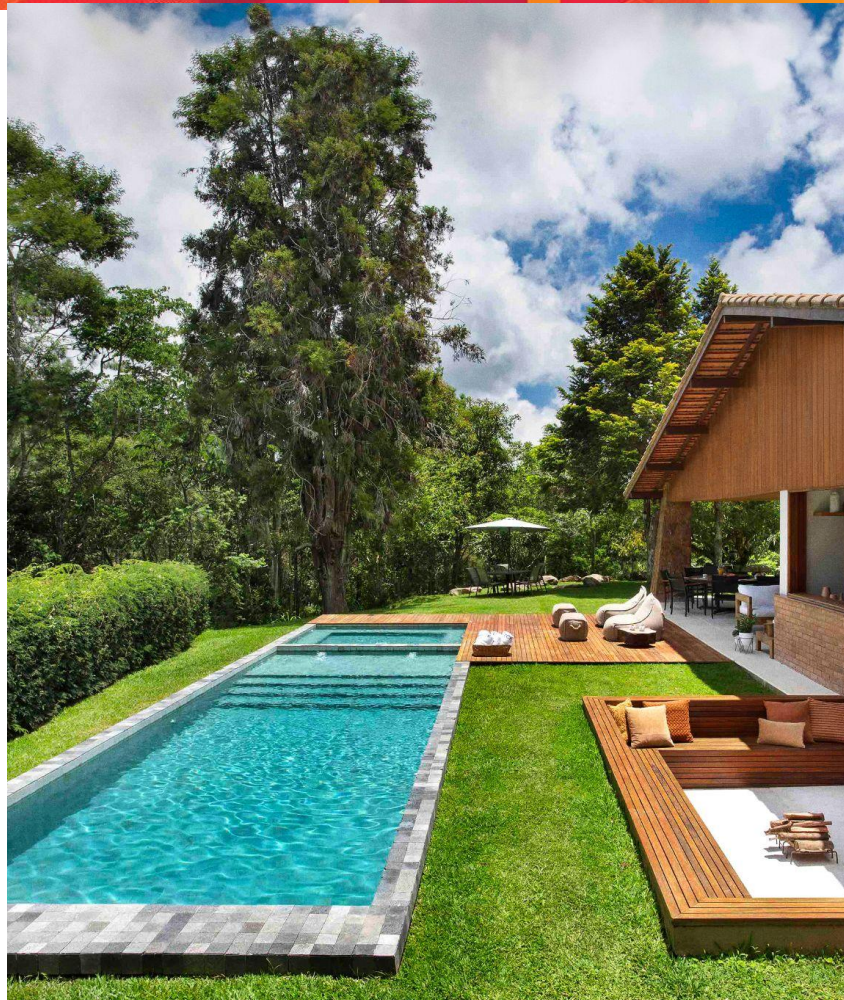
FEATURE STORE SUMMIT 2024

DATA FOR AI:
REAL-TIME, BATCH, AND LLMS



Chronon Way

- Single Python config file
- Same Scala compute engine for both envs
- Take care of the infra orchestration



Chronon Architecture

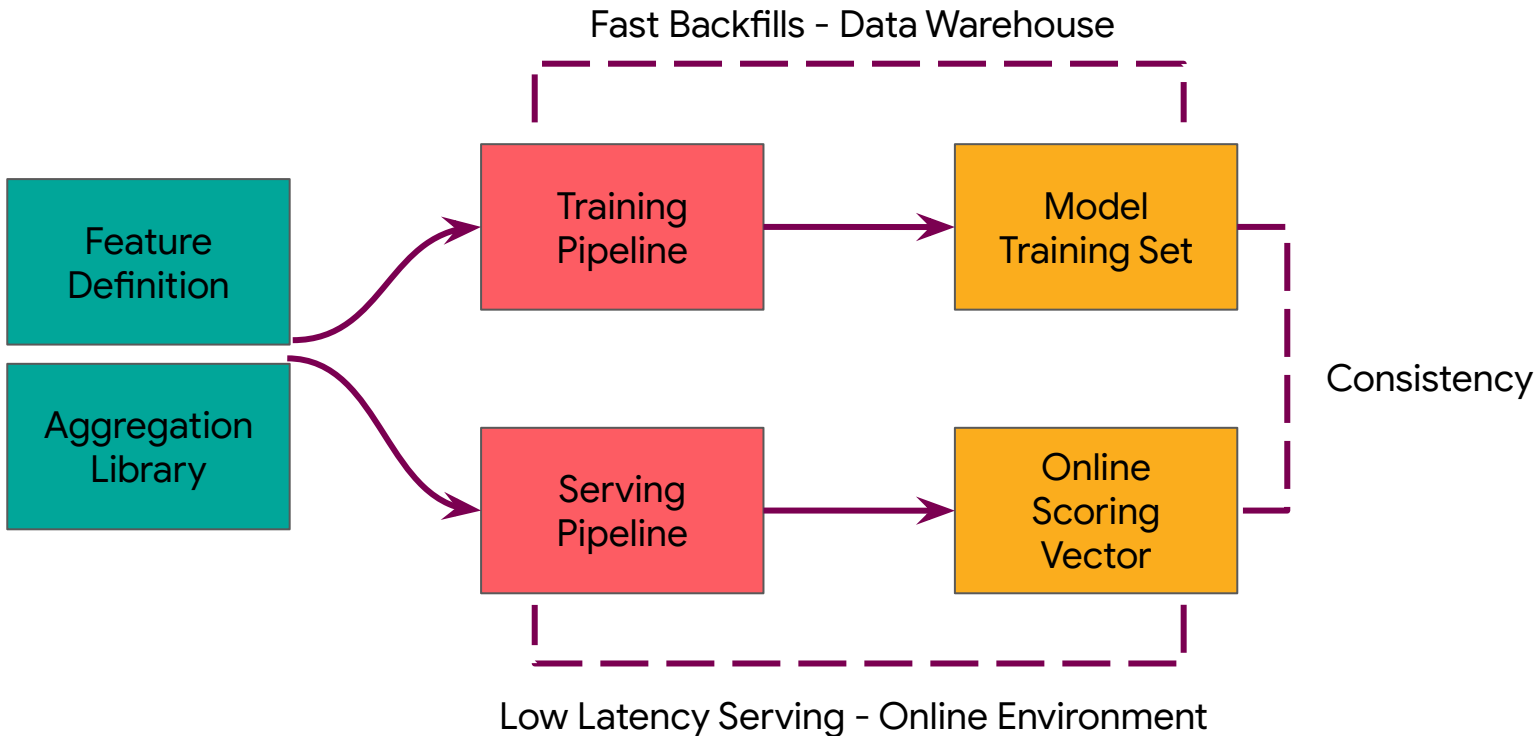


FEATURE STORE SUMMIT 2024

DATA FOR AI:
REAL-TIME, BATCH, AND LLMS



Feature Management System



Examples



FEATURE STORE SUMMIT 2024

DATA FOR AI:
REAL-TIME, BATCH, AND LLMS



Feature definition (in Python!)

```
v1 = GroupBy( —————> GroupBy: Aggregation over a single Source  
    source = source,  
    keys = ["user_id"],  
    aggregations = [Aggregation( —————> Aggregation: Operation + Windows  
        name = "page_view_sum",  
        operation = Operation.SUM, —————> Operation: functions i.e. SUM, AVG, LAST  
        windows = [  
            Window(length = 7, timeUnit = TimeUnit.DAYS), ——> Windows: for time-series data  
            Window(length = 14, timeUnit = TimeUnit.DAYS),  
            Window(length = 1, timeUnit = TimeUnit.MONTH),  
            Window(length = 1, timeUnit = TimeUnit.YEAR),  
        ],  
        inputColumn = "page_view_count"  
    )  
]
```



Joins

- **Join**

- What is the <aggregation statistic> for <key> from <time A> to <timestamp>?

- **Challenges**

- Data skew e.g. bots
- Computing with midnight accuracy can be faster

```
v1 = LeftOuterJoin(  
    left=leftSource, # (user_id, timesteps)  
    rightParts=[  
        JoinPart(  
            group_by=page_views.v1,  
            keyMapping={'user_id': 'host_id'},  
        ),  
        JoinPart(  
            group_by=profile_change.v2,  
        ),  
    ],  
)
```


Fun fact

One user achieved:
250 lines of code → 1000 features



FEATURE STORE SUMMIT 2024

DATA FOR AI:
REAL-TIME, BATCH, AND LLMS

New Functionalities



FEATURE STORE SUMMIT 2024

DATA FOR AI:
REAL-TIME, BATCH, AND LLMS



Derivations

- **Why?**
 - Computing acceptance ratio = $\text{accept_sum} / \text{requests_sum}$
- **Challenges**
 - Uniform API - online & offline
 - Online Latency

```
1  v1 = Join(  
2      # it supports group by level as well  
3      left=...,  
4      right_parts=[...],  
5      derivations=[  
6          Derivation(  
7              name="acceptance_ratio",  
8              expression="accept_sum /  
9                  requests_sum"  
10             ),  
11             ...  
12         ]  
13     )
```



Real time chaining features

- Why?

- Computing average price of ice-creams a user viewed in the last 14 days
- Require transformations and aggregations along with denormalization

- Challenges

- online & offline consistency
- Online Latency

```
1  group_by = GroupBy(  
2      name="enriched_ice_cream_info",  
3      sources = join_source,  
4      keys = ["user"],  
5      aggregations = [  
6          # creates an aggregate of last 100 ice  
7          # creams viewed in a 14 day window  
8          Aggregation(  
9              operation = Operation.LAST_K(100),  
10             input_column = "price",  
11             window = [Window(14, TimeUnit.DAYS)]  
12         )  
13     ],  
14     derivations = [  
15         Derivation(  
16             name = "avg_price_last_14d",  
17             expression = "aggregate  
18                 (price_last100_14d, 0, (acc, x) -> acc +  
19                 x) / size(price_last100_14d)"  
20         )  
21     ]  
22 )
```


Key Takeaways



FEATURE STORE SUMMIT 2024

DATA FOR AI:
REAL-TIME, BATCH, AND LLMS



Key takeaways:

- Provides declarative language to define features once for both online and offline features.
- Lambda architecture to achieve online offline consistency and low latency serving
- Windowed operation supported
- Operations: currently 15 aggregation operation
 - Can be extended easily

Call to actions



FEATURE STORE SUMMIT 2024

DATA FOR AI:
REAL-TIME, BATCH, AND LLMS



Call to actions:

- Star the repo and fork it:
<https://github.com/airbnb/chronon>
- Kick starter: report issues or ideas to git issues
- Discord: <https://discord.gg/GbmGATNqqP>

Additional Resources



FEATURE STORE SUMMIT 2024

DATA FOR AI:
REAL-TIME, BATCH, AND LLMS



Additional Resources:

- [Chronon — A Declarative Feature Engineering Framework](#)
- [Chronon, Airbnb's ML Feature Platform, Is Now Open Source](#)
- [Shepherd: How Stripe adapted Chronon to scale ML feature development](#)
- <https://www.chronon.ai/>
- <https://github.com/airbnb/chronon/>

Thank you



FEATURE STORE SUMMIT 2024

DATA FOR AI:
REAL-TIME, BATCH, AND LLMS