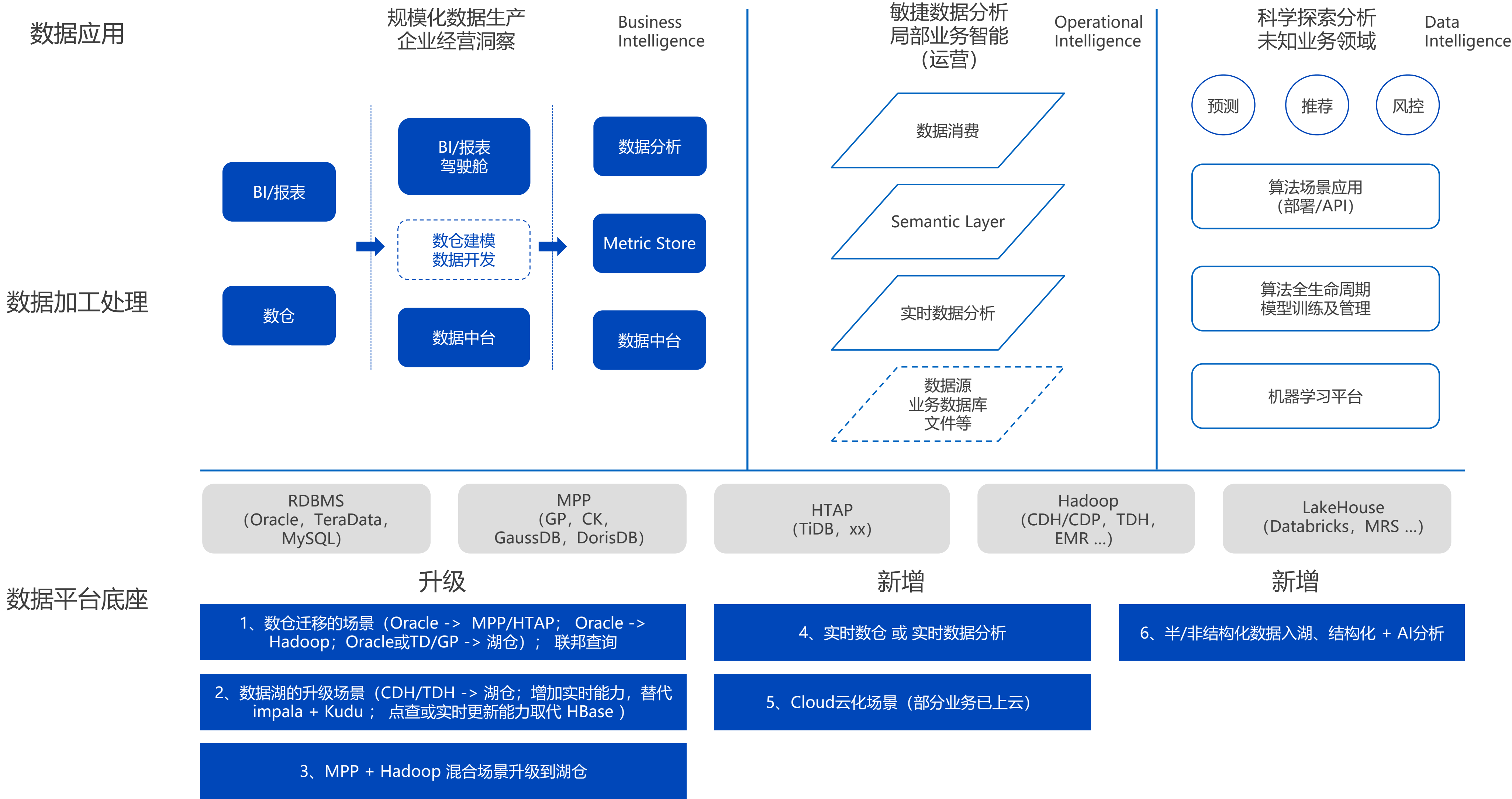


# 湖仓数据平台的技术核心和价值探索

杨磊

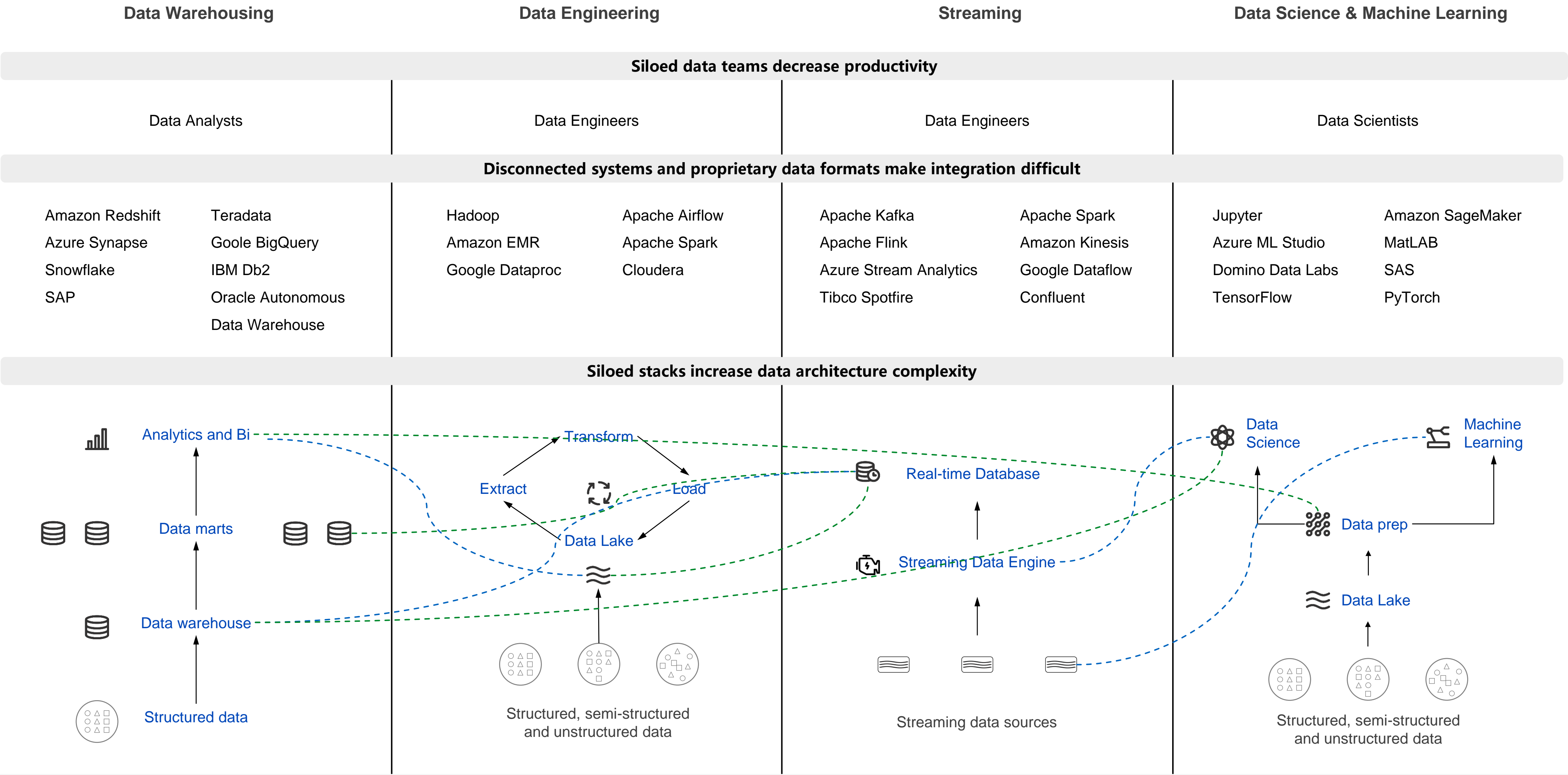


# 企业数据平台场景





# 企业内数据平台的现状 (复杂、低效)





# 目标和挑战

## 数据特点

时效性: T+1, T+10min, T+0

数据类型: 半结构、非结构

数据存量和数据增量

## 技术债务

烟囱开发的积累

临时方案债务 (无统一服务...)

历史原因导致多技术架构...

## 平台能力

满足在线业务和分析SLA

计算能力服务化、多负载

统一存储、存算分离

## 扩展和演进

过重的数仓模式

无ACID能力

无法应对业务对数据体系要求



# 湖仓数据平台架构

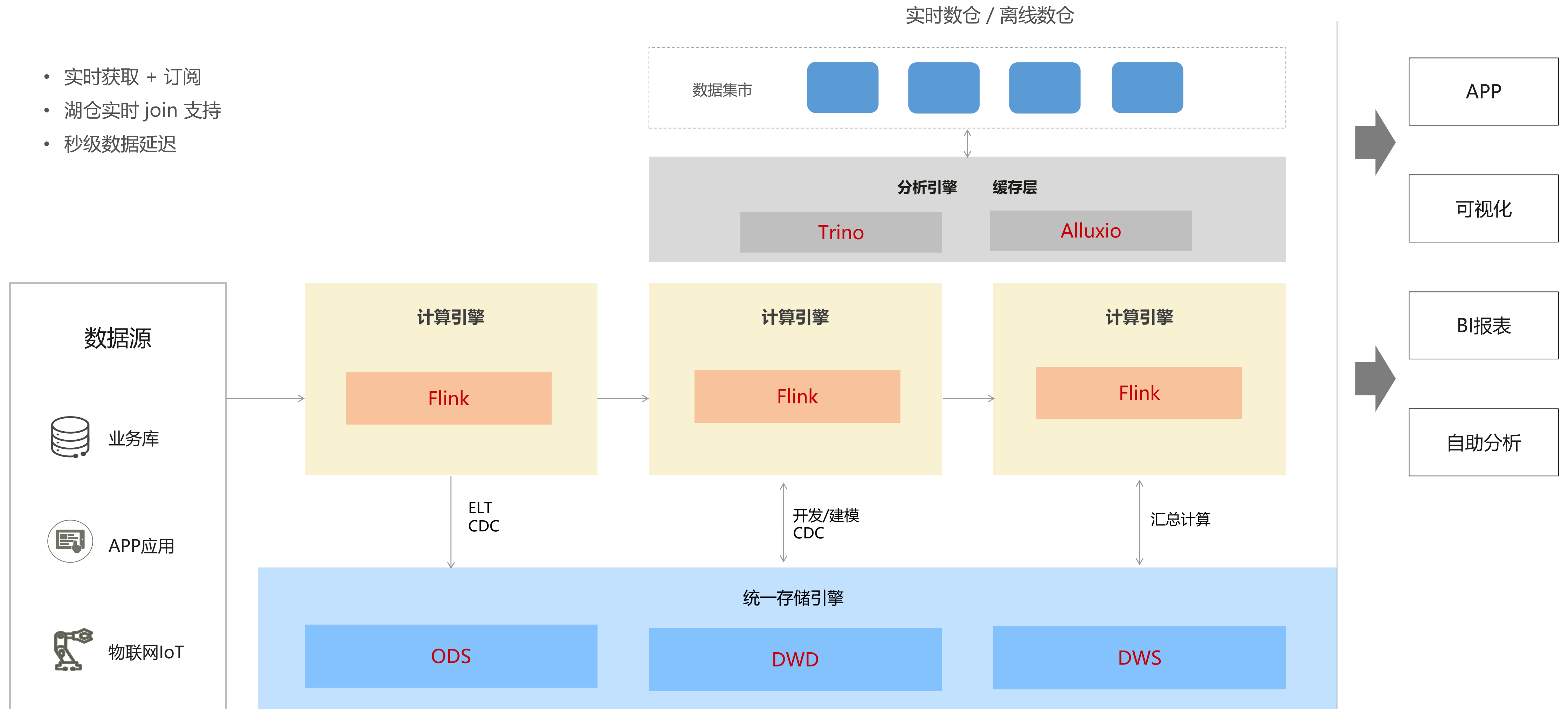


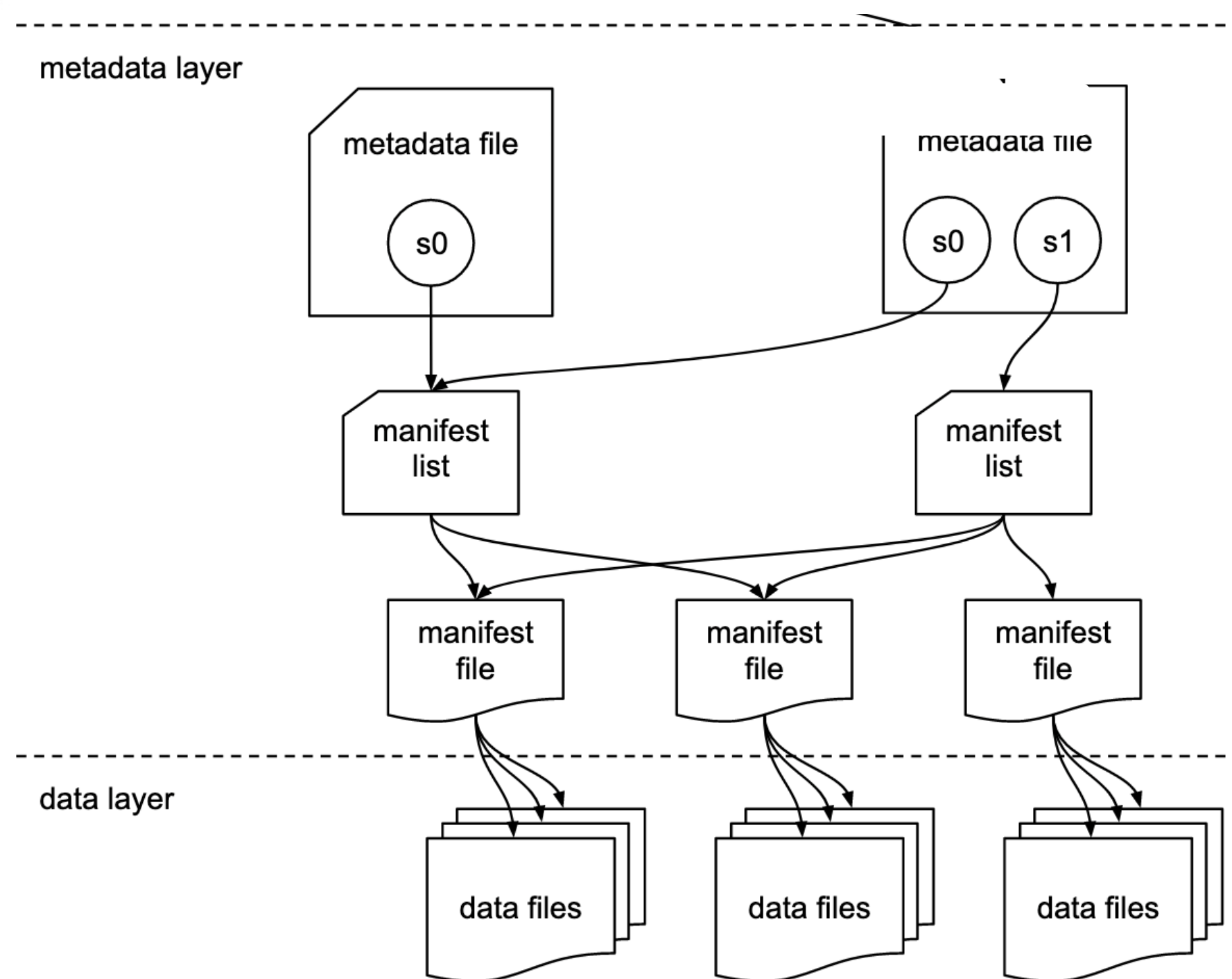




# FastData实践：批流融合 & 全链路CDC

- 实时获取 + 订阅
- 湖仓实时 join 支持
- 秒级数据延迟





Metadata file: 元数据文件（存储某个时间点的表元数据）

Manifest list file: Manifest列表文件（文件列表）

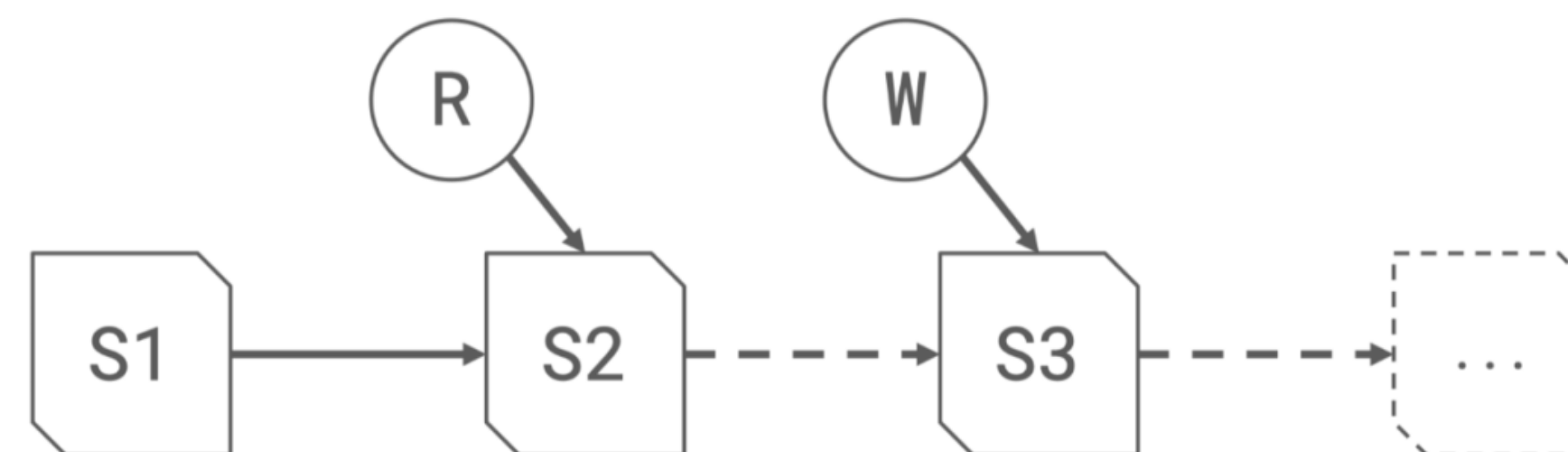
Manifest file: Manifest文件（数据文件列表，以及关于每个数据文件的详细信息和统计信息）

Data files: 数据文件（数据文件对象存储）

每次写入都会成一个snapshot, 每个snapshot包含着一系列的文件列表

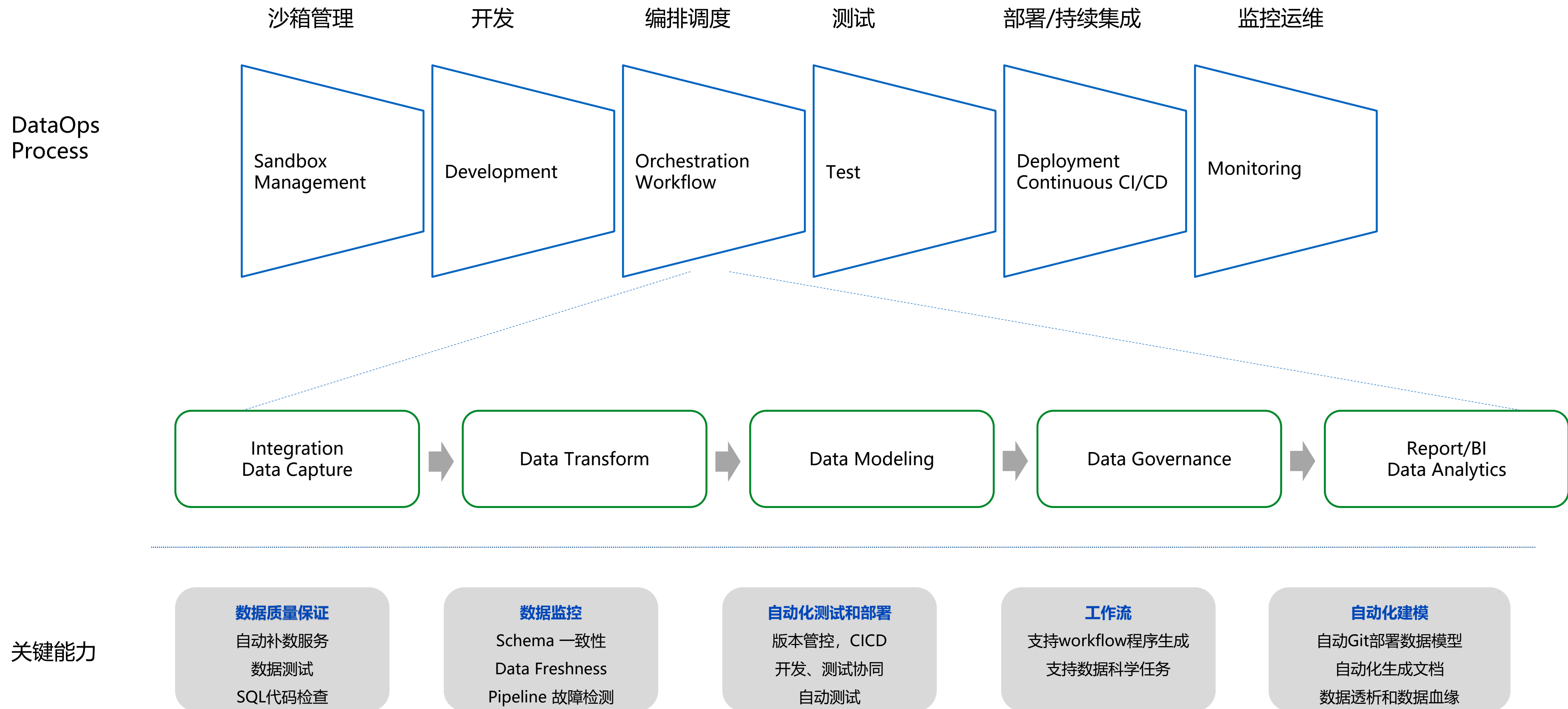


基于MVCC(Multi Version Concurrency Control)的机制,默认读取文件会从最新的版本, 每次写入都会产生一个新的snapshot, 读写相互不干扰





# FastData实践: DataOps







# FastData实践：统一元数据管理

## 统一数据标准

数据架构

数据规范

数据安全

数据质量

生命周期管理

任务治理

## 统一数据源管理

租户隔离

唯一性校验

解析/检验

帐号管理

数据源应用

配置信息

## 统一元数据

Schema-Mapping

流表MetaData构建

统一物理转换

租户及项目级别Catalog

## 统一存储

统一存储类型(HDFS/对象)

存储格式(Parquet or ORC)

统一表索引引擎

多模态存储(结构/非结构)

## 基础数据源

RDBMS  
(Oracle, MySQL,  
GoldenDB,  
OceanBase)

MPP  
(GP, CK,  
GaussDB, GBase)

HTAP  
(TiDB, xx)

Hadoop  
(CDH/CDP, TDH, ,  
FI, EMR ...)

KV & MQ

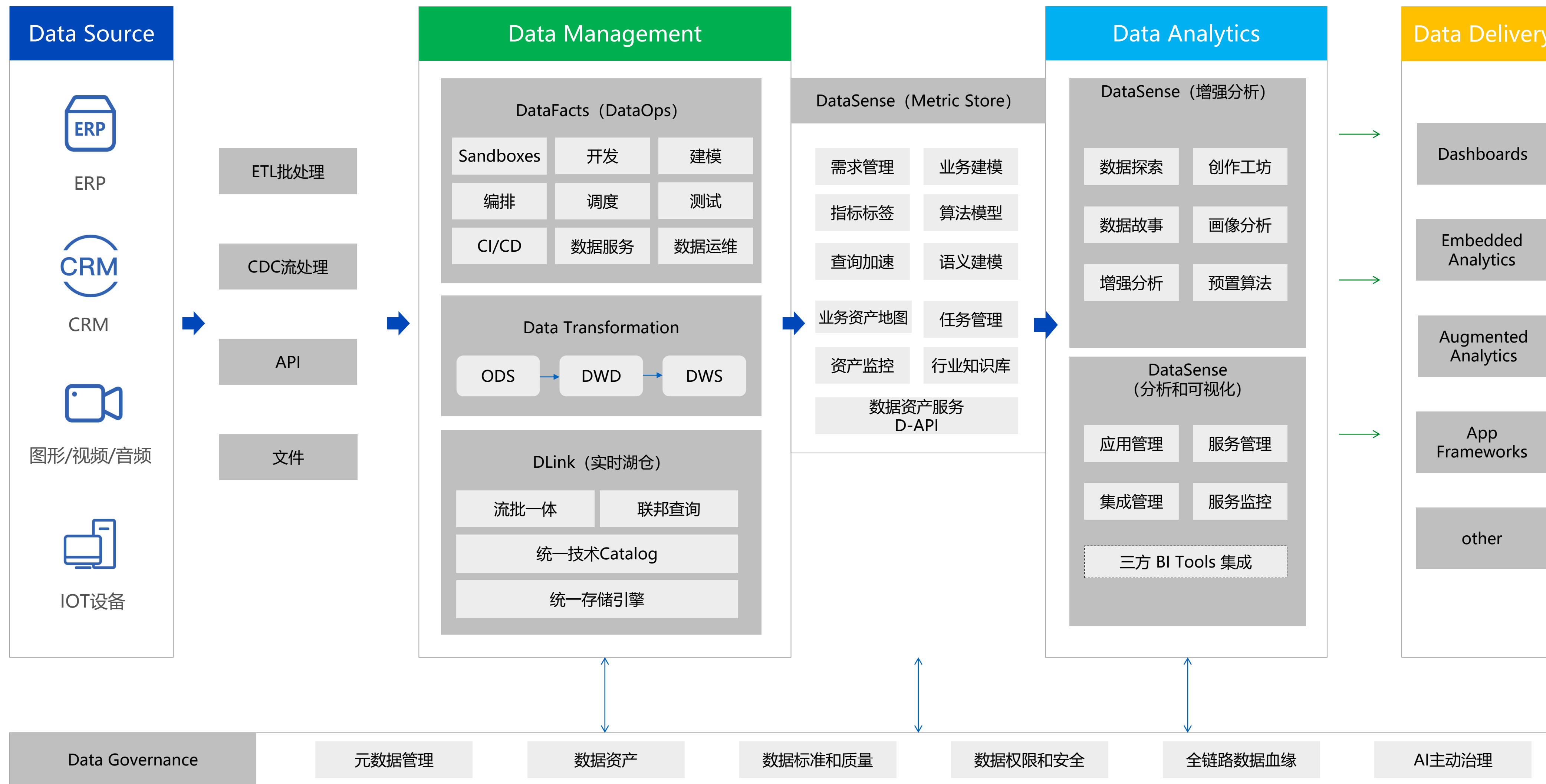
### 湖仓一体的底座核心能力：

- ✓数据存储的类型和格式尽量统一，可以在分散在不同的物理机房或节点，逻辑统一；
- ✓结构、半或非结构化数据的多模态一体化存储；
- ✓统一基础技术元数据；
- ✓统一各板块数据标准体系；



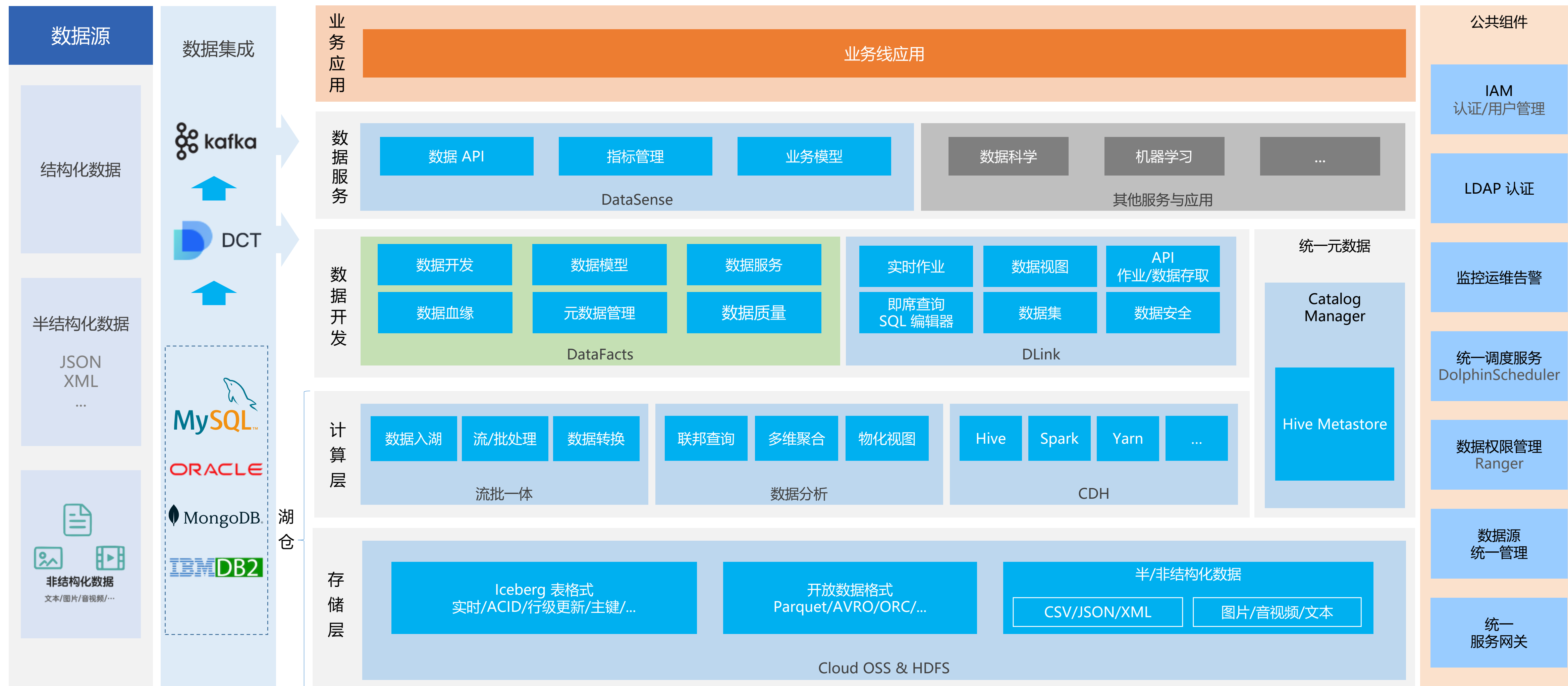
## FastData实时数据平台架构







# FastData价值实践

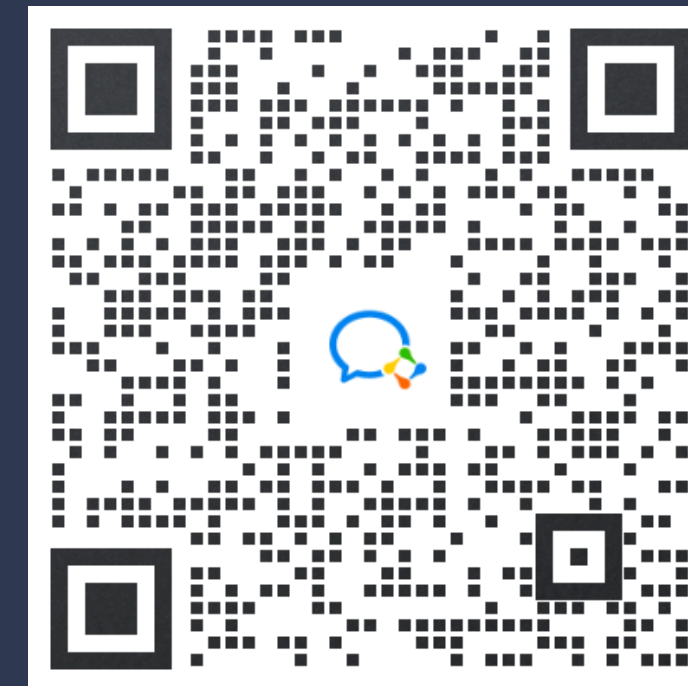




想一想,  
我该如何把这些  
技术应用在工作实践中?

---

# THANKS



技术交流群



社区公众号