

讲师介绍

王小静

- · 北京滴滴无限科技发展有限公司(2015.01~NOW)
 - 一 基础平台-大数据架构部-Kylin引擎负责人/专家工程师
 - 一 曾负责了滴滴大数据任务调度系统、数据资产、数梦底层调度执行引擎的架构和落地
 - 一 作为主力成员之一多次参与滴滴hadoop集群异地迁移工作
- · 亚信科技中国有限公司(2012.12~2015.01)
 - 一参与了中国移动云审计项目架构设计和核心模块开发工作





Kylin 在滴滴的应用&架构

议程 Agenda 02 滴滴全局字典最新版本介绍

13 Kylin RT OLAP探索经验分享

04 Q&A

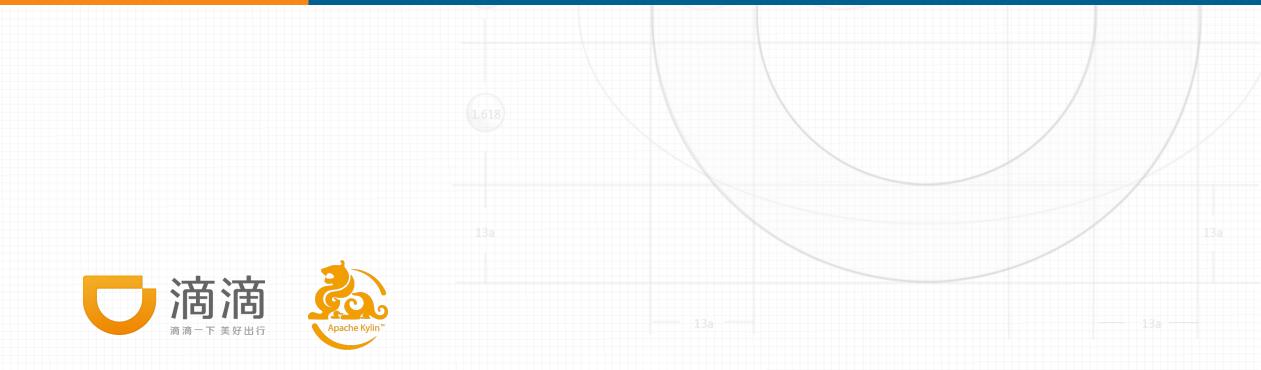






01

Kylin 在滴滴的应用&架构



Requirements/ Background



产品类需求

大屏展示、企业用车分析等明确的产品

SQL固定

用户要求查询响应时间快且稳定,SLA要求高 尽量对产品RD使用简单,要求接入成本低

报表类需求

自定义报表
SQL 复杂且不确定,各种Join/函数
数量众多
用户大多数为运营分析人员

活动营销

活动前范围圈定 活动后效果分析 对精准去重需求大 分析方式灵活,生命周期多为1-2个月



Data Display

• **6** 个Kylin集群

国内:4个Kylin 2.0 ,1个kylin3x

国外:1个Kylin 2.0

• Cube数量 **3.2K**+

日构建任务数 4k+

• Hbase表数量 **5W**+







Data Display



TOTAL CUBE

ots

8

More Details

AVG CUBE EXPANSION

1.83

QUERY COUNT

2,457

More Details

AVG QUERY LATENCY

0.29 sec

More Details

JOB COUNT

6

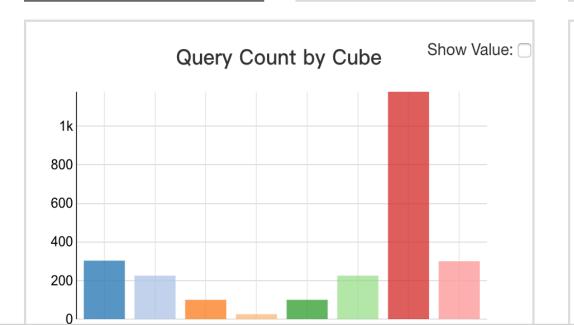
More Details

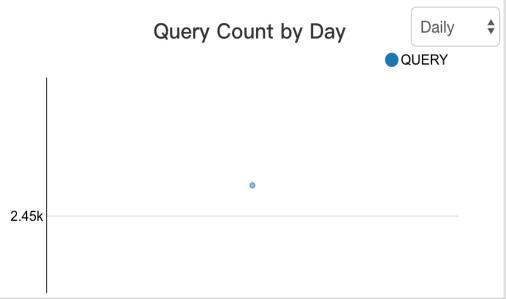
AVG BUILD TIME PER MB

2019-11-13 - 2019-11-13 **▼**

27.49 sec

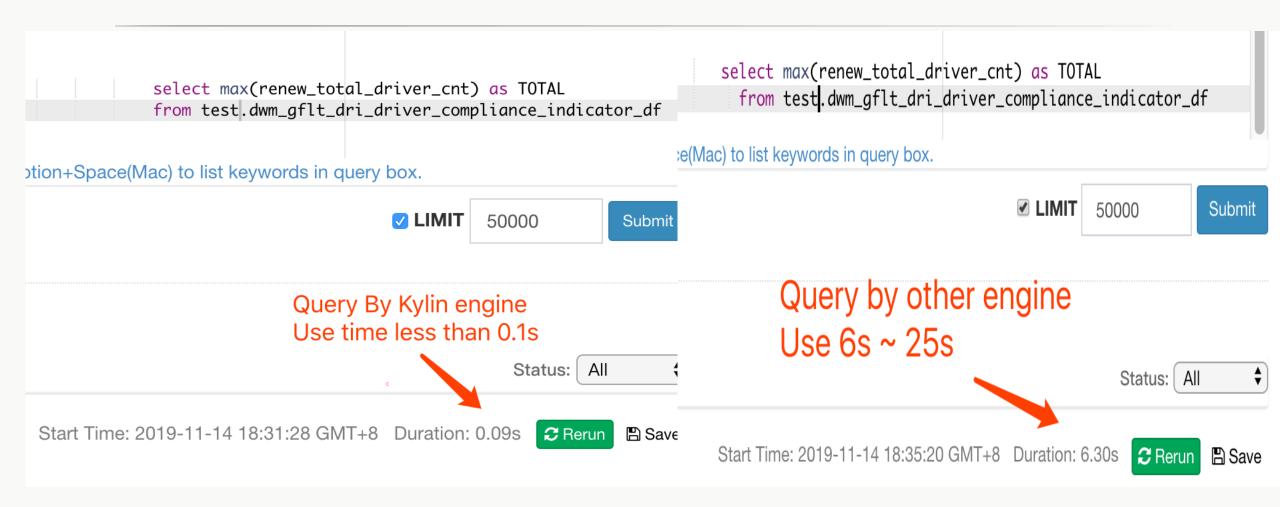
More Details





Data Display





探索架构路上经验分享之易用性、灵活性:



● 复杂SQL

- 用户来源表关联关系复杂,多种join
- 某些SQL预发不支持,函数不支持
- 维度表数据变化要求严格:如全局唯一等

● 字段预处理

- 对原始表进行字段简单处理如substr,concat等
- 分区字段格式限定

● Model需管理但复用率底

• 每创建一个Cube均需要创建一个Model,在我们的场景复用率底,维护成本大

● 未建模SQL查询

· 用户提交一个SQL查询,如果未建模报错





建模转换



Kylin建模



Kylin|下推

建模平台:

- 工作:负责封装临时表, Kylin建模、调度任务创建
- 好处:多表join支持,字段转换,复杂sql支持,屏蔽load table, model, cube等

报表平台

其他

建模平台

Kylin Cluster

HBase

Pushdown

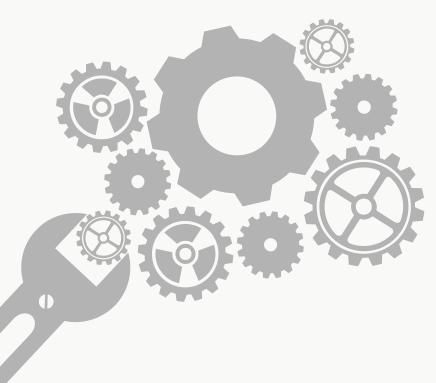




广播延时或者分布式事物等导致的一致性问题







- · 虚拟State 角色, Standby/Active
- · 元数据补全API

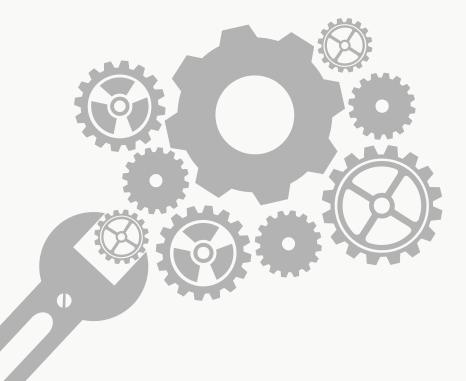




集群治理







- · 对于长时间运行的任务自动触发Discard
- · 只调度最近3天的job,减少Job轮询
- · Api与数据清理分离,提高Api相应时间
- 历史数据清理,计算集群&存储集群清理
- · Meta元数据清理
- · Rs group 拆分

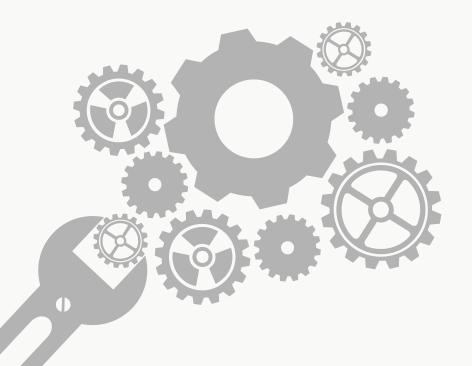




集群运维







- · 配置与代码分离
- 多集群负载均衡,对于版本升级,流量管控等
- 关键日志监控,如待运行任务数,错误任务数
- 添加探活服务,实时监测每个节点的查询情况
- 远程操作,在线日志查看,远程起停





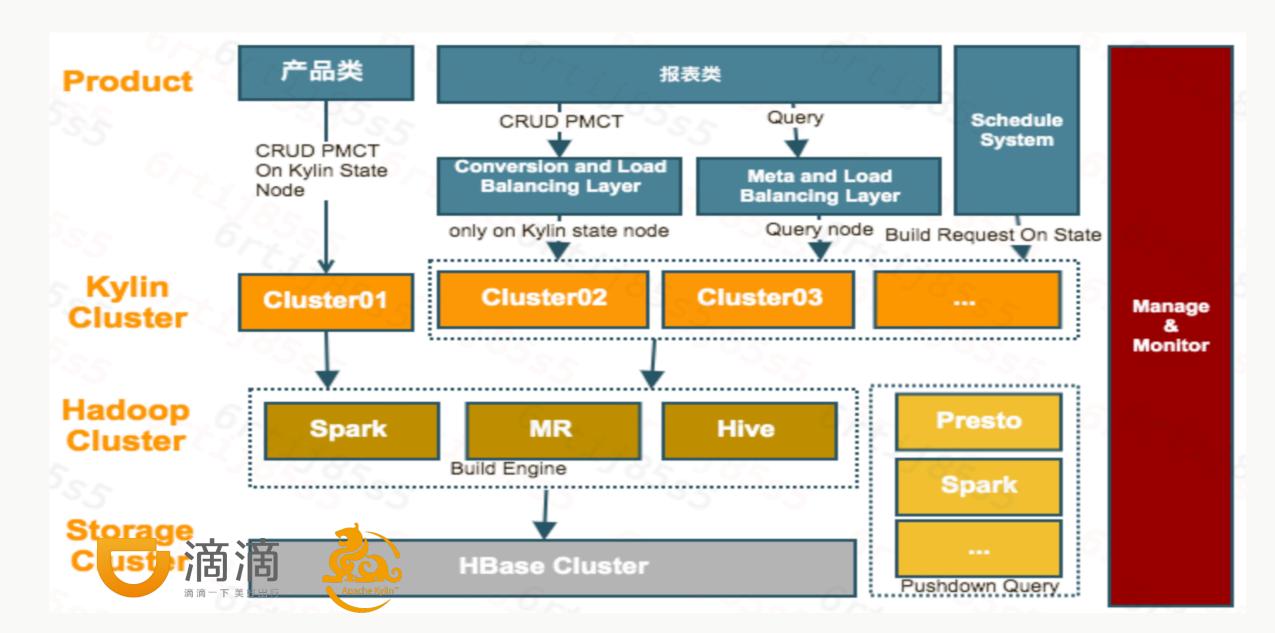
构建速度







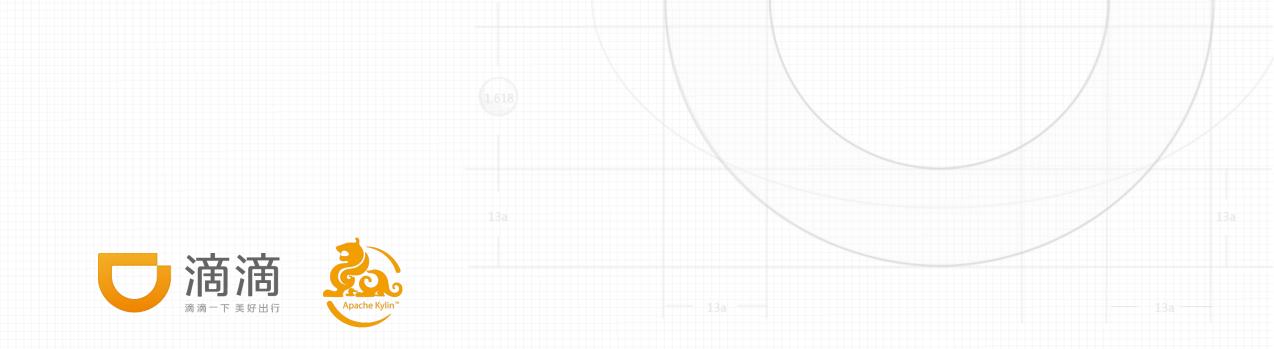
- Spark Livy
- · Cube 优化
- 全新全局字典构建
- 全域字典复用





02

滴滴全局字典最新版本介绍



WHY

突破超大数据量字典构建瓶颈

Tiretree方式数据量达到亿级别规模时随时数据的增加逐渐凸显瓶颈,需要支持更大基数的精准去重

MR 字典

提升构建速度

原tiretree方式当数据量超大时构建速度也变得不可控,分钟级别变为几个小时,甚至无法出结果的情况

减少重复构建

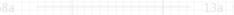
通常多个cube会有对相同字段同时求UV场等景,每个Cube均需要构建 复用

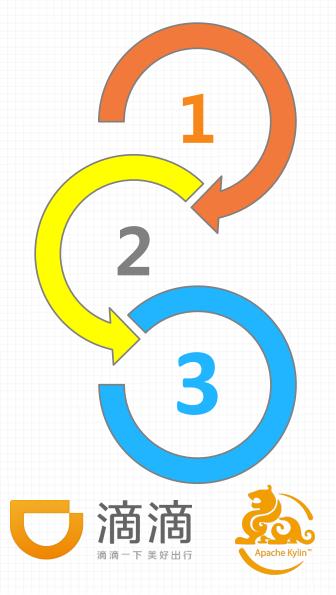




Version Iteration







V1 Hive

第一版本采用Hive实现,可较好完成千万到四亿级别字典规模构建,数据量继续扩大构建时长和需要的单任务内存无法满足需求

V2 MR TotalOrderPartition

突破单列构建数量瓶颈,单列理论数十亿规模均能恒定在15分钟内完成,但同cube有多列构建全局字典需求时,需要串行构建

V3 MR Multi Cols Parallel

在V2基础上进一步优化,由cube多列串行构建改为并行构建,达到理论同cube无论多少列单列字典数十亿均能恒定在15分钟左右完成

__13a

uname

先二

李三

Extract

• 通过hive获取本次需要字典编码的原始值;

• 获取需要字典编码列之前最大编码;

Dict

根据hive order by / row number 获得每个值编码

key	value	
李三	2	
张三	3	

Merge

通过hive与每列之前的字典合并

key	value
王芳	1
李三	2
张三	3

Replace

• 通过Hive用编码后的值替换flat table中原始值

滴滴一下 美好出行



uname 2

L



13a

uname

张三

李三

Extract

• 通过hive获取本次需要字典编码的原始值;

• 获取需要字典编码列之前最大编码:

Dict

- 1、MR totalOrderPartition局部编码
- 2、MR 全局编码

key	value	
李三	2	
张三	3	

Merge

通过hive与每列之前的字典合并

key	value
王芳	1
李三	2
张三	3

Replace

• 通过Hive用编码后的值替换flat table中原始值

滴滴一下 美好出行



13a

uname

张三

李三

• 通过hive获取本次需要字典编码的原始值;

Extract

• 获取需要字典编码列之前最大编码;

Dict

- 1、MR 多列并行局部编码,MultipleOutputs,SelfPartitoner
- 2、MR 全局编码, Multiple Outputs

key	value
张三	2
李三	3

Merge

通过hive与每列之前的字典合并

key	value
王芳	1
张三	2
李三	3

Replace

• 通过Hive用编码后的值替换flat table中原始值

滴滴一下 美好出



- 单个字典内部构建全部并行化,理论上在Kylin允许的字典基数范围(Integer)内均可恒定在15分钟内完成
- 多列全局字典在字典编码步骤同样也采用并行化,增加全局字典列基本构建时长不变



② 2019-07-25 10:30:18 UTC

#2 Step Name: Build Global Dict - extract distinct value from data

Duration: 5.53 mins Waiting: 0 seconds





② 2019-07-25 10:35:56 UTC

#3 Step Name: Build Global Dict - parallel part build

Data Size: 13.06 MB

Duration: 1.51 mins Waiting: 45 seconds





② 2019-07-25 10:37:33 UTC

#4 Step Name: Build Global Dict - parallel total build

Data Size: 14 68 MB

Duration: 1.29 mins Waiting: 52 seconds





More Additional

68a

13a

精准去重

- UV
- 留存

数据公海

- 全域字典
- One ID
- One Service





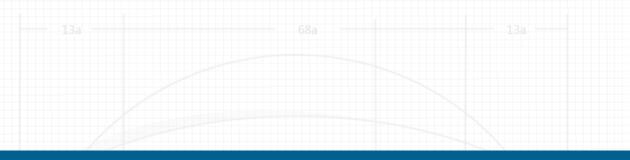


明细钻取

• 营销分析

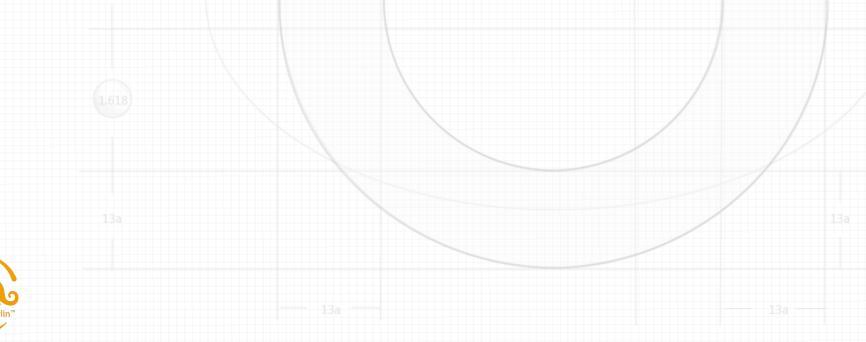
优化

- 小数据量构建提速
- Spark函数



03

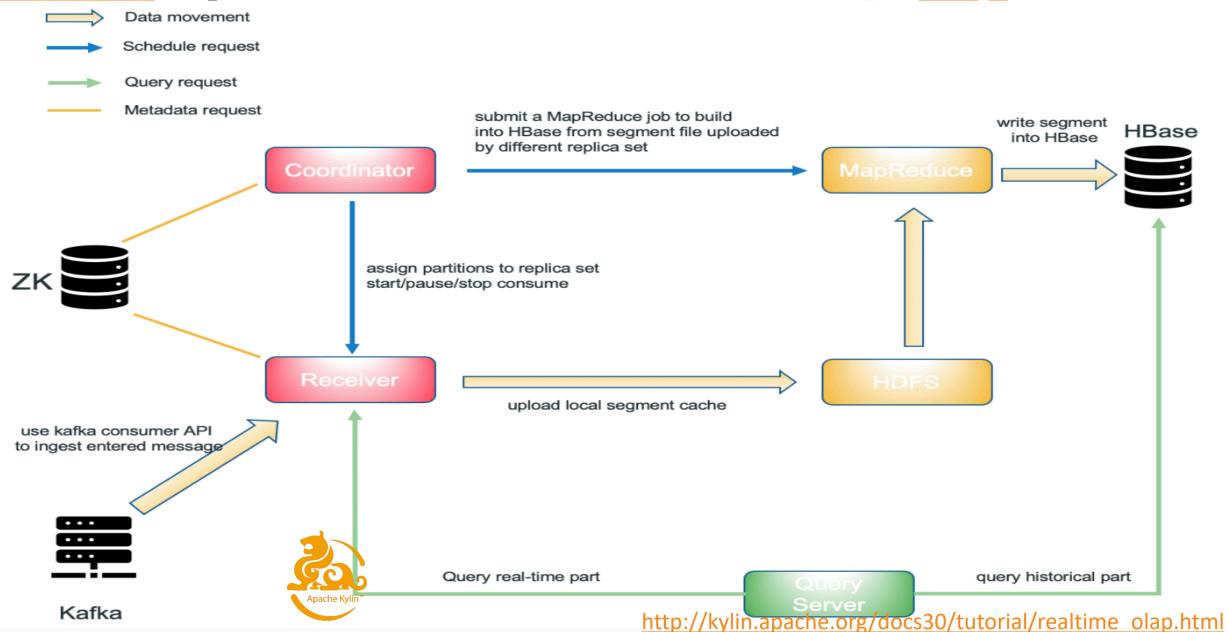
Kylin RT OLAP探索经验分享



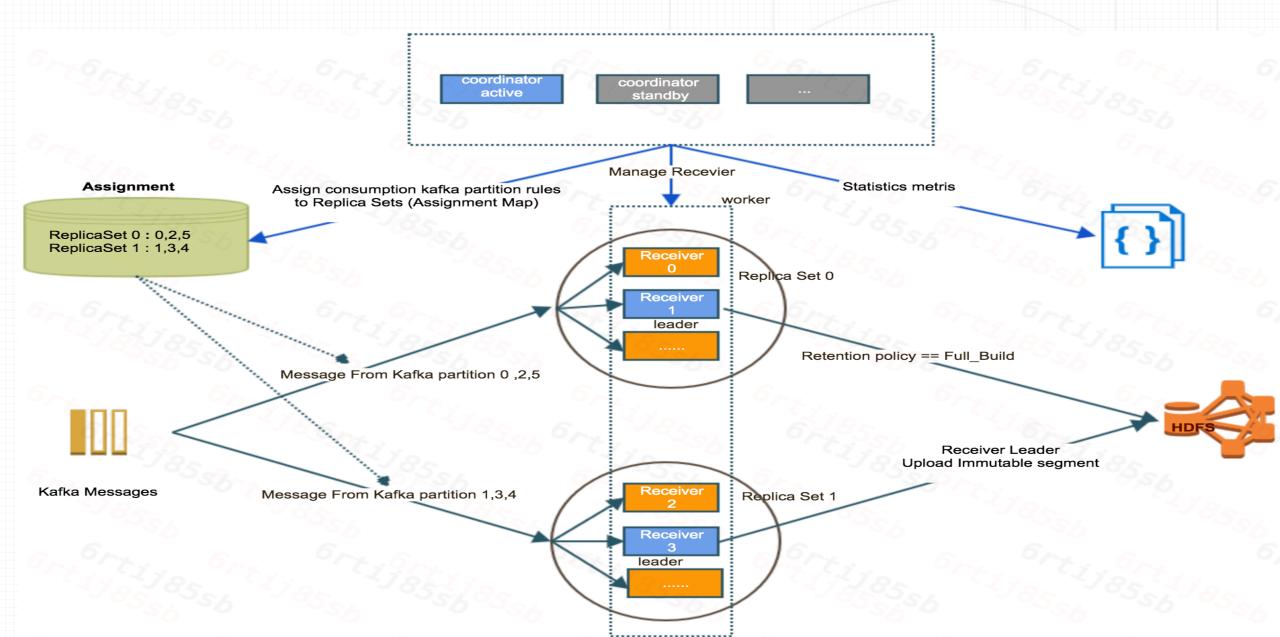




Kylin RT OLAP Architecture



Kylin RT OLAP Architecture



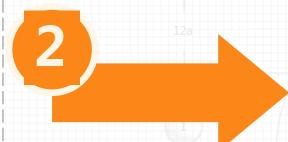
State Flow





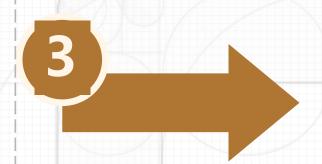
- Reciever负责创建
- EnableCube | Start Receiver 消费Kafka数据时

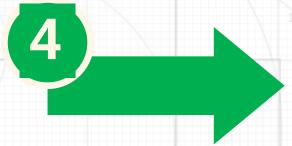






- Segment不再接收新数据
- Receiver负责状态变更
- 消费新数据同时判断达到 Immutable条件的active Segment | reassign





Remote Persisted

- Segment数据上传到HDFS
- Receiver 启动定时后来线程 轮训将满足条件的segement (receiver leader && FULL_BUILD_POLICY) One Build Finish |Coordina 上传到HDFS并触发一次build检查 api

Ready

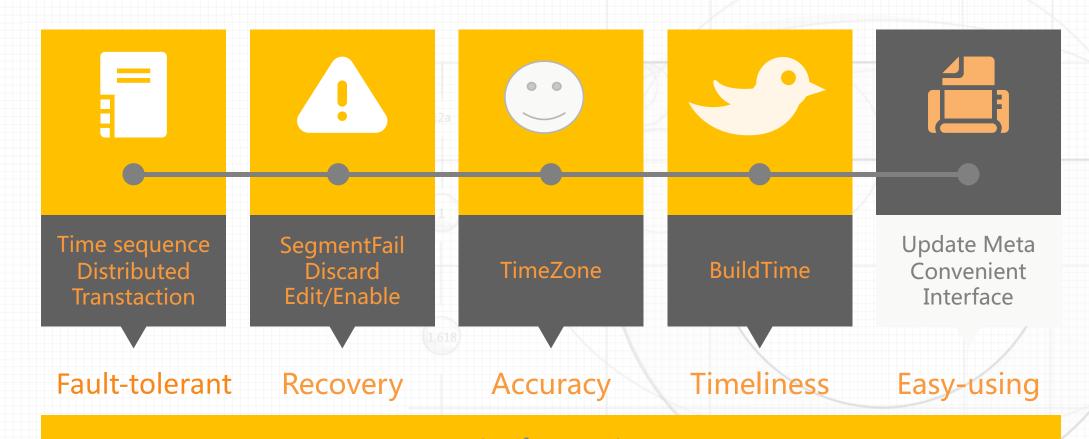
- SegmentBuild完成数据已经
- Coordinator负责提交build
- Receiver uploda HDFS后触

StreamingBuildJobStatusChecker of



Points for Attention











04

Q&A

