



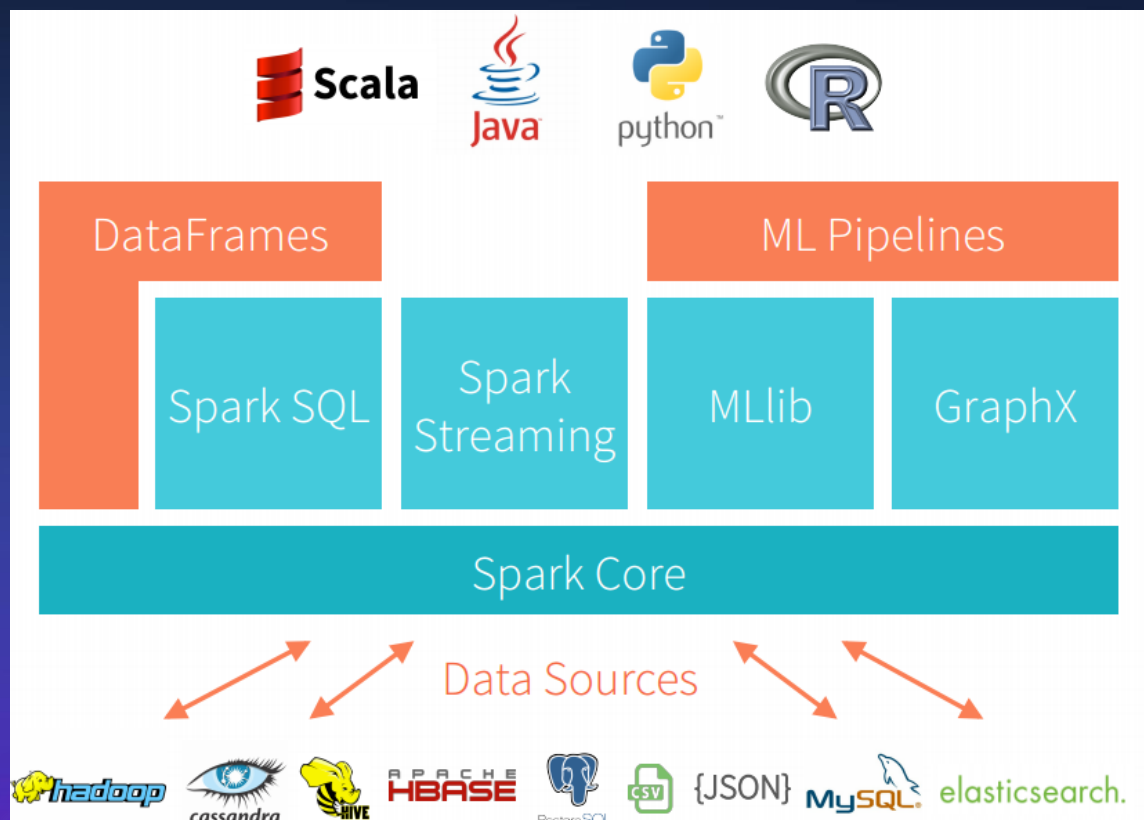
华为云
技术
私享会

华为云Spark技术创新和实践

王飞



Spark：统一大数据分析处理引擎

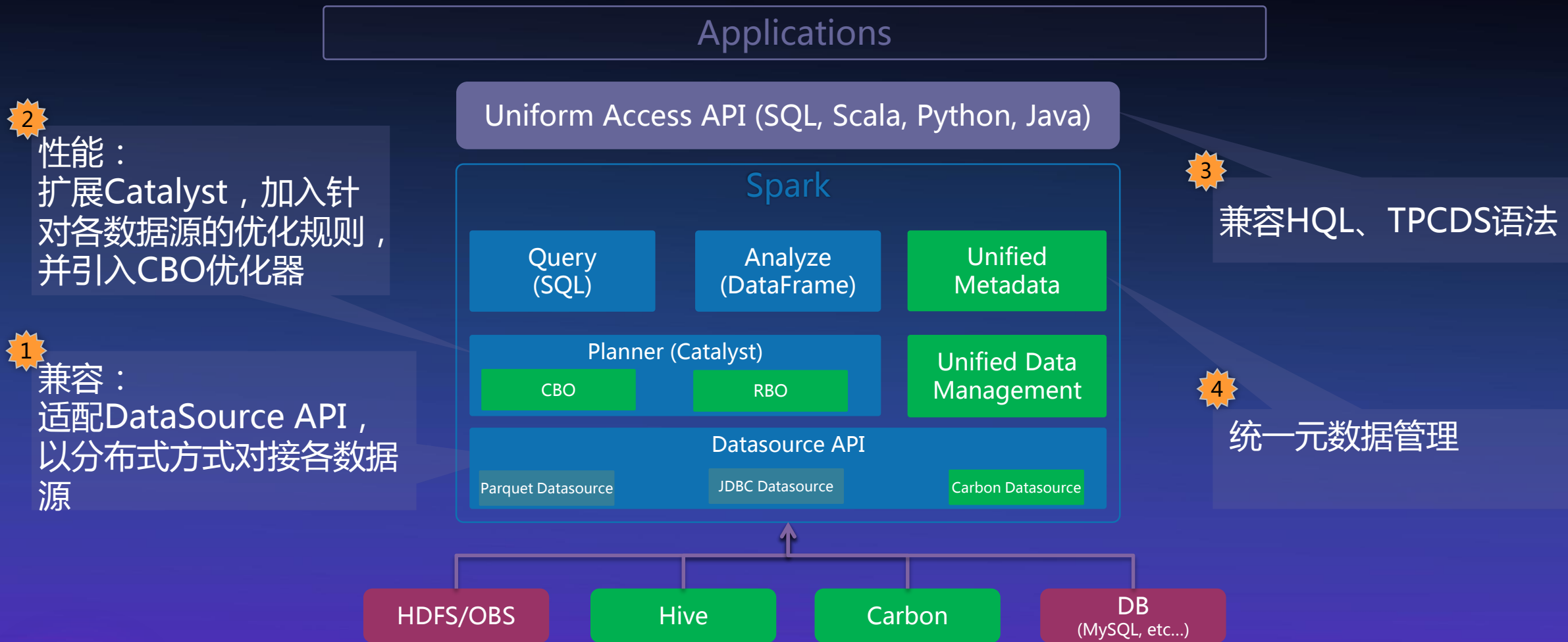


业界主流大数据分析引擎

- Apache 顶级项目，拥有最强大的大数据开源社区，业界公认的大数据分析引擎翘楚
- 支持流、SQL、机器学习、图多种处理范式，满足企业多样化分析模式需求
- 大数据领域最好的生态圈，多种 Datasource 对接支持
- 针对数据分析师和数据科学家提供最简洁易用的分析API

```
sqlCtx.table("people") \
    .groupBy("name") \
    .agg("name", avg("age")) \
    .collect()
```

华为Spark核心能力构建



围绕Spark构建高稳定、高安全、兼容性、极致性能大数据统一分析能力

Spark 查询优化器业界主导者

典型场景	业务特点
关联文件可过滤下压 (大表变小表)	broadcast 和hash join算法路径选择
压缩文件关联查询 (小表变大表)	文件压缩比大，实际大小和表数据量规格不符

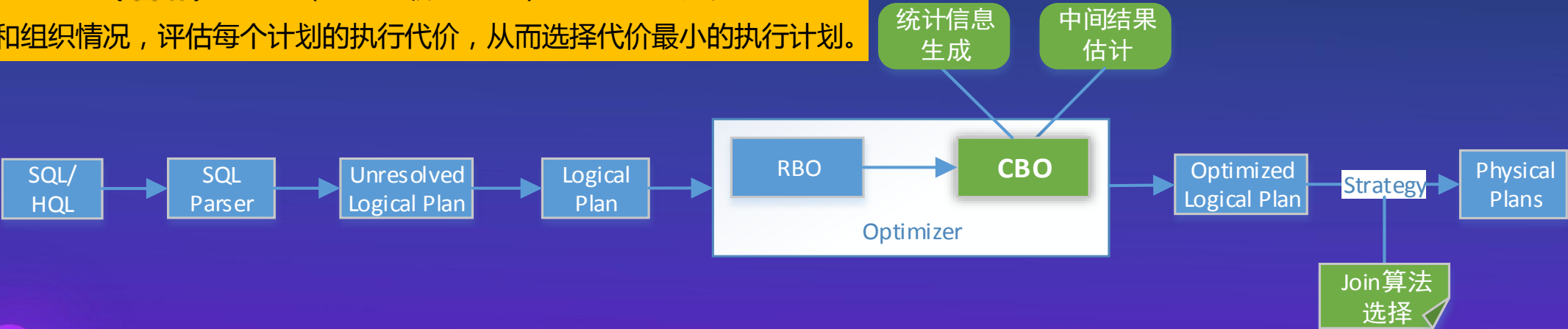
问题痛点：

- ✓ 多表关联下，表规格无法估计，导致逻辑计划优化不合理，执行效率低；
- ✓ 高度压缩的表文件表面size很小，实际数据量很大，BroadCast Join后导致Driver内存被撑爆；

现有能力（开源）：RBO（基于规则的优化）是根据经验形成的固定规则，有效但僵硬。Spark中已有数十条优化规则，如谓词下推、投影裁剪等。



华为构建（自研）：CBO（基于代价的优化）则可以根据实际数据分布和组织情况，评估每个计划的执行代价，从而选择代价最小的执行计划。



华为主导Spark社区查询优化器框架
Spark社区Committer 1名contributor 20+
效果：TPC-DS 99个Query 查询性能端到端提升20倍

当前大数据分析探索面临的挑战

原始数据类型多样，格式不一



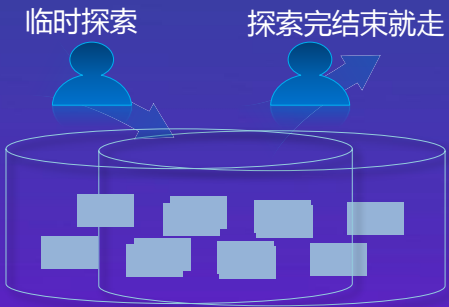
依赖平台技术，门槛高，ETL工作量大，不同系统数据倒换



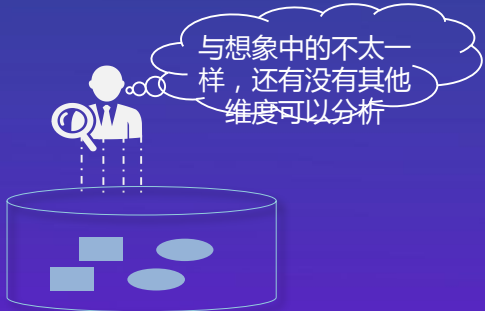
数据共享难操作，代价大



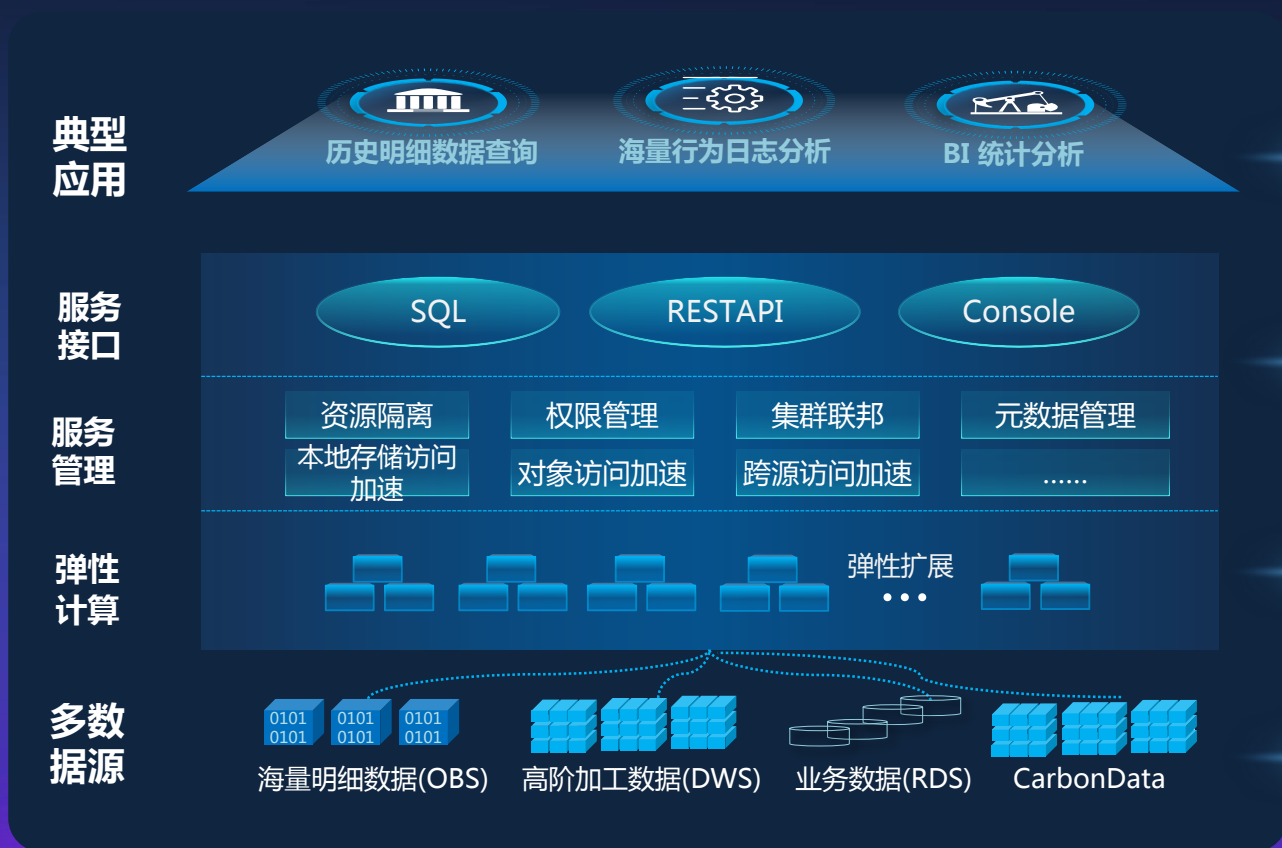
临时性探索查询，需要即来即用



分析维度不全，不灵活
分析模式不足，无法充分发挥数据价值



UQuery，完全托管的云上Serverless Spark统一分析服务



便捷易用

服务完全托管，无需管理任何基础设施，零维护成本，使用SQL接口/原生Spark（Dataframe/Dataset/ML/StructStreaming）接口进行数据查询分析，即来即用。

多数源联合分析

支持文本/CSV/Parquet/ORC/JSON多类型数据格式，支持OBS/服务本地存储加速/CloudTable/DWS/CarbonData跨源联邦查询，用户无需额外数据搬迁，可直接通过SQL，轻松实现跨源数据探索分析

企业多租户

支持按租户申请资源队列，指定独占的计算资源，用户作业运行时资源隔离，有效应对业务峰值，保障用户作业SLA，支持表/列数据权限控制，方便企业内部安全访问管控。

无限扩展

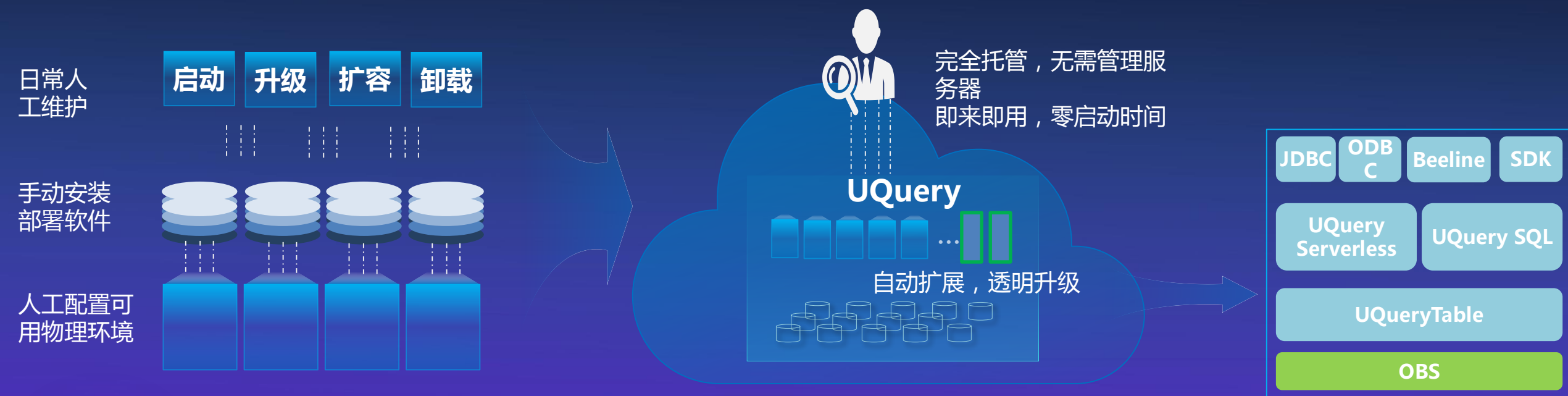
支持集群联邦技术，资源、数据自动弹性扩展，用户不感知，无需关心资源是否够用。

UQuery特点1：云上Serverless Spark，统一多种分析范式

华为云社区
bbs.huaweicloud.com

传统数据分析过程

Serverless：云上托管，用户透明



华为云
技术
私享会

HUAWEI

UQuery SQL：简单易用，低门槛上手

The screenshot displays the UQuery SQL interface. On the left, there's a sidebar with '资源列 (2)', '数据库 (11)', and '表 (8)'. The main area shows a SQL query: `SELECT * FROM dimensions_adgame.join_agg1 LIMIT 1000`. Below the query, there are buttons for '执行' (Execute), '语义校验' (Semantic Check), '格式化SQL' (Format SQL), '设为模板' (Set as Template), '选择模板' (Select Template), and '执行历史' (Execution History). The execution results show '查询耗时2.08s, 已扫描235.37 KB'. On the right, there's a '选择参数' (Select Parameters) section with '图形类型' (Chart Type) set to '扇形图' (Pie Chart), '图例' (Legend) set to 'GID', '指标' (Metric) set to 'download', and '结果数目' (Result Count) set to '20'. A pie chart is displayed with various segments labeled with values like 9744.0, 977.0, 895.0, etc. The text 'SQL兼容标准2003，完全兼容Spark接口' is overlaid on the query editor. The text '结果图形化 / 表格呈现' is overlaid on the pie chart.

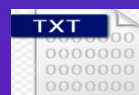
三步上手应用

登陆控制台

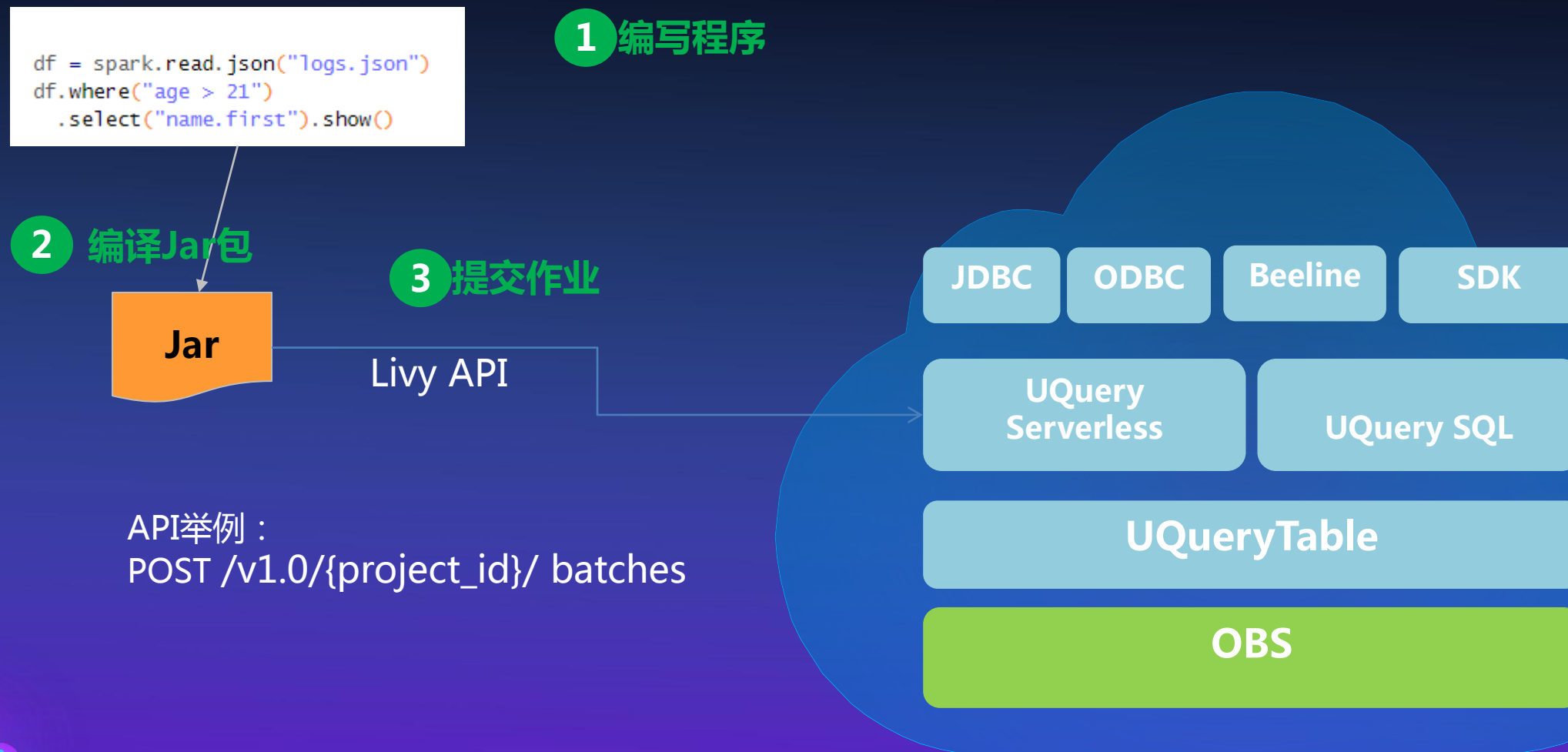
建库建表

开始查询

支持丰富数据格式
对开源的共享，体现能力：
语法，核心CBO优化器commiter



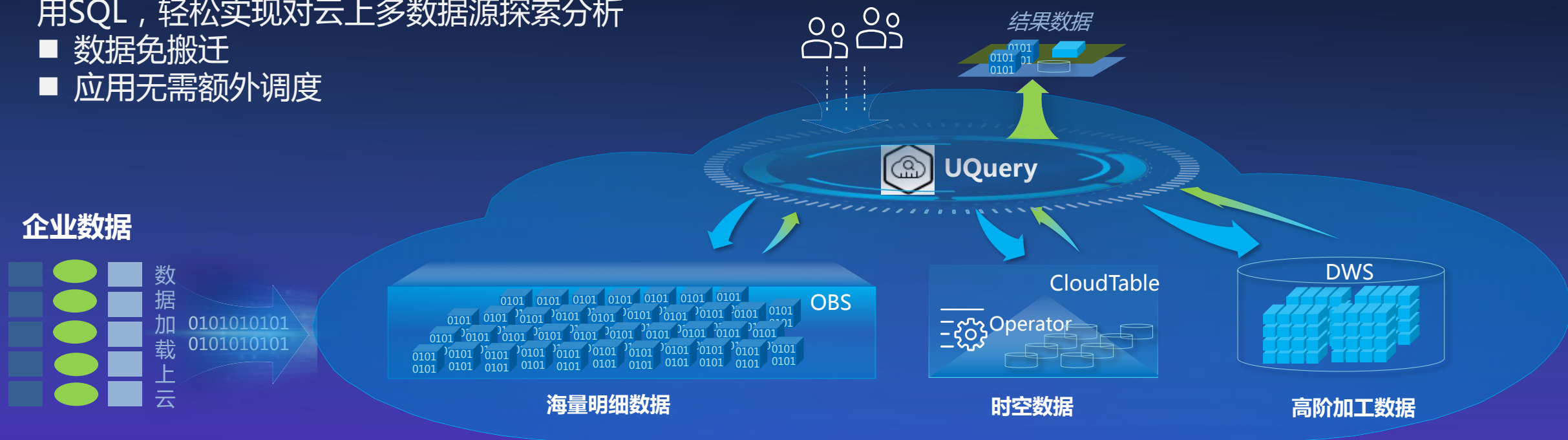
UQuery Serverless : 原生编程接口支持



UQuery特点2：跨多数据源联邦分析，数据免搬迁

UQuery支持**跨源联邦查询**能力，用户直接使用SQL，轻松实现对云上多数据源探索分析

- 数据免搬迁
- 应用无需额外调度

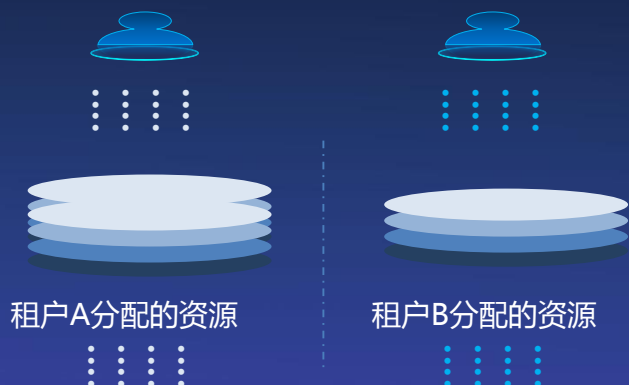


UQuery特点3：完整的企业多租户管理

资源隔离，保障用户作业SLA

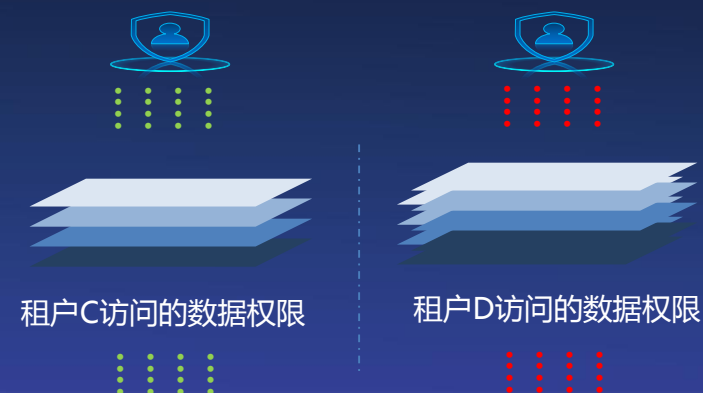
用户可选择对资源进行配置，包括CPU，内存，不同用户对应不同资源队列，保障作业相互不受影响

用户作业提交
不同队列，队
列间相互隔离



权限控制，提供数据精细化管理

用户可以选择表或者部分列进行权限设置，以保障企业内部各部门之间的数据安全访问，实现数据精细化管理



数据管理员可
以对库，表，
列进行授权，
实现数据控制
与分享

UQuery特点4：提供本地存储加速，提供秒级查询性能



索引：快速定位数据，降低I/O

多维索引、高维索引、Min/Max、倒排索引



预聚合：预先组织数据，秒级查询响应

字典编码，内存缓存，延迟解码



引擎加速：执行优化，DAG调度提升2倍

Codegen，堆内存，序列化器优化



优化器：计算存储联合优化

CBO，计算下压，基于ML自调优

UQuery采用业界先进的CarbonData存储技术，结合分区表，缓存加速，索引等技术加速海量数据查询性能

BI分析

场景特点：

统计计算，对比分析，交互查询

典型应用：

固定报表，多维OLAP分析，交互式报表等

海量行为日志分析

场景特点：

半结构化，海量数据，部分涉及模糊查询

典型应用：

学习习惯分析，系统操作日志查询等

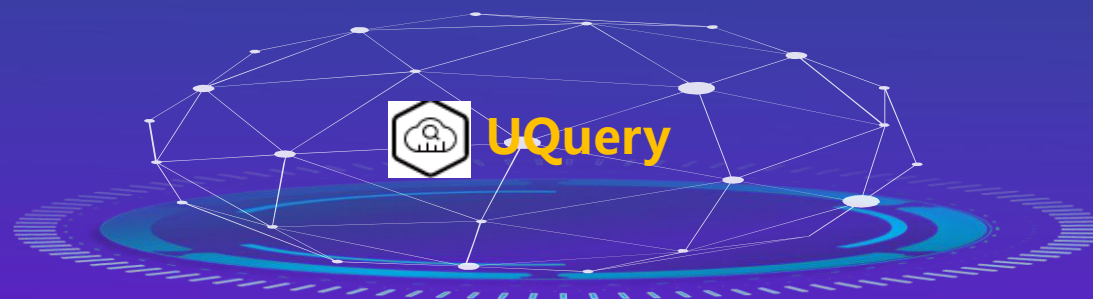
历史数据明细查询

场景特点：

直接查原始海量数据，带过滤条件，查询维度不固定，灵活多变

典型应用：

流水审计，设备历史能耗分析，轨迹回放，车辆驾驶行为分析等



为您推荐数据查询服务UQuery几个常用链接

数据查询服务UQuery官网：华为云→EI企业智能→EI大数据→数据查询服务
<http://www.huaweicloud.com/product/uquery.html>

数据查询服务UQuery论坛：华为云→云社区→论坛→EI企业智能→数据查询服务
<http://forum.huaweicloud.com/forum-599-1.html>

数据查询服务UQuery帮助资料：
http://support.huaweicloud.com/uquery_dld/index.html

Q&A



THANK YOU

