



Flink: 构建下一代大数据处理引擎

巴真 @FlinkChina Meetup





Agenda

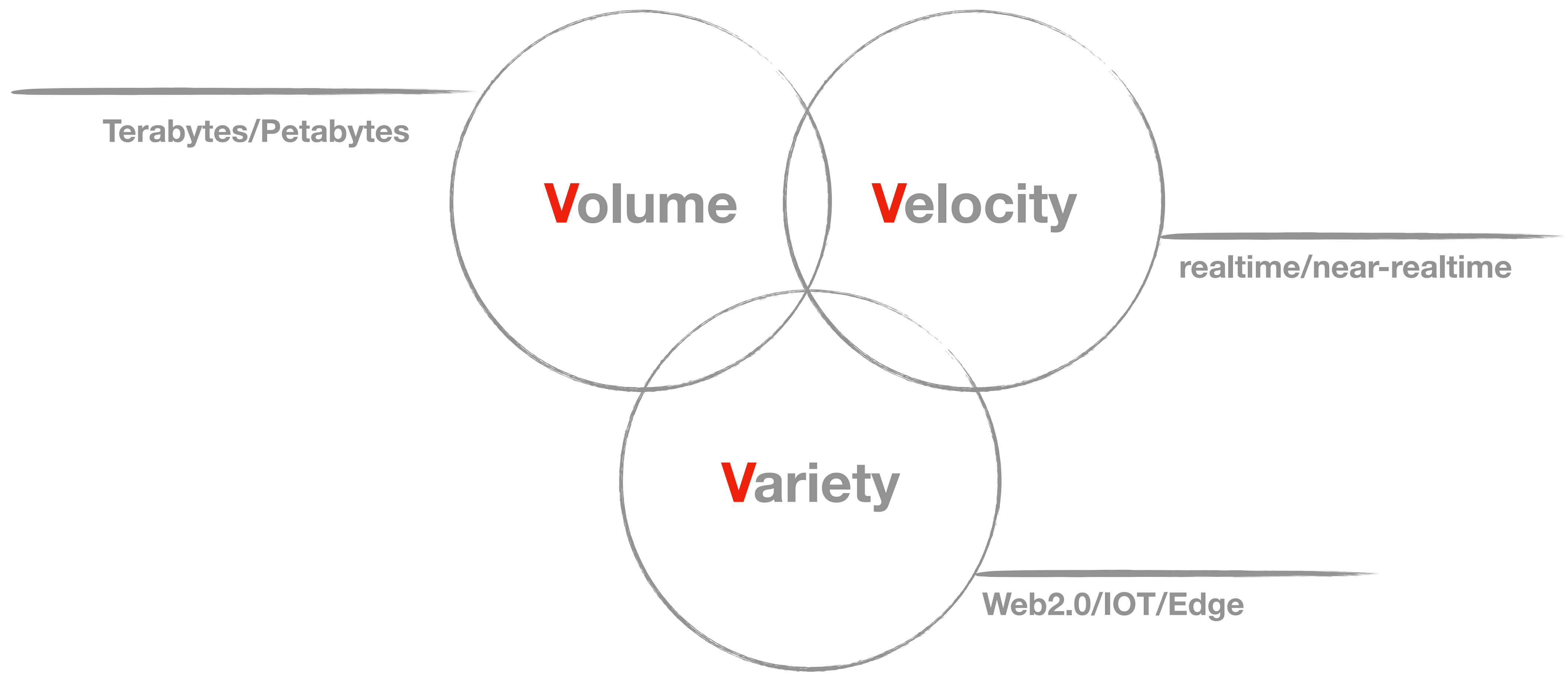
数据趋势

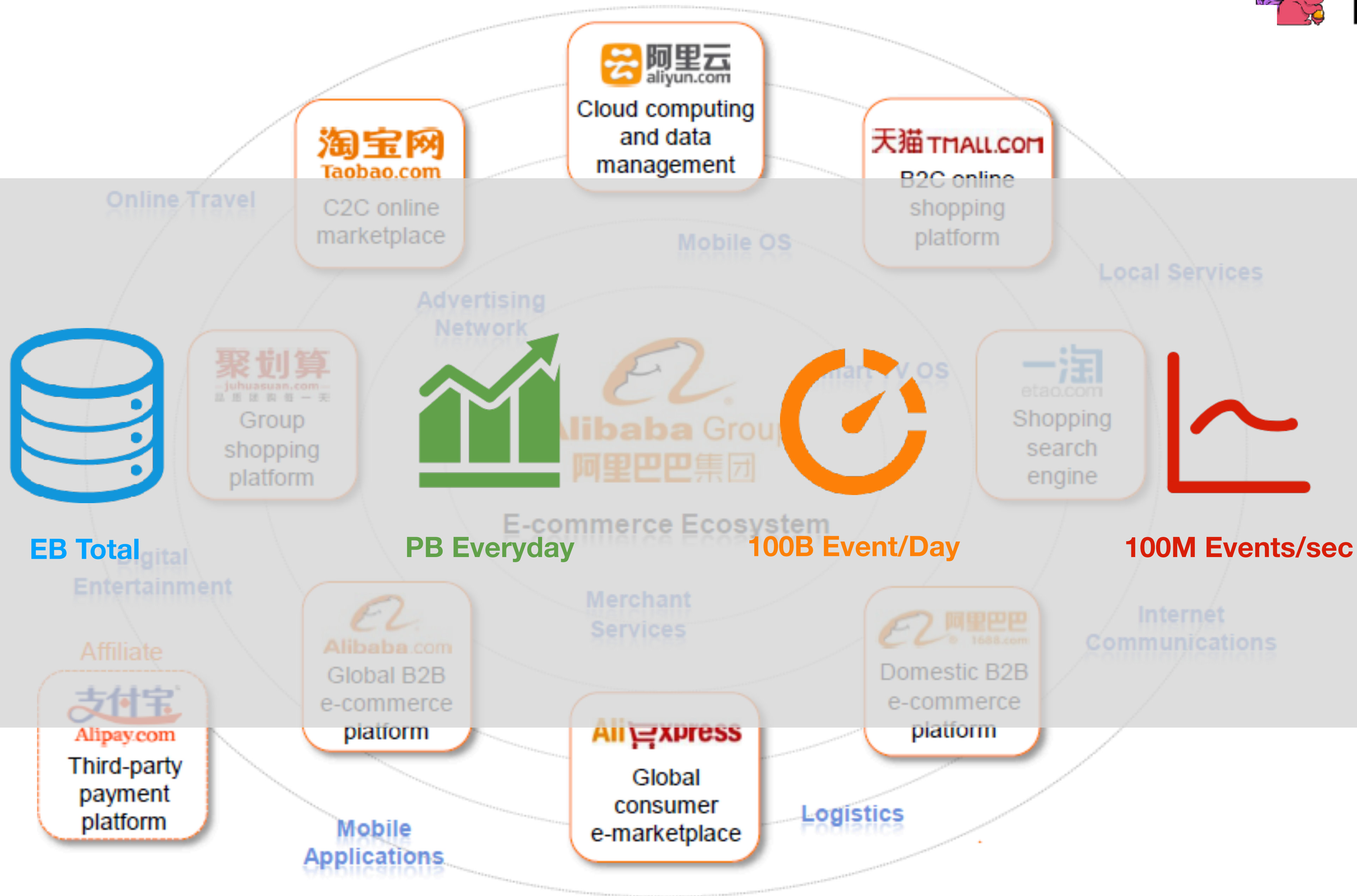
业界案例

阿里思考

Flink@阿里









3亿用户



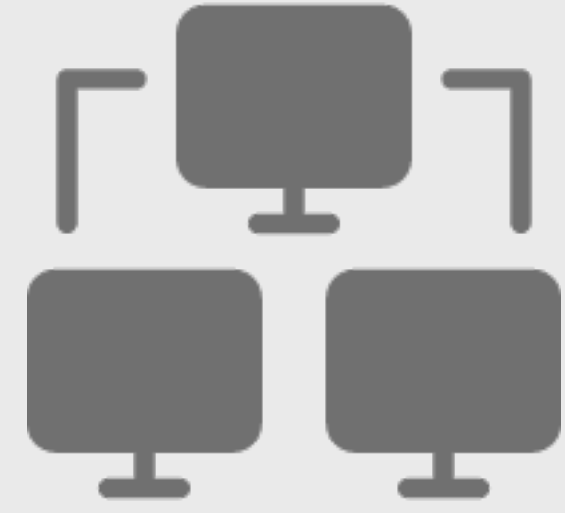
百万商家



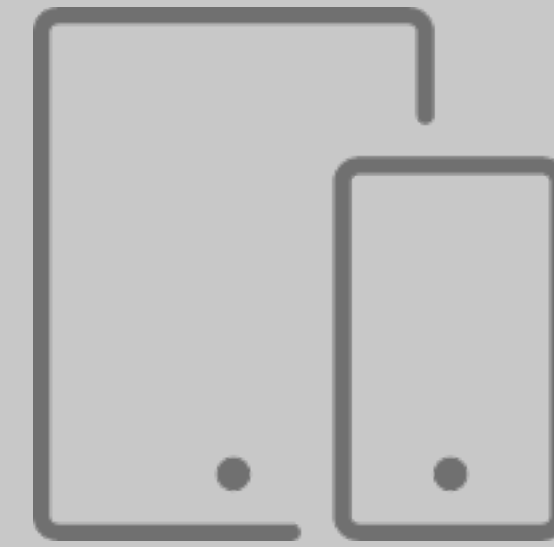
千亿成交



IT化



网络化



移动化



万物互联



AWS S3



阿里云OSS

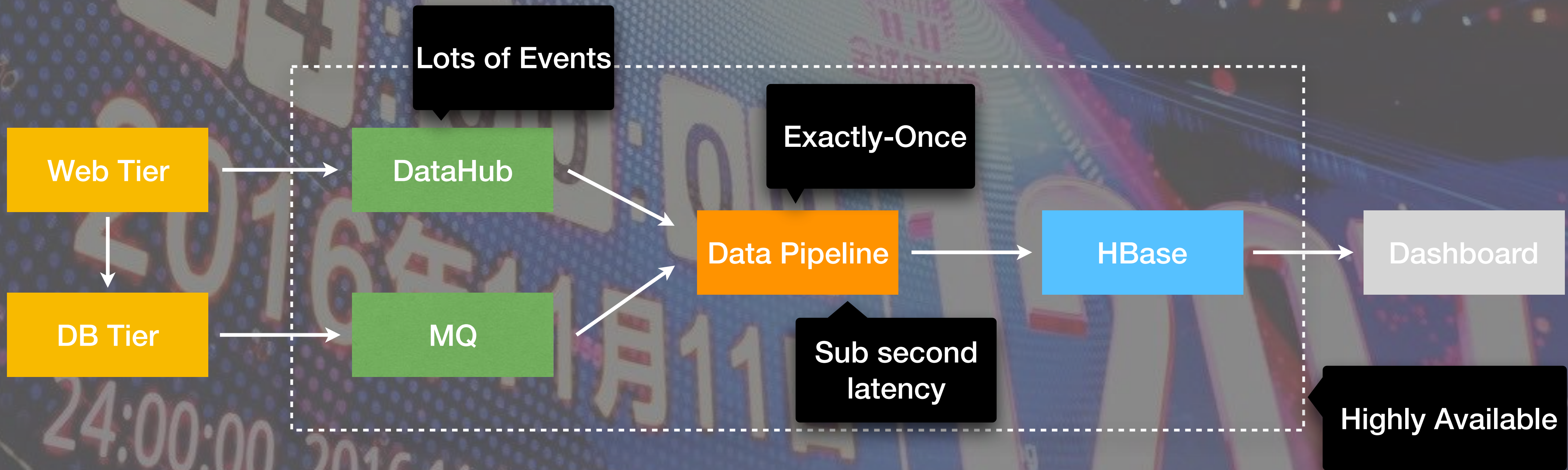
Semi-Structured Unstructured

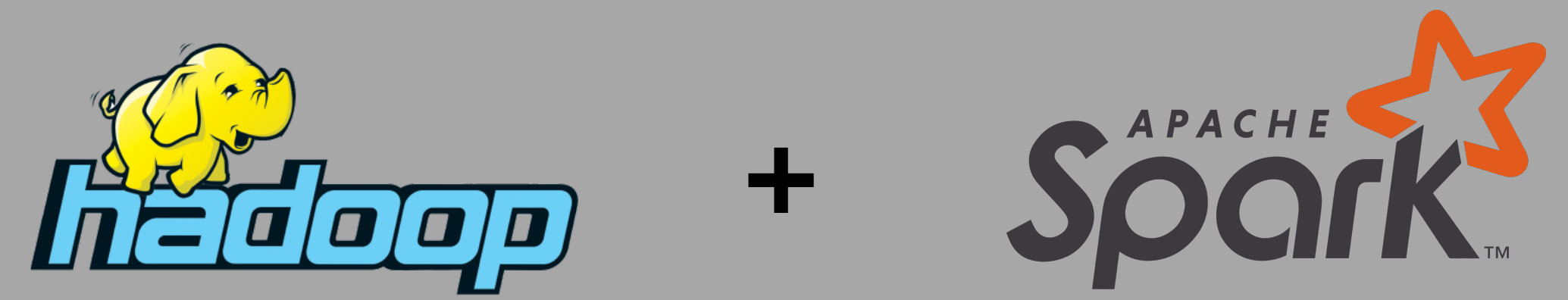
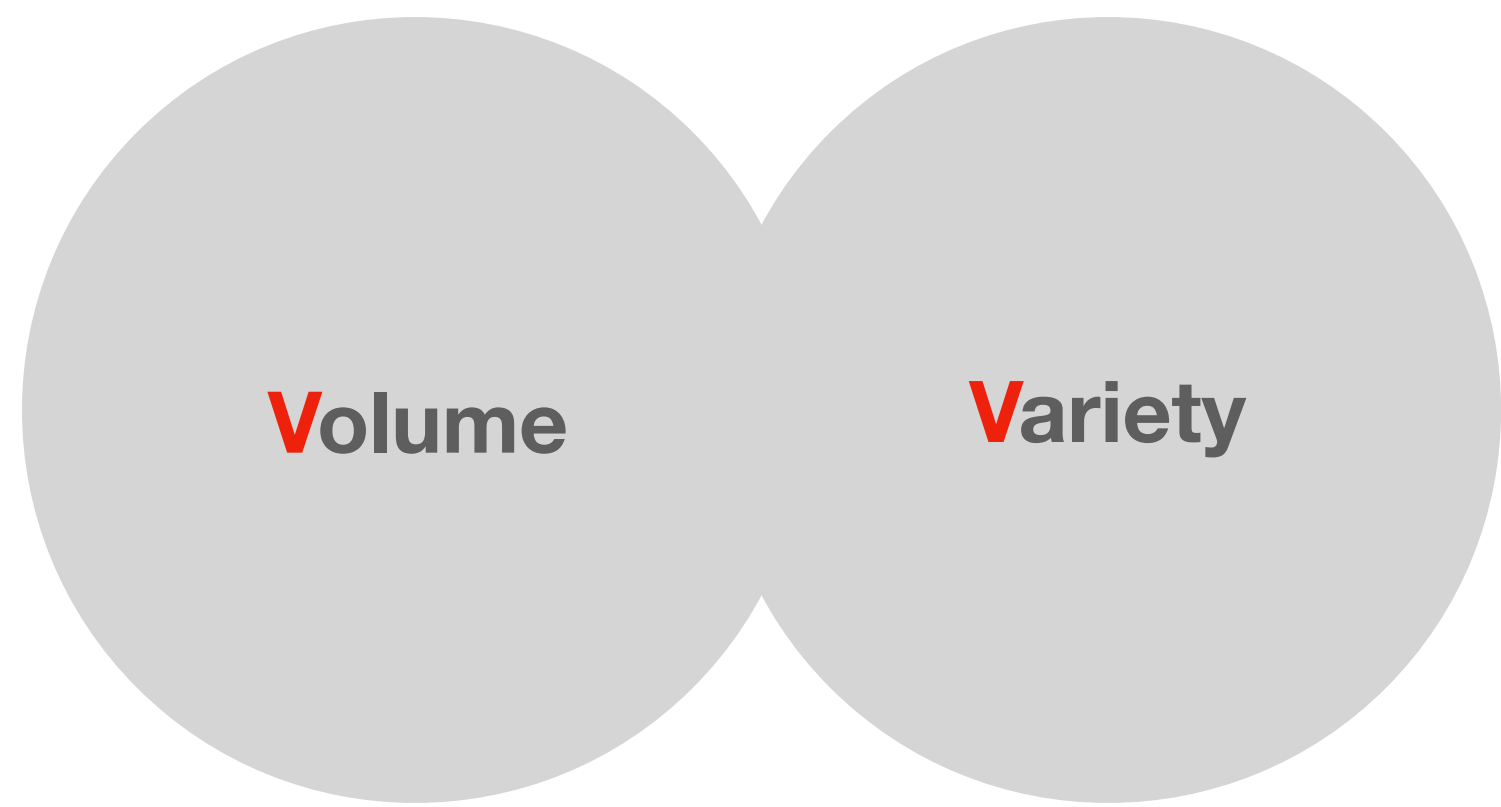


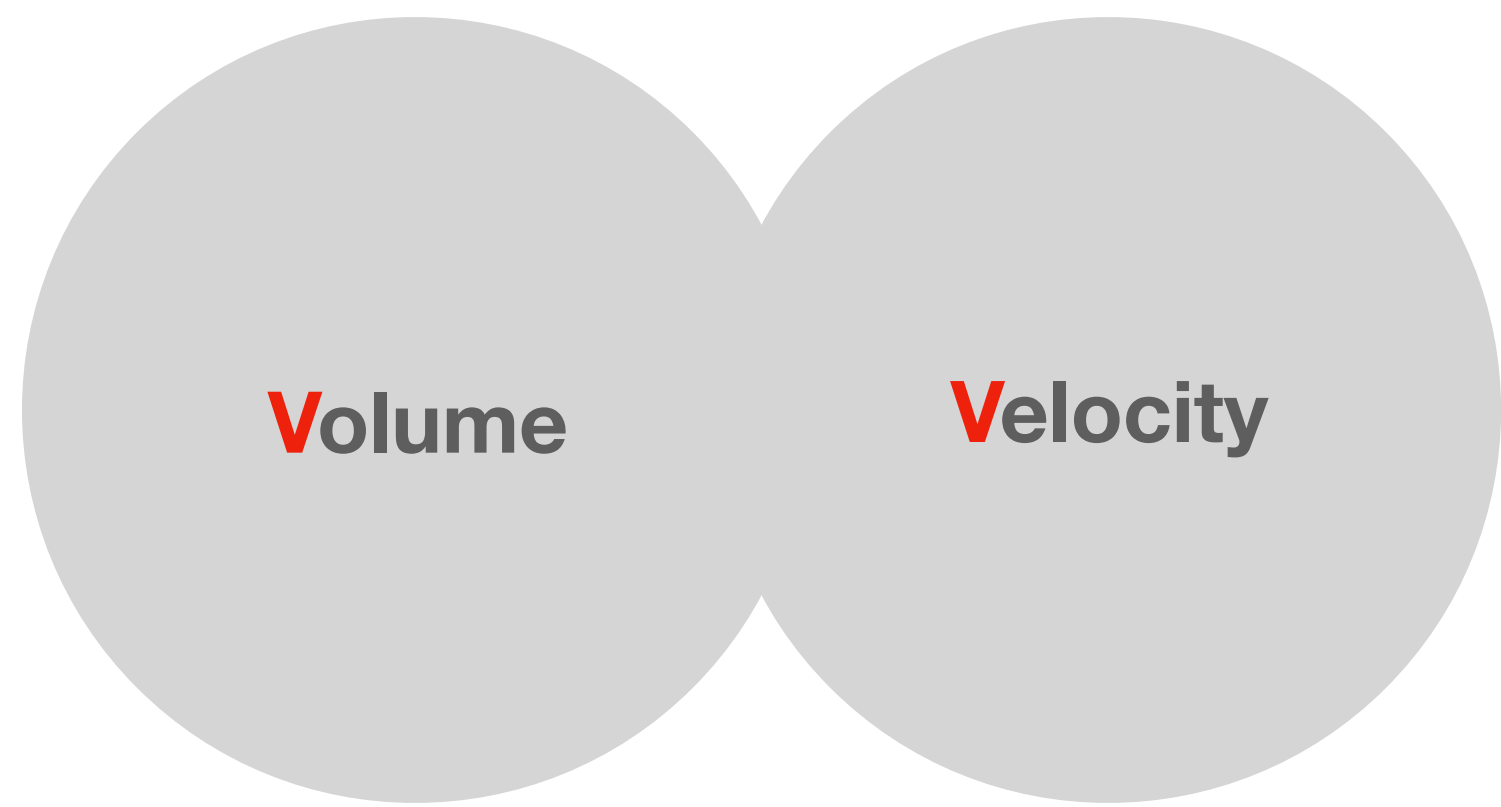
472M events/sec

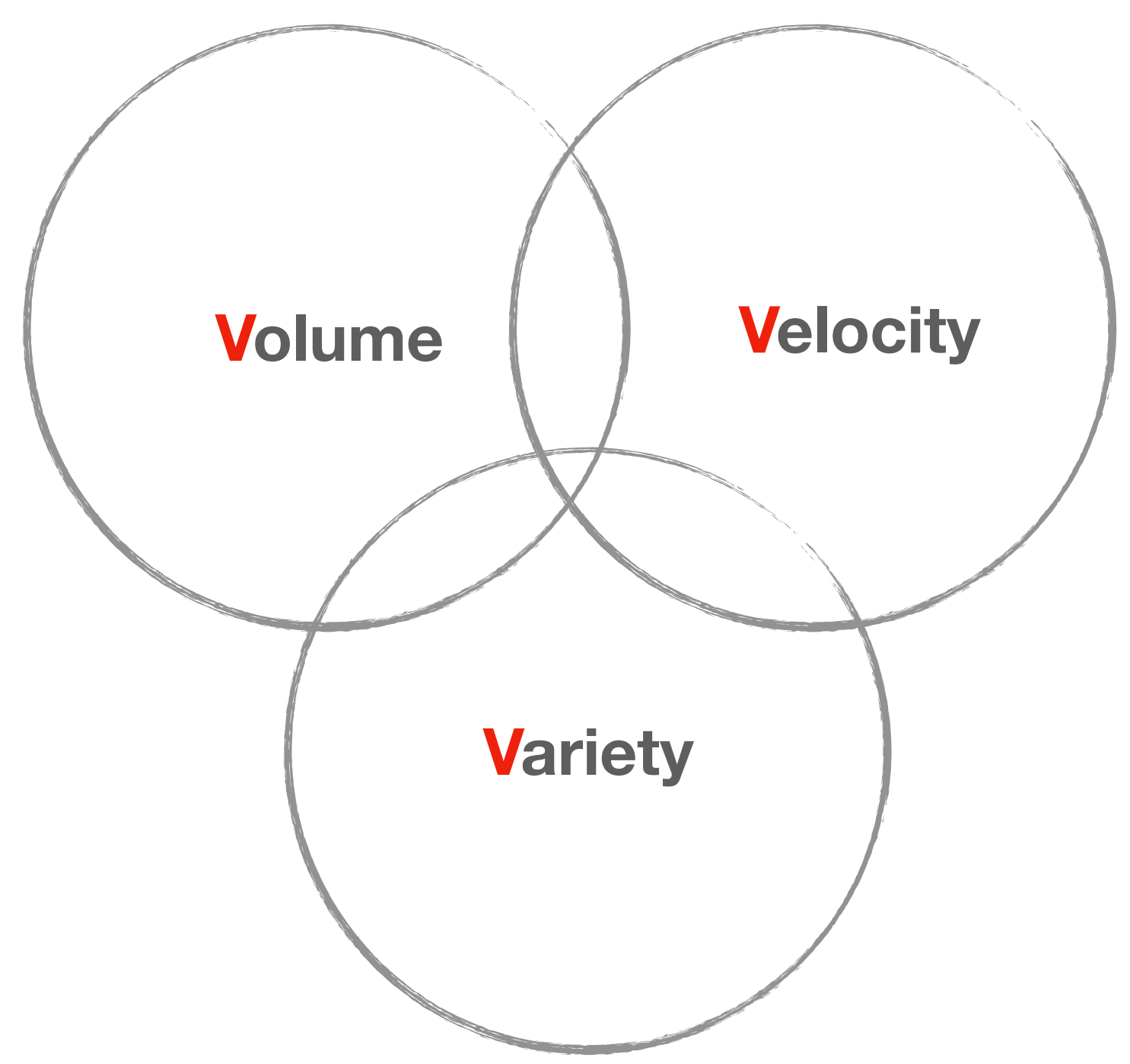


Sub-second Latency









?



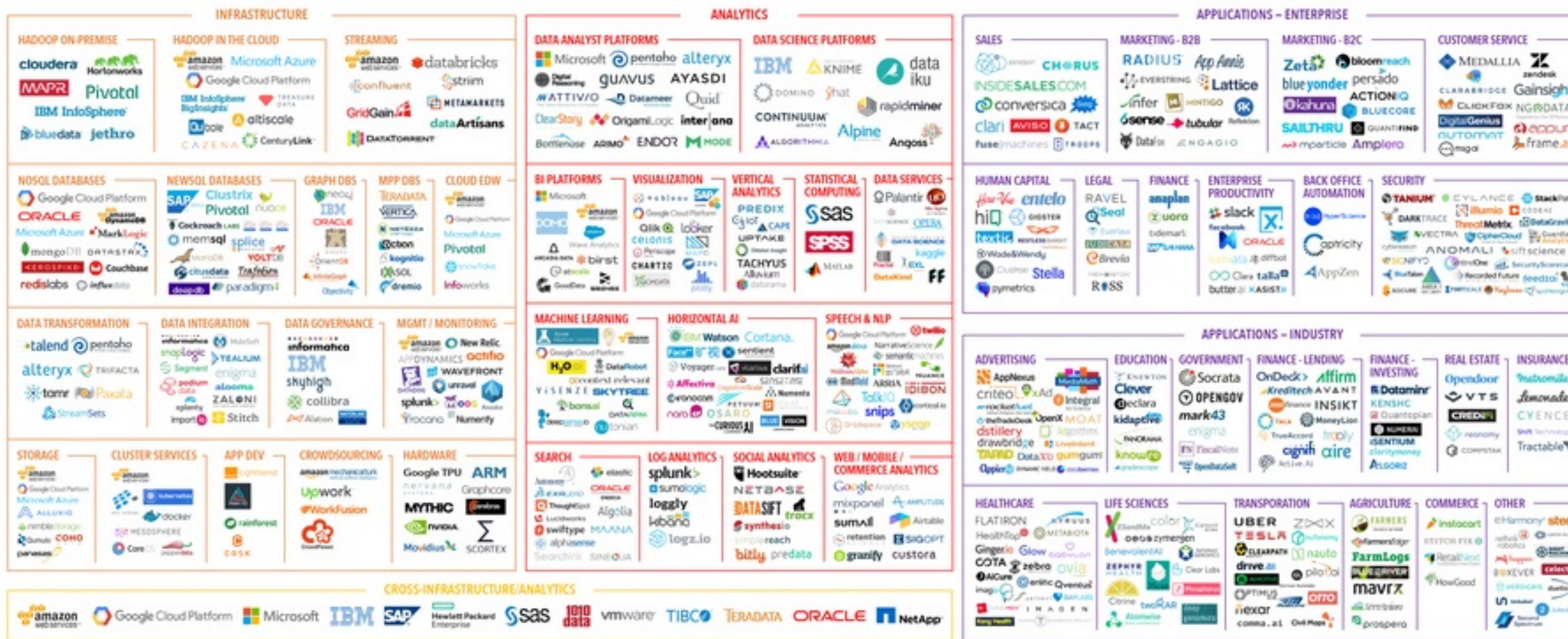
计算碎片化



模型多样化



BIG DATA LANDSCAPE 2017





数万Dataer



百万Job

在阿里，人人都是数据分析师





BI Engineer



AI Engineer

在阿里，BI工程师在转型AI工程师





典型案例分析

阿里业务

超大规模、超级复杂、计算多样、数据多样、用户众多

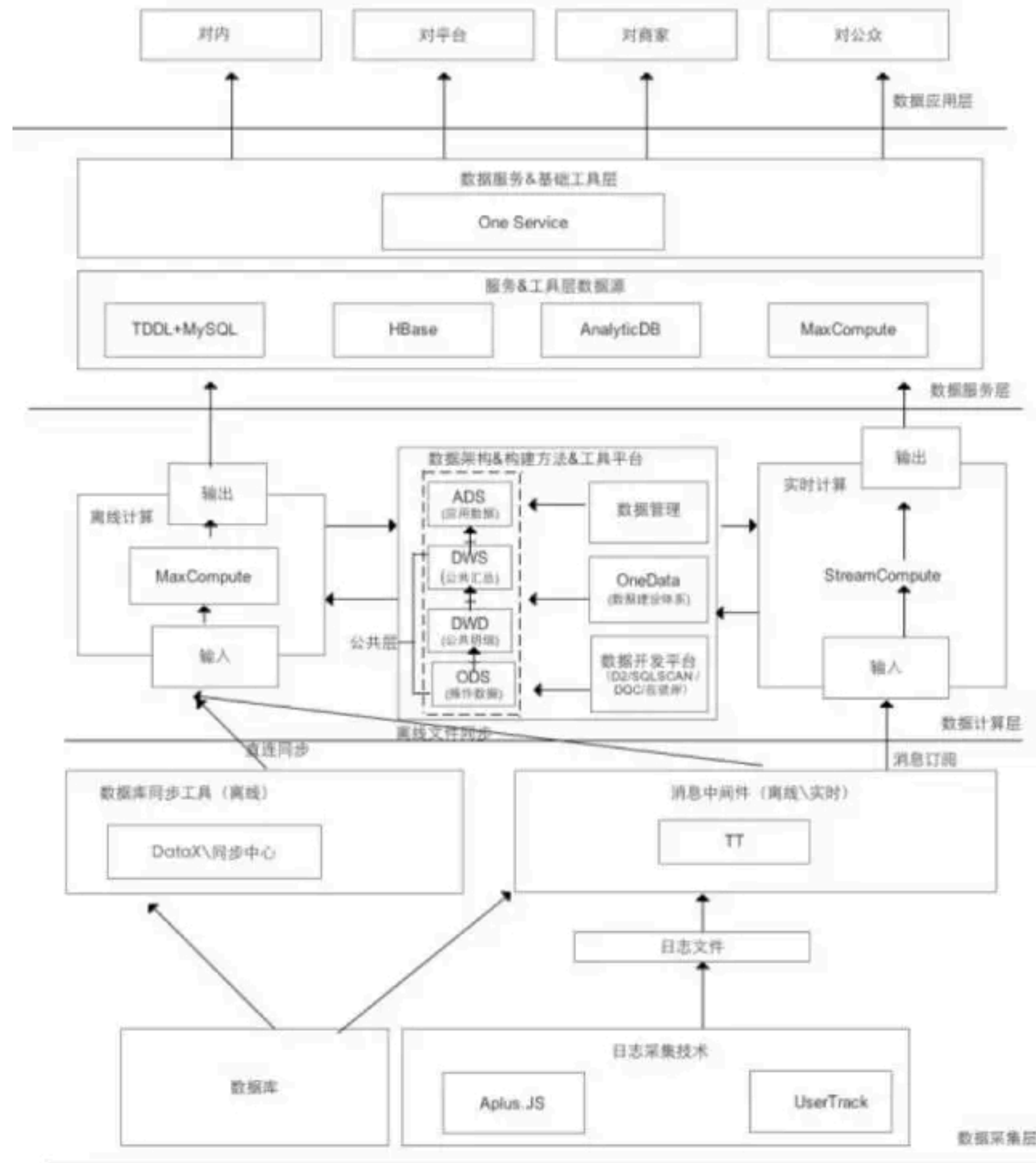
滴滴业务

独角兽公司，爆发式增长、计算时效性强

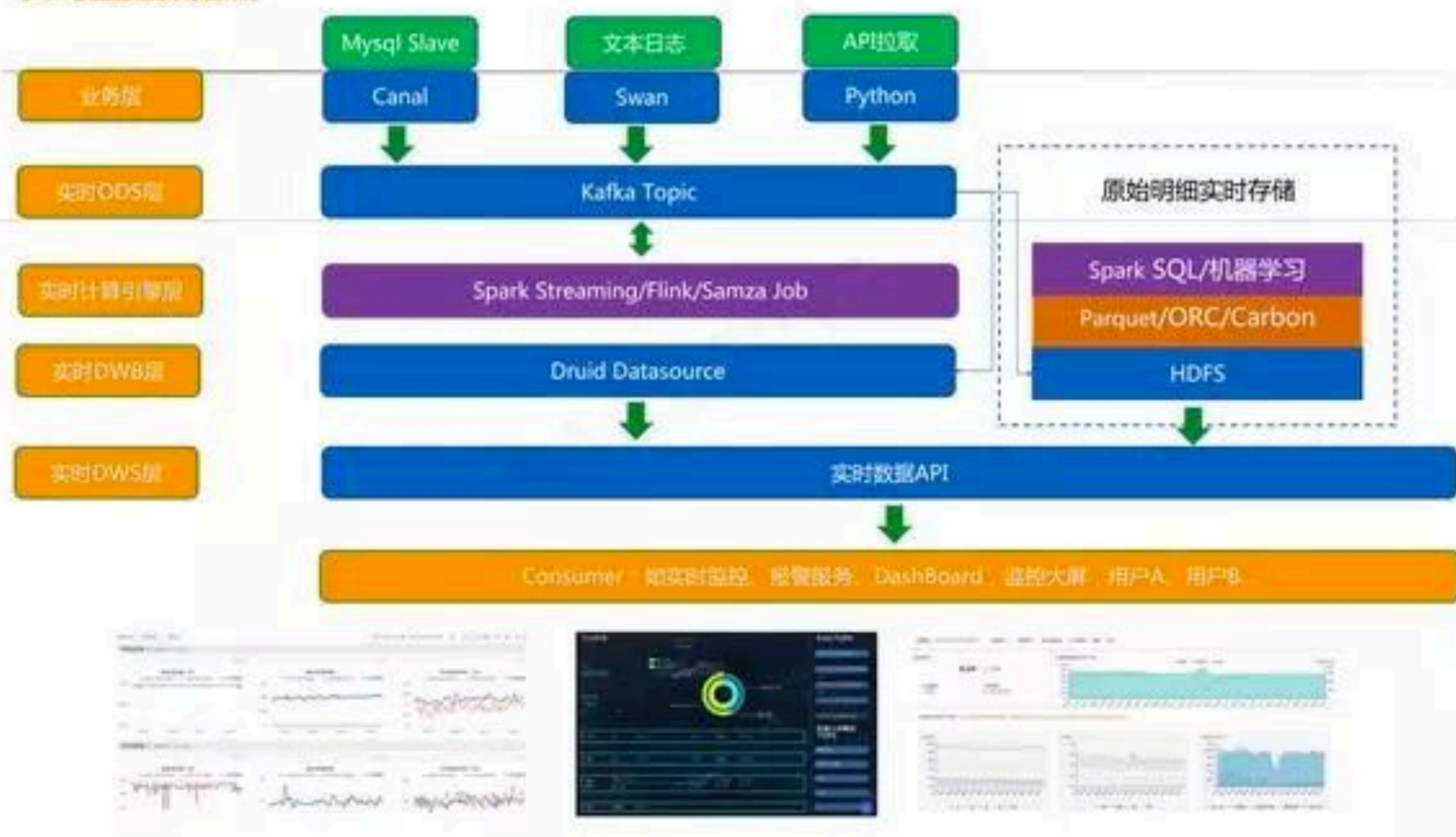
知乎业务

典型UGC、Web2.0业务

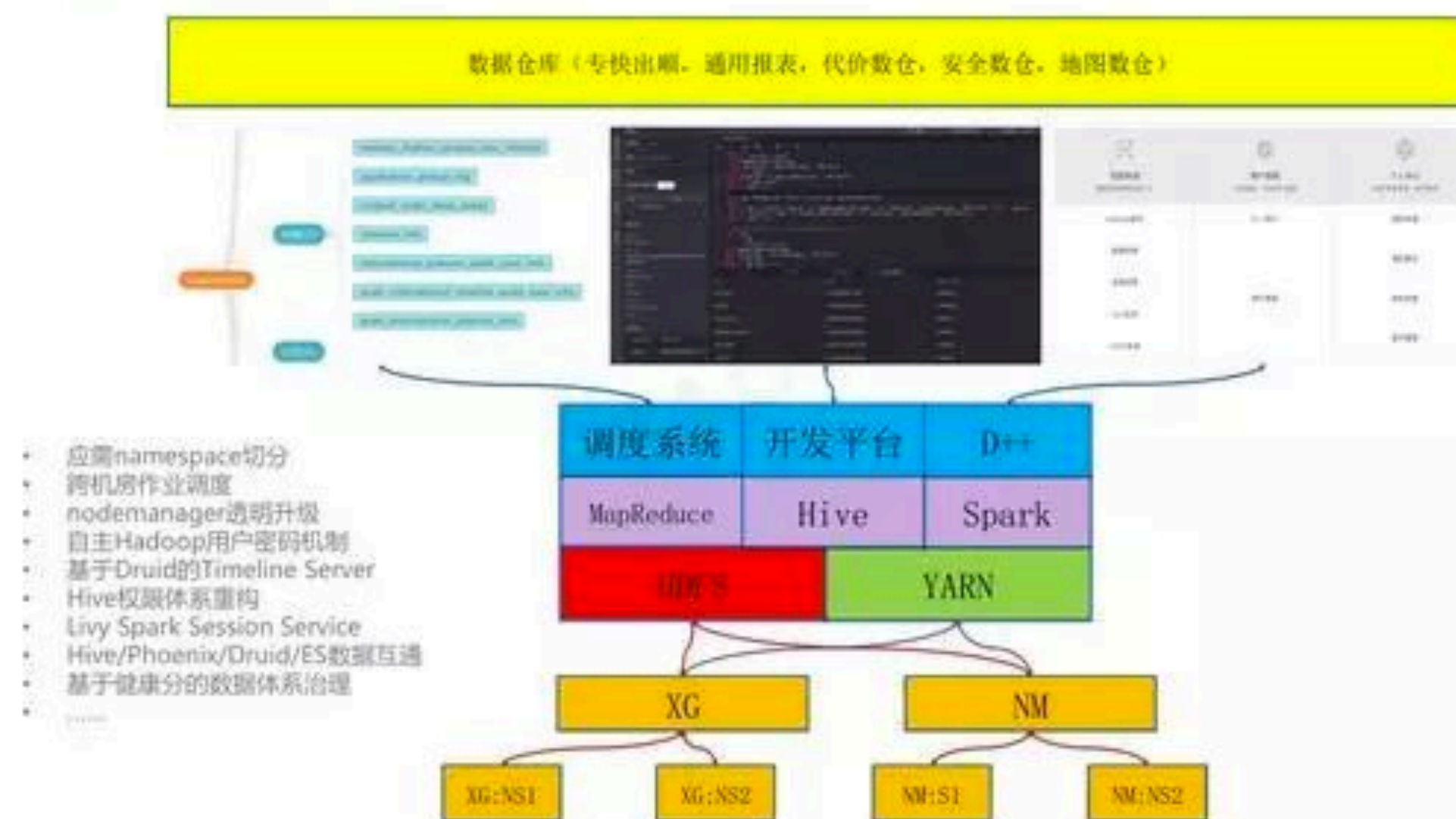
阿里业务



实时监控数据流

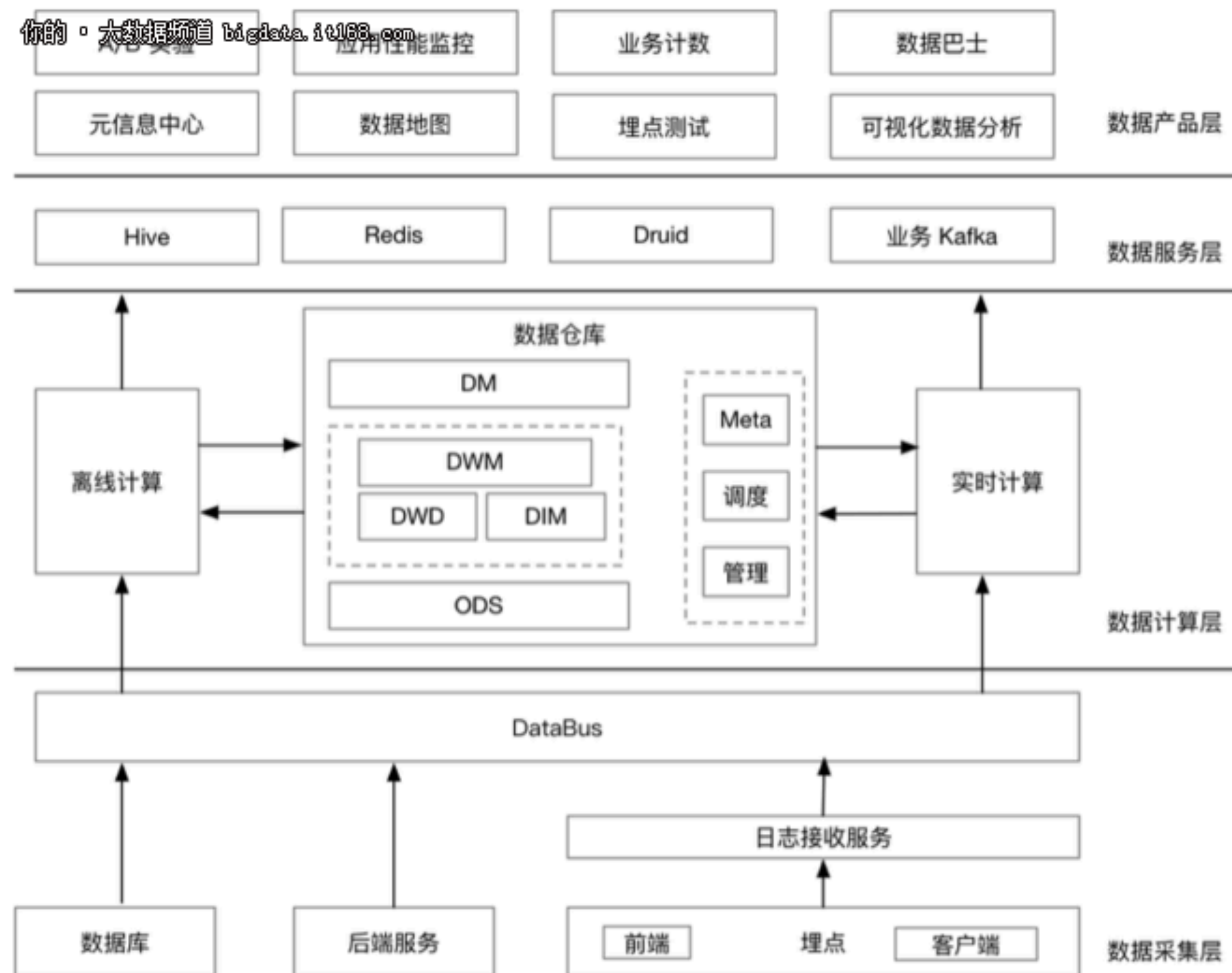


离线计算平台



滴滴业务

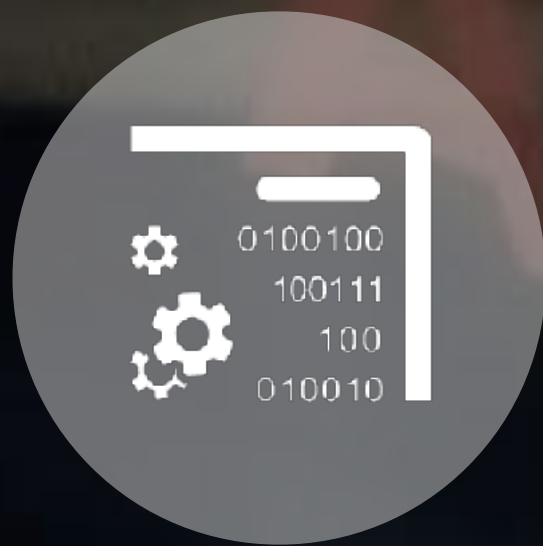
知乎业务



阿里在思考



| 计算类型多，但核心模型可列举



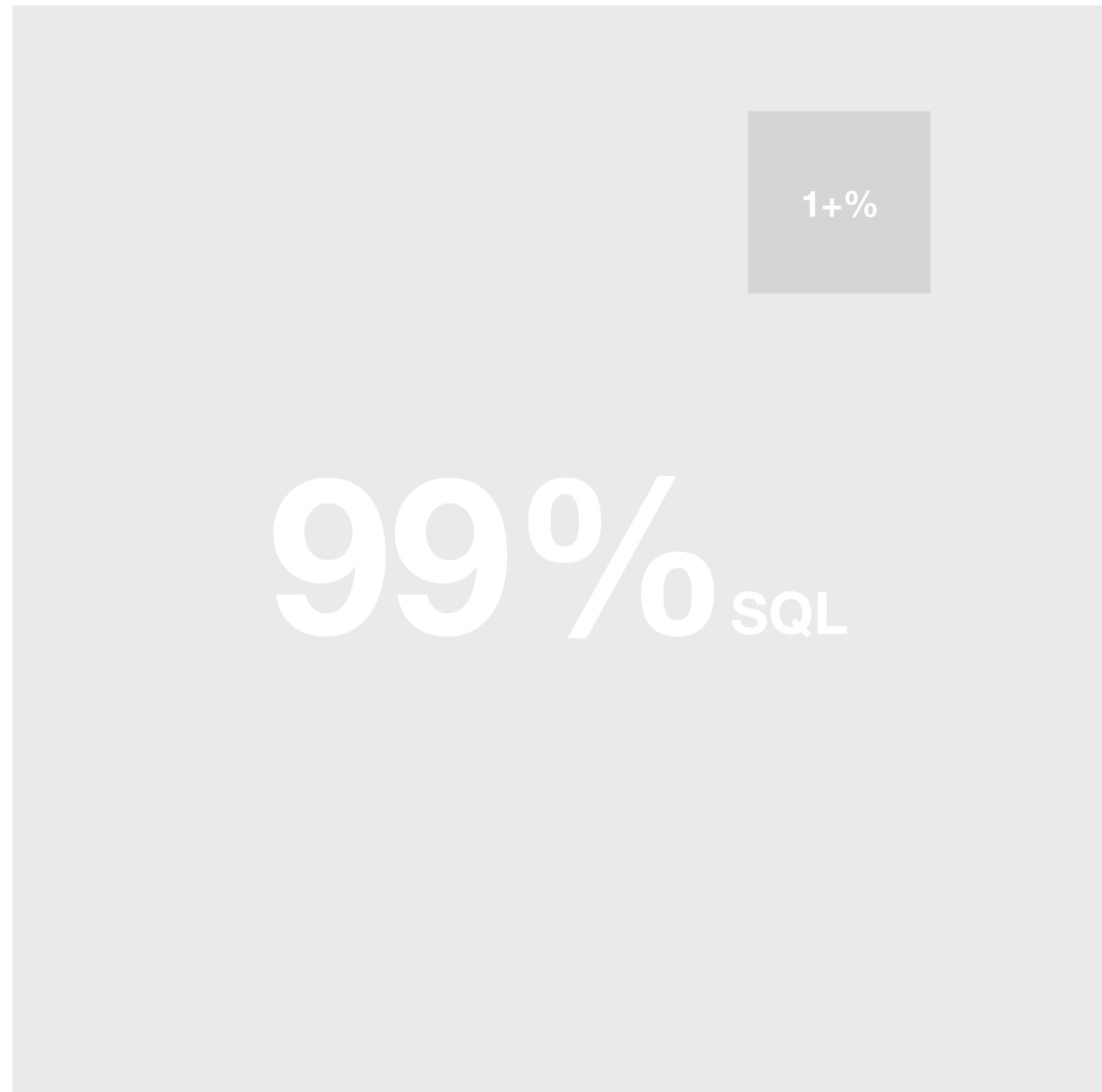
| 描述方式多，但核心抽象可列举



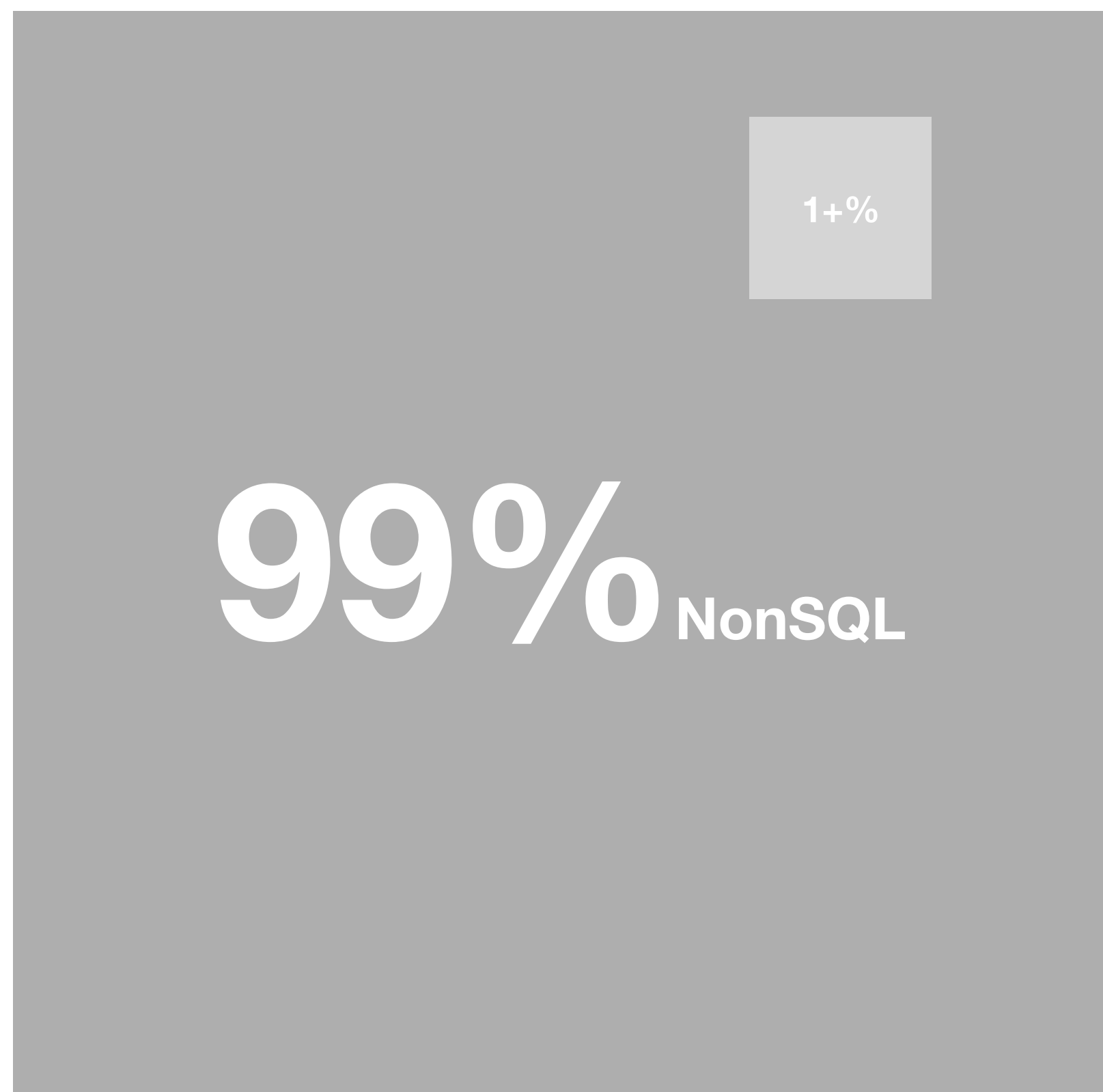
| 人参考数据决策 -> 机器计算数据决策



计算类型多种多样，但入口业务仅此数种



Batch

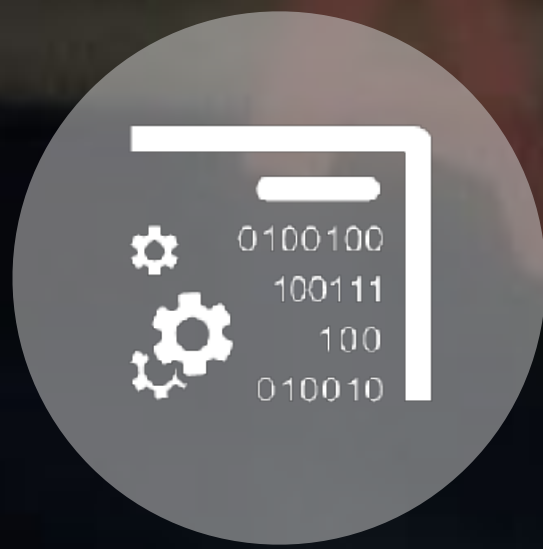


Stream

阿里的思考



| 统一计算引擎: Batch/Stream/OLAP/ML/Graph...



| 统一抽象方式: Unified SQL & API



| 统一BI+AI引擎: 数据清洗&数据训练

阿里的思考



如何构建下一代大数据处理引擎？



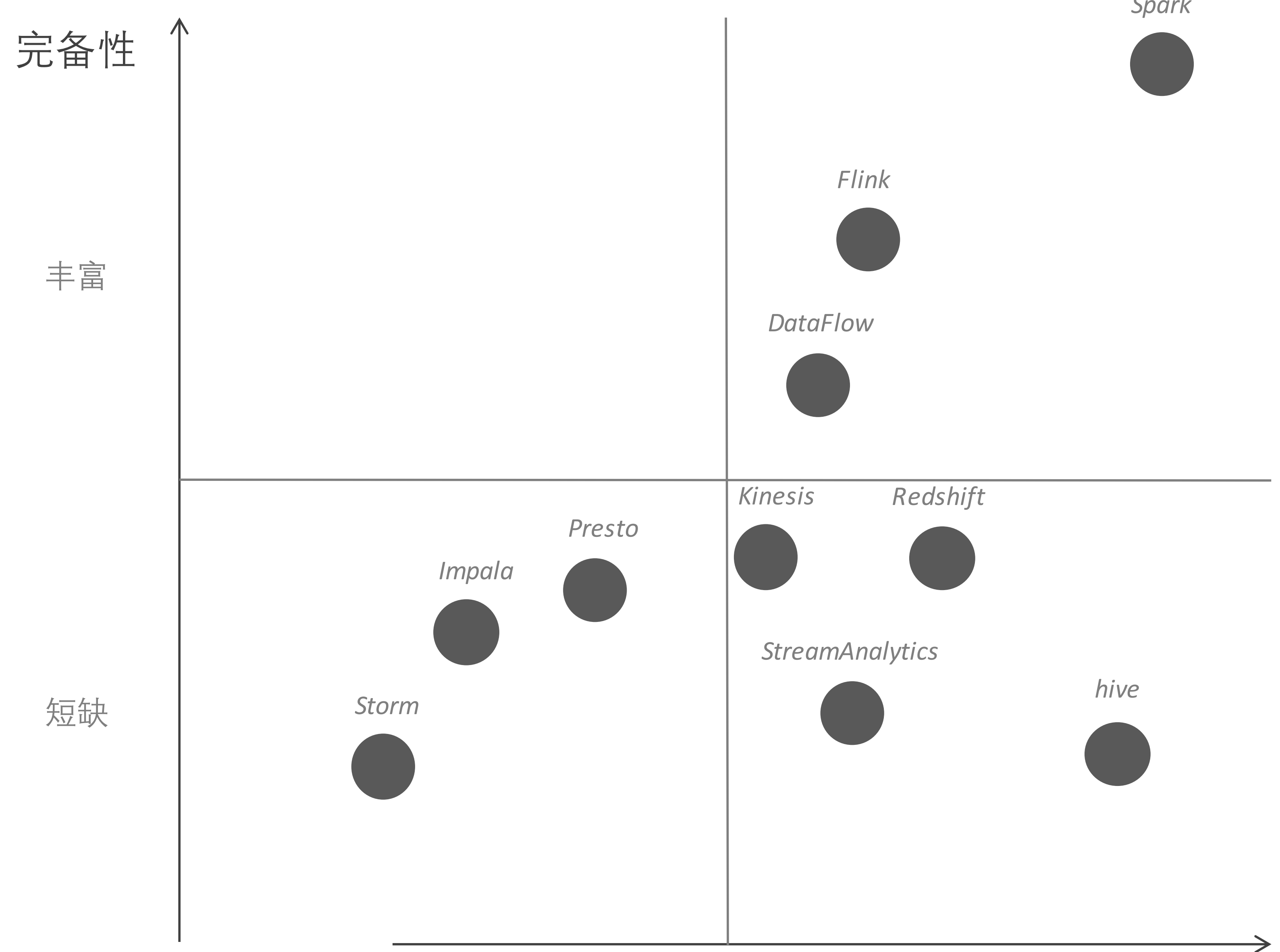
第一代



第二代



第三代



Spark优势明显

功能完备: Batch/Stream/ML
简单易用: API丰富/文档较丰富/生态对接
社区活跃: 运营出彩、大数据最活跃社区

云产品系列

产品化: 产品化成熟, 易于上手
稳定性: 稳定性好, 商业托管
社区小: 产品封闭, 不易形成社区

开源软件

社区活跃: 产品开放, 易于形成社区
稳定性差: Bug多, 不稳定
产品粗糙: 产品化初级, 上手较难



Spark期望一套软件覆盖主要计算模型，但实际覆盖不完整

21%用户认为Spark Streaming在功能(集中在窗口)和时延(亚秒)等比不上Flink，增量流式业务考虑使用Flink

17%用户认为Spark ML部分落后，包括提供更多算法、对接TF，部分业务迁移到TF框架运行

25%的用户认为当前缺乏好用、内置的上层平台，包括开发界面、工作流调度，用户使用Spark同样需要重新搭建平台系统

稳定性/调优/排错 仍未解决

31% 用户吐槽Spark集群不稳定，经常性OOM导致业务产出不稳定

另外，几乎同样客户群体(说明都是深入生产使用Spark用户)都认为Spark作业排错、调优困难，易用性不够

中文资料/社区严重缺乏，未能形成有效组织

30%用户吐槽当前文档、案例过少，特别在调优、排错方面，用户往往不知所措
用户同时认为相关中文资料相比更少，时效性也落后英文社区太多
中美语言差异导致中国市场更加空白



Flink: 下一代流式处理系统

——为什么Flink要比Storm/Spark更加优秀？



+



=





Scalability

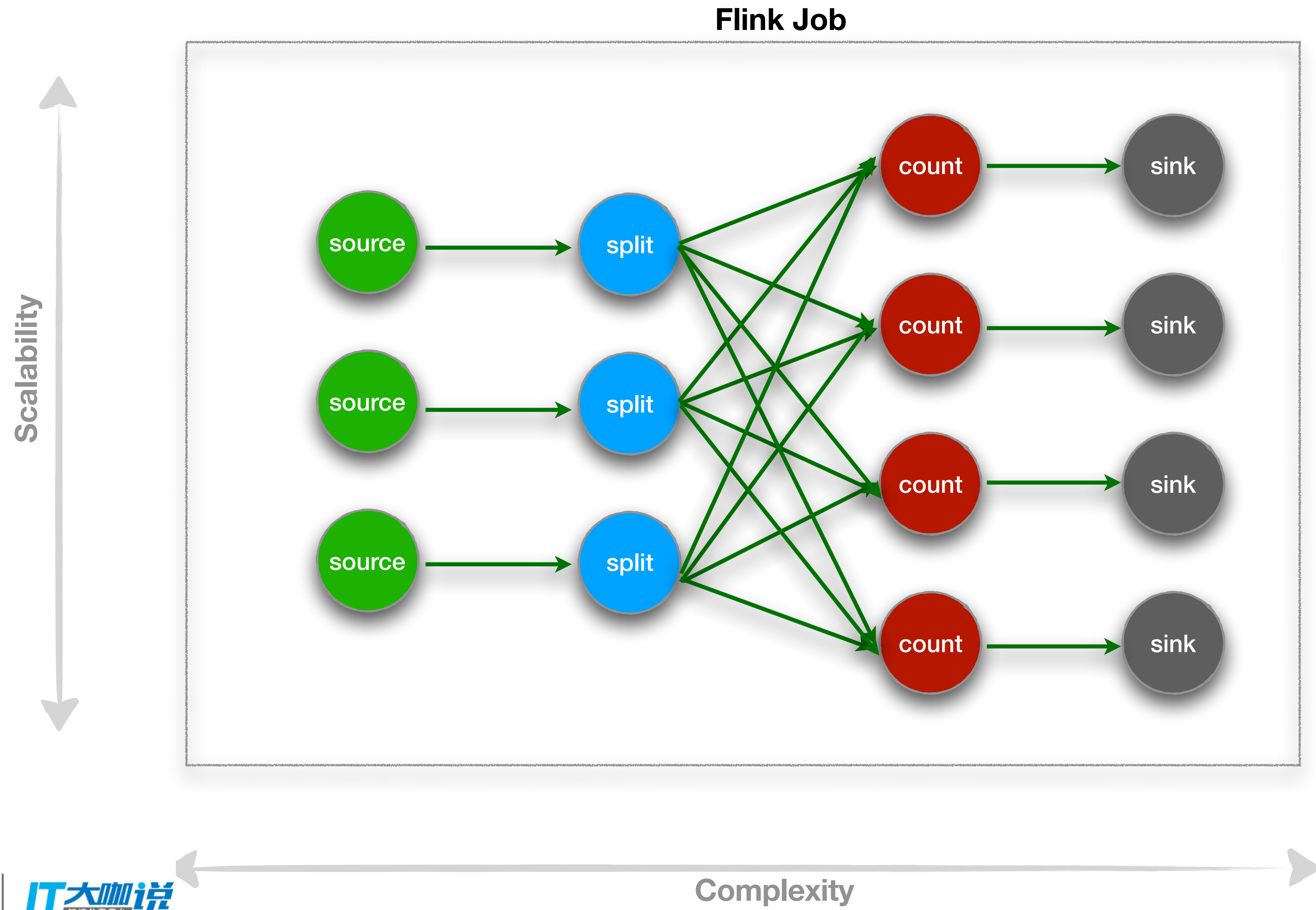


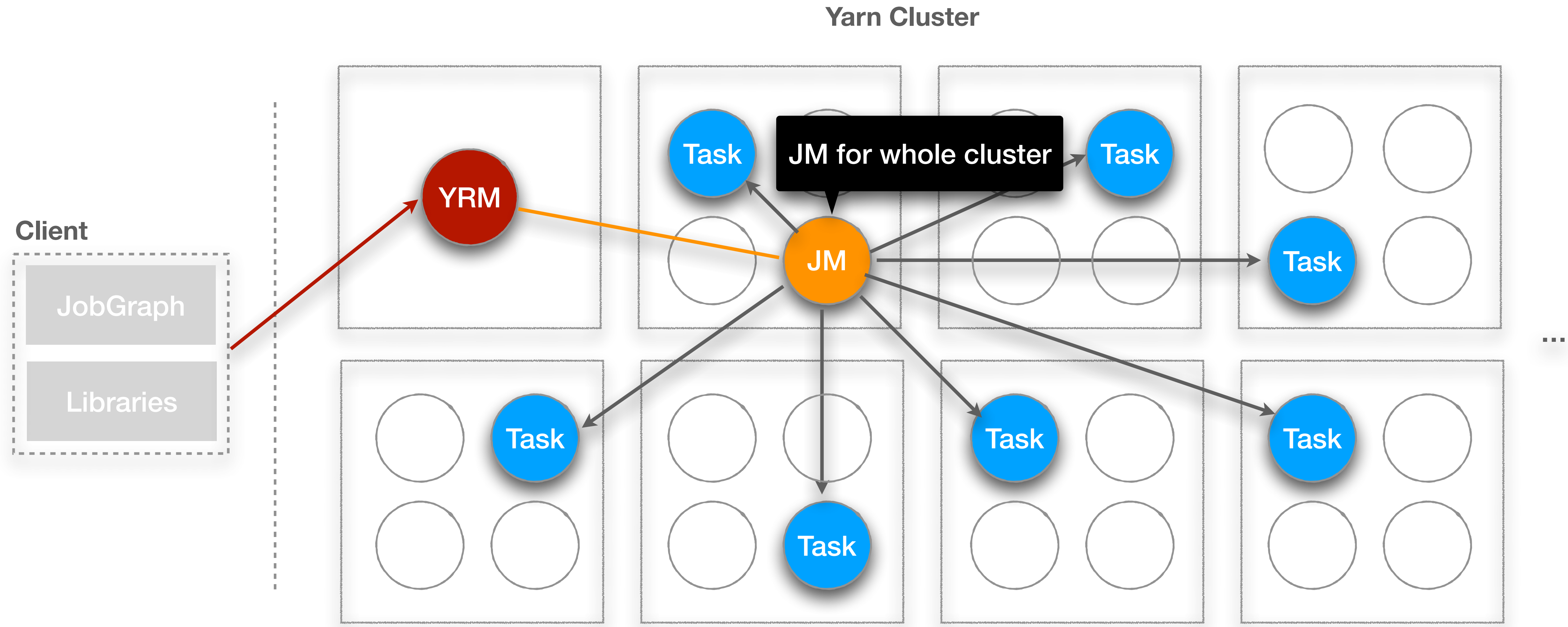
Unified

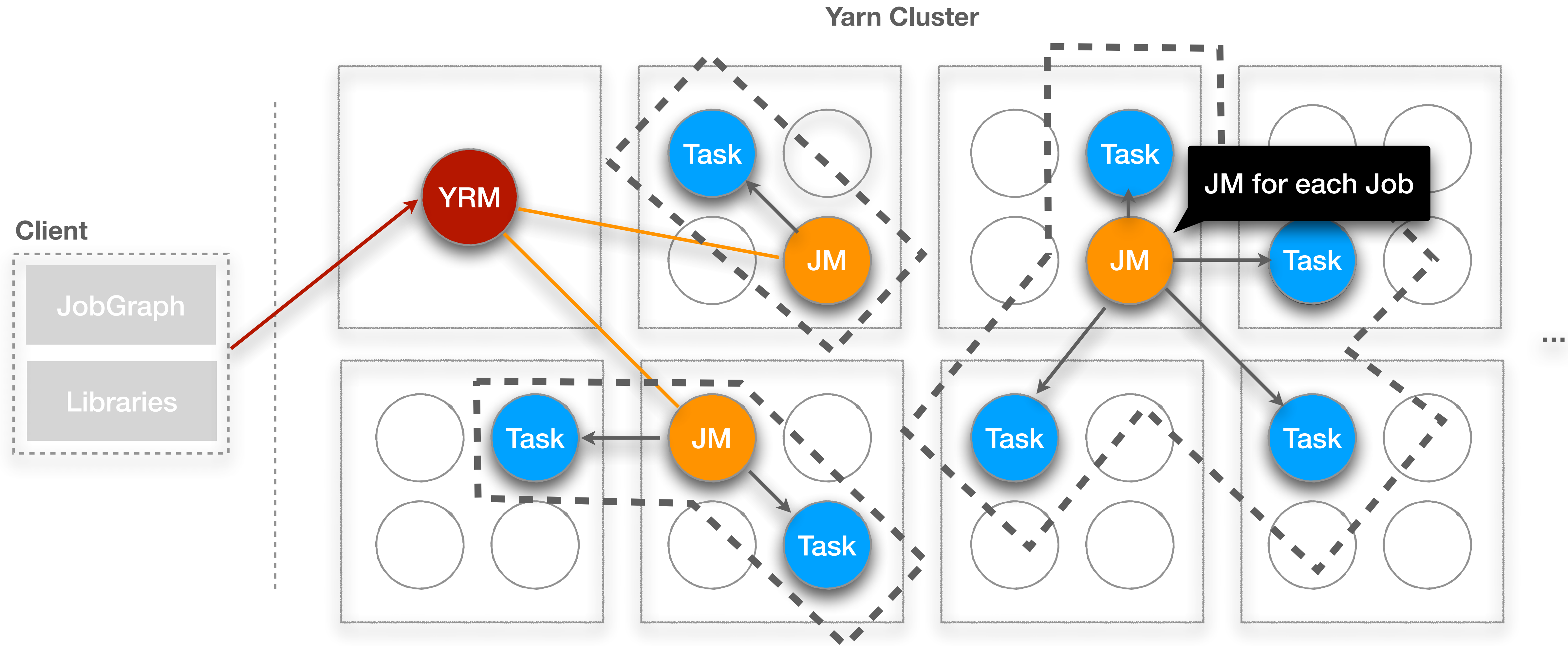


API&SQL











Scalability&Complex

FLIP-6: Deployment and Process Model

解决大规模部署+云计算隔离问题

FLIP-12: Asynchronous IO

解决流计算时延问题

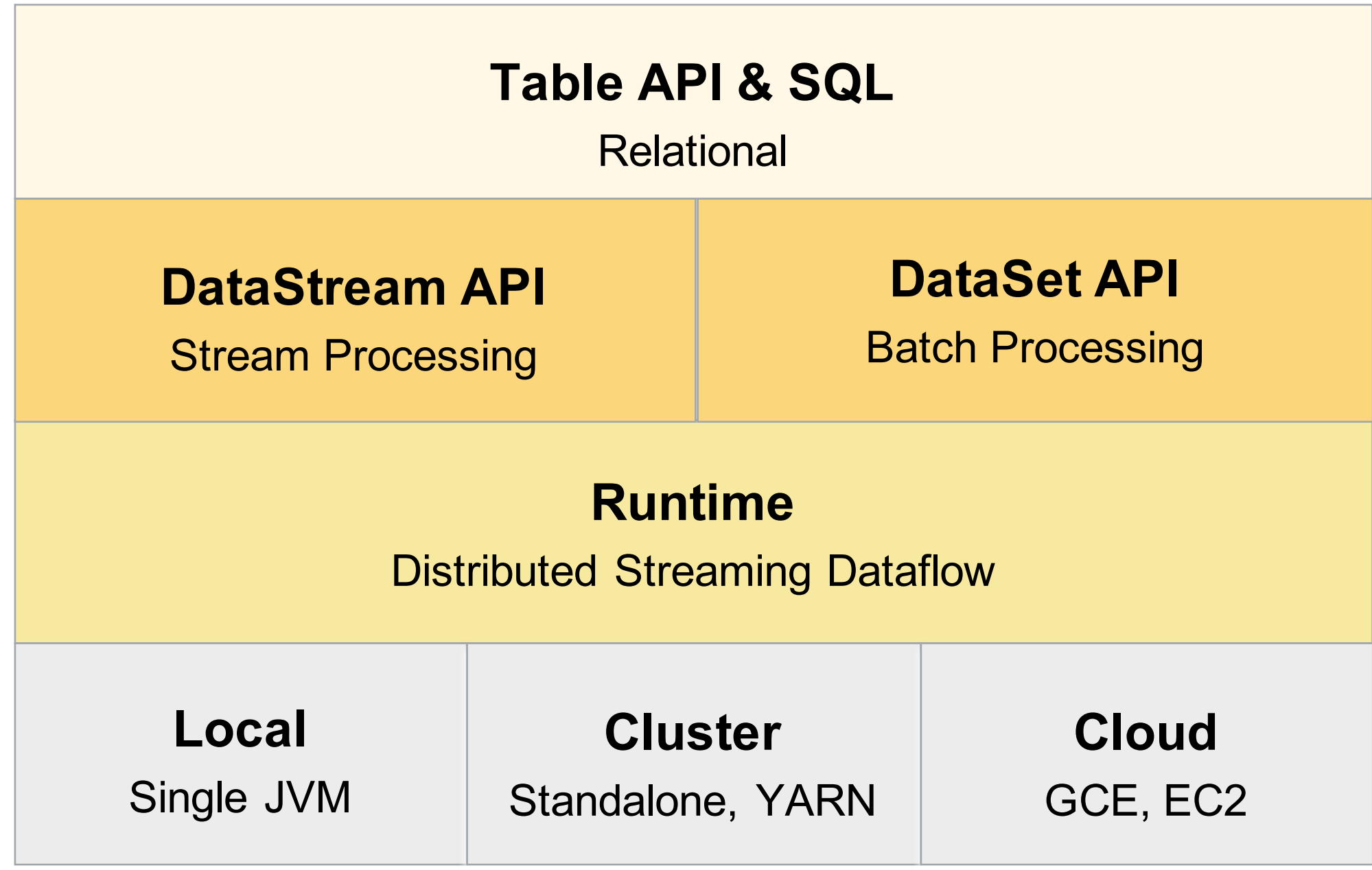
Incremental Checkpoints

减少CP磁盘IO，解决流计算时延问题

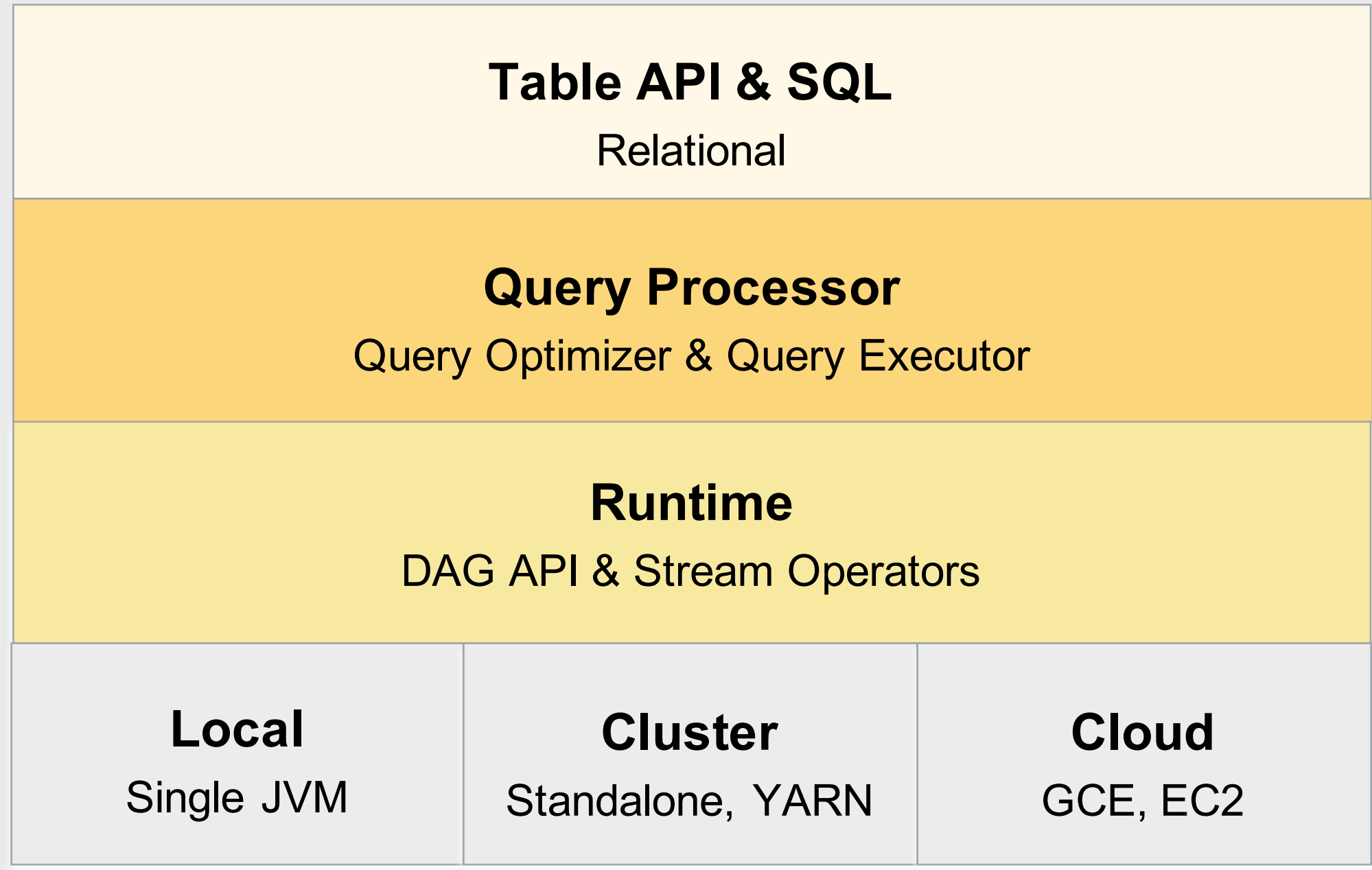
Dynamic Balance&Local-Global AGG

减少Data Skew，避免木桶效应



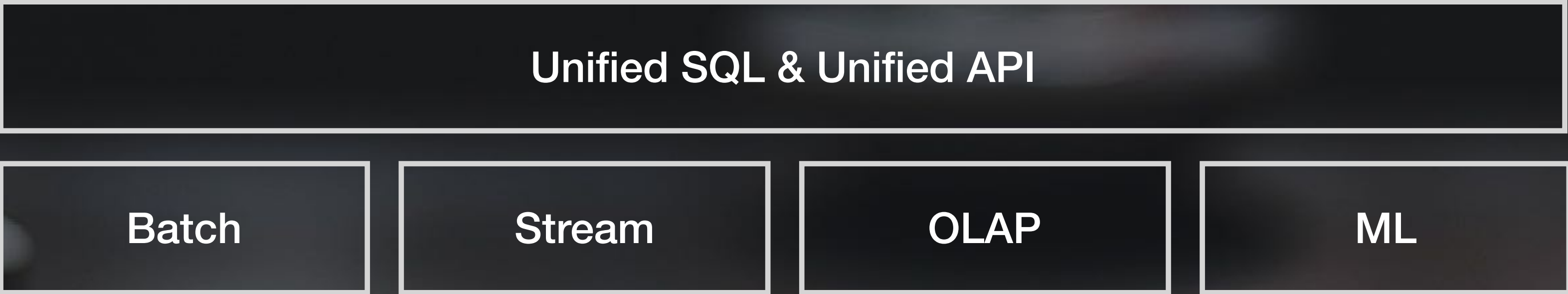


Flink Stack



Blink Stack

Unified Engine Unified API





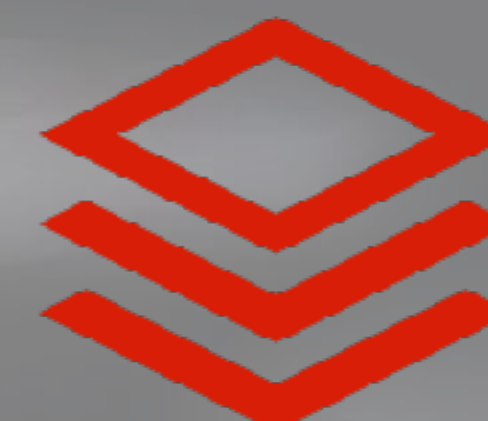
Declarative



Optimizable



Understandable



Stable



Unify

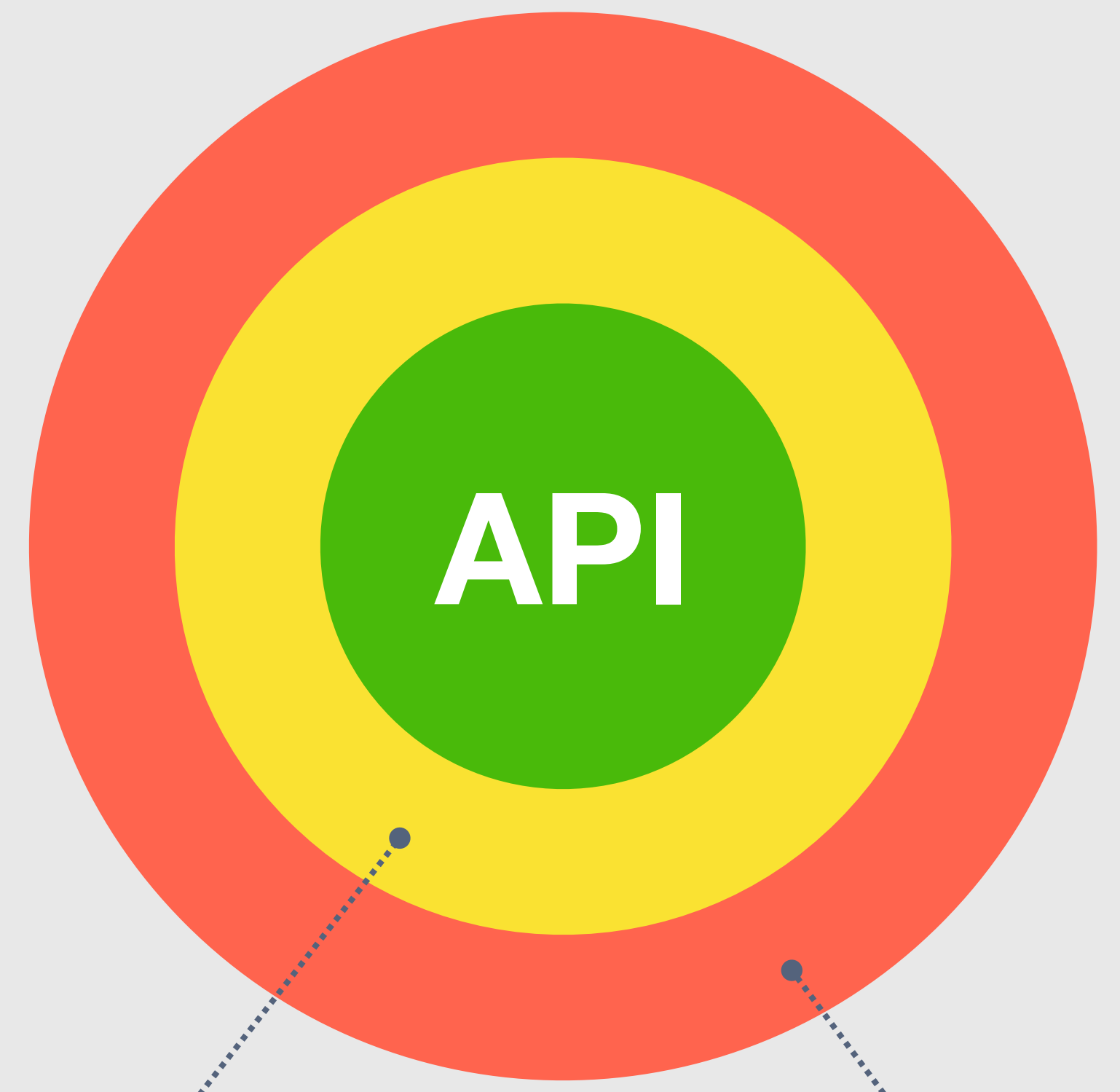


Just API?



SQL/API Only Solves Part of the Problem

We Provide More!



Blink Platform

- Development
- Debugging
- Deployment
- WebConsole
- ...

Ecosystem

- Metadata Management
- Data Quality
- ...

开源计划



(约)一个季度内Blink分支贡献给Apache Flink



Flink Forward China

12月20-21日 @北京 国家会议中心

亚洲区第一场Flink官方大会，预计3000+参会人员

Flink中文社区大牛&公司悉数到场

华为、腾讯、滴滴、饿了么、阿里等共同举办

Call For Keynotes&Sponsors

召集keynotes&Sponsors





Thanks

