

基于 Apache Doris 的 数据湖仓与数据治理实践

沈杨

Cisco WebEx 数据治理平台架构师



目录

01 Webex 业务及数据平台介绍

02 Webex 的数据湖仓架构演化

03 基于 Doris 的数据湖仓实践

04 Doris 与 Webex 数据治理平台的融合

Apache Doris x 讯飞开源联合 Meetur

Apache Doris x 讯飞开源联合 Meetur

联合 Meetur

Apache Doris x 讯飞开源联合 Meetur

Apache Doris

Apache Doris x 讯飞开源联合 Meetur

Apache Doris x 讯飞开源联合 Meetur

Apache Doris x 讯飞开源联合 Meetur

Apache Doris x 讯飞开源联合 Meetur

Apache Doris x 讯飞开源联合 Meetur

Apache Doris x 讯飞开源联合 Meetur

Apache Doris x 讯飞开源联合 Meetur

Apache Doris x 讯飞开源联合 Meetur

Apache Doris x 讯飞开源联合 Meetur

Apache Doris x 讯飞开源联合 Meetur

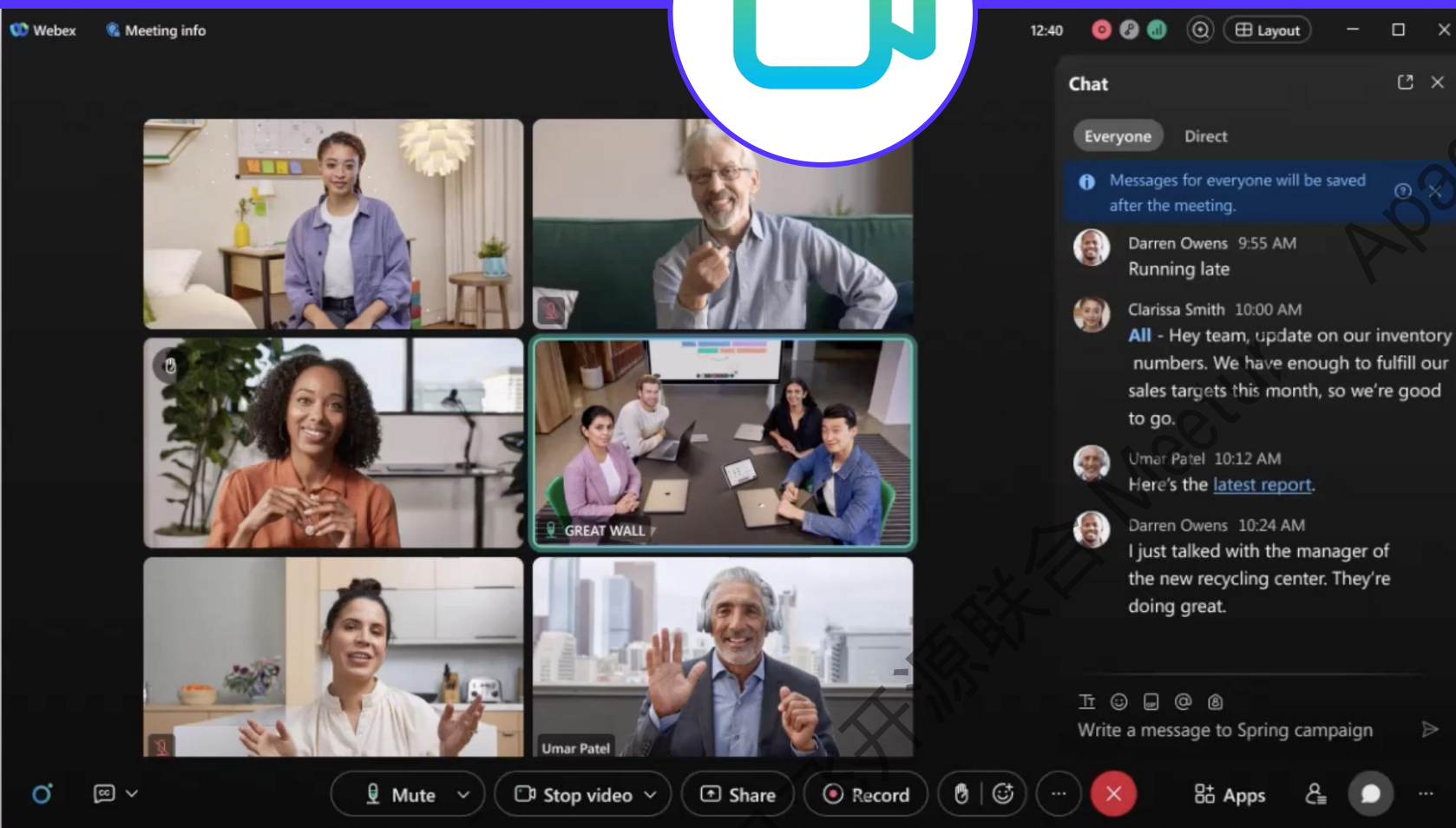
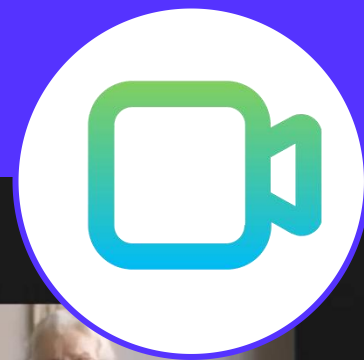
Apache Doris x 讯飞开源联合 Meetur

Apache Doris x 讯飞开源联合 Meetur

WebEx 业务介绍

Meeting

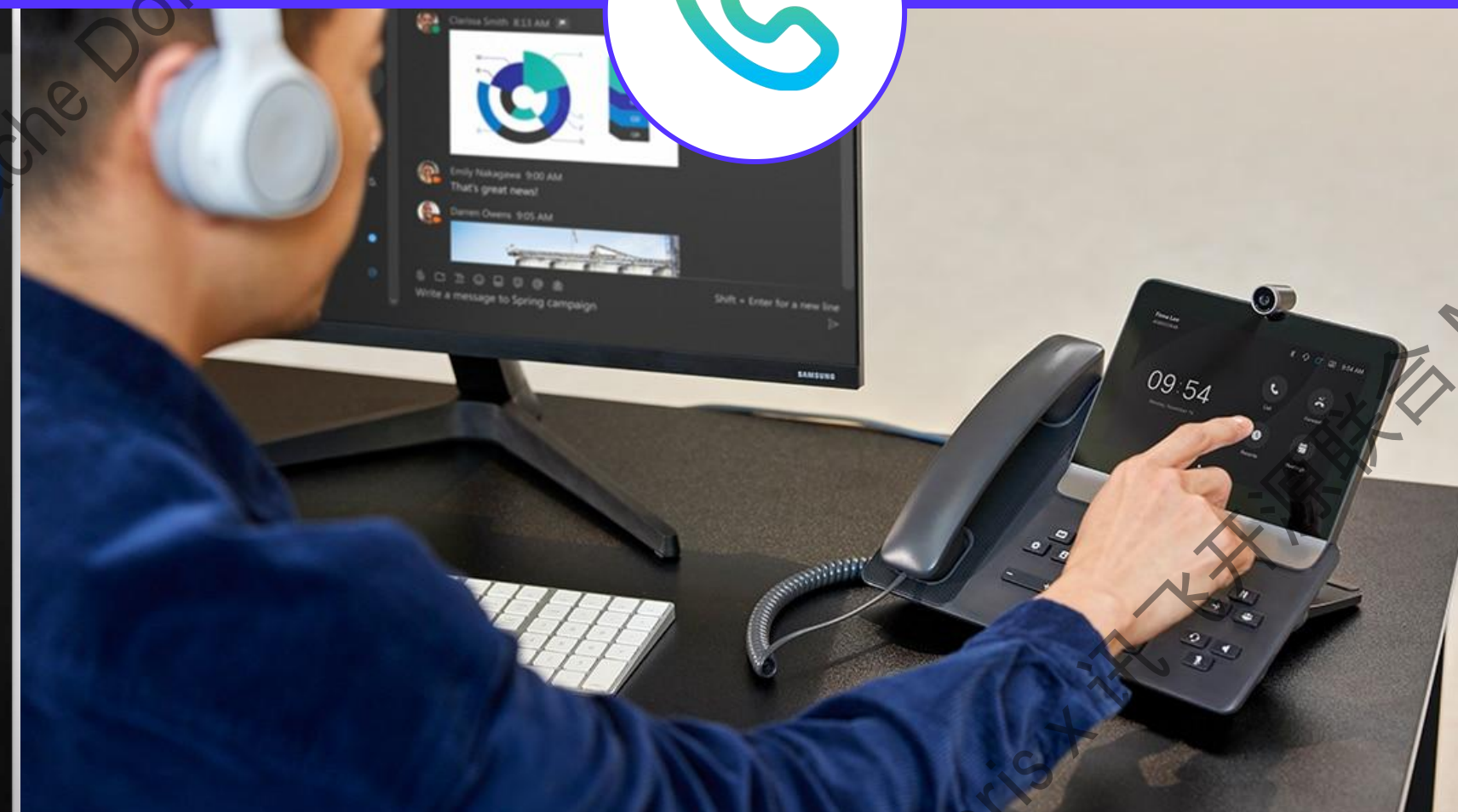
AI powered video conferencing and screen sharing platform.



95%
of Fortune 500 companies

Calling

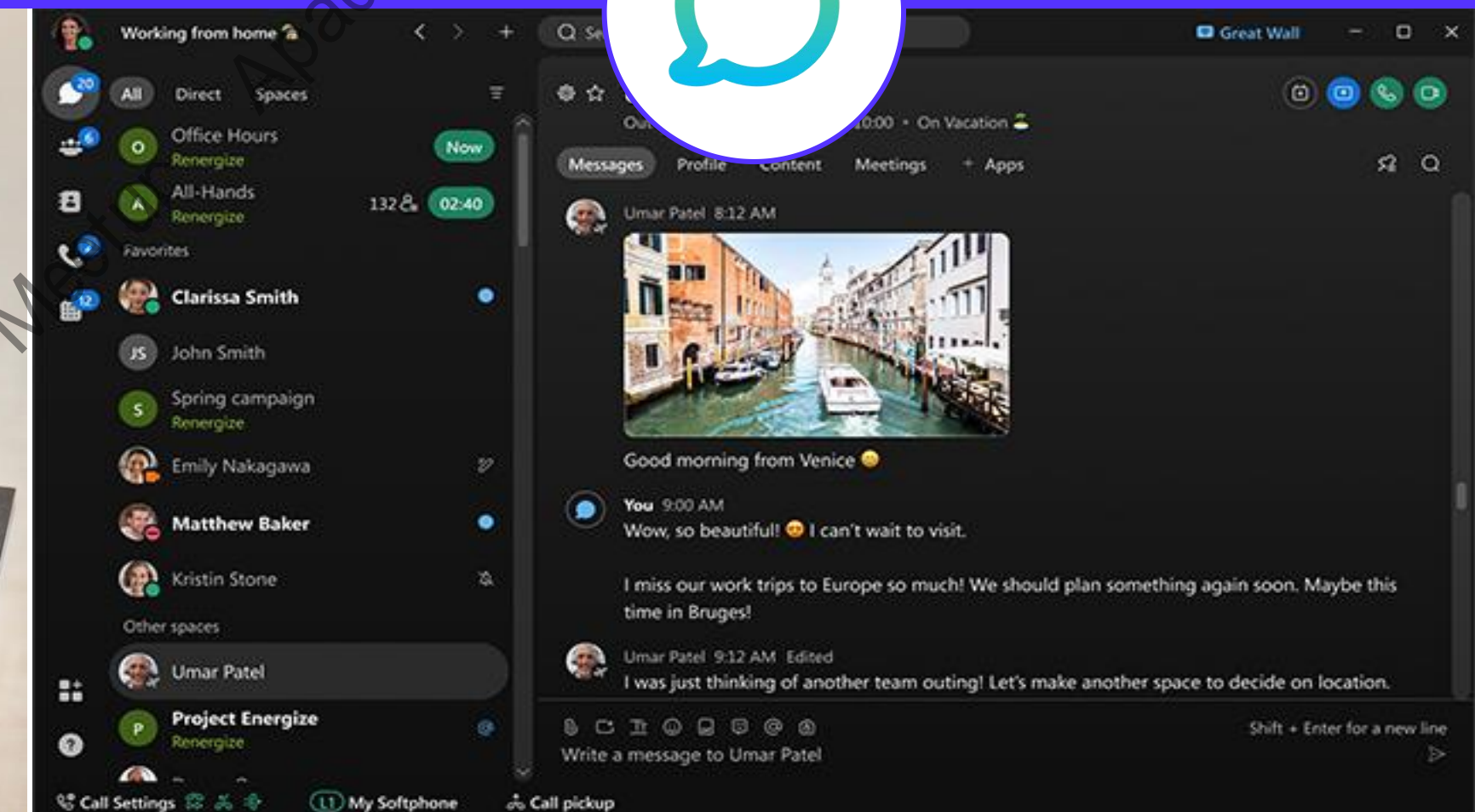
A modern and complete business cloud calling and phone system.



15000000
meetings per day

Messaging

Group messaging, chat, and file sharing that make collaboration engaging and effortless.



160
markets worldwide.

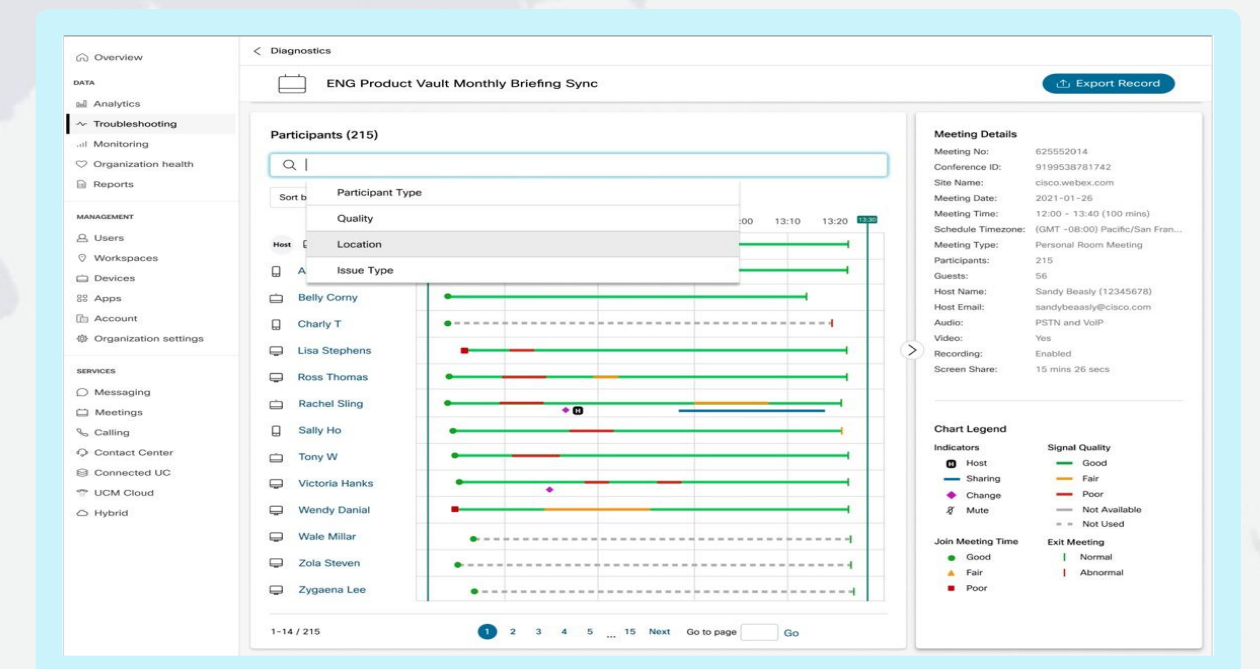
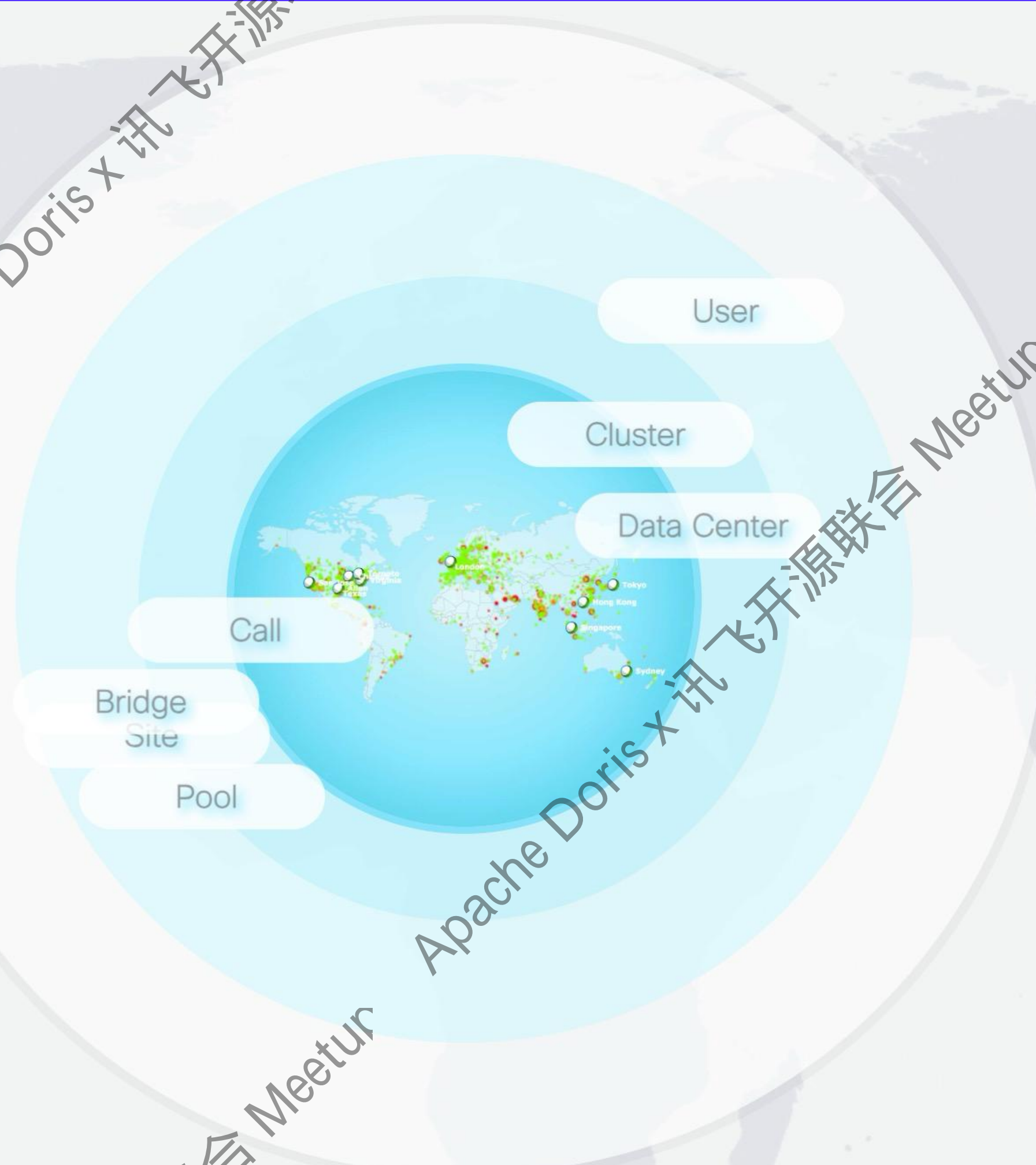
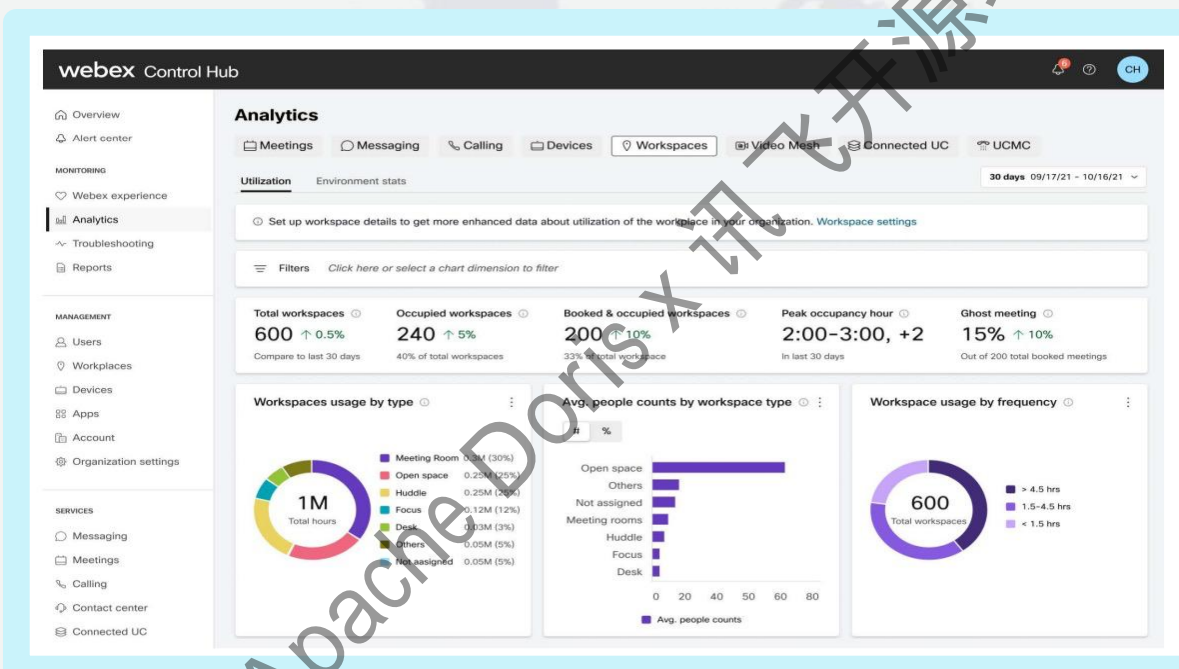
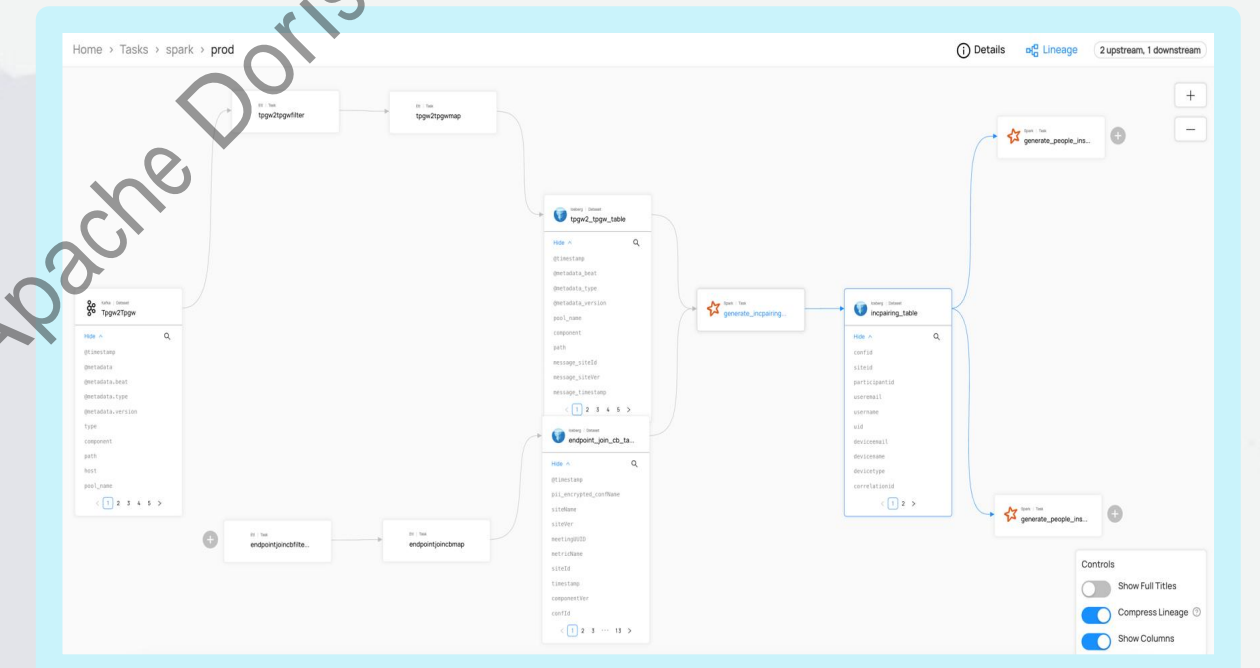
WebEx 数据平台介绍

Troubleshooting & Diagnostic

Batch & Real-time Analytics

BI & LLM Applications

Observability & Data Governance



WebEx 数据平台架构


Business People


Data Engineers


Data Analyst


Data Scientist

WAP Data Workbench

Sources

Webex App
(Metrics-A)

Client
(Telemetry)

Application Logs

Calling

Contact Center

More ...

Data Ingestion and
Processing

Kafka

ETL Jobs

Flink Jobs

Spark Jobs

Kyuubi

UDP

Data Lakehousing

Pinot

TiDB

ADS

Iceberg/Paimon

ODS/DWD/DWS

HDFS/S3

Storage

Data Governance

Data Catalog

Metadata

Lineage and
Monitoring

Quality

Security and
Access

Retention
Management

MetaHub

Data Applications

ControlHub

ROMA

Notebook

Report

Dashboard

Governed Data Platform

平台规模

Cluster

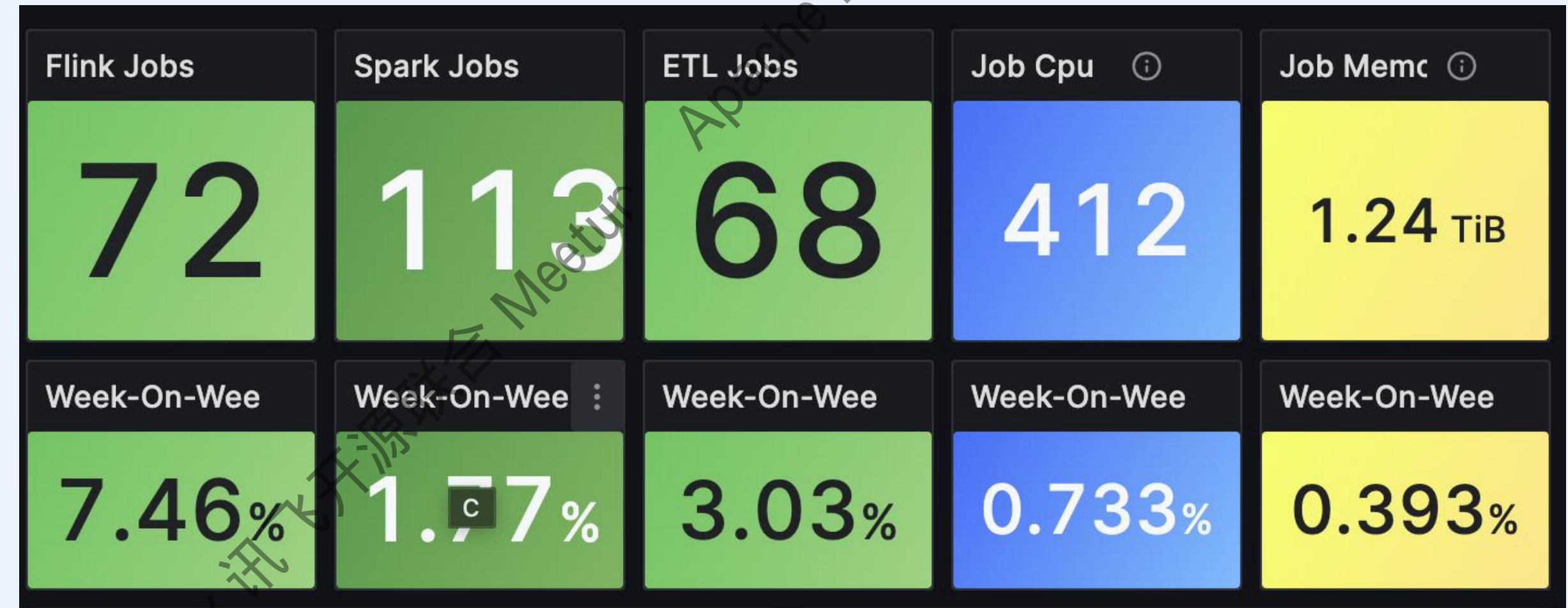
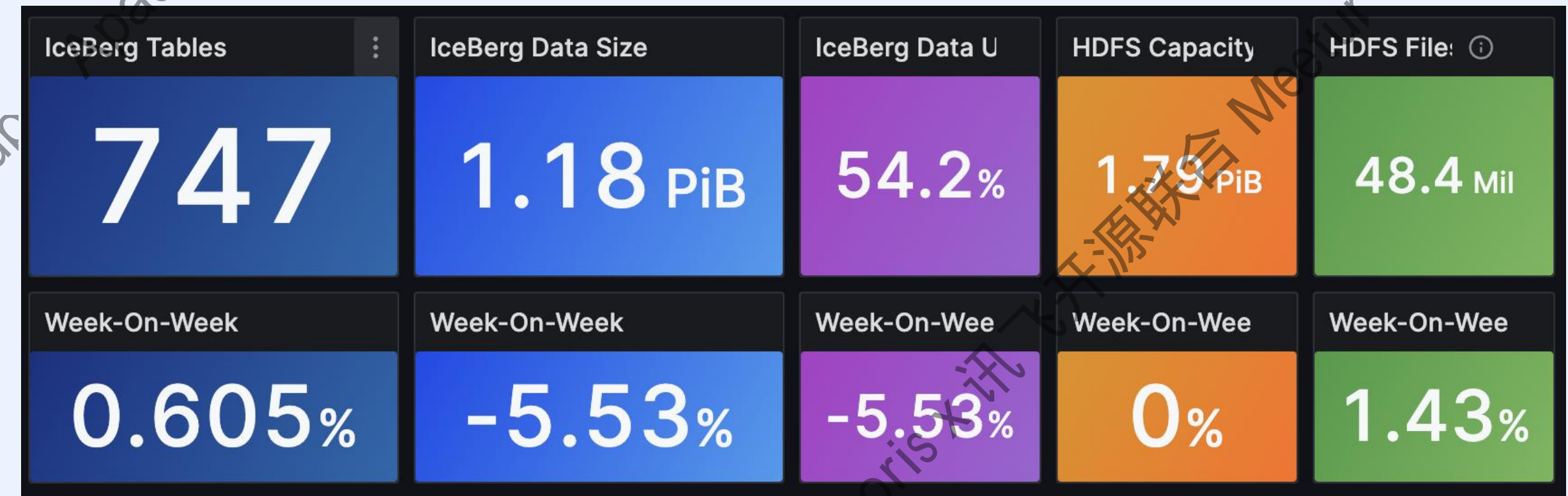
- 全球 15 个数据中心
- 每个数据中心下若干 cluster
- 自有机房与 AWS 公有云 hybrid 架构

Jobs

- 离线批处理任务（Spark）各集群累计仅 1000 余个
- 实时处理任务（Flink/KStreams）各集群累计仅 1000 余个

Data

- 每天超过 150 万会议
- 每天超过 10 亿分钟视频会议时长
- 每天 PB 级别原始数据规模
- 每天超过 10TB 数据接入数据湖仓



目录

01 Webex 业务及数据平台介绍

02 Webex 的数据湖仓架构演化

03 基于 Doris 的数据湖仓实践

04 Doris 与 Webex 数据治理平台的融合

Apache Doris x 讯飞开源联合 Meetur

Apache Doris x 讯飞开源联合 Meetur

联合 Meetur

Apache Doris x 讯飞开源联合 Meetur

Apache Doris x 讯飞开源联合 Meetur

WebEx数据湖仓架构演化

传统架构阶段
存储: RDBMS + NoSQL
计算: 存储过程、脚本、定时程序

转型阶段

Governed Data Platform

探索 Apache Doris 在数据湖仓中的应用

2018

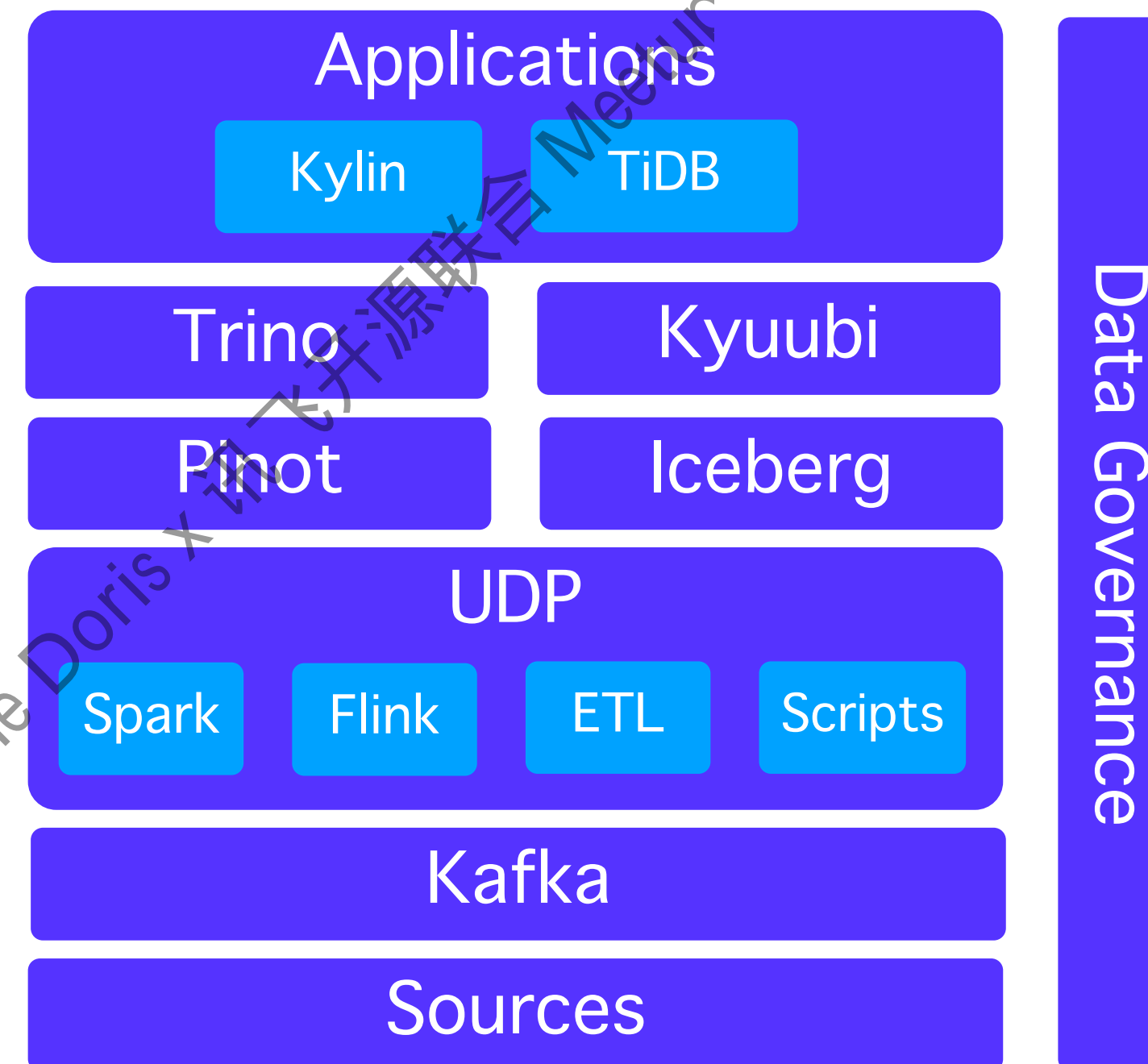
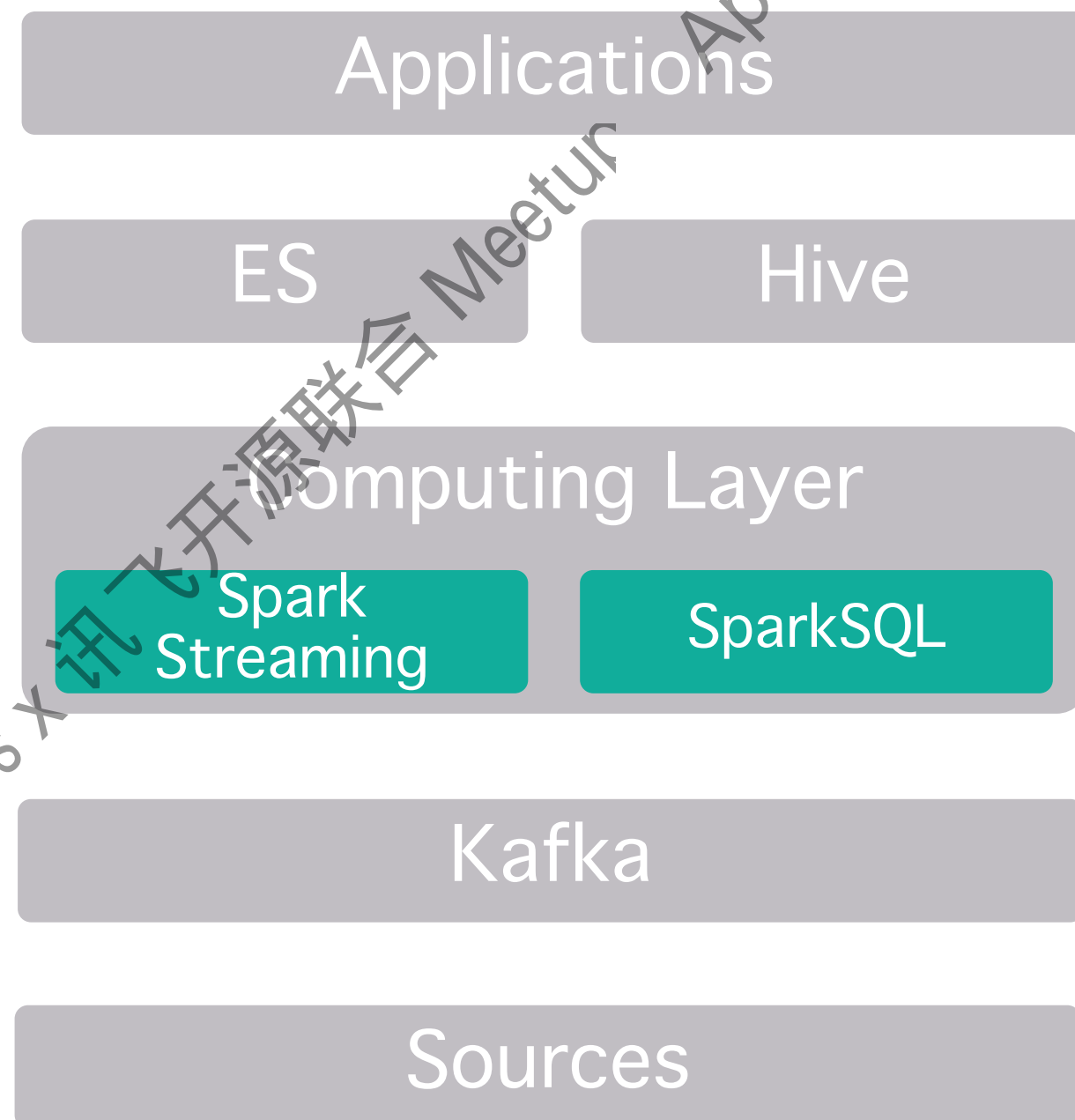
2019

2020

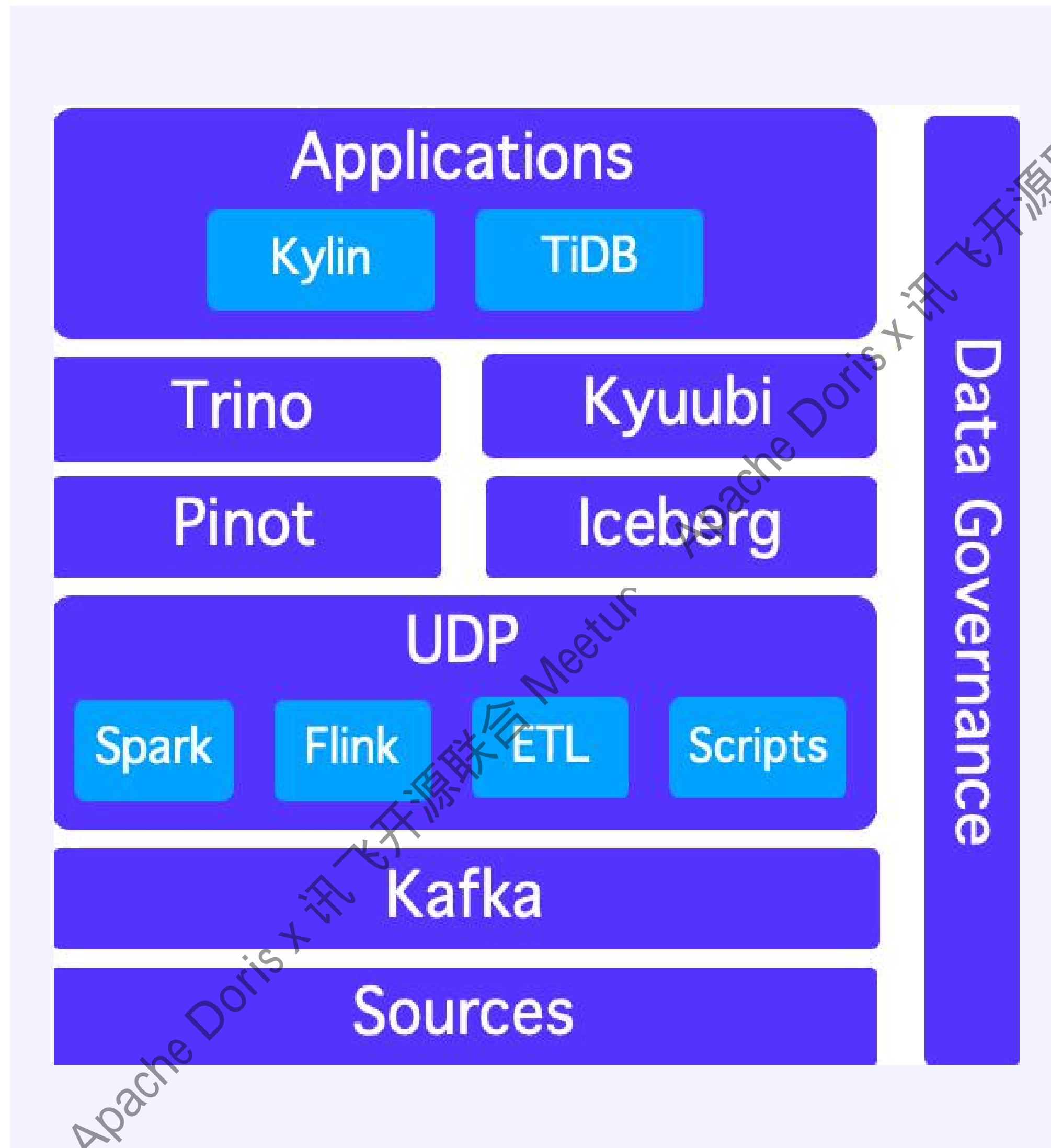
2021

2022

2024



引入 Doris 前的架构



Sources

- 数据源种类多，客户端、服务端metrics，业务数据库等
- 需要进入数据湖中的数据配置Kafka通道进行发送

Kafka

- 贴源层数据摄取统一走Kafka
- 多数据中心之间通过MirrorMaker将必要数据集中化

UDP

- 计算资源管理、隔离与任务调度
- 支持多种类型的任务

数据湖仓存储与建模

- Iceberg
- Pinot

查询与分析

- Trino
- Kyuubi
- SparkSQL

应用层

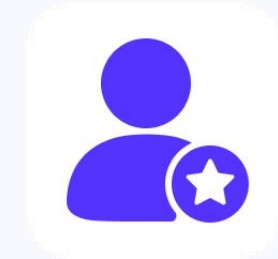
- ADS层可能会应用自建（Kylin/TiDB）

痛点与需求



架构复杂

- 多种 OLAP 型数据库方案并存
- 运维难度大



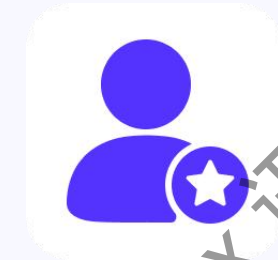
数据治理难度大

- 元数据来源多
- 数据质量问题易扩散难管控



资源利用率低

- 数据冗余程度高 (Storage 成本高)
- 查询入口不统一, 无法做到统一调度 (CPU/Mem)



数据一致性问题频发

- 各种数据仓库中的数据单独计算, 计算口径不一致
- 数据用户抱怨多

目录

01 Webex 业务及数据平台介绍

02 Webex 的数据湖仓架构演化

03 基于 Doris 的数据湖仓实践

04 Doris 与 Webex 数据治理平台的融合

Apache Doris x 讯飞开源联合 Meetur

Apache Doris x 讯飞开源联合 Meetur

Apache Doris x 讯飞开源联合 Meetur

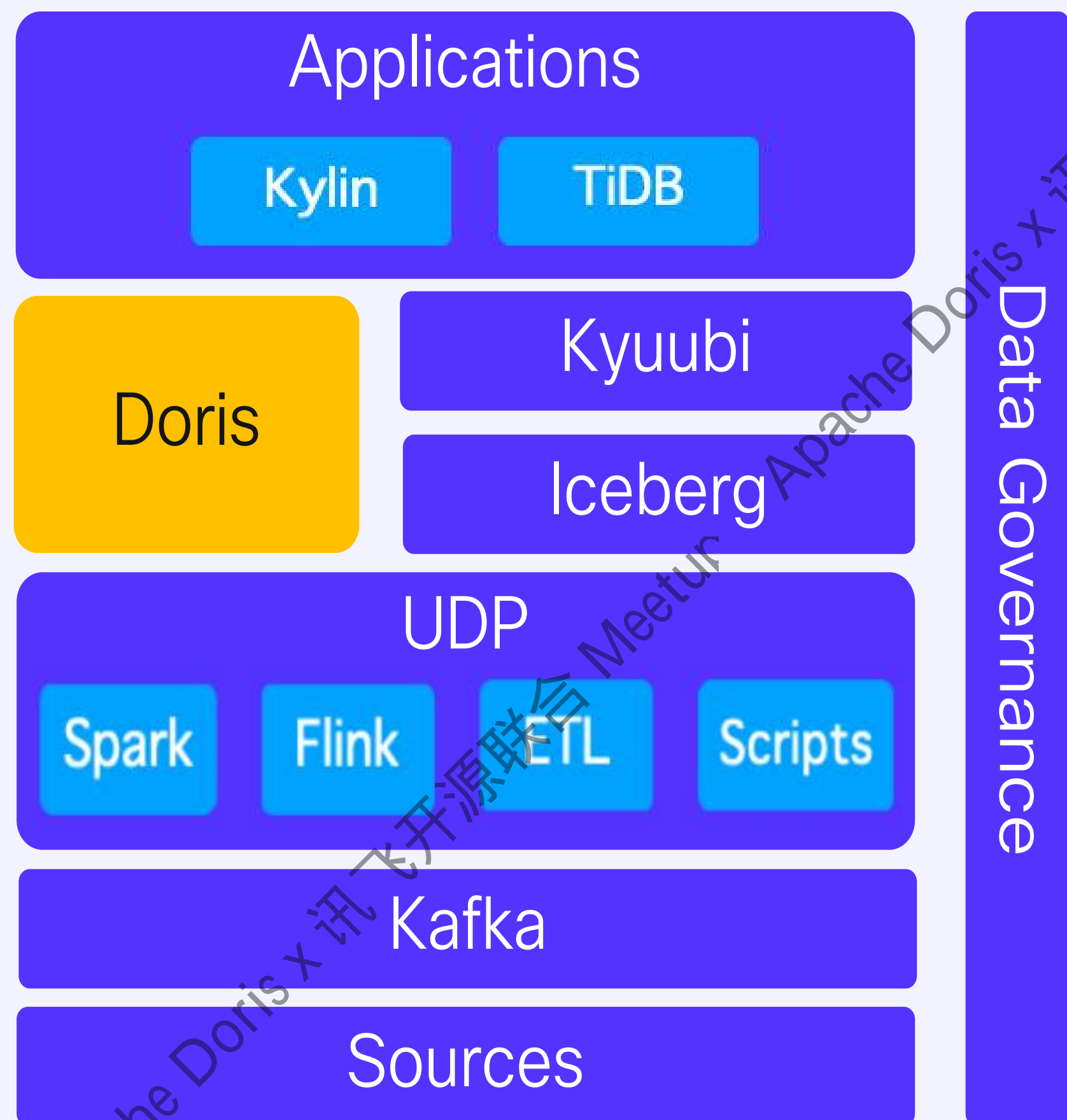
Apache Doris x 讯飞开源联合 Meetur

Apache Doris 初识

| 关键指标 | Apache Doris |
|------|--|
| 查询性能 | 存储层性能优秀，TPC、SSB 测试及业务场景测试均满足需求，1.1 之后支持向量化执行，支持查询优化器等，有更多的性能提升 |
| 时效性 | 支持 Routine Load、Stream Load，并且有事务支持，让实时导入数据更高效、安全 |
| 功能 | 支持 Aggregate、Unique、Duplicate 数据模型，让业务有更多选择空间。 并且支持物化视图，让数据聚合更轻松 2.0 支持倒排索引，让全文检索更容易 2.1 版本后支持 VARIANT 类型，对 JSON 操作更简单、高效 FE 可扩展性强 |
| 使用方式 | 兼容 MySQL 协议，学习成本基本为 0；接入和迁移成本低 |
| 运维部署 | 架构简单，不依赖第三方组件，扩缩容简单 数据自动平衡到各 BE，运维更省心 支持 zstd 等高效压缩方式，更节省存储资源 |
| 开源方式 | Apache License 开源协议，安全性高、灵活性强 |
| 其它 | 文档、资料丰富；社区活跃度高 |

基于 Apache Doris 的新架构

新架构：



变化：

Kafka

- Doris直接摄取Kafka topic数据

UDP

- 部分计算任务重构
- 退役一部分任务

数据湖仓存储与建模

- Doris取代Pinot
- 一部分数据模型迁移到Doris

查询与分析

- Doris 使用 Doris 取代 Trino

应用层

- 部分业务迁至Doris

收益：

简化架构（降本）

- 原Pinot和Trino两个集群由Doris替代
- 提高了资源利用率节省了30%资源
- 减少了系统运维的工作量

提升效率（增效）

- 提升了新集群部署效率
- 部分业务的数据流水线得到了简化

进化创新

- 重新思考业务数据建模方式
- 为更大范围的精简架构带来可能

Doris 集群情况

4

集群总数
各机房做冗余备份

30+

总节点数
独立部署 FE/BE

100000+

在线服务
平均每日查询总数量

40TB+

存储总量
按照业务设置不同的 retention

5TB+

实时
导入的日增数据量

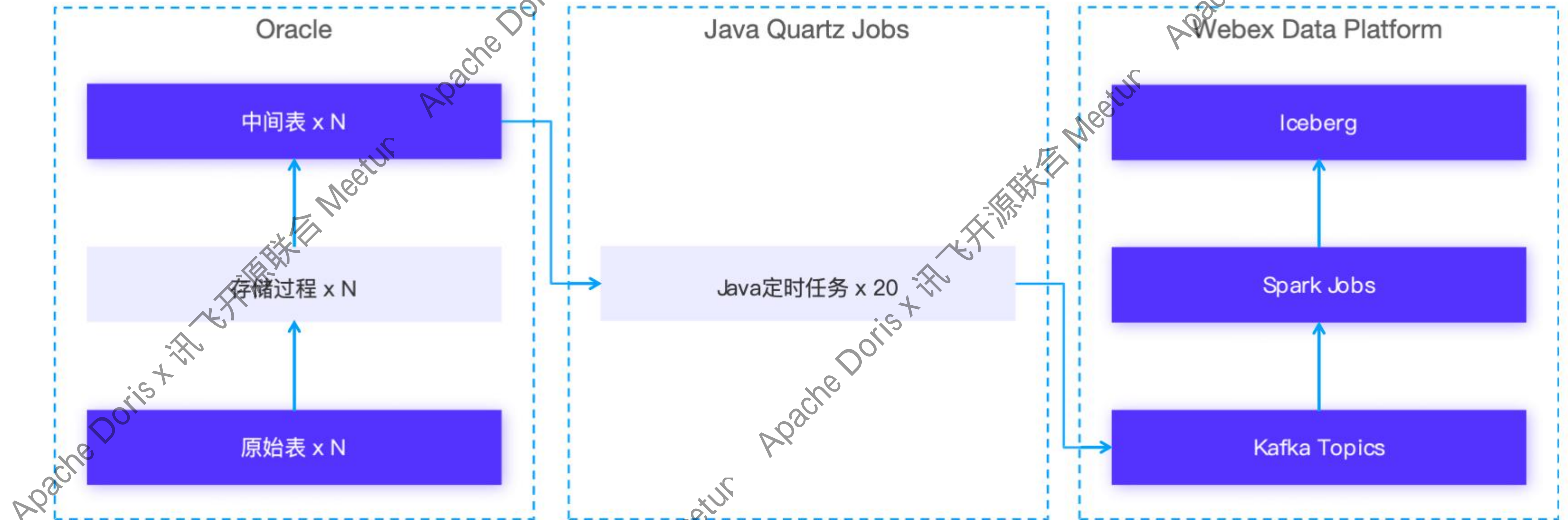
Apache Doris x 讯飞开源联合 Meetur

基于 Doris 的业务场景

基于 Doris 的 CCA Peak Ports 报表业务：

CCA Peak Ports 项目，用于在基于 Peak Ports 套餐的计费模式下，计算 Webex 与 Partner 之间的对账报表。该项目历史久远，并且长期以来未进行过大的技术革新，所采用的技术栈较为传统。此项目大小适中，数据源清晰，被选为第一批基于 Doris 改造的项目之一。

原方案：



基于 Doris 的业务场景

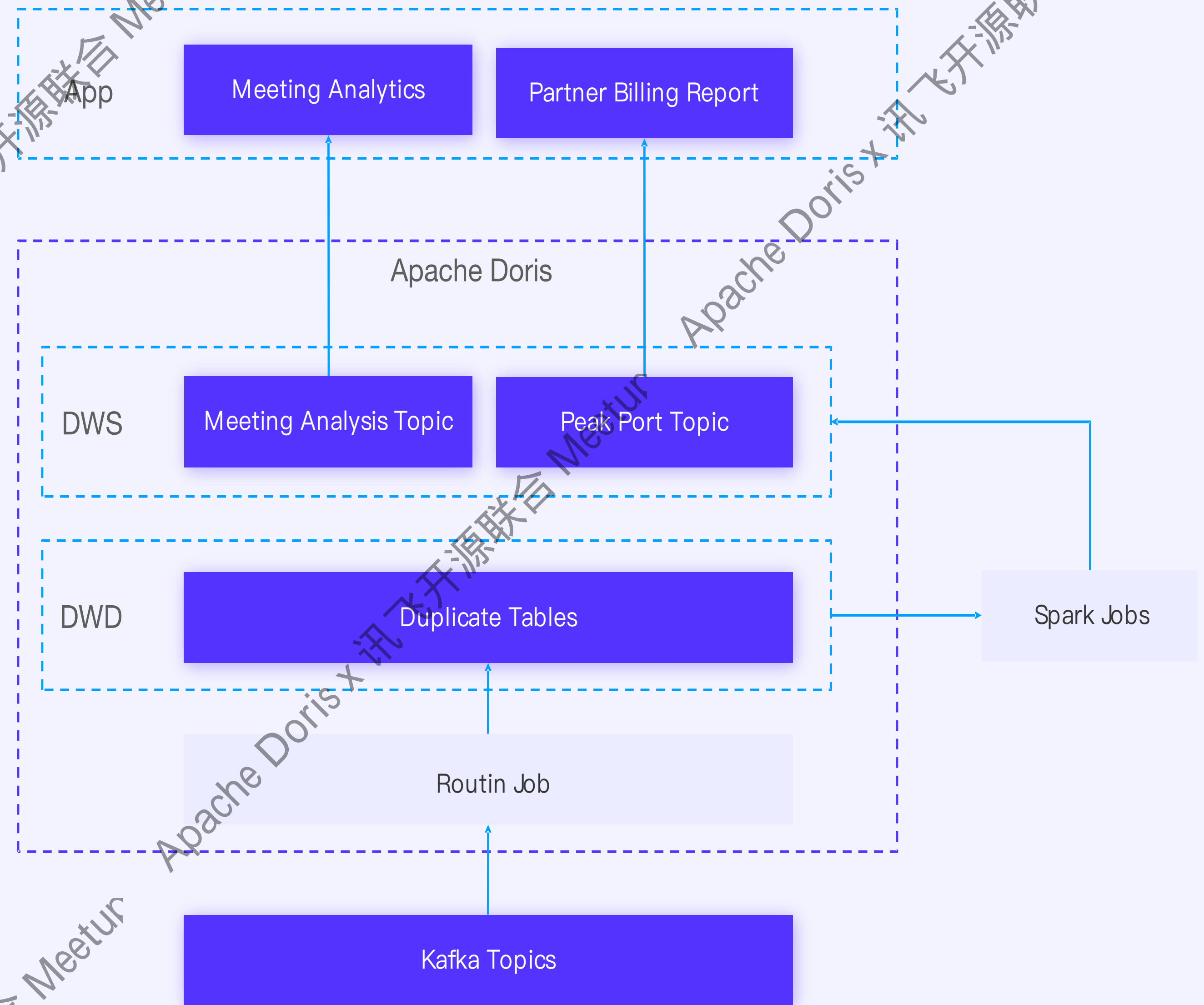
模型：

| peak_ports (Unique/Aggregate) | |
|-------------------------------|---------------------------|
| date | key/partition |
| confID | key / DISTRIBUTED BY HASH |
| siteID | key |
| orgID | key |
| peakTime | |
| peakValue | |

| meetings (Unique/Aggregate) | |
|-----------------------------|---------------------------|
| confID | key / DISTRIBUTED BY HASH |
| meetingID | key |
| siteID | key |
| orgID | key |
| startTime | partition |

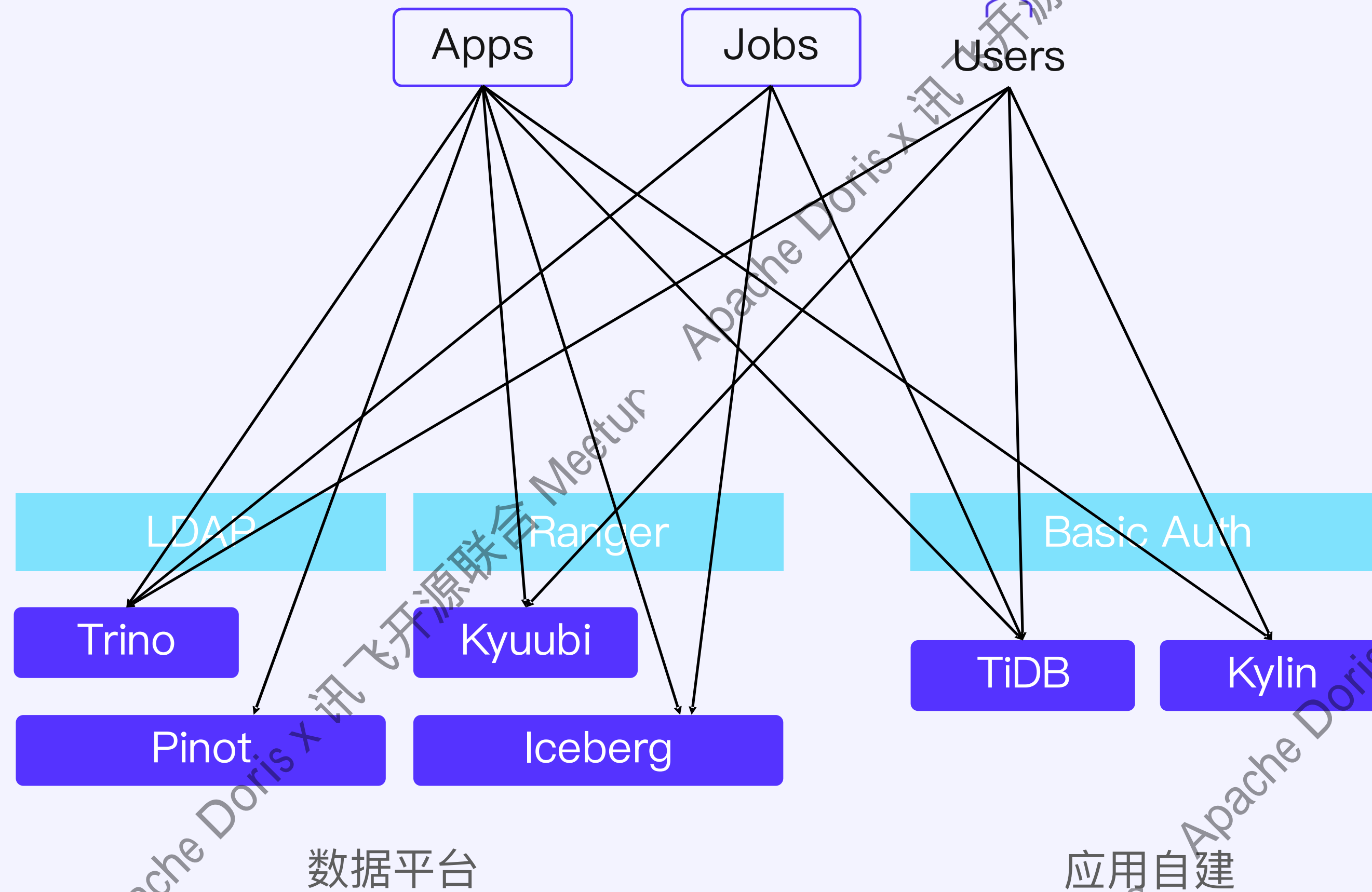
| *_raw_events (Duplicate) | |
|--------------------------|---------------------------|
| timestamp | key/partition |
| confID | key / DISTRIBUTED BY HASH |
| eventType | |
| userID | |
| ... | |

新方案：基于Doris的方案更加现代化，免去了存储过程和Java定时任务的维护工作，整体技术栈的技术风格更为统一。



基于 Doris 的统一查询引擎

引入 Doris 前的数据访问



存在的问题

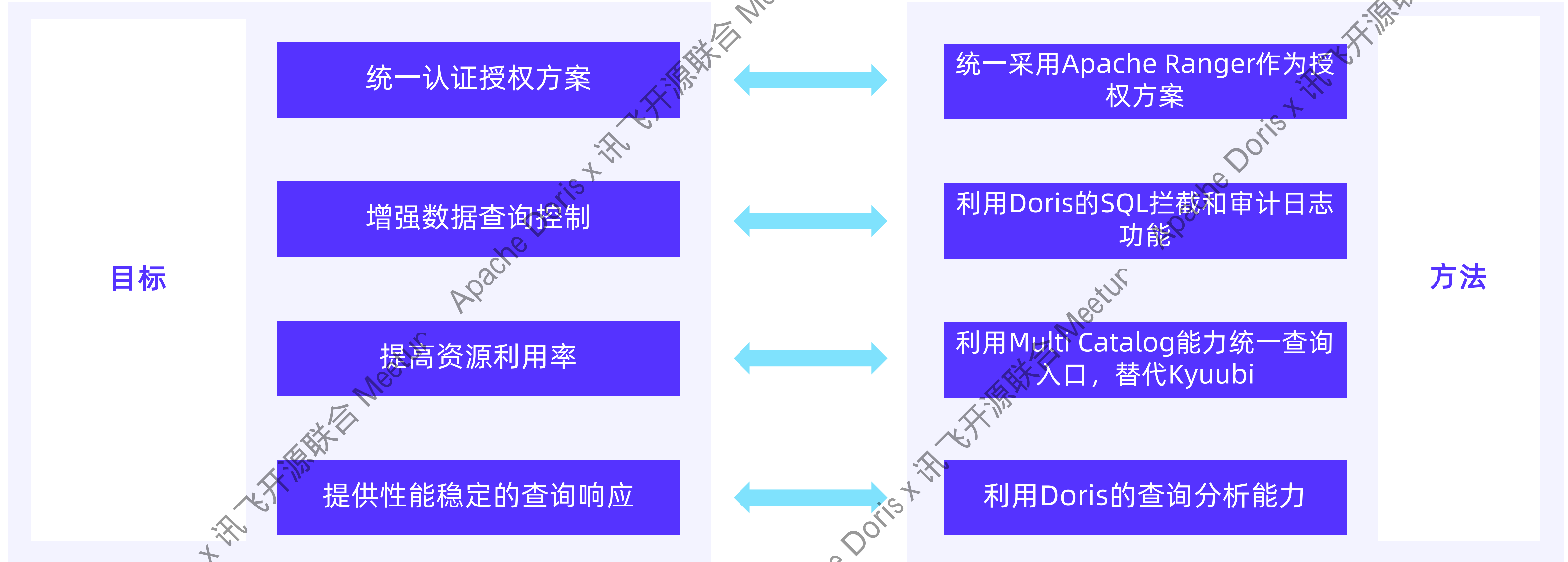
用户角度

- 链接管理复杂，需要管理多个系统的链接信息
- 密码更新频率不统一且密码更新时需要应用停机修改
- 认证授权申请繁琐，不同系统需要分别申请权限
- 查询语法有差异
- 性能体验差异大
- 应用稳定性难以保证

平台角度

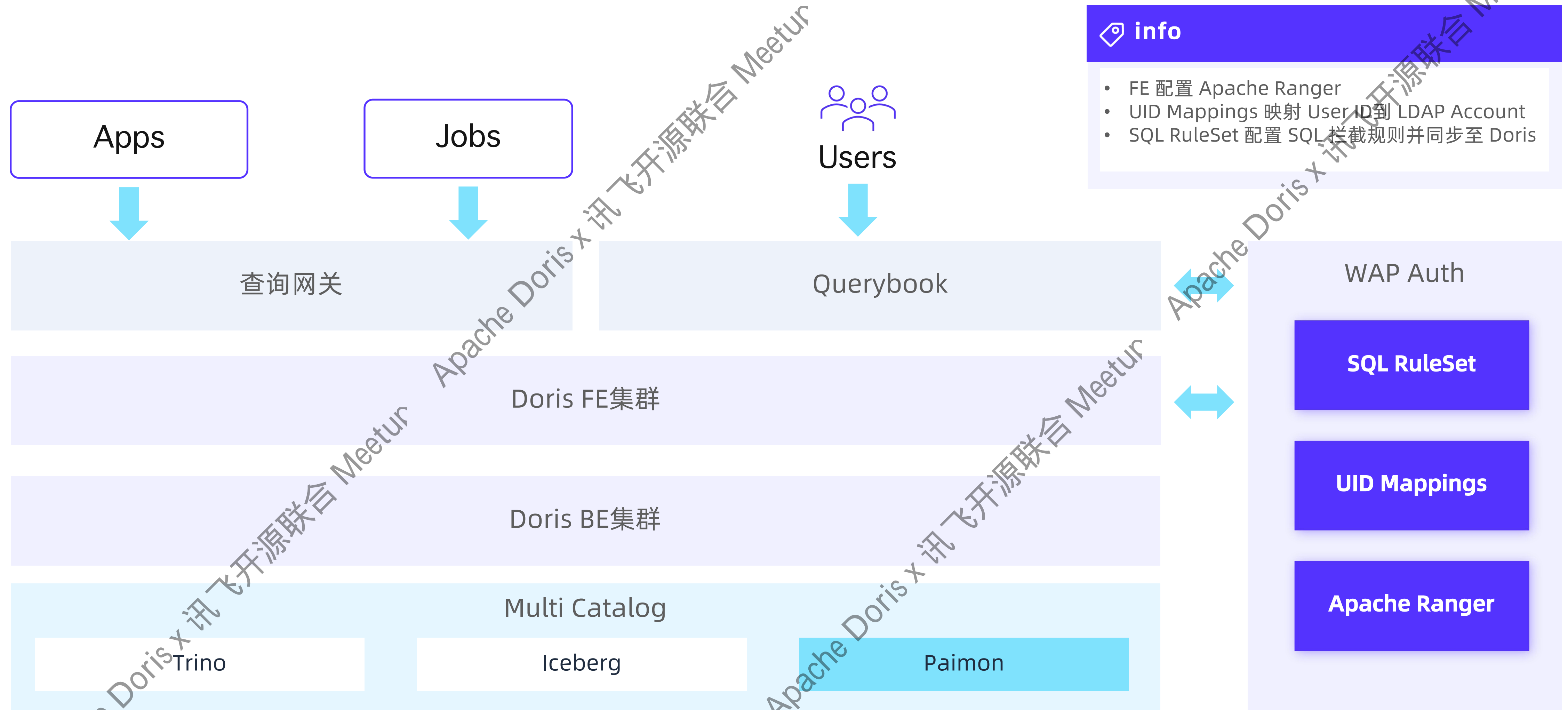
- 一个应用多个账号，账号管理难度大
- 认证授权不统一，授权管控难度大
- 不同用户查询复杂度不同，Kyuubi资源利用率低
- 对于执行的查询的管控能力低，如，无法识别风险查询

解决方案



Apache Doris x 讯飞开源联合 Meetur

基于 Doris 的统一查询服务



目录

01 Webex 业务及数据平台介绍

02 Webex 的数据湖仓架构演化

03 基于 Doris 的数据湖仓实践

04 Doris 与 Webex 数据治理平台的融合

Apache Doris x 讯飞开源联合 Meetur

Apache Doris x 讯飞开源联合 Meetur

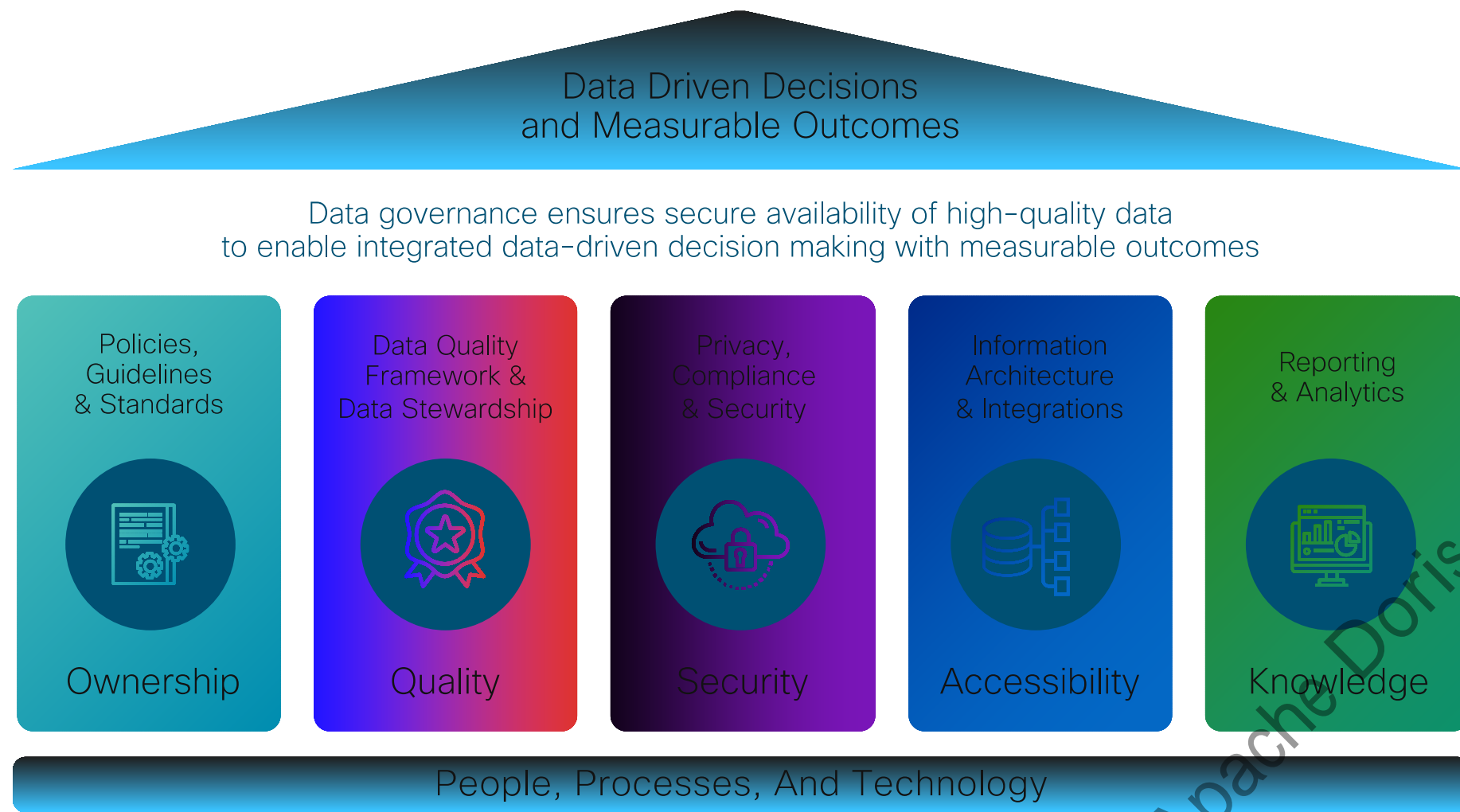
Apache Doris

Apache Doris x 讯飞开源联合 Meetur

Apache Doris x 讯飞开源联合 Meetur

联合 Meetur

数据治理平台融合Doris



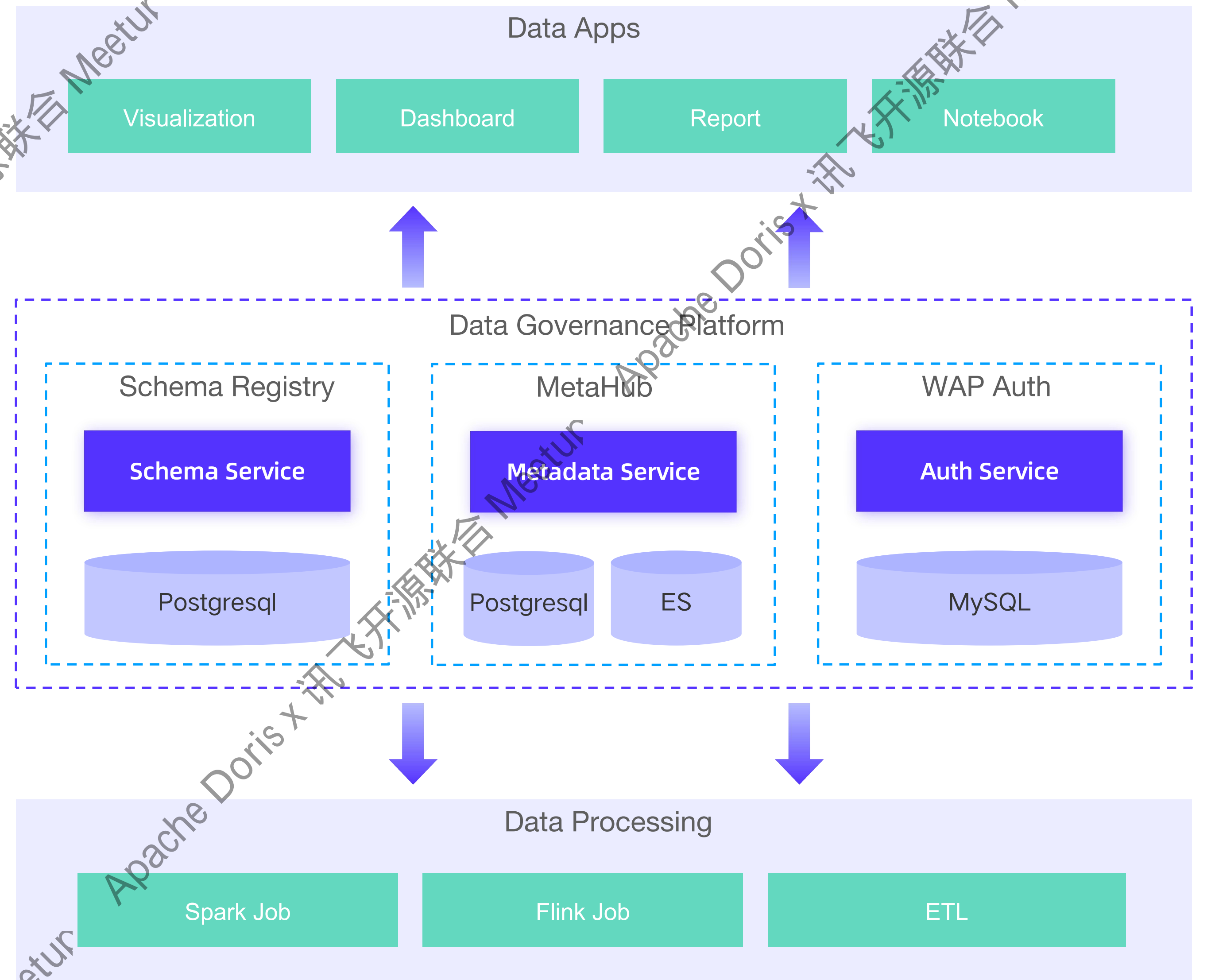
Effective governance is an ongoing effort-executed by people, enabled by processes and supported by technology

目标

为数据管理建立一个组织结构，使其与业务目标、法规遵从性和有效的数据使用保持一致。

已实现

- 由策略、流程和技术工具组成的数据治理框架。
- 丰富数据属性，包括目录、元数据、血缘、所有权等。
- 向数据生产者 and 消费者开放元数据发现、验证能力以了解和实现数据标准。
- 实现嵌入式数据治理，以协调所有数据服务/产品遵守治理策略。



数据治理平台融合 Doris

Webex Data Platform 是受数据治理的平台，引入 Apache Doris 后需要将其纳入数据治理体系中

元数据融合

确定元数据采集方法
明确元数据采集范围
数据血缘集成

授权与审计

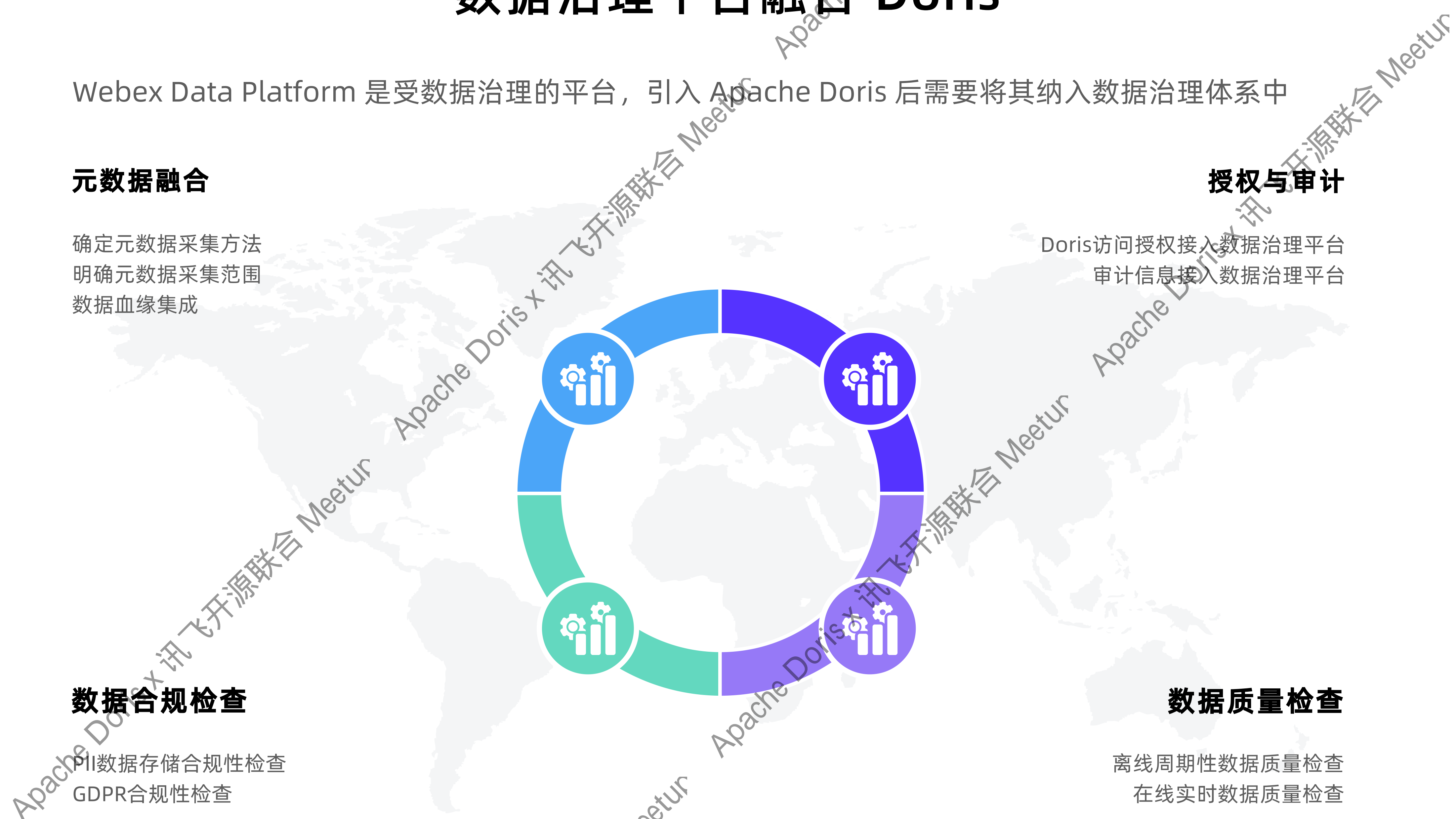
Doris访问授权接入数据治理平台
审计信息接入数据治理平台

数据合规检查

PII数据存储合规性检查
GDPR合规性检查

数据质量检查

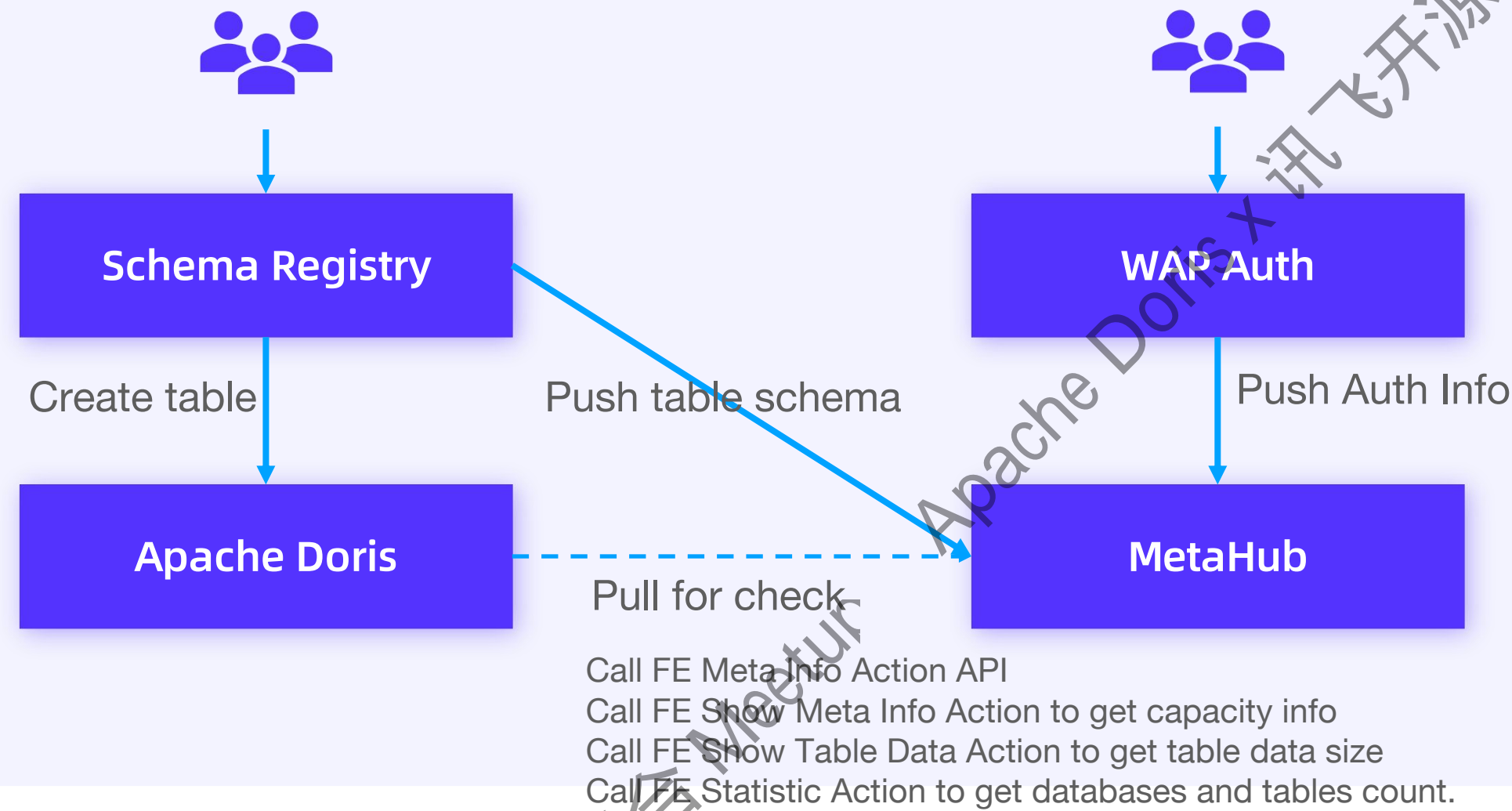
离线周期性数据质量检查
在线实时数据质量检查



数据治理平台融合 Doris

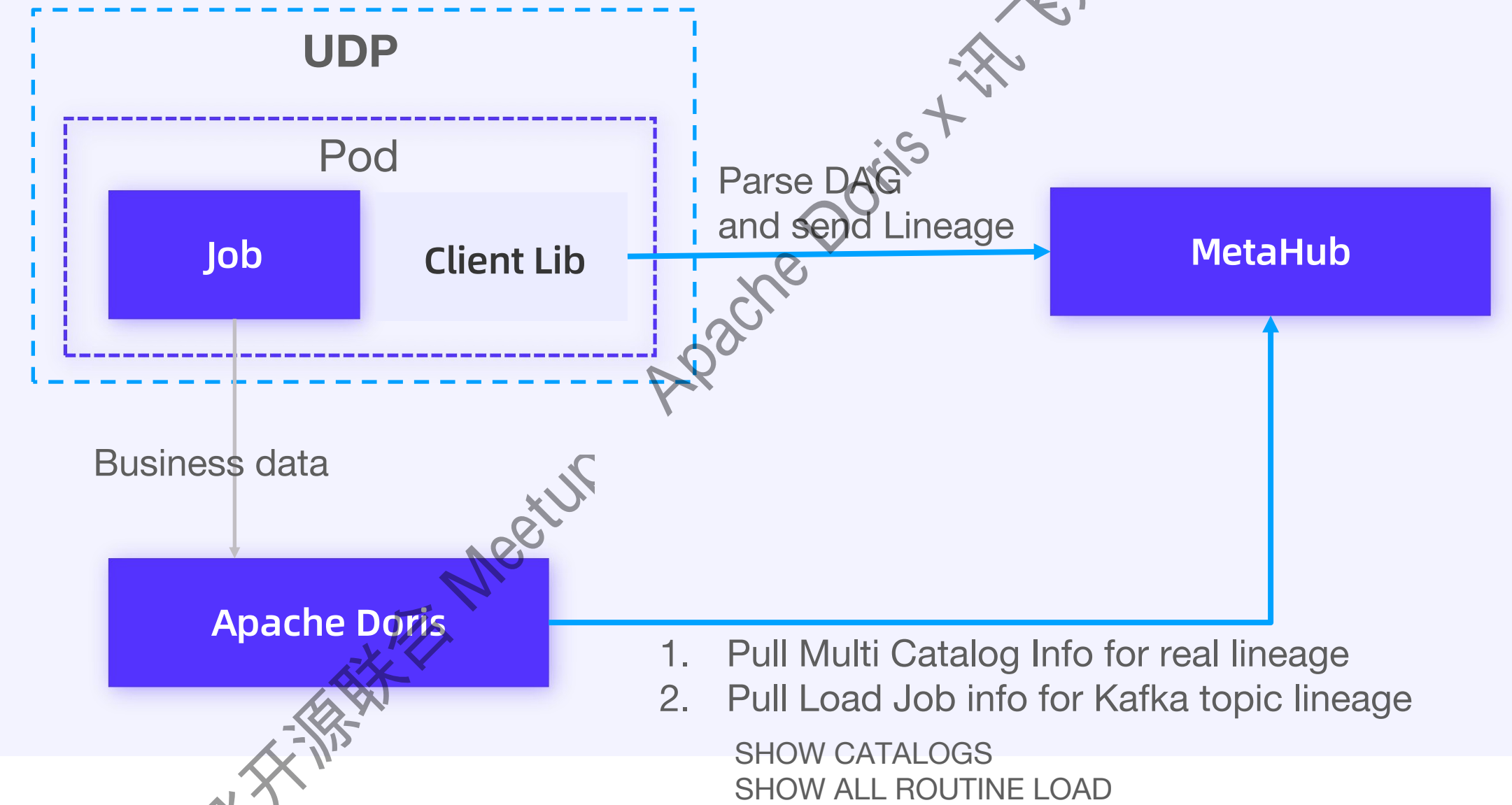
元数据流

(Table Schema + Ownership + Attributions)



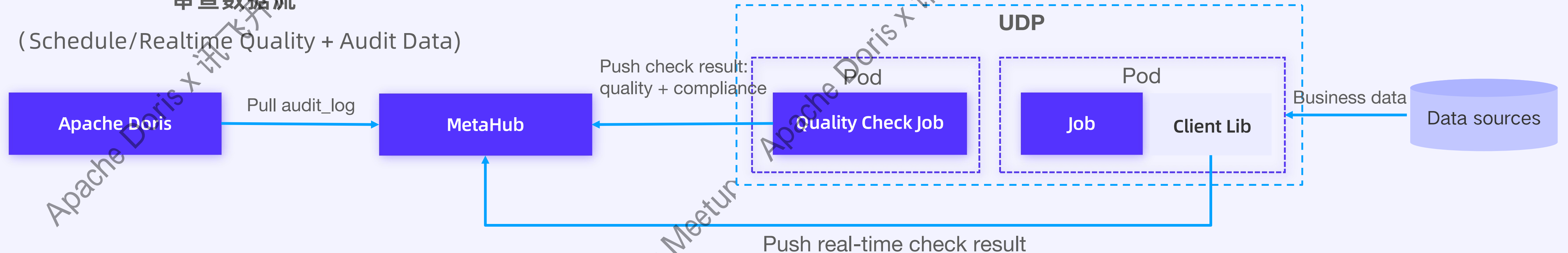
血缘数据流

(Lineage + Catalog + Load Job Metadata)

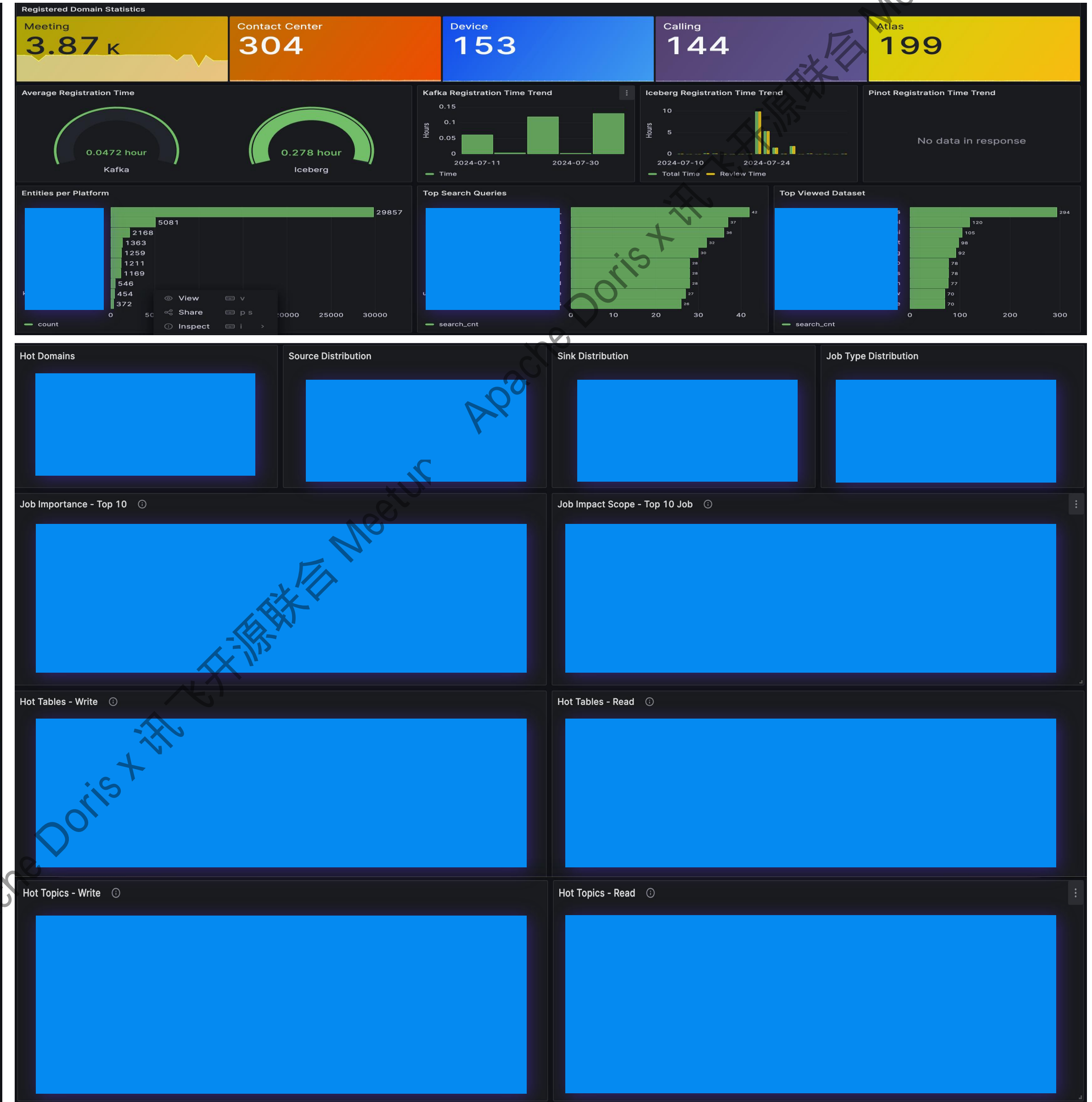
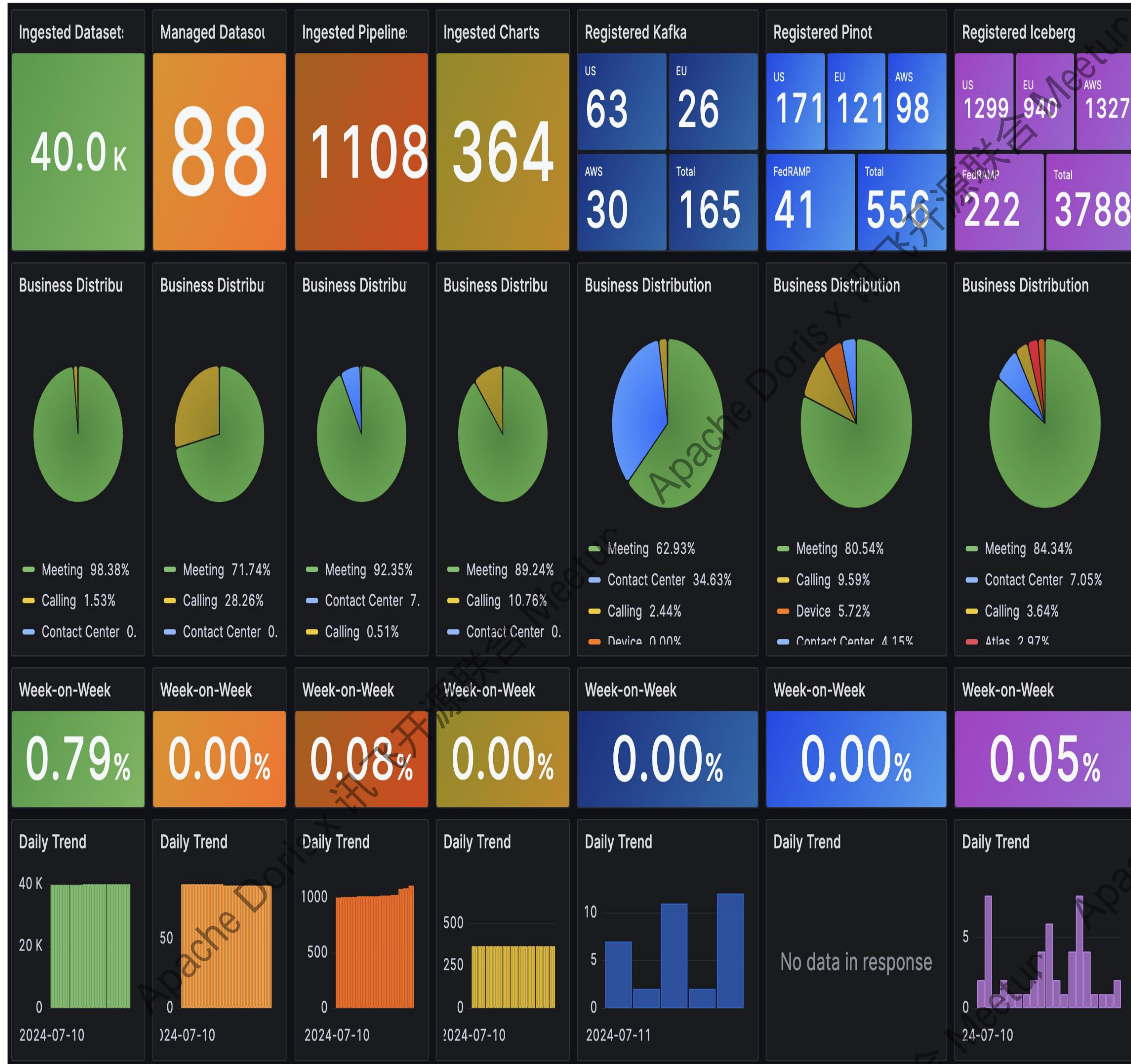


审查数据流

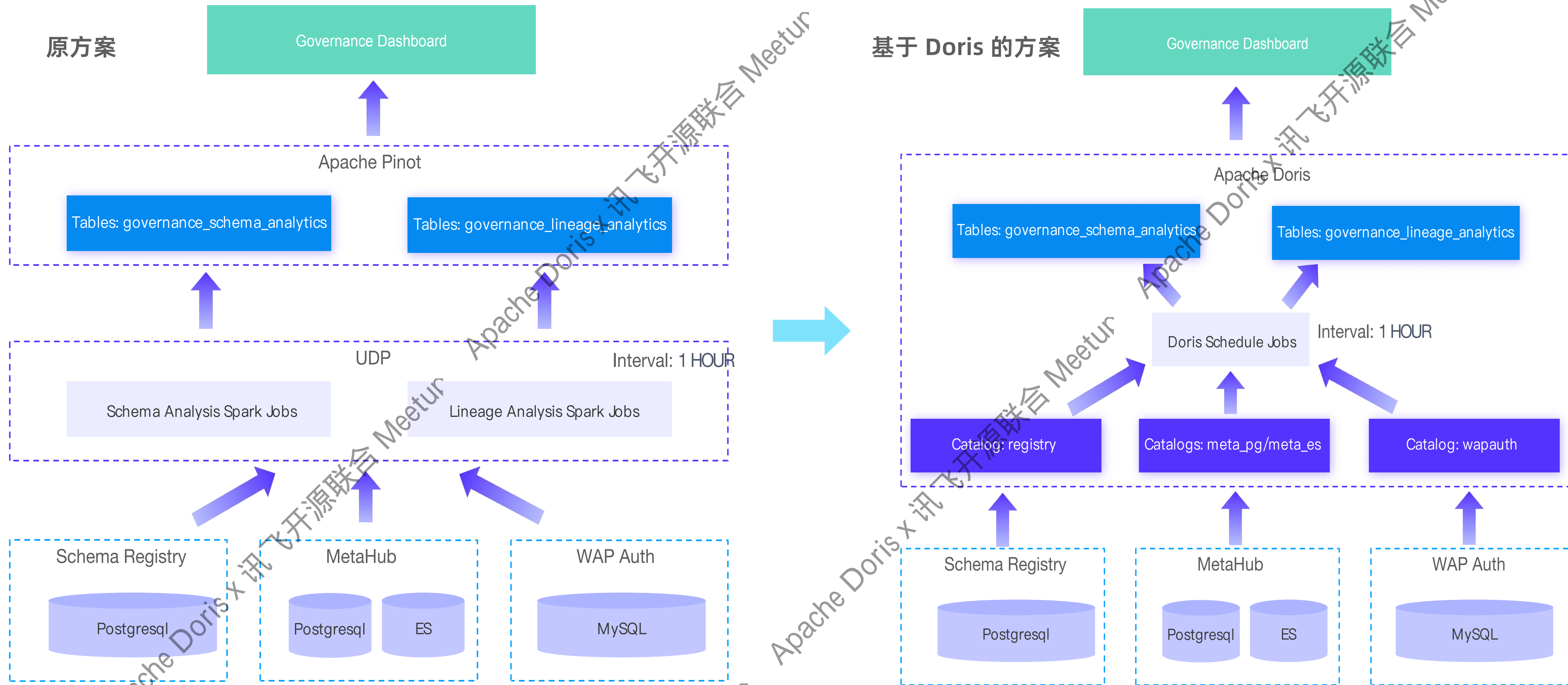
(Schedule/Realtime Quality + Audit Data)



基于 Doris 的元数据分析

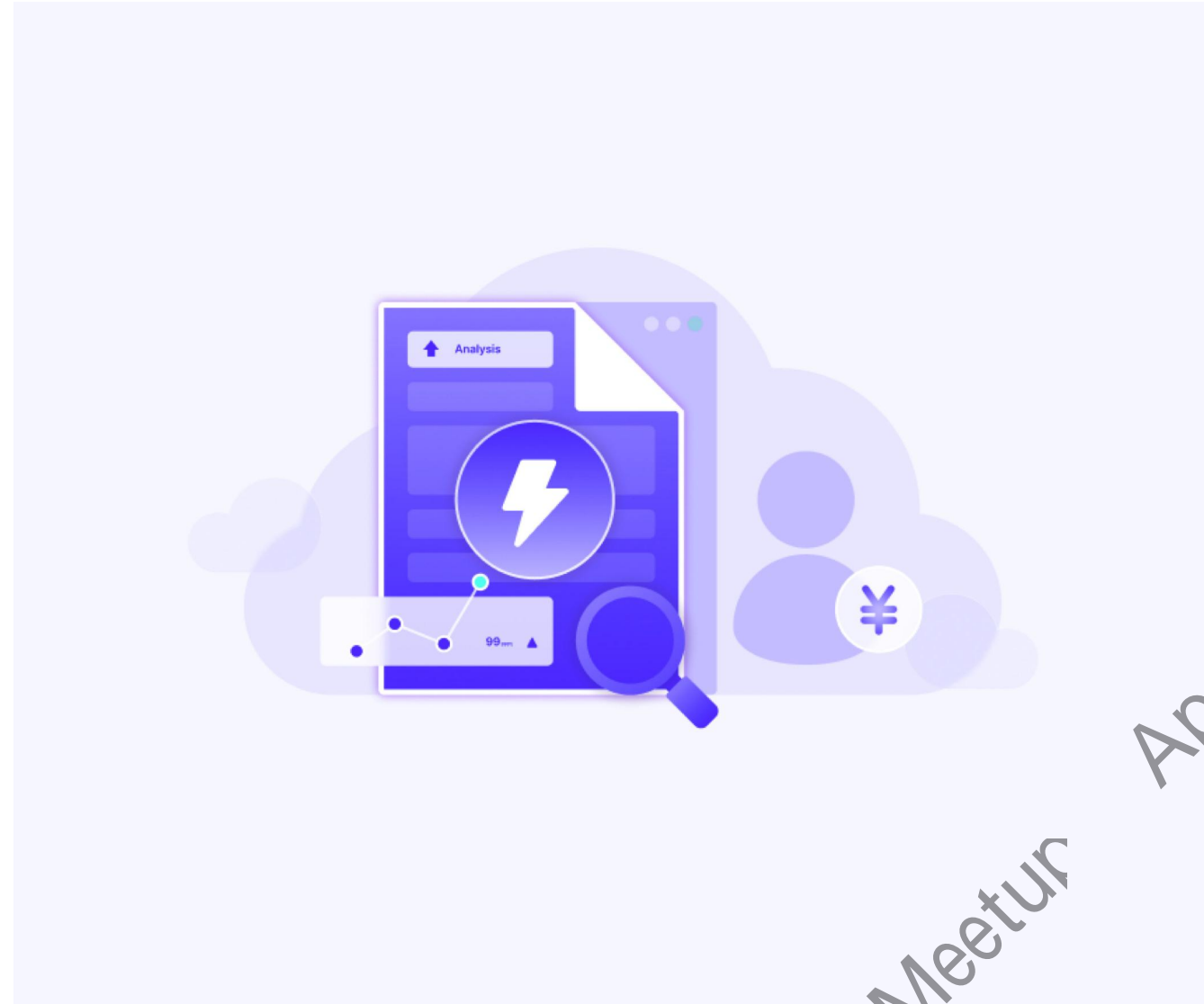


基于Doris的元数据分析



基于 Doris 的方案省去了 Spark Job 的维护工作，整个 Data pipeline 都在 Doris 中进行创建和管理，过程更加清晰

未来计划



扩大Doris应用范围

逐步迁移更多数据湖仓中的业务至 Doris
逐步迁移更多数据应用层 App 至 Doris



增强数据平台功能与性能

基于 Doris 构建高性能数据分析平台，
逐步替代应用自建的分析性存储方案，
如 TiDB，Kylin 等。



探索更多应用场景

探索 AI on Doris 场景
探索 Doris on Paimon 场景

Thanks !



Apache Doris x 讯飞开源联合 Meetur

Apache Doris x 讯飞开源联合 Meetur

Apache Doris

联合 Meetur

Apache Doris x 讯飞开源联合 Meetur

Apache Doris x 讯飞开源联合 Meetur