

Kylin在小米大数据中的应用

陈学辉

小米大数据部数据平台组

2019.4.13

- **Contents**

01

小米业务场景

02

数据分析演进

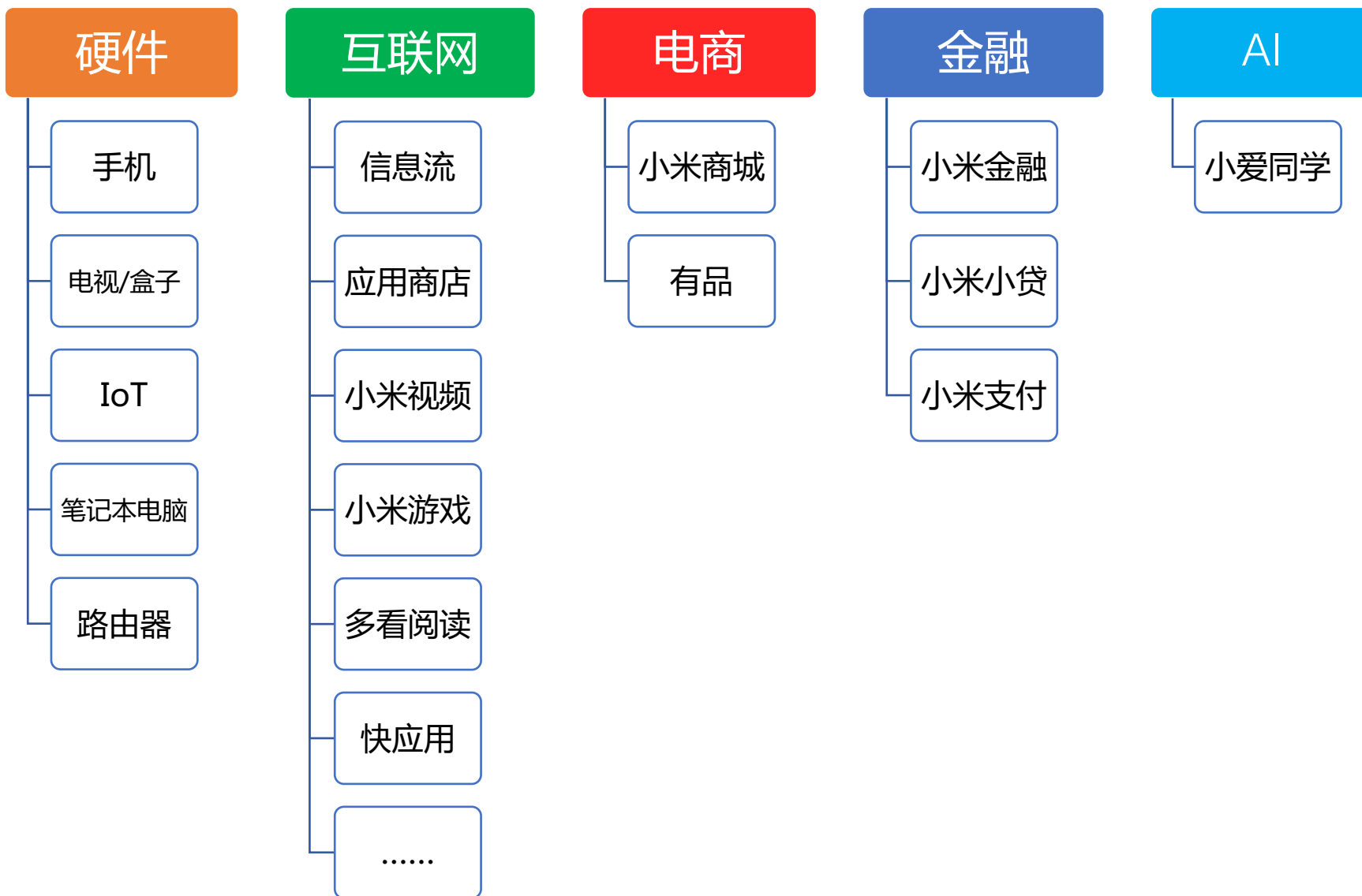
03

Kylin使用情况

01 小米业务场景

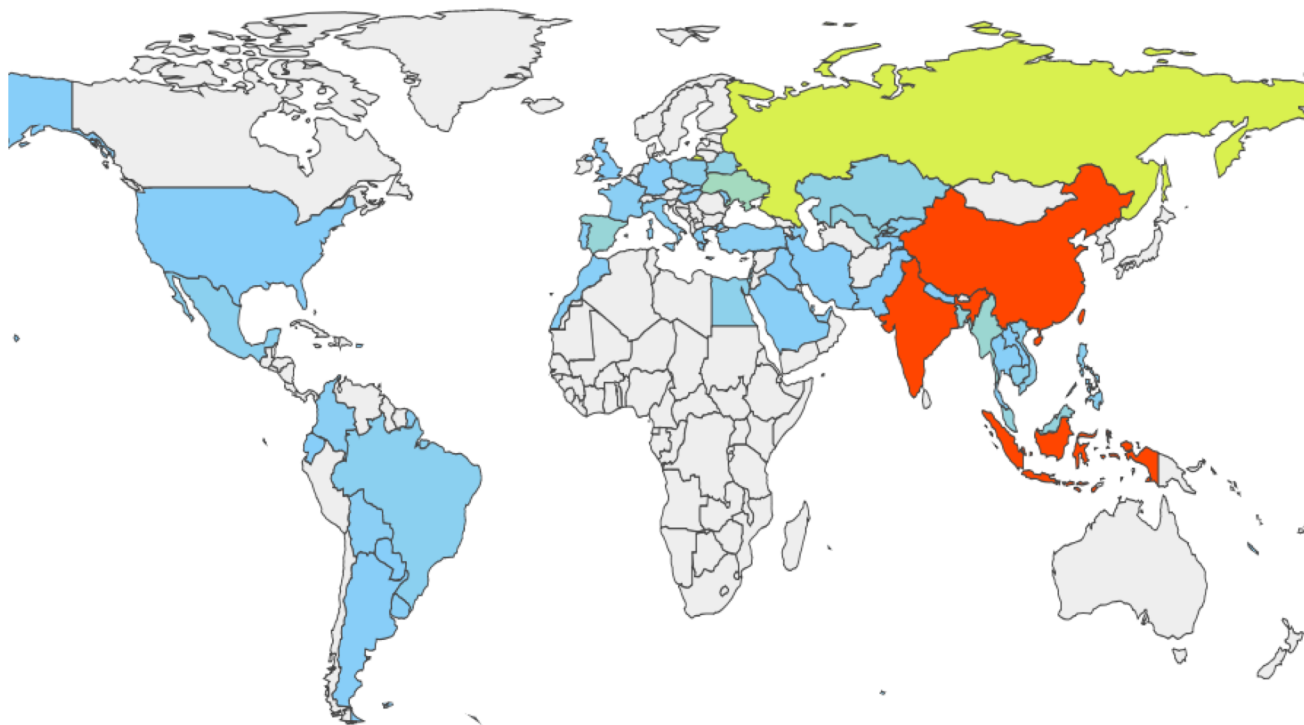
1. 业务产品多样化
2. 全球化业务布局
3. 业务数据分析场景
4. 业务数据需求

● 小米多样化的业务产品线



● 全球化的商业布局

- 中国区
- 印度
- 东南亚
- 独联体（俄罗斯）
- 欧洲
- 美洲



业务数据分析场景

核心指标

- **数据需求**：日活、新增、次数、时长等
- **使用系统**：小米统计、MIUI BI



用户行为

- **数据需求**：前台用户的行为，如事件、留存、漏斗、访问路径等
- **使用系统**：小米统计、GA



竞品对比

- **数据需求**：市场占有率、竞品指标等
- **使用系统**：MIUI BI



后台数据

- **数据需求**：后台业务数据，如内容、活动、订单、搜索、信息流、收入等
- **使用系统**：Big BI、Xdata、业务自建数据系统、建光统计、数据工厂任务、云平台ABC平台



● 业务的数据需求

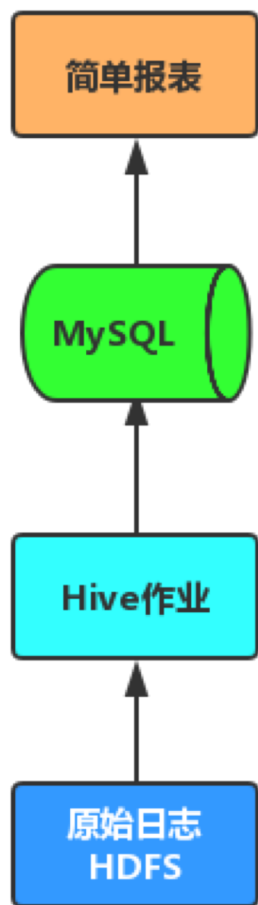
	业务需求点	重要程度
性能	大数据规模（单表每日10~100亿条记录）	高
	毫秒级/秒级查询响应	高
系统	数据实时性	低（天级/小时级）
	稳定性	高
	高并发查询	中
功能	多表关联	中
	Schema灵活	高/中
	精确去重	高
运维	扩展性	高
	多租户	高

02

数据分析演进

1. 多种数据系统
2. 小米内部的数据分析OLAP方案
3. Big BI 数据流程
4. Big BI 技术架构

● 数据工场统计平台



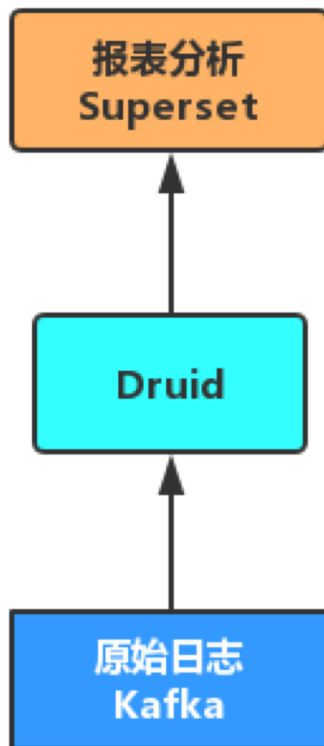
优点：

- 满足天级离线统计需求
- 接入、使用简单

不足：

- 报表展现简单，无法满足复杂可视化需求
- 一个任务对应一个图表，需求和工作量成正比

● 简单实时统计系统



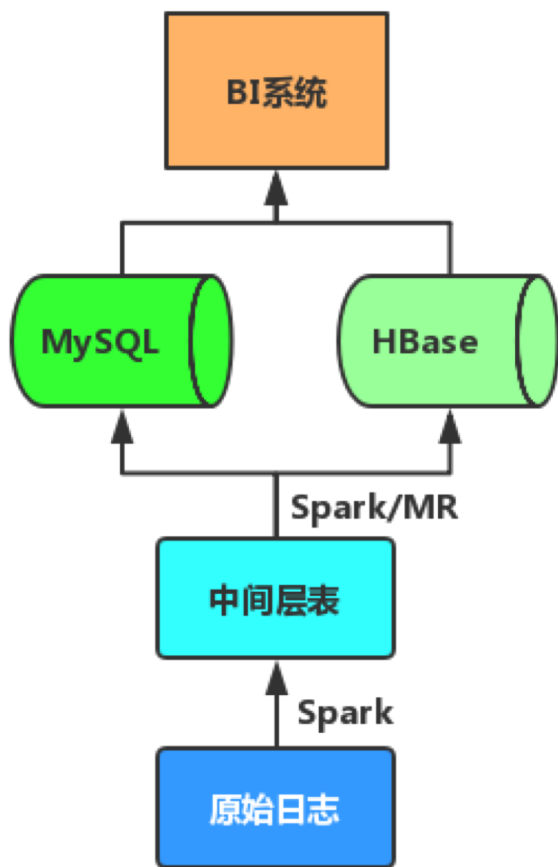
优点：

- 满足实时统计需求
- 开源方案，维护成本低

不足：

- 历史数据查询性能无法保障
- 不支持多表JOIN

● 传统 BI 系统



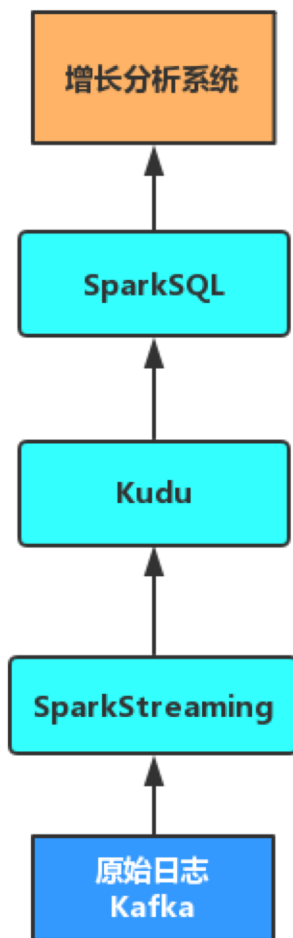
优点：

- 丰富的可视化图表，统计分析功能强大
- 能够满足多种数据需求

不足：

- 难以支撑海量数据多维的灵活分析
- 系统复杂，开发测试维护成本较高

● 增长分析系统



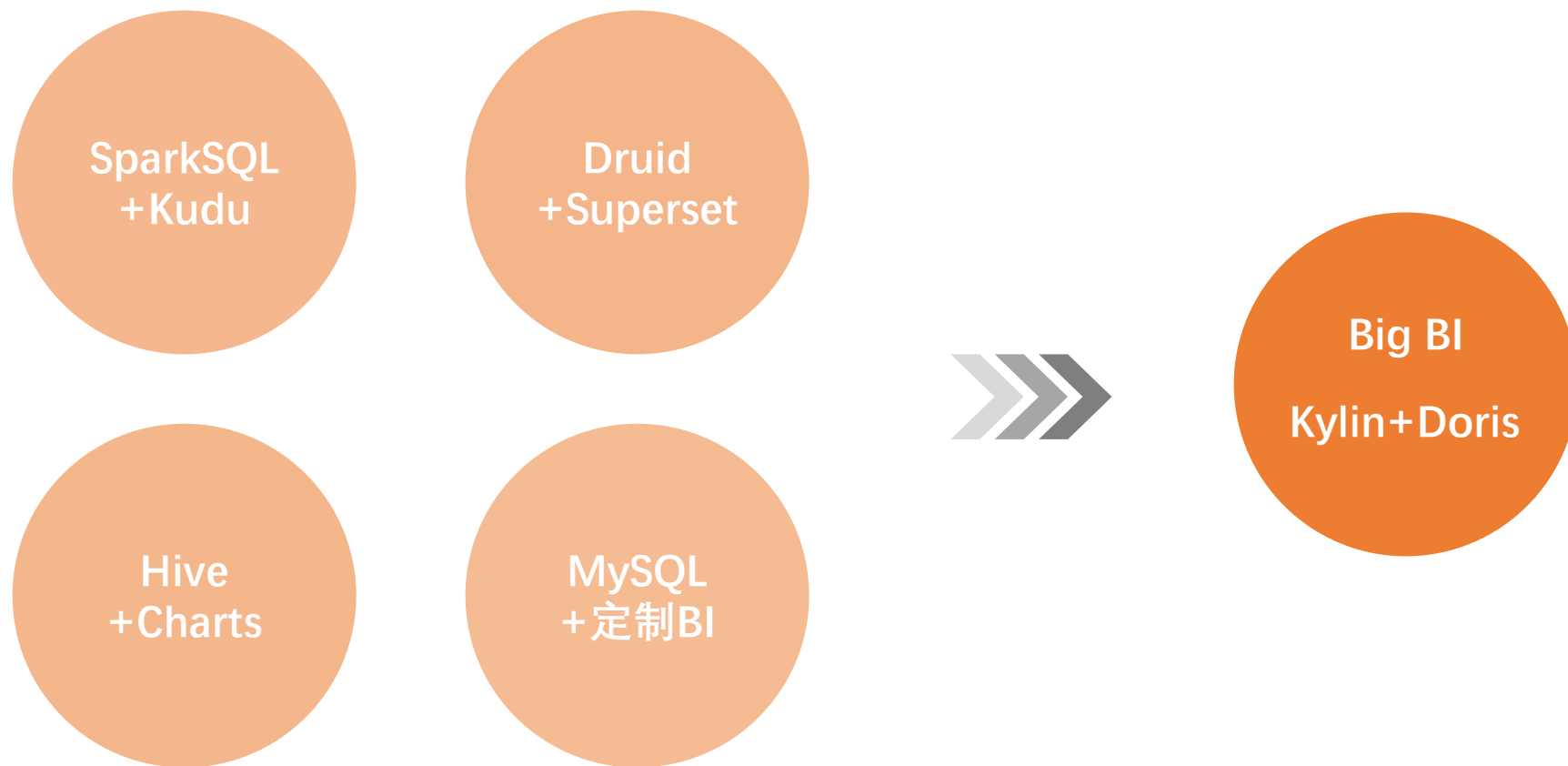
优点：

- 支持实时数据多维即席查询分析
- 支持灵活的Schema
- 支持UDF实现用户行为分析

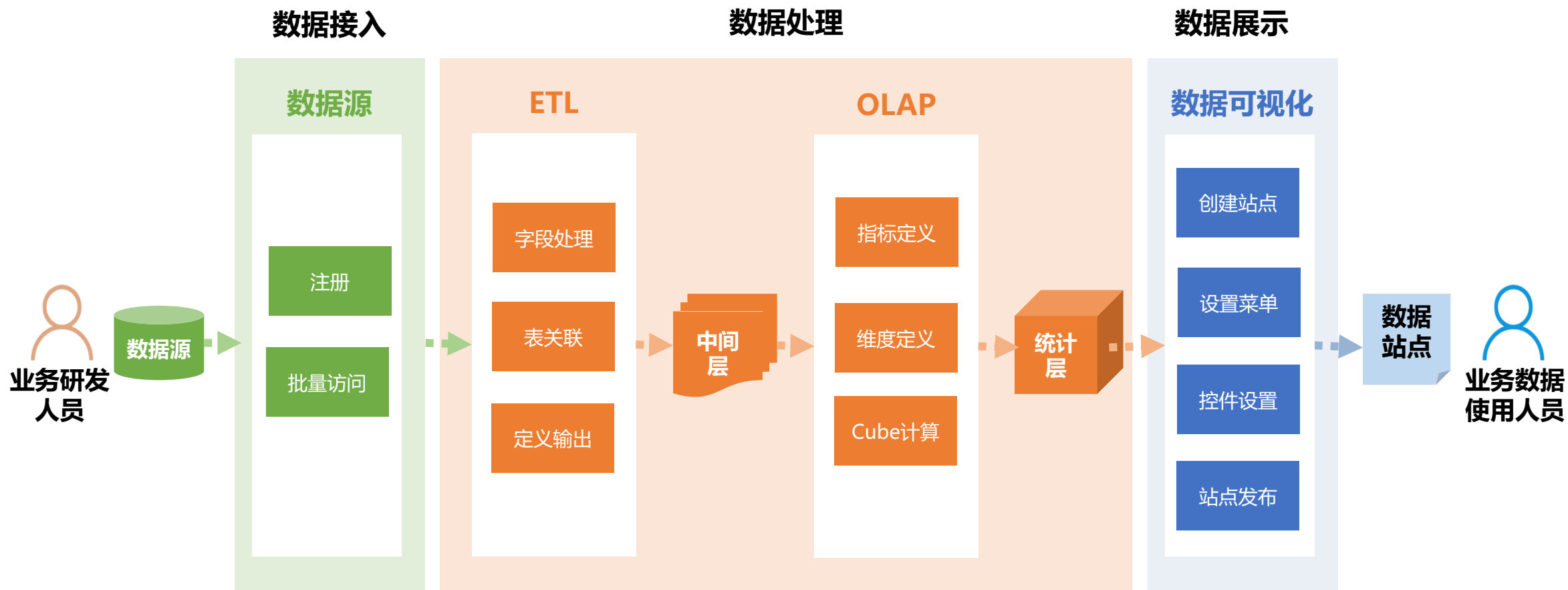
不足：

- 较大数据规模下查询时间较长
- 系统硬件成本较高

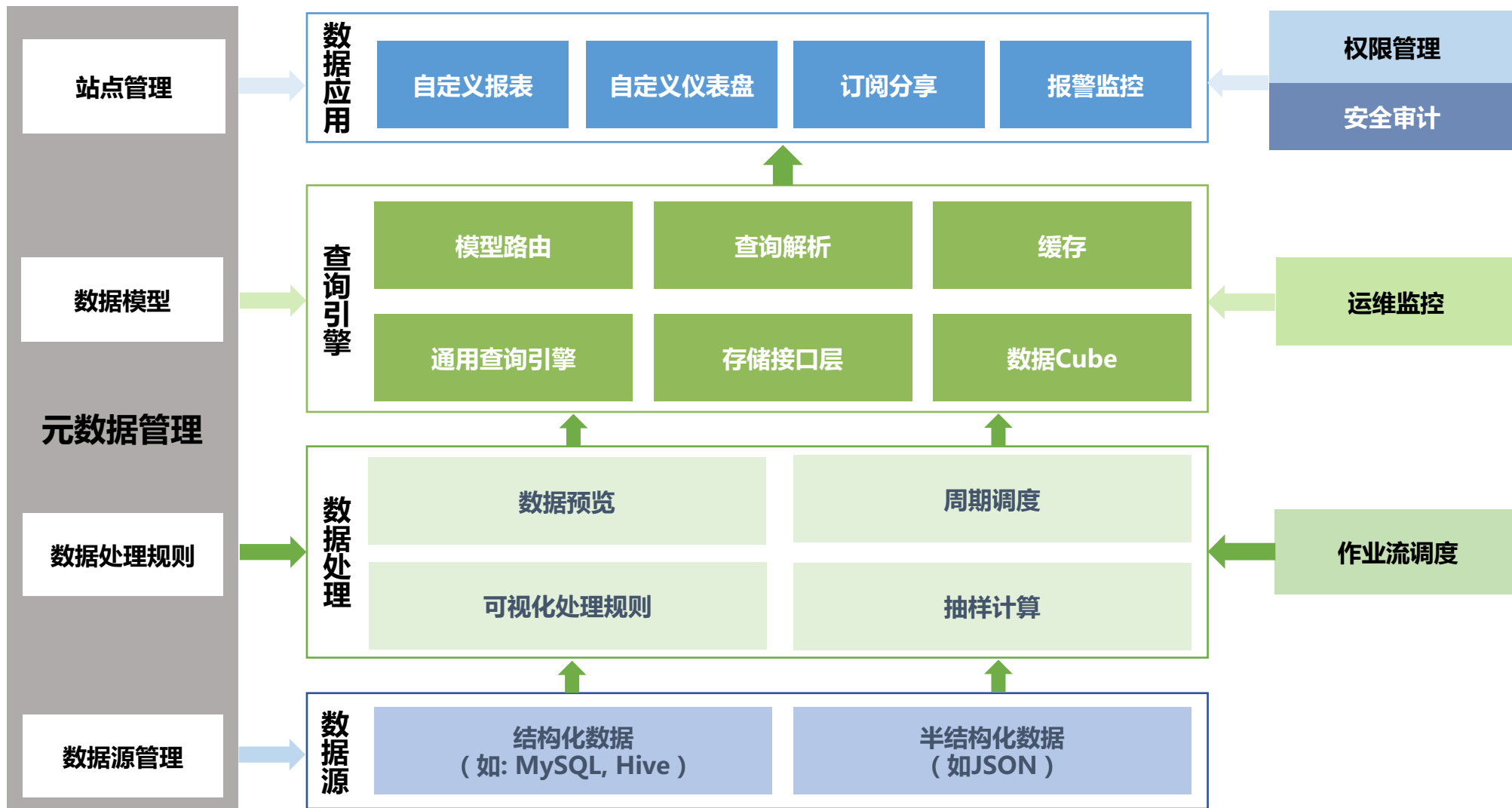
- 小米内部的数据分析OLAP方案



● Big BI 数据流程



● Big BI 技术架构



03 Kylin使用情况

1. 使用现状
2. 版本维护
3. Kylin集群
4. 易用性优化和改动
5. 使用技巧分享

● Kylin使用现状

20+人/月

活跃管理员用户数

~100TB

数据量

30w+次

日查询量

0.8s

平均查询时间

150+个

Cube数

重度用户
(2个)

- MIUI BI
- 小爱 (AI)

一般用户
(6个)

- 有品 (电商)
- Push
- 手机质量系统
- 笔记本
- 安全中心
- MiTV

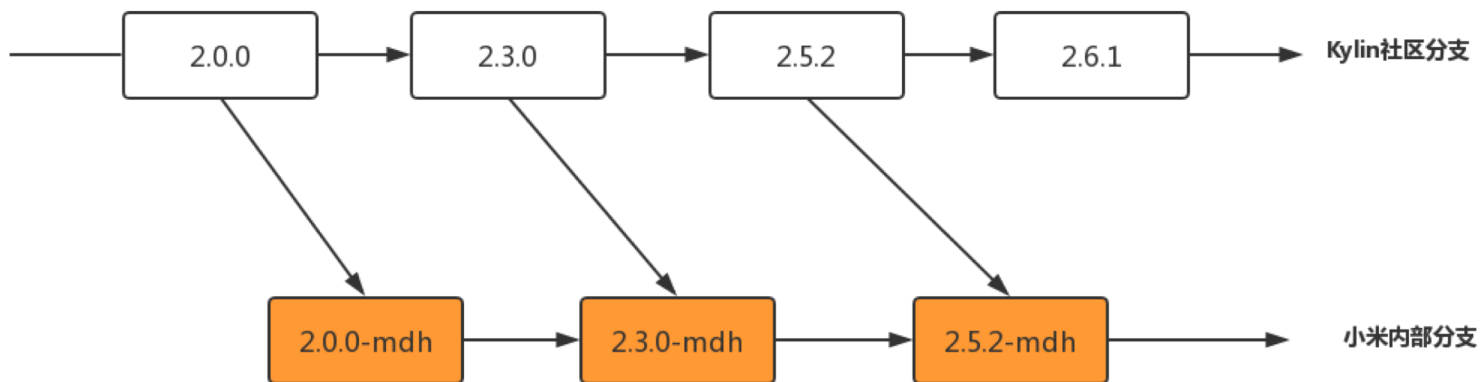
业务试用
(10+个)

- 小米金融
- 应用商店
- 路由器
- 用户反馈
- 游戏中心
-

● 版本维护

维护小米内部版本，定期与社区同步，快速开发迭代

- 2019-3 kylin-2.6.1-mdh 测试预览
- 2018-12 kylin-2.5.2-mdh 线上版本 (Cube Planner)
- 2018-5 kylin-2.3.0-mdh(用户权限)
- 2017-11 kylin-2.1.0-mdh
- 2017-8 kylin-2.0.0-mdh

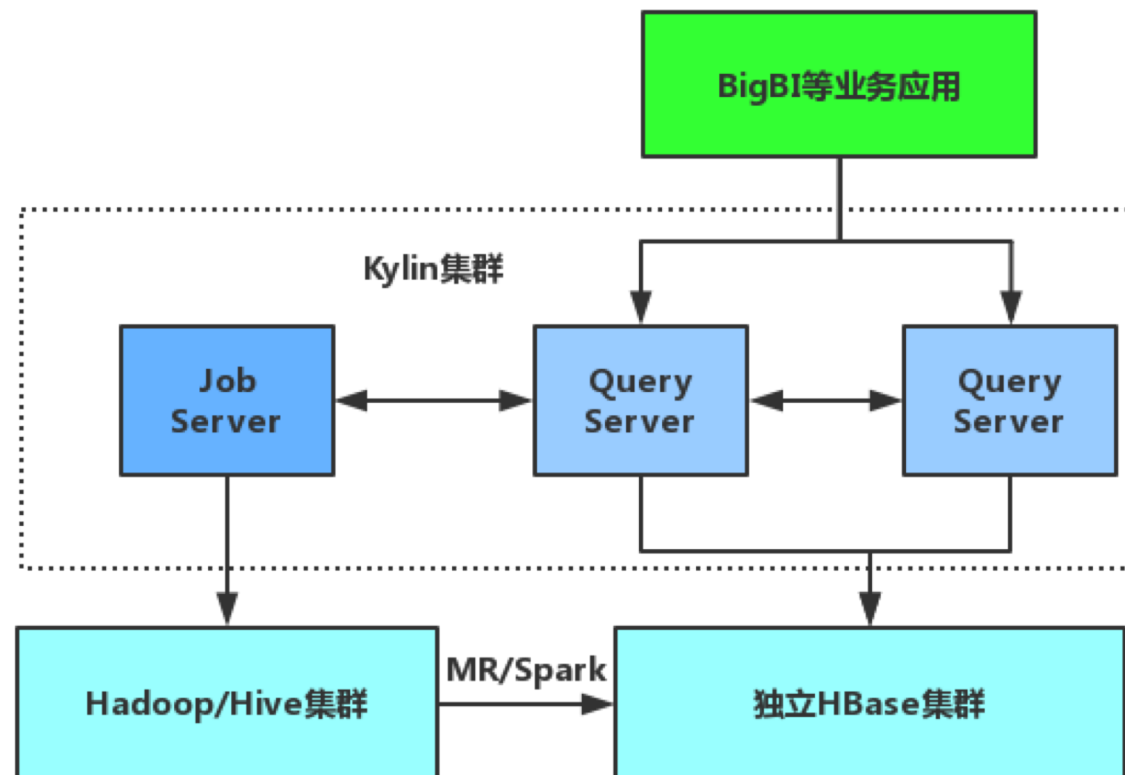


● Kylin集群

Staging集群/Production集群

- 集群包含一个JobServer, 2个QueryServer
- 使用公司共享的离线HBase集群, 独立的Namespace
- Project/Cube级别配置独立计算队列, 保证计算资源隔离

痛点: HBase数据迁移较复杂, 无法一键复制
Staging集群数据至Production集群



易用性优化

Cube Designer

1 2 3 4 5

Cube Info Dimensions Measures Refresh Setting Advanced Setting Configur

Model Name * miui_app_retain_preinstall_model

Cube Name ⓘ * miui_app_retain_preinstall_cube

Notification Email List litong5@xiaomi.com, wushuangxi1@xiaomi.com, yuechongyuan@xiaomi.com

Notification [Sms Group List](#) miui-bi-kylin

Build success callback callback url: http://dev.kylin.com

Schedule time 目前仅支持按日调度，例：17:30

Notification Events ⓘ ERROR ✕ DISCARDED ✕

- 基于时间的周期性调度
- 增加Cube构建回调
- 支持Kylin事件的短信通知

Kylin [Search] Insight Model Monitor System User Dashboard Help Welcome, [Avatar]

+ User

Users			
用户名	用户权限	状态	操作
[Avatar]	ROLE_ADMIN	已激活	[Edit] [Delete]
[Avatar]	ROLE_ANALYST ROLE_MODELER	已禁用	[Edit] [Delete]
[Avatar]	ROLE_ANALYST	已激活	[Edit] [Delete]
[Avatar]	ROLE_ANALYST ROLE_MODELER	已禁用	[Edit] [Delete]

- 增加用户管理
- 将Kylin metrics推送到监控系统falcon
- 定期检查和修复segment空洞和重叠问题
- 基于ZK的服务发现，避免繁琐的服务列表配置

● 其他改进和修复

- Support max segment merge span
- HBase 0.98适配和支持
- KYLIN-3780: Add built instance in Job info
- KYLIN-3918: Add project name in cube and job pages
- KYLIN-3912: BeelineHiveClient支持Project/Cube级别的资源隔离
- KYLIN-3880: DataType is incompatible in Kylin HBase coprocessor
- KYLIN-3882: kylin master build failed for pom issues
- KYLIN-3884: loading hfile to HBase failed for temporary dir in output path
- KYLIN-3886: Missing argument for options for yarn command
- KYLIN-3909: kylin job failed for MappeableRunContainer is not registered
- KYLIN-3913: Remove getAllOutputs api in ExecutableManager to avoid OOM for large metadata

● 使用技巧

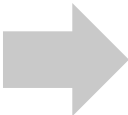
1. 如果表中包含date、sum、day、hour等关键字，查询时需要用 "" 括起来，双引号表示列名，单引号表示字符串
2. 查询时，表名即是hive表名，cube/model是数据表之上的抽象，不是查询实体
3. RowKey顺序
 - ✓ 查询中作为过滤条件的维度放在其他维度前面
 - ✓ 查询中将经常出现在的维度放在不经常查询的维度的前面
 - ✓ 对于基数较高的维度，如果查询有该维度上的过滤条件，那么将它往前调整，如果没有，则向后调整
4. UV计算，生成宽表时，确保每个ID只有一行数据，同时字段uv值为1表示该ID是活跃的，通过sum(uv)来计算UV，大幅节省计算时间和存储空间



● 维度互斥问题

维度	非互斥原因
国家	中国，海外国家(Global)，印度，新加坡
机型	小米机型、小米国内机型、小米国际机型
MIUI版本	CnStableAll、GlobalStableAll

机型
小米国际机型
小米国内机型
小米Pad
第三方机型
未知



机型	是否国际机型
小米6	是
小米6	否
小米Pad	未知
第三方机型	未知
未知	未知

MIUI版本
CnStable
GlobalStable
V8
Dev&Alpha
未知



MIUI版本	是否GlobalStable
V8	是
V8	否
Dev&Alpha	未知
未知	未知

国家	国家	是否Global国家
Global	中国	否
中国	印度	是
未知	未知	未知



拆分维度导致查询复杂：

查询某时间段内某App的多个国家的指标值
(select sum(dayInstallUser) as metricsValue, 'All' as dimensionChoice, "DATE" from table where "DATE " between 20171110 and 20171115 and packageId = 93386 and category = 'A_A' group by "DATE ")/* -- 国家all -- */
Union all
(select sum(dayInstallUser) as metricsValue, 'Global' as dimensionChoice, "DATE " from table where "DATE " between 20171110 and 20171115 and packageId = 93386 and category = 'A_A' and isGlobal = 'Y' group by "DATE ")
/* -- Global国家 -- */
Union all
(select sum(dayInstallUser) as metricsValue, country as dimensionChoice, "DATE " from table where "DATE " between 20171110 and 20171115 and packageId = 93386 and category = 'A_A' and country in (国家列表) group by country, "DATE ")

● Kylin服务限制

添加一些必要的限制, 保证服务质量

- 限制每次构建的Cube膨胀率
- 限制查询SQL语句的长度
- 限制SQL语句中in值的个数
- 限制一个查询里面设置segment个数
- 限制每次构建cube数据量大小
- 限制度量使用内存大小(2M), 避免HBase集群压力
- ...

● 未来规划

OLAP解决方案

整合Kylin、Doris、SparkSQL，实现查询智能路由，充分发挥每种引擎的长处，形成统一的OLAP解决方案

优化性能

测试Kylin on Parquet方案，提升查询性能

弹性运维

支持Kylin直接部署在Kubernetes里

Q&A

Thank You
