

OpenMLDB: An Open-Source Enterprise-Grade Real-Time Feature Platform

Mian Lu, PhD
OpenMLDB PMC Member
System Architect at 4Paradigm



Organized by  **HOPSWORKS**

Project Background

4Paradigm Inc.: offering platform-centric AI solutions

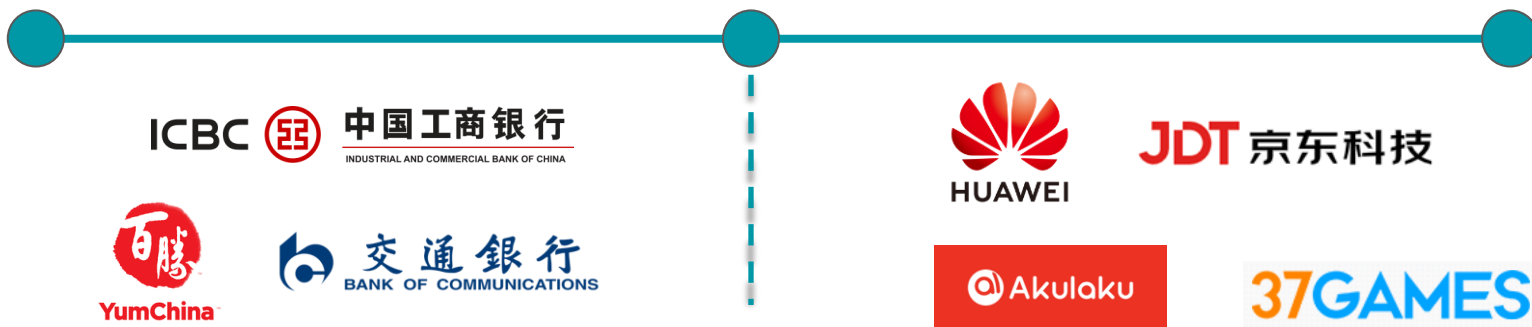
Sage: 4Paradigm's Enterprise-Grade MLOps Platform

OpenMLDB: used by Sage to define and compute real-time features; open-source in Jun 2021

2017.2: the first commit

2021.6: open-source v0.1

2022.10: v0.6 release

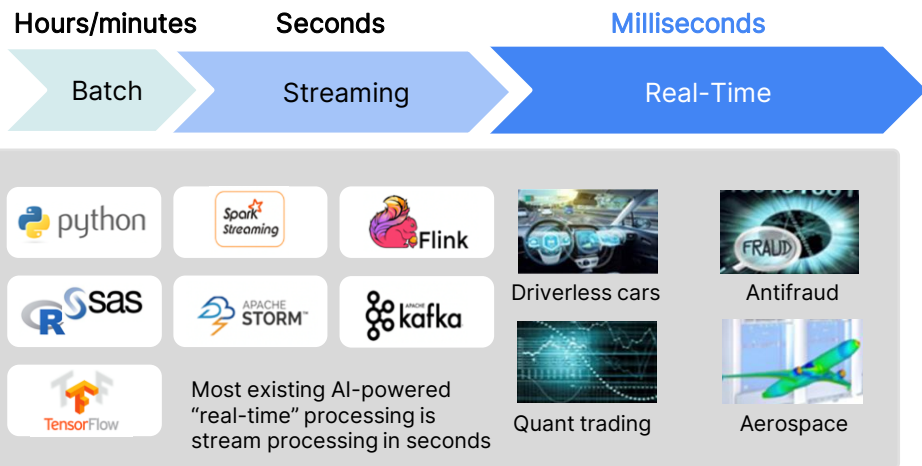


Real-Time Features for Real-Time Decision Making

REAL-TIME FEATURES: **real-time data** + **real-time processing**

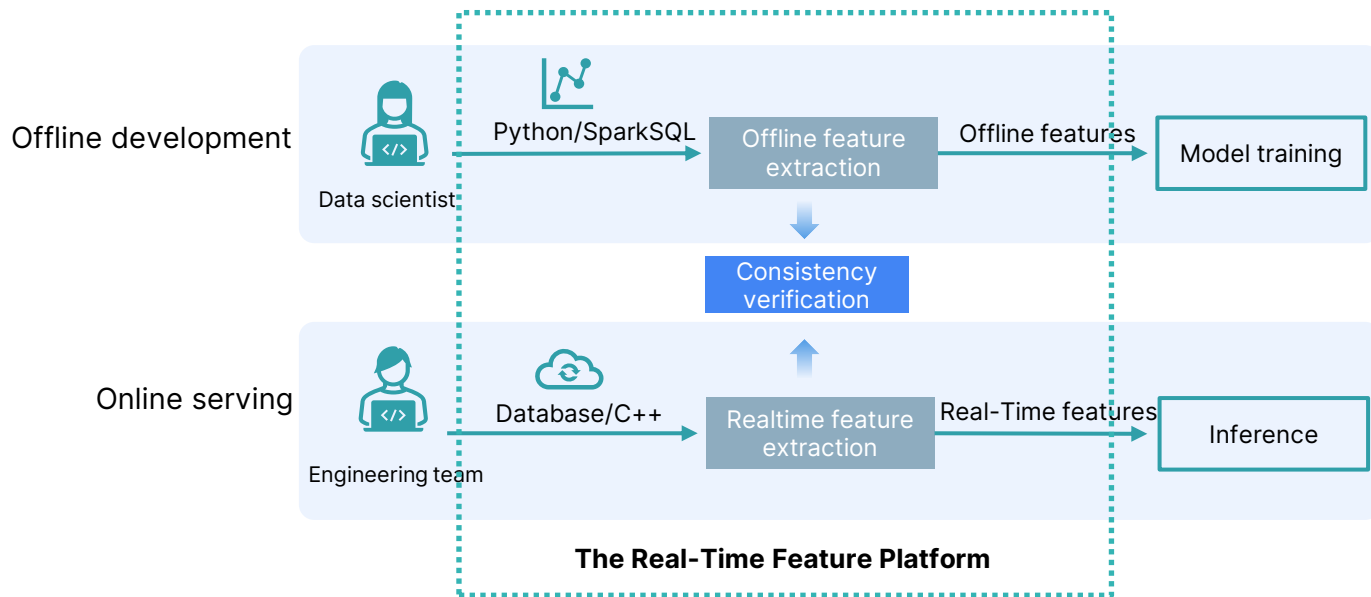
Real-time antifraud: computing real-time features in **milliseconds**

Solutions	Latency	Recall
Rule-based	~200ms	Low
In-house platform	~50ms	Medium
OpenMLDB	<20ms	High



Real-time decision-making in milliseconds is useful for a wide range of applications

Challenges of Computing Real-Time Features



Correctness:

Consistent computation logic between offline and online features, no data leakage

Efficiency:

Computing real-time features in milliseconds

OpenMLDB: An Open-Source Enterprise-Grade Real-Time Feature Platform

Three Steps

From development to deployment



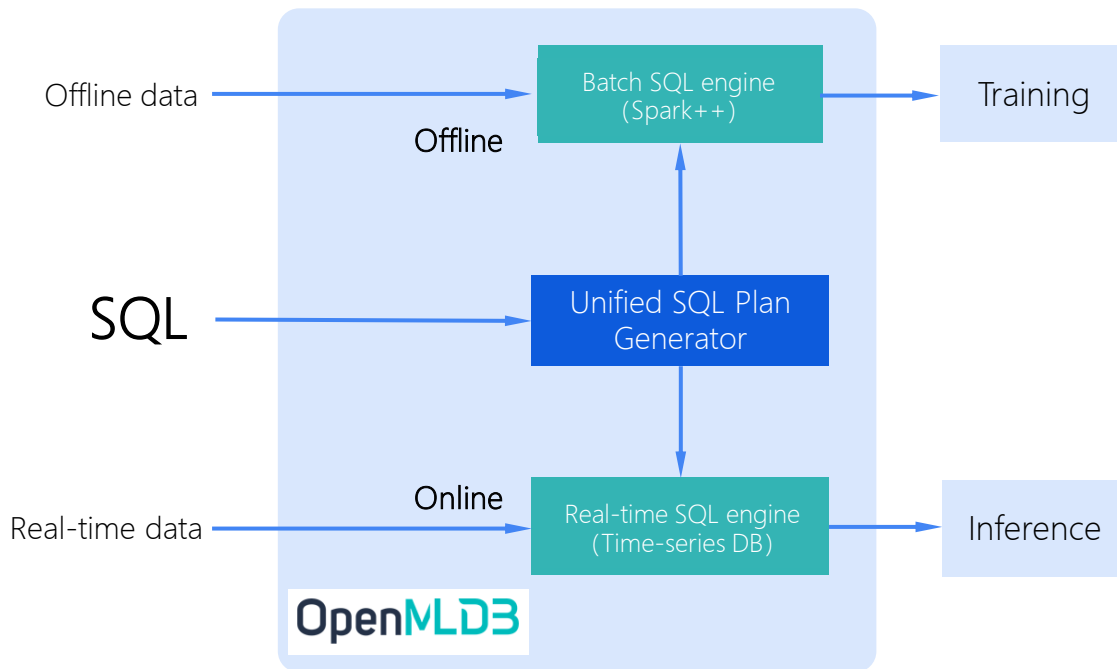
Offline development



Deployment to online



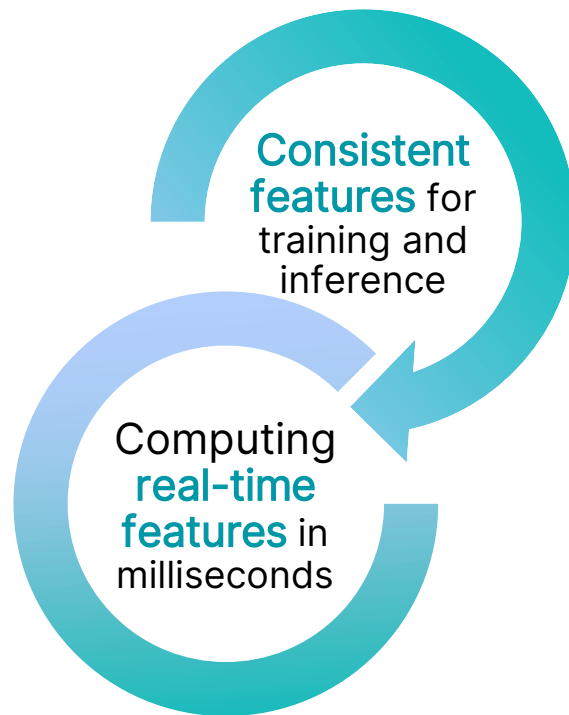
Real-time data input



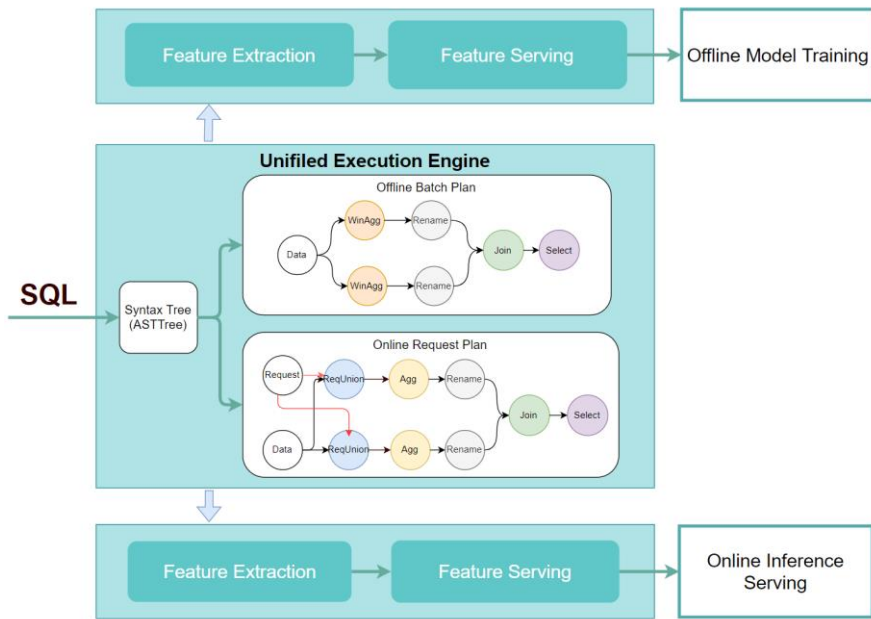
Tackling One Challenge, with One Distinguishing Feature

Challenge

Feature



Consistent Features for Training and Inference

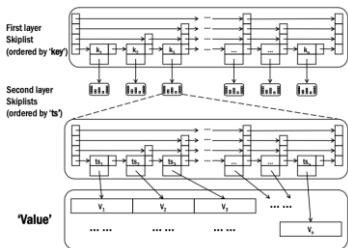


- Adaptable online and offline SQL translation from the logical plan to physical plan
- Unified computing libraries

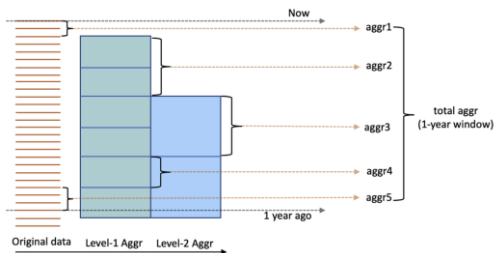


Inherent guarantee of
online-offline consistency,
no data leakage

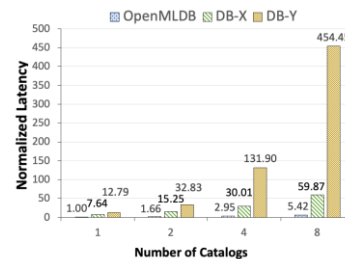
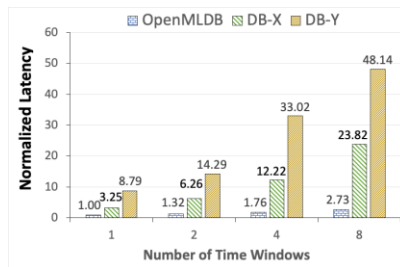
Computing Real-Time Features in Milliseconds



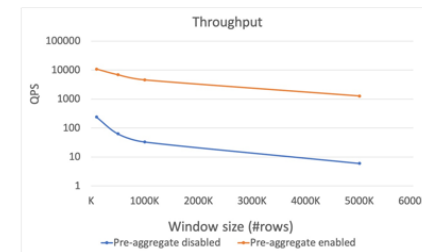
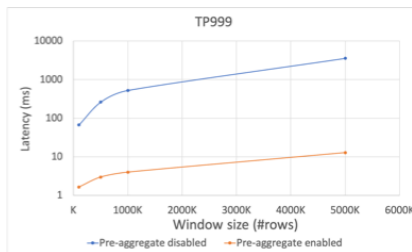
Key opt.: double-layer skiplist



Key opt.: pre-aggregation



Real-time perf: OpenMLDB vs commercial in-mem. DB

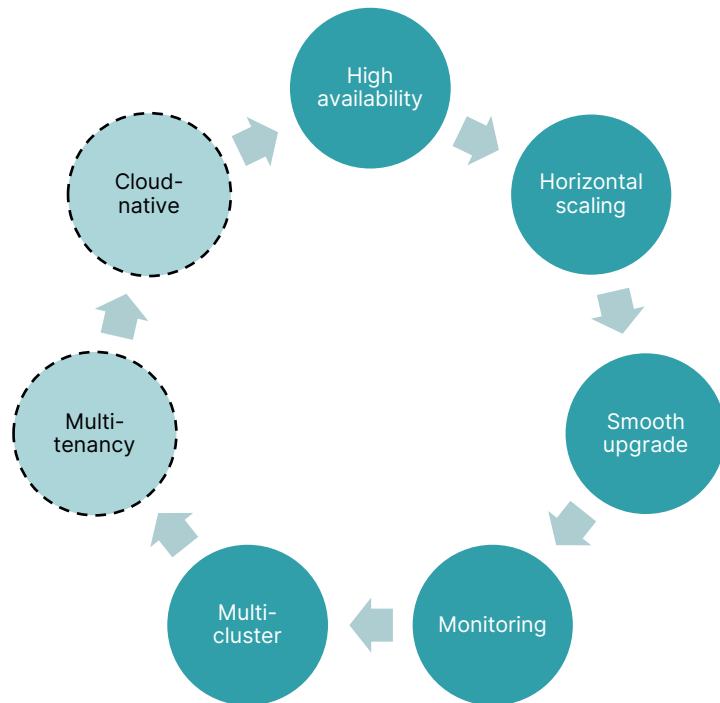


Performance improvement by pre-aggregation

Enterprise-Grade Features

Born for large-scale
enterprise applications

Deployment **in hundreds**
of real-world applications



Use Case: Akulaku's ML Platform to Produce Real-Time Features in 4 Milliseconds

Challenges of Features

- **Online serving:** low latency, high concurrency, features based on fresh data
- **Offline batch:** high throughput
- **Consistency:** consistency for online and offline

The OpenMLDB Solution

- **Scenario-driven:** Realtime features driven by scenarios
- **Solution:** 1) bridge the online and offline by SQL; 2) realtime features based on sliding window

```
(product_id:1, price:2, .update_time:100L )
(product_id:1, price:3, .update_time:200L )

(product_id:2, price:3, .update_time:1000L )
(product_id:2, price:4, .update_time:1201L )
(product_id:2, price:3.5, .update_time:1100L )
```



```
SELECT COUNT(*)
FROM w100ms
```

```
WINDOW w100ms AS
(PARTITION BY product_id ORDER BY update_time
ROWS_RANGE BETWEEN 100ms PRECEDING
AND
CURRENT ROW)
```

Business Scenarios

- **Scenario:** realtime processing for orders within a day
- **Data:** 1 billion orders per day
- **Requirement:** realtime features with sliding window
- **Latency:** 4 milliseconds



Welcome to Join the OpenMLDB Community!

OpenMLDB

<https://github.com/4paradigm/OpenMLDB>



Slack Workspace
[contact@openmldb.ai](https://openmldb.slack.com)

Contact me: (LU Mian) lumian@4paradigm.com