

大数据和人工智能计算

王绍翯（大沙）

2018.11

2018携程技术峰会 上海



王绍翯（花名“大沙”）

阿里巴巴资深技术专家

shaoxuan.wsx@alibaba-inc.com

北京大学

EECS

美国加州大学圣地亚哥分校

Computer Engineering

博通（Broadcom）

High-Perf Platform

脸书（Facebook）

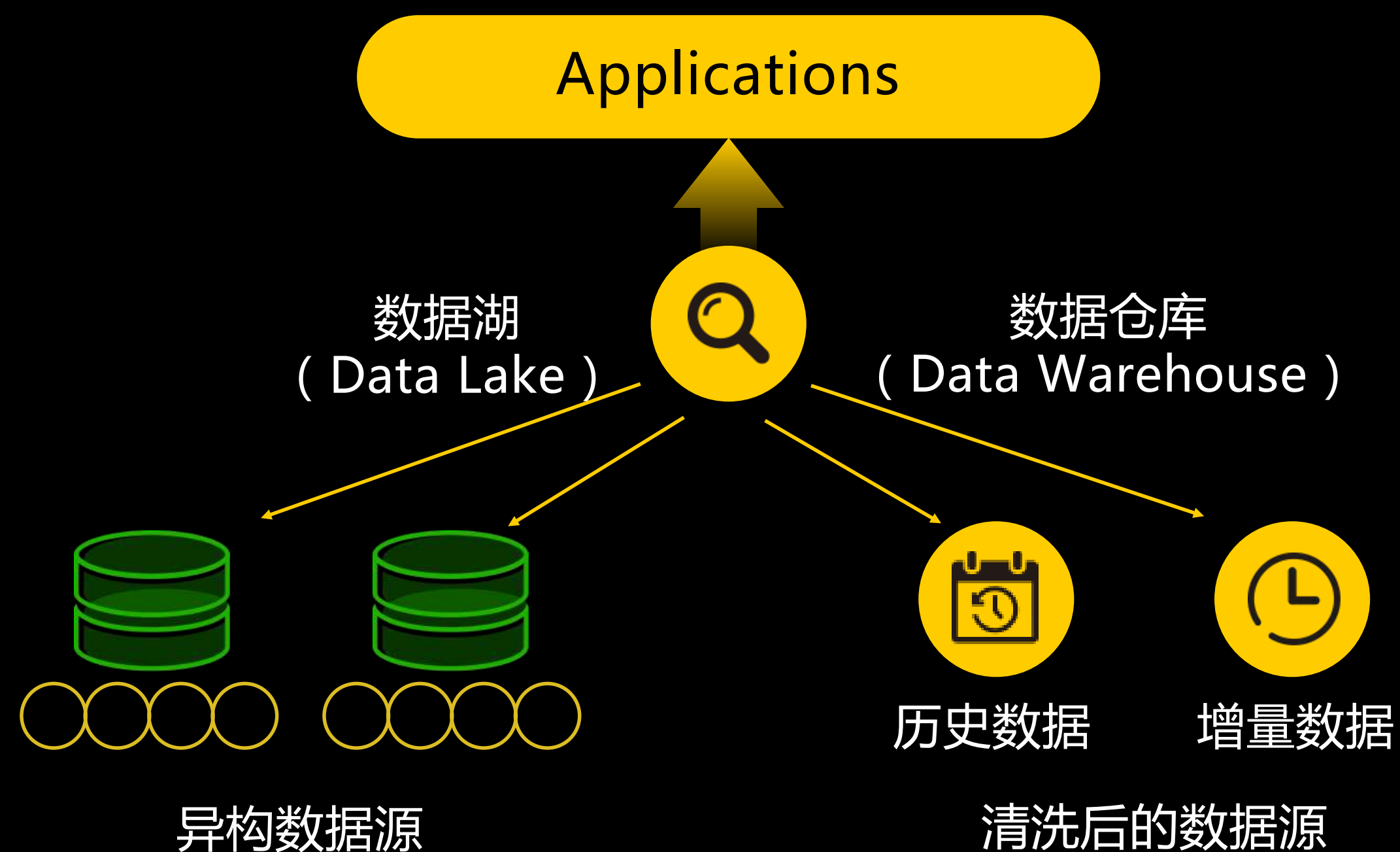
Social Graph Storage

阿里巴巴

Real-Time Data Infra

负责大数据实时计算平台
和算法平台

大数据计算的类型



大数据计算的类型



批计算

Changing Query
Fixed Data



流计算

Fixed Query
Changing Data

双十一大屏



大数据计算的类型



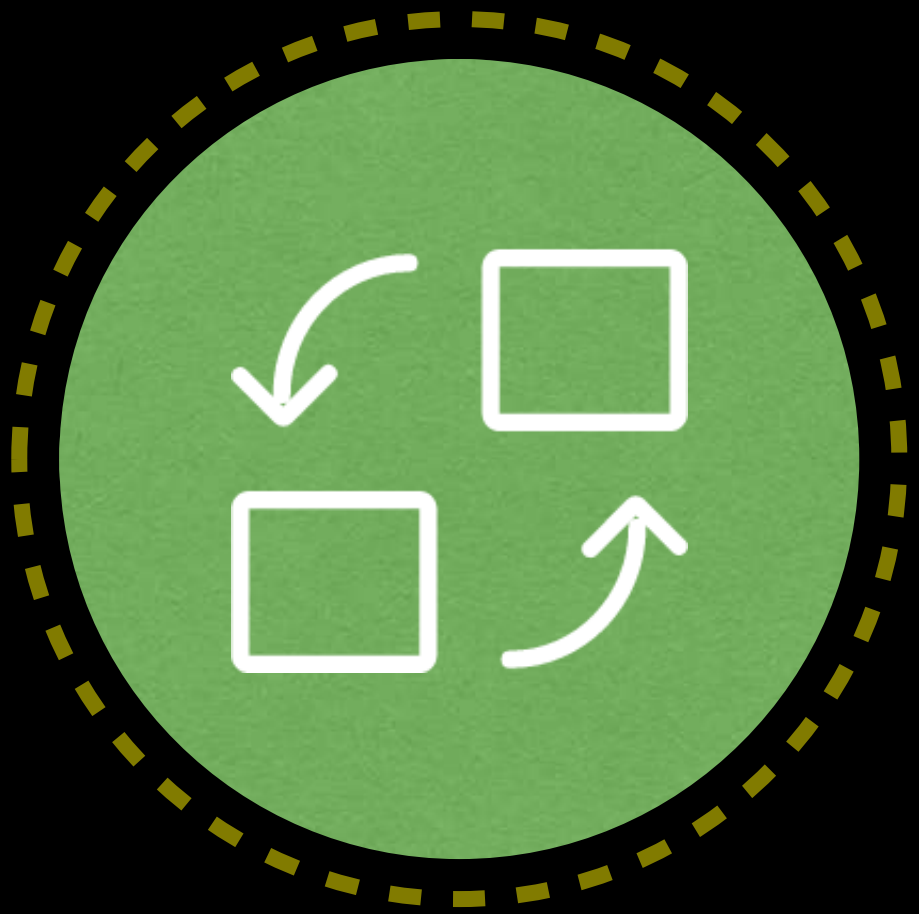
批计算

Changing Query
Fixed Data



流计算

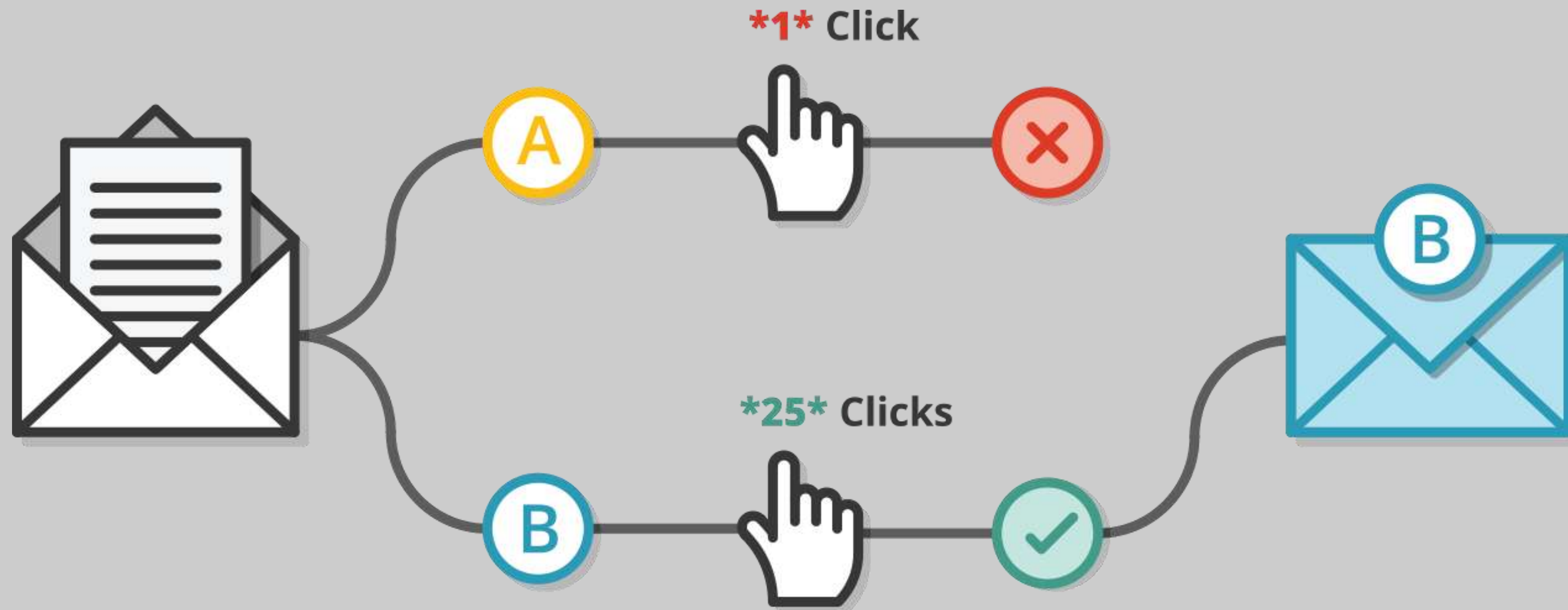
Fixed Query
Changing Data



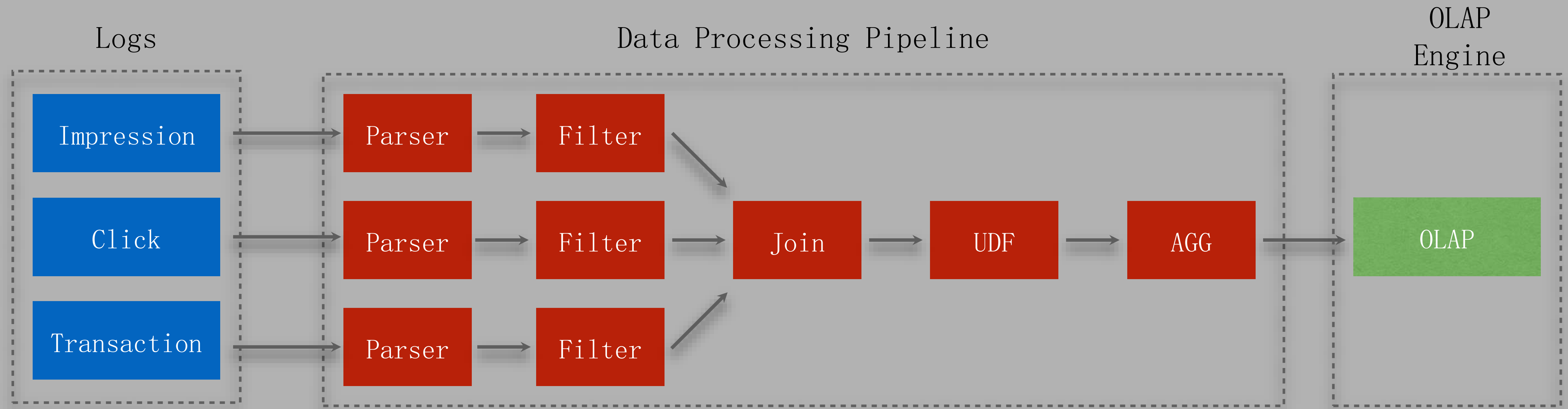
交互式分析

Changing Query
Changing Data

A/B Testing



A/B Testing



大数据计算的应用场景



数据计算 Data Lake/Warehouse Platform

提供完整的数据湖 / 数仓解决方案

Data Lake/Warehouse Solutions



城市大脑 ET City Brain

支持视频流解析&结构化处理赋能城市大脑数据处理

Support video stream and structuring process

Empower ET City Brain data processing



工业大脑 ET Industrial Brain

工业大数据实时监控&预测赋能工业大脑解决方案
Industrial big data real time monitoring and
prediction

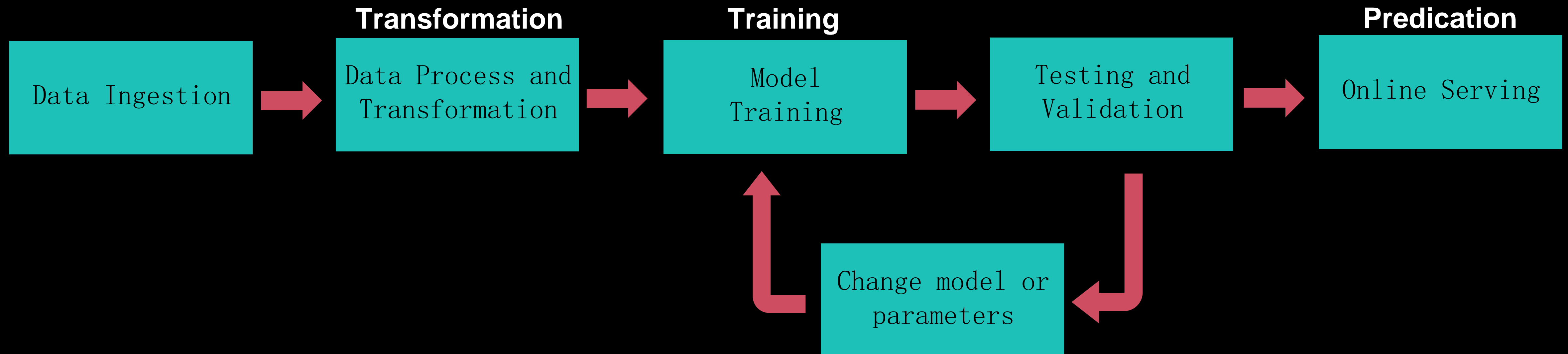
Empower ET Industrial Brain Solutions

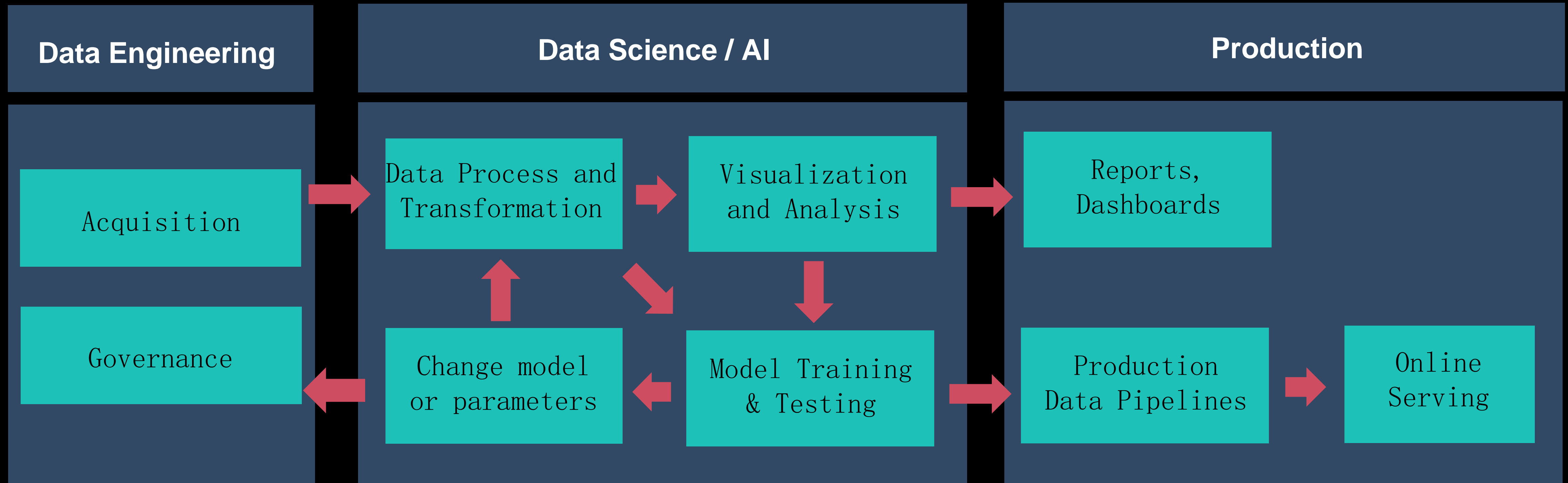


IOT

支持云上IOT数据处理实现云端计算一体化

Support IOT data process on Cloud





杭州大脑 交通数据实时监控

Traffic Data Real-Time Monitoring

🕒 10-28 19:21:05

车辆监控

Vehicle Monitoring

视频源01 视频源10 视频源12 视频源23

浙A3002R	浙AT1225	浙A61EFF
浙J3SD23	浙C6JF40	浙G7999M
浙ADM608	浙AZ1369	蒙JG1768
浙AES048	浙A22ZV5	鄂A1Z501
浙AMB39A	浙GX1886	浙ADD772

30分钟内统计: 9

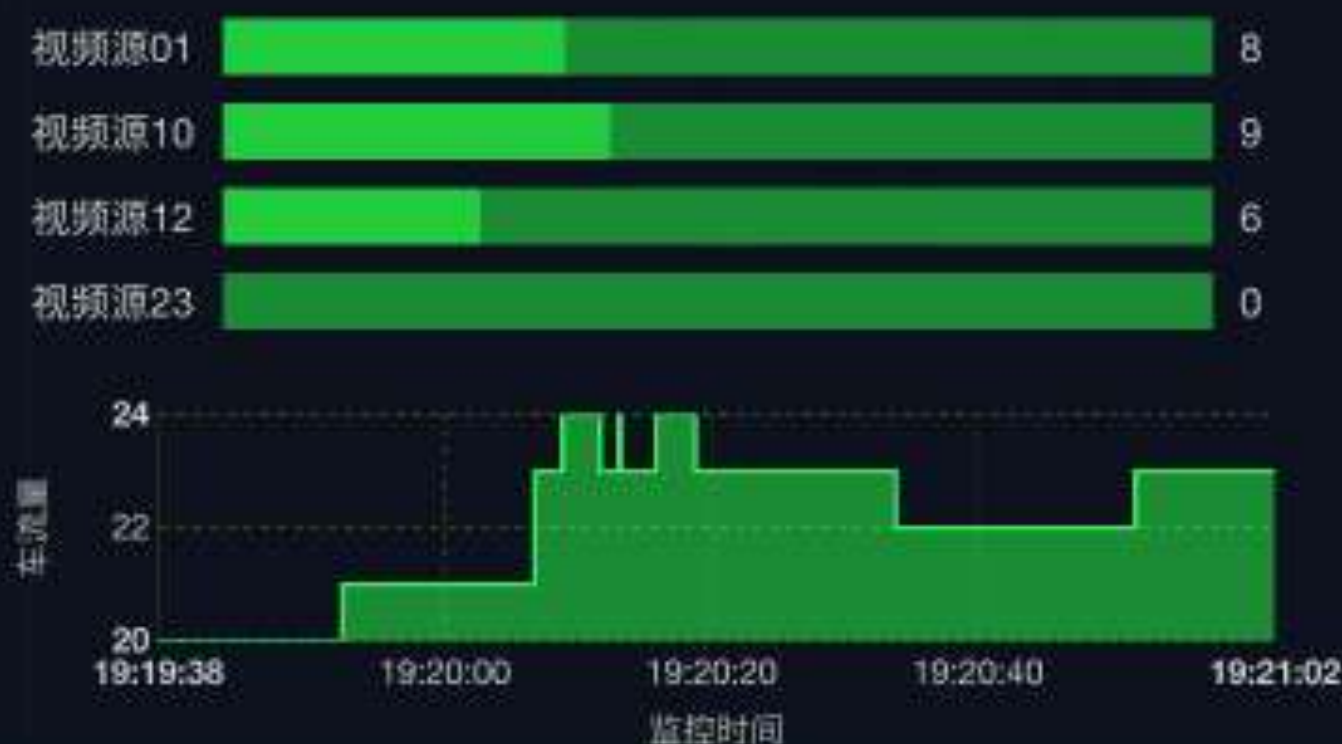
地域信息

Area Information



流量监控

Traffic Volume Monitoring



监控数据

Vehicle Information

城市大脑

杭州城市大脑使用阿里云流计算，实时监控交通状况，通过智能调度道路红绿灯、实时更新导航系统，优化城市道路拥堵情况。





传媒视频广告植入效果分析

Intelligent access control solution

使用多模态AI技术，实时分析视频
植入广告的露出节点和效果



智能客服解决方案

Intelligent customer service solution

辅助坐席客服人员了解客户意图，
提供金牌话术应答，促进销售成单



外贸客服实时同声传译

Real-time simultaneous
interpretation of foreign trade
customer service

不同语言的语音和文本实时互译，
比如外贸客服、出境旅游



传媒视频直播解决方案

Live broadcast solution

实时监控敏感图像，给出等级区分，
准确率高达99.9%



舆情监控与分析

Live broadcast solution

全方位分析互联网舆论，实时监控
舆论动态，追溯事件脉络



媒体内容安全解决方案

Smart metro solution

图片、视频，文字等多媒体的内容
风险智能识别服务



传媒视频多模态关键人物识别

Key character recognition in video

使用多模态AI技术，可以快速的锁
定目标人物出现的时间轴



智慧地铁解决方案

Smart metro solution

首个AI地铁之城，买票动动嘴、闸
机能看脸、摄像头会数数



智能门禁/闸机解决方案

Public opinion monitoring and
analysis

软硬件结合门禁方案，提升人员系
统化管理的安全性及便捷性



智能司法解决方案

Intelligent judicial solution

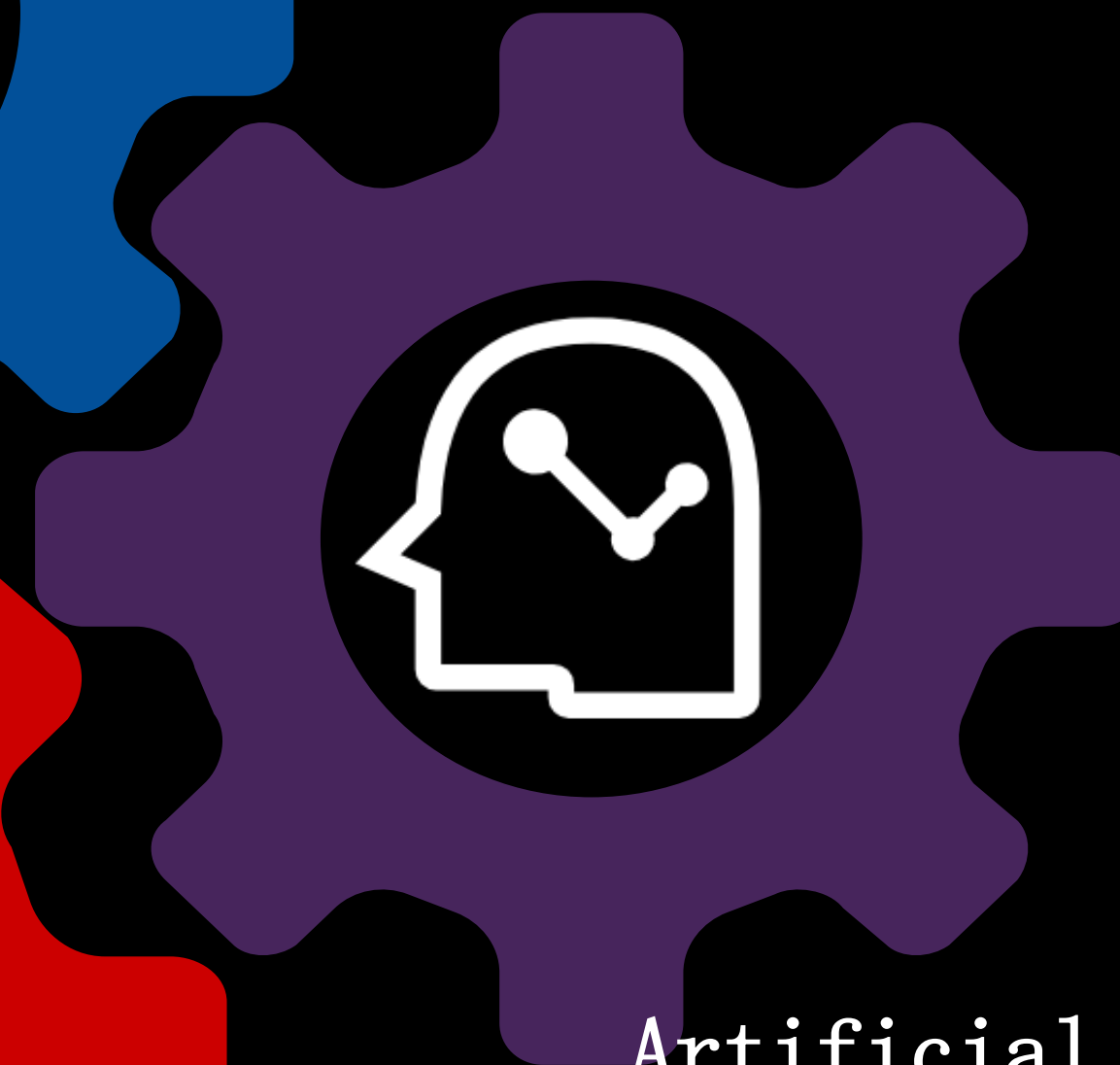
将语音识别、OCR文本识别、同案
类推和司法大数据融合

大数据智能计算

Business Intelligence



Artificial Intelligence

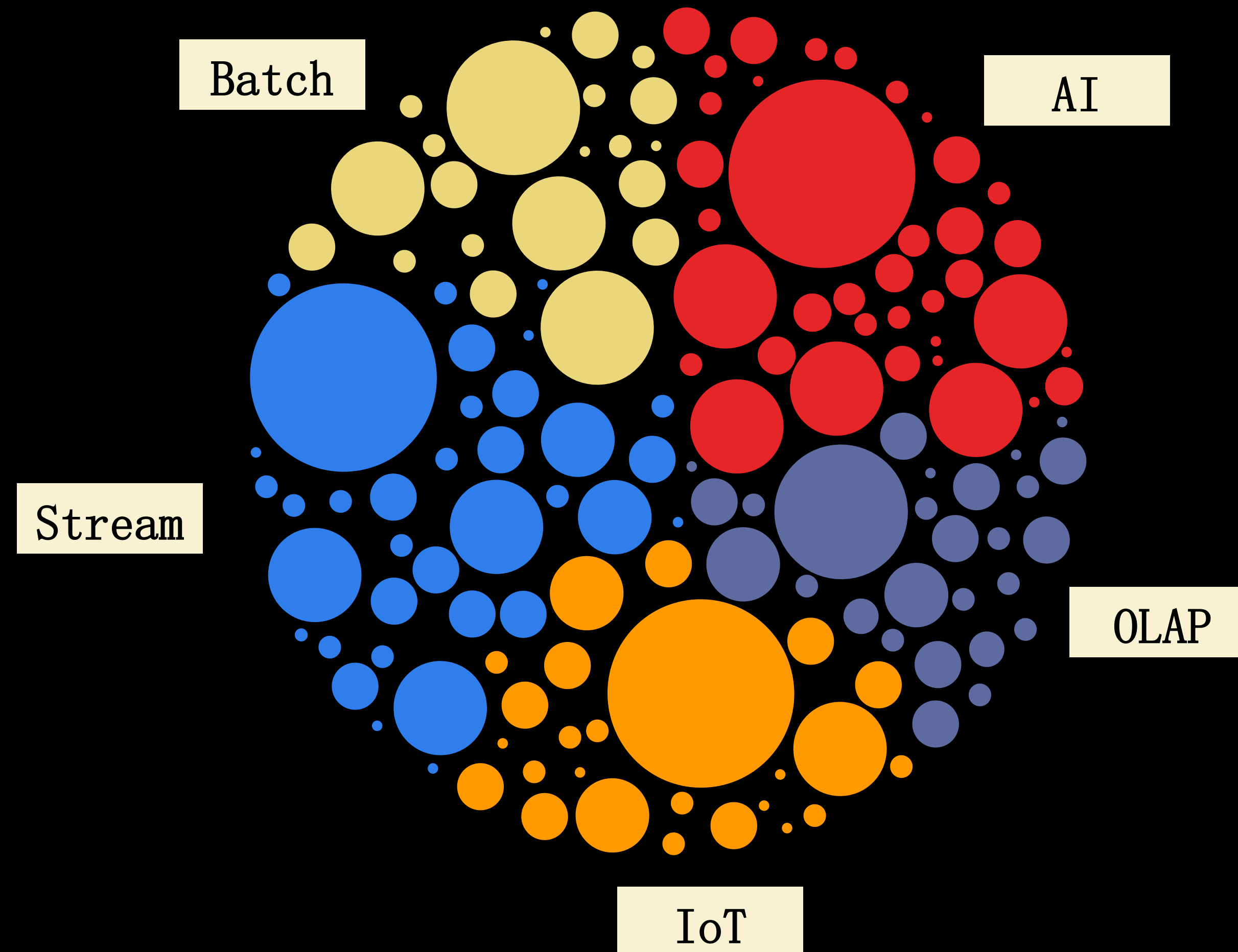


Big Data Infrastructure



BIG DATA & AI LANDSCAPE 2018

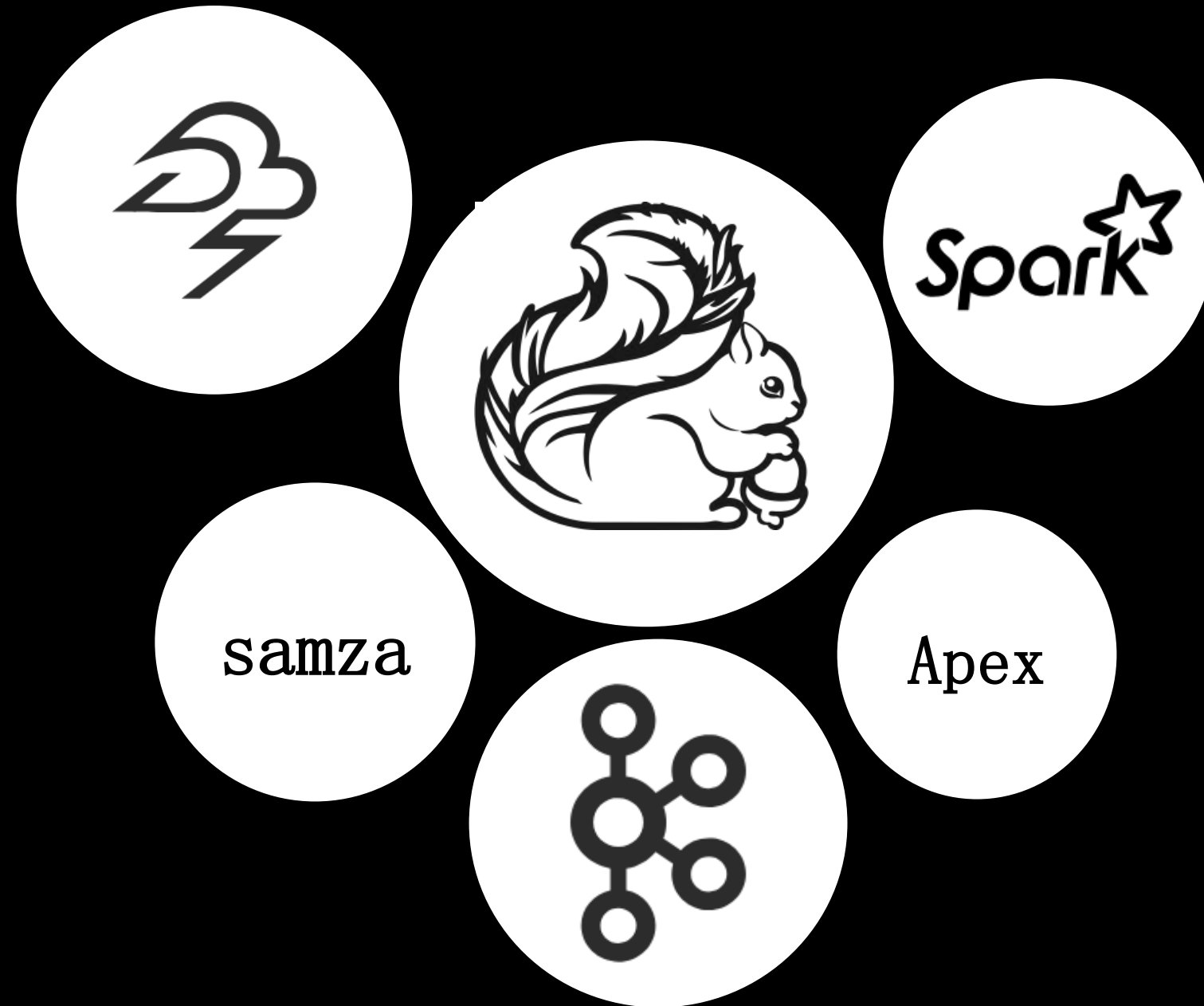




能否用一套计算引擎解决所有问题

用Apache Flink构建大数据智能平台

流计算引擎

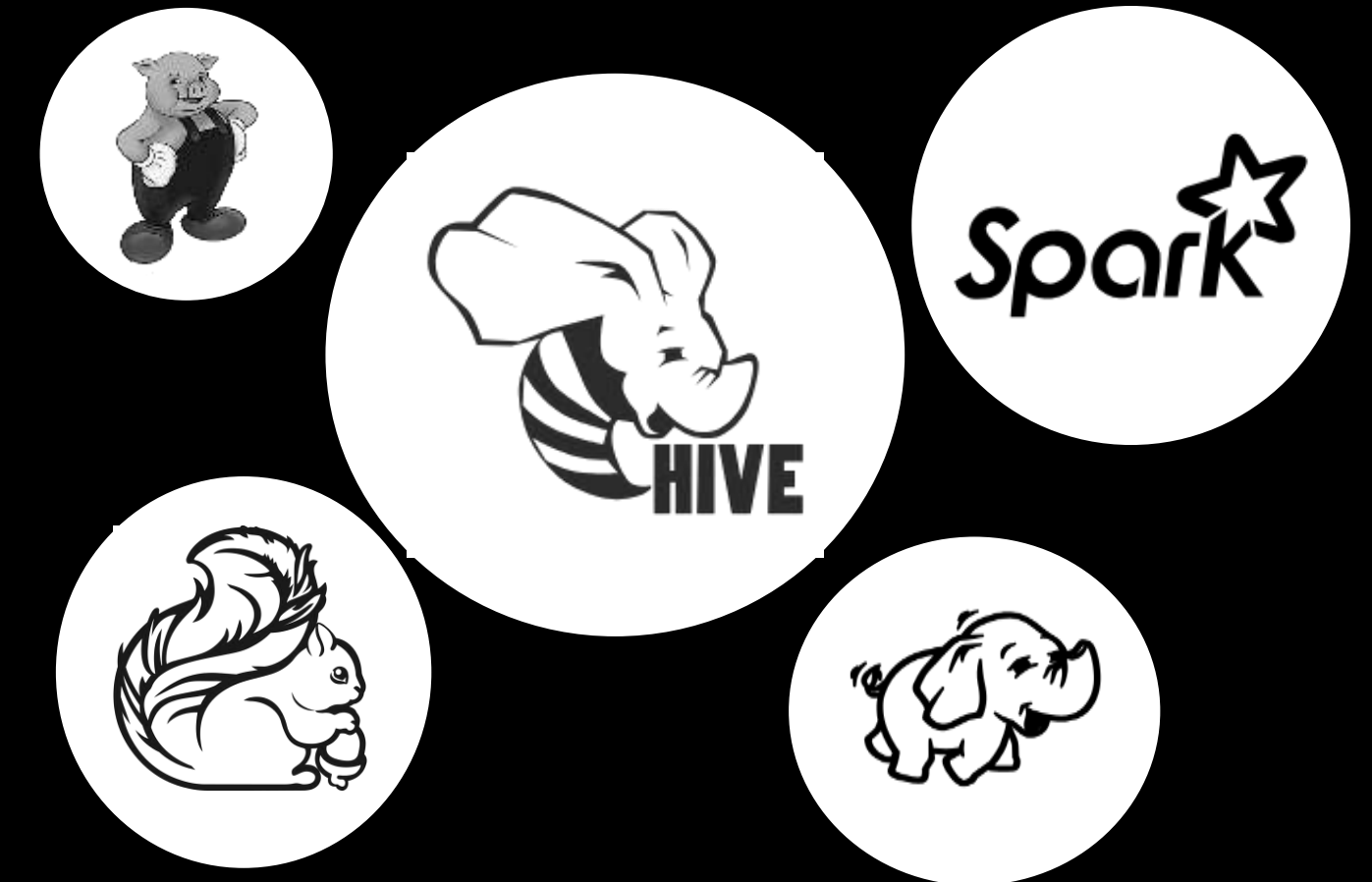


统一的大数据
智能计算引擎



Apache Flink

批计算引擎



OLAP分析引擎



AI计算引擎





Apache Flink



Alibaba's Improvements

Blink1.0:
high performance stream
compute engine

Blink2.0:
a new unified high performance compute engine for
complete data applications (including batch, stream,
AI, and IoT)



天然的Flink纯流式计算的低延迟



大规模高并发部署的优化 2016

贡献社区



快速的容错

- Incremental Checkpoint 2016
- Fine-grained Recovery 2016
- Barrier Alignment Improvement 2017

贡献社区



性能提升

- Async Operator 2016
- Credit Based Flow Control 2017
- Load Auto Balance 2017

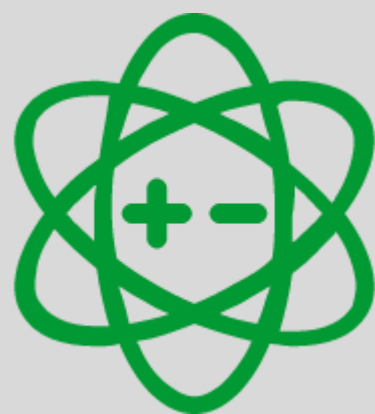
部分贡献社区



主导制定 Flink SQL 语义

- Dynamic Table 2016-2017
- Retraction 2016-2017

贡献社区



完善 Flink SQL 功能

- Aggregation, Join, Window 2017
- 跑通全部TPCH/TPCDS Query 2018

贡献社区



性能提升

- 大量的Query Optimization 2017-2018

部分贡献社区



资源配置自动化 2018

Blink2.0: unified compute engine for complete data applications

API的统一:

same query, same results

Query Processing的统一:

unified query optimization and query execution framework

应用的统一:

switch between batch processing and streaming seamlessly

批计算和流计算API的统一

Batch Processing

correctness



return one final result

VS

Stream Processing

real-time



emit results as early as possible

VS

in stream processing, it **emits** intermediate results,
and keeps **refining** the results to ensure **correctness**

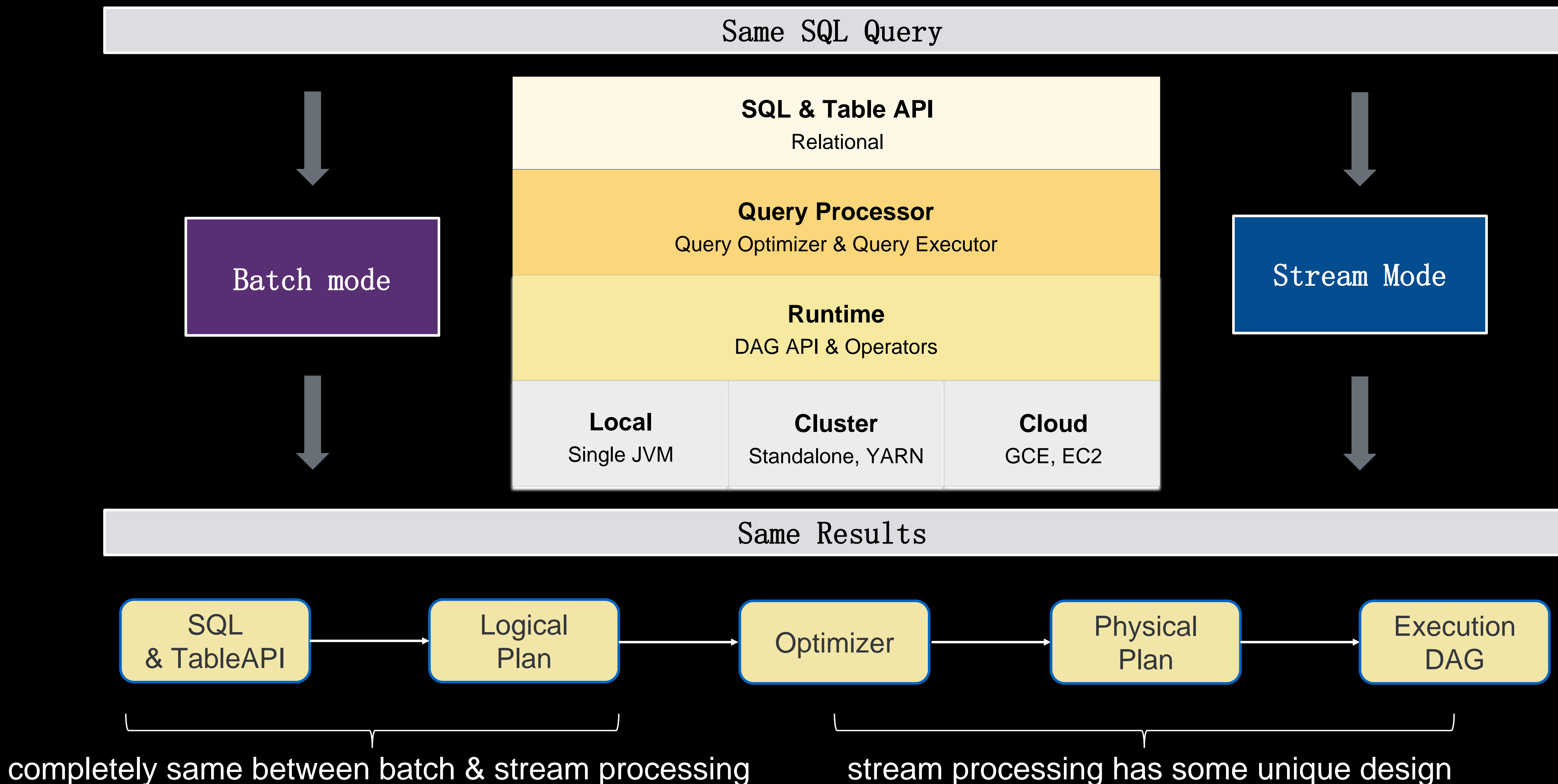
批计算和流计算API的统一

WHAT & HOW: results are calculated	-----	Can be fully described by SQL
WHEN: to emit a (intermedia) result	-----	Does not affect business logic
HOW: to refine the results	-----	Can be solved by SQL engine

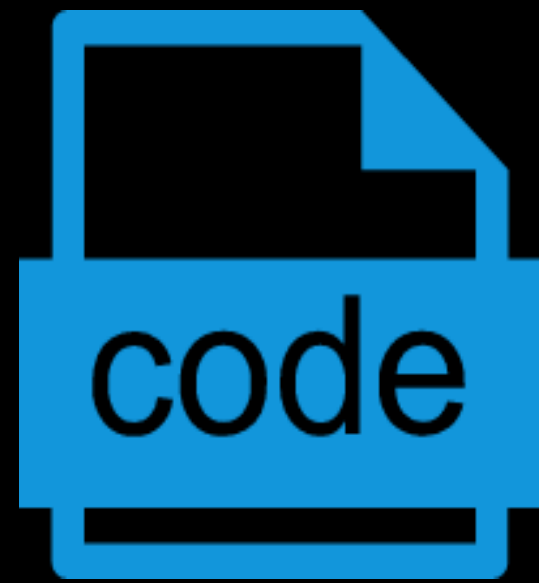
ANSI SQL can Describe Stream Processing

可以用SQL统一流和批的计算

Blink2.0 统一了SQL Engine的QP



Optimizations for Batch Processing in Blink2.0



Execution
Optimizations



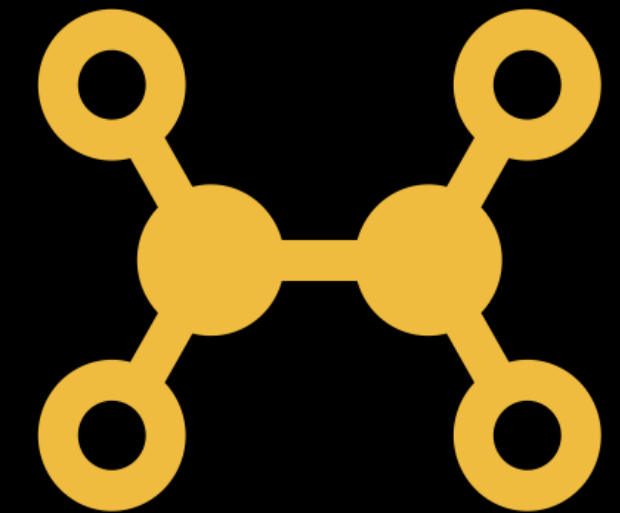
Query
Optimizations



Customizable
Job Scheduling



Batch Friendly
Failover



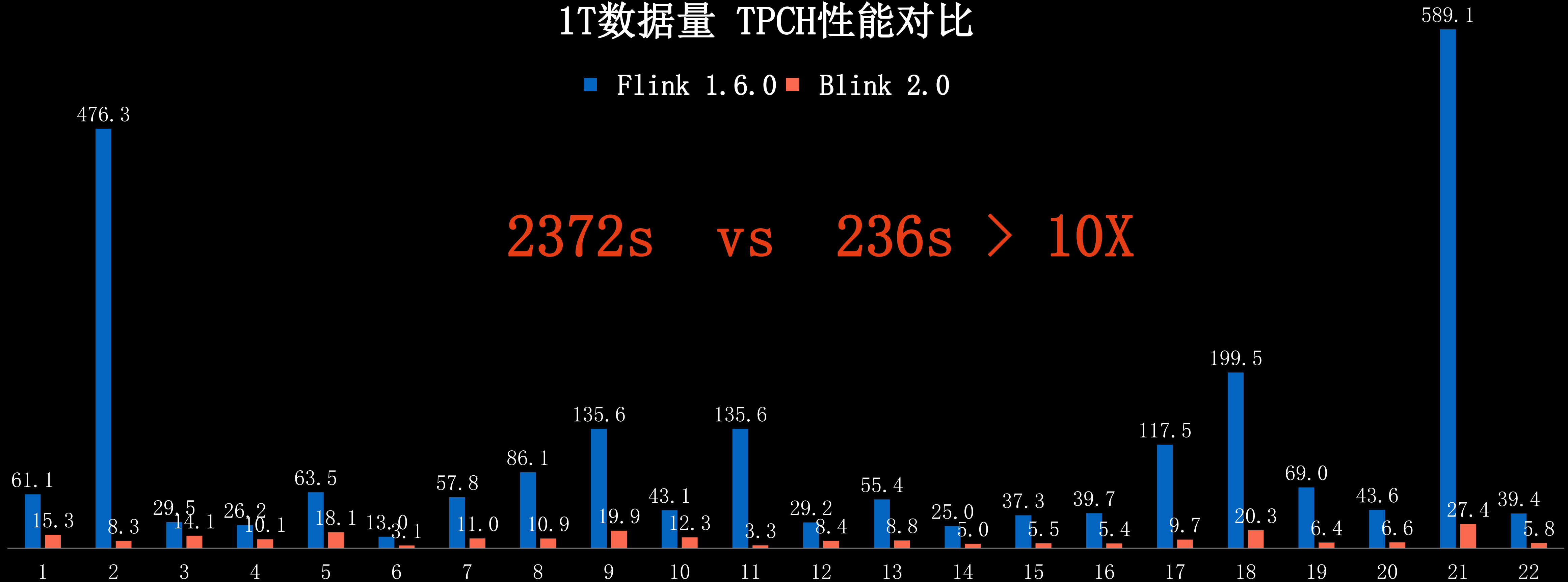
Various
Shuffle Service

Optimizations for Batch Processing in Blink2.0

1T数据量 TPCCH性能对比

■ Flink 1.6.0 ■ Blink 2.0

2372s vs 236s > 10X





Next Steps

Grand Unification of Data Processing

Switch between batch processing and streaming
seamlessly

Flink Machine Learning/AI

PyFlink, TableAPI, TensorFlow, Julia, Flink ML

Improvements



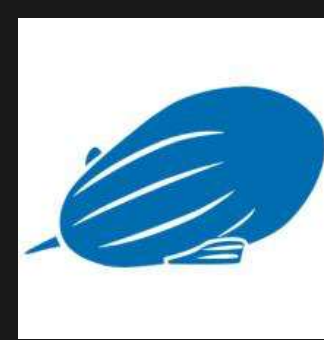
Ecosystem

Hive,

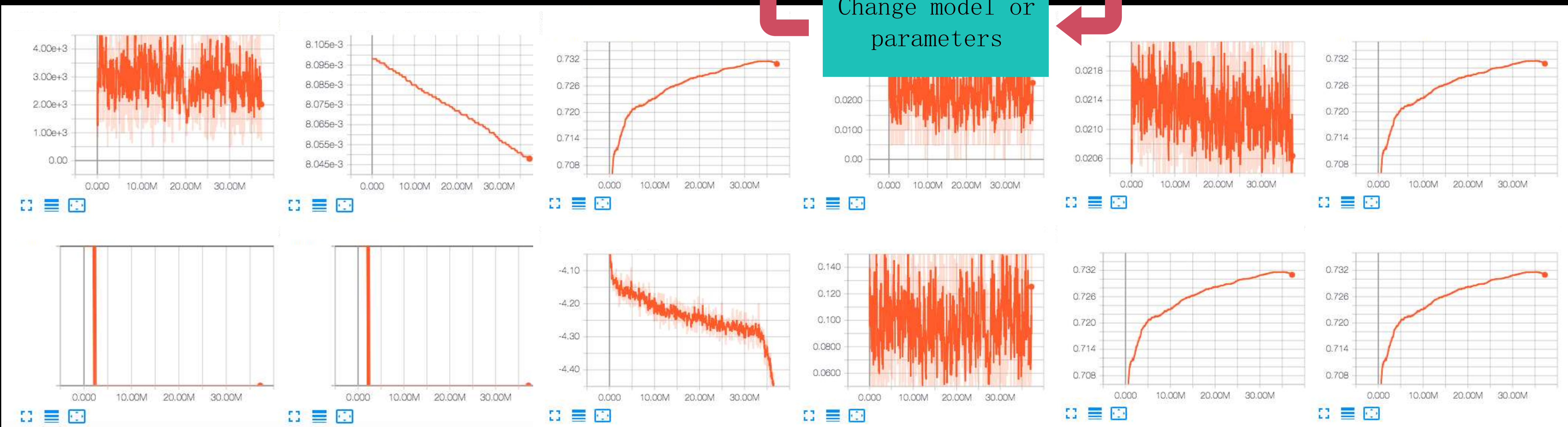
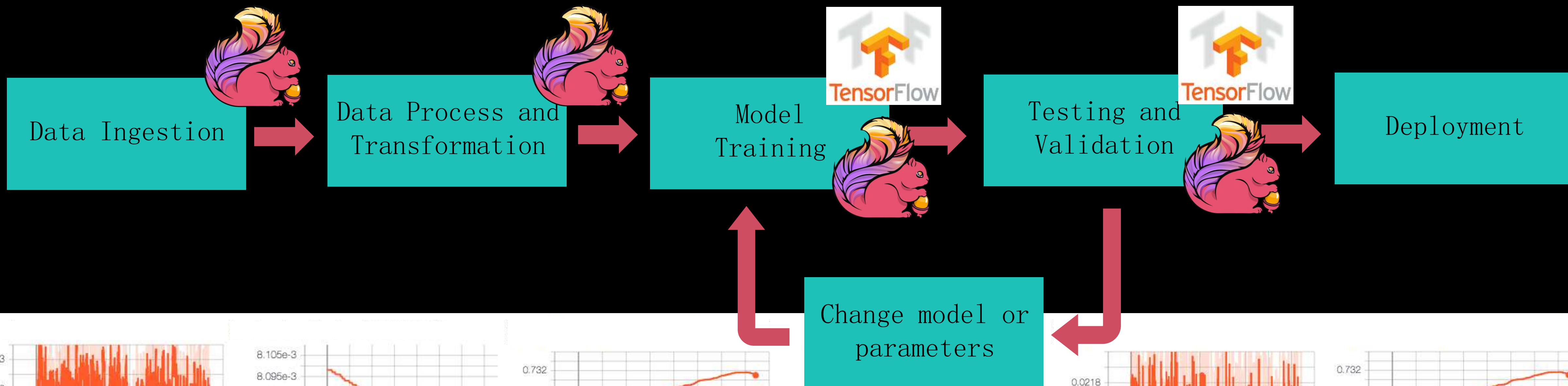
Zeppelin,

Jupyter,

Livy



TensorFlow on Flink





Home | Location ▼

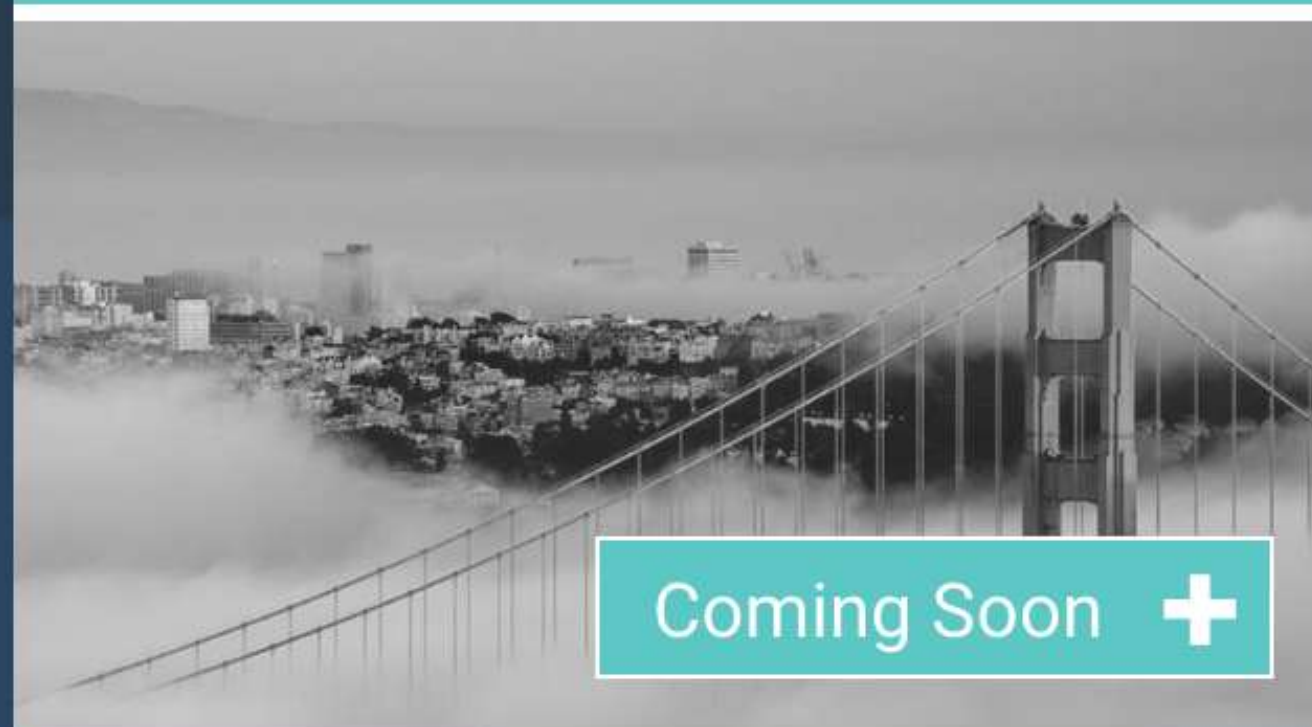
The Apache Flink® Conference

Stream Processing | Event Driven | Real Time

Flink Forward - Beijing
December 20-21, 2018



Flink Forward - San Francisco
April 1-2, 2019



Flink China社区大群



 扫一扫群二维码，立刻加入该群。

本PPT来自2018携程技术峰会

更多技术干货，请关注“携程技术中心”微信公众号

