



大数据计算 产品文档





文档目录

产品简介

流计算产品概述

相关概念

产品优势

使用与约束

应用场景

快速入门

查看服务

创建项目

操作指南

进入项目

创建数据源和topic (以Kafka举例)

新建数据源

设计态-新建topic

发布到测试态

测试态审批

发布到生产态

申请数据权限 (跨项目的数据库表)

申请数据权限

申请权限

权限审批

数据授权

流计算作业开发

新建作业

在线开发模式

在线开发模式

作业参数配置

数据预览

添加数据源

SQL Operator

添加目标表

连接source、Operator、sink插件

单元测试

版本管理

自定义开发说明

自定义开发说明

ETL Operator

ETL Operator

JAR开发

在线开发

Custom Operator

自定义UDF

JAR开发

在线开发

自定义UDAF

JAR开发

在线开发

维表配置

功能说明

与ETL结合方式

维表参数说明

JAR包上传模式



测试及运行

发布测试

测试运行

发布生产

上线启动

生产运行

最佳实践

场景实践

背景

业务

业务

业务目标

数据格式

场景

技术实践

故障指南

常见问题

产品介绍常见问题

产品使用常见问题



产品简介

流计算产品概述

最近更新时间: 2019-11-26 14:47:19

流计算产品致力于给用户提供更高效、更稳定、更易用的流计算处理服务，解决实时数据处理需求，如网站点击流量分析、电商精准推荐、金融实时风控等流计算应用。

流计算组件通过借助数据开发、智能调度等组件的能力，提供流计算开发测试环境到部署运维等一站式服务。在使用流计算服务时，需要依赖数据采集组件将大数据云外的流式数据采集到Kafka中统一存储。同时，通过数据管理组件获取元数据信息后，在数据开发组件提供的IDE环境中完成流计算应用的开发，最终通过流计算引擎完成实时数据处理。

本文档主要描述流计算服务的功能简介、快速入门、操作指南、常见问题等内容。



相关概念

最近更新时间: 2019-11-26 14:47:19

使用流计算产品时，通常会涉及到以下概念：

- 实时流数据：实时、持续生成的数据，如业务日志、系统日志等各类日志信息。
- 主题（topic）：流计算服务订阅和发布的最小单位，用户可以用topic表示一种流数据，类似与数据库中的表（table）。在大数据云服务中，一个流计算服务的数据源对应一个topic，单个topic可以存储一个或多个日志中的流式数据。
- 流连接分区（partition）：topic存储数据的最小单元，对于吞吐较高的topic，可以创建多个分区。
- 状态保存点（SavePoint）：SavePoint由用户手动触发，可以支持程序升级后，继续从升级前的那个点开始执行计算，保证数据不中断。
- CU（Compute Unit）：流计算所提供计算资源的最小单位，1CU的具体含义为：CPU 1核、内存 4GB。



产品优势

最近更新时间: 2019-11-26 14:47:19

流计算产品具有以下优势：

- 优良的计算性能 流计算支持每秒千万级数据处理，响应延迟达到毫秒级。
- 简单易用的计算服务 上层大数据加工提供这些服务的可视化开发界面，降低使用和开发门槛。借助数据管理服务，能迅速获取元数据信息，解耦与数据源的依赖，实现端到端的数据处理。
- 灵活的资源规划 用户可根据自己的需求，灵活划分资源组，满足个性需求，不浪费。



使用与约束

最近更新时间: 2021-09-15 15:59:08

支持的流计算引擎：流计算服务支持Flink、Spark Streaming两种计算引擎，其中Flink为全功能支持，Spark Streaming只支持上传JAR包方式使用，推荐用户优先使用Flink引擎，两种引擎的对比如下：

流计算引擎对比	Flink	Spark Streaming
引擎版本	Flink 1.11.3	Spark 2.3.0
产品生态	Apache顶级项目，Flink 产品发展速度快，在流处理方面处于业界最前沿水平	归属于Spark的子项目，Spark社区规模大，产品相对更加成熟
运行模式	事件驱动	时间驱动（微批处理，运行的时候需要指定批处理的时间，每次运行 任务时处理一个批次的数据）
时间机制	支持三种时间机制：事件时间，注入时间，处理时间	只支持处理时间
支持程度	在大数据云服务中，支持可视化拖拽、上传JAR包两种开发方式	在大数据云服务中，只支持JAR包的开发方式

支持的source数据源、sink数据源：流计算服务支持的source源、sink源如下：

	开源组件	组件版本
支持的source源	Kafka	0.10及以上版本
支持的sink源	Kafka	0.10及以上版本
	Oracle	11.2.0.1.0
	MySQL	5.6及以上版本
	HBase	2.0
	Redis	4.0
	Elasticsearch	6.4.2
	InfluxDB	2.1.5
	Phoenix	5.0.0
	Hive	3.1.0



应用场景

最近更新时间: 2019-11-26 14:47:19

流计算产品具有以下应用场景:

- 互联网点击流分析: 用户在网站浏览时会产生许多的点击行为, 对这些点击行为加以分析可精确把握热点趋势。借助于云端流计算服务, 构建实时分析可达到分钟级, 对用户行为数据进行实时汇聚分析, 持续地挖掘出有价值信息, 帮助更好地做出运营决策, 改进用户体验。
- 金融实时风控: 在众多金融风险中, 及早探测到风险往往能有效地减少损失, 将金融交易大数据与流计算服务相结合, 引入特征模型算法, 及早地过滤出诸如盗刷卡等异常交易行为, 实施风险控制, 提升金融安全性。
- 物联网监控: 在工业设备的运转过程中, 及早发现潜在故障会极大降低维修成本。借助于云端流计算服务, 及时收集设备传感器数据, 并进行聚合、分析筛选, 可实现秒级设备异常告警, 提升设备利用率。
- 电商精准推荐: 在电商交易中, 借助于云端流计算服务, 实时提取特征变量, 及时跟踪用户关注品类, 预测用户消费趋势, 为精准推荐提供基础能力。从而提升用户购物体验, 促进消费行为。



快速入门 查看服务

最近更新时间: 2019-11-11 07:47:14

登录后, 点击左侧的“资源管理-产品与服务”菜单, 或点击上方的“产品与服务”, 可查看所有开通的服务。

建行大数据云 产品与服务 帮助与文档 uat_te

资源管理 > 已购产品

产品名称	产品分类 (大数据开发套件产品) v	资源配额(剩余/总量)	购买时间	所在项目数	状态	操作
数据采集	大数据开发套件产品	DCU:8/8 个 DIU:0/0 个	2019-08-16 15:42:29	51	● 正常可用	立即使用 资源配置
数据集成	大数据开发套件产品	CU:48/48 个 DCU:60/60 个	2019-08-16 15:42:29	52	● 正常可用	立即使用 资源配置
离线计算	大数据开发套件产品	CU:32/32 个 DCU:40/40 个	2019-08-16 15:42:29	53	● 正常可用	立即使用 资源配置
流计算	大数据开发套件产品	CU:12/12 个	2019-08-16 15:42:29	46	● 正常可用	立即使用 资源配置
数据挖掘	大数据开发套件产品	CU:40/40 个 DCU:60/60 个	2019-08-16 15:42:29	46	● 正常可用	立即使用 资源配置
数据服务	大数据开发套件产品	CU:6/6 个 DMU:0/0 个	2019-08-27 11:49:38	36	● 正常可用	立即使用 资源配置

建行大数据云 产品与服务 帮助与文档

平台概览

资源管理

已购产品

项目空间

数据管理

运维中心

组织管理

个人中心

产品详情 立即购买

产品详情 立即使用

产品详情 立即使用

产品详情 立即使用

流计算

流计算

建行大数据云流计算服务提供面向高速流式数据进行实时快速计算的解决方案, 开通后可满足您实时风控、...

产品详情 立即使用

数据挖掘

数据挖掘

建行大数据云数据挖掘服务提供可视化建模和Notebook建模能力, 内含多种数据挖掘算法, 开通后您可以...

产品详情 立即使用

数据服务

数据服务

建行大数据云数据服务产品是构建数据中台的基础, 可提供统一的数据访问能力, 提供可视化的API开发和服...

产品详情 立即使用

创建项目

最近更新时间: 2019-11-11 07:47:14

点击页面左侧“项目空间-我的项目”选项卡，可以看到当前所有的项目列表。

项目名称	创建时间	负责人	开通的服务	最近更新时间	状态	操作
Test	2019-07-15 15:07:29	ccb_test_4	数据挖掘	2019-07-15 15:07:29	正常	配置项目 修改服务 更多
test_xcy	2019-07-15 15:01:48	ccb_test_4	流计算 数据采集 数据集成 批计算 数据挖掘 数据服务	2019-07-15 15:01:48	正常	配置项目 修改服务 更多
test_1	2019-07-09 09:47:15	wangdasheng_xm	流计算 数据采集 数据集成 批计算 数据挖掘 数据服务	2019-07-09 09:47:15	正常	配置项目 修改服务 更多
litman	2019-07-08 19:11:49	ccb_test_4	流计算 数据采集 数据集成 批计算 数据挖掘 数据服务	2019-07-08 19:11:49	正常	配置项目 修改服务 更多
test20190704	2019-07-04 17:39:30	ccb_test_4	数据采集 数据集成 批计算 数据挖掘 数据服务	2019-07-04 17:39:30	正常	配置项目 修改服务 更多
test0705	2019-07-04 16:55:18	ccb_test_4	流计算 数据采集 数据集成 批计算 数据挖掘 数据服务	2019-07-04 16:55:18	正常	配置项目 修改服务 更多

点击蓝色的“新建项目”按钮，弹出“新建项目”对话框。

< 创建项目

1 基本信息 ———— 2 选择服务 ———— 3 配置资源 ———— 4 信息确认

* 项目名称: test

项目描述:

取消 上一步 下一步

这里可以根据需要选择各项服务，本次新建项目需包含“流计算”，其余服务按需选择。点击“下一步”配置项目的各项属性参数。



- ① 基本信息
- ② 选择服务
- ③ 配置资源
- ④ 信息确认

请选择服务 (至少选择一个)

- 数据采集 已购买
数据采集是一种面向开发者提供的端到端数据采集服务，是数据进入大数据平台的第一道关卡。支持日志文件、数据库、报文接口等多种数据源的的流式、批量采集，提供向导式的采集配置、在线Agent管理、实时采集任务监控等服务能力。
- 数据集成 已购买
建行大数据云数据集成服务支持金融级强监管要求下的数据集成功能，开通后您可以轻松实现数据在不同系统间的整合和流通，实现向业务快速供数。
- 离线计算 已购买
离线计算是一种经济并高效的分析和处理海量数据的开发平台，可提供快速、完全托管的PB级数据仓库解决方案，支持以SQL代码、shell脚本、拖拽式等多种开发模式构建金融企业级数仓。
- 流计算 已购买
建行大数据云流计算服务提供面向高速流式数据进行实时快速计算的解决方案，开通后可满足您实时风控、实时推荐等业务场景数据开发。
- 数据挖掘 已购买
建行大数据云数据挖掘服务提供可视化建模和Notebook建模能力，内含多种数据挖掘算法，开通后您可以轻松构建算法模型，快速服务风险控制、精准营销等业务。
- 数据服务 已购买
建行大数据云数据服务产品是构建数据中台的基础，可提供统一的数据访问能力，提供可视化的API开发和服务治理能力。

默认包含的公共基础服务

- 基础版数据管理
- 运维中心 (调度系统)
- 项目管理

点击下一步，在项目资源组页面，通过下拉框中，为所选的服务选择一个资源分组。配置完成后，点击“下一步”。

- ① 基本信息
- ② 选择服务
- ③ 配置资源

项目资源组

流计算: 默认资源分组 (Yarn) ▾

默认资源分组 (Yarn)

对以上配置



信息进行确认，无误后点击“确定”按钮。

- ① 基本信息
- ② 选择服务
- ③ 配置资源
- ④ 信息确认

基本信息

项目名称: test

项目描述:

选择的服务

流计算

配置资源

流计算资源组: 默认资源分组

取消

上一步

确定



操作指南

进入项目

最近更新时间: 2019-11-11 07:47:14

创建成功的项目会显示在项目列表的最上面。点击“修改服务”可以对项目中的服务进行新建、修改配置、删除等操作。点击“开通的服务”列中的流计算，进入此项目下流计算的开发页面。

项目空间 > 我的项目

新建项目

项目名称	创建时间	负责人	开通的服务	最近更新时间	状态	操作
test	2019-09-03 21:29:26	uat_test	流计算	2019-09-03 21:29:26	● 正常	配置项目 修改服务
my_stream_test_proj	2019-09-03 17:55:08	uat_test	流计算 数据采集 数据集成 离线计算 数据挖掘 数据服务	2019-09-03 18:00:46	● 正常	配置项目 修改服务
atk_proj08	2019-09-02 13:14:08	uat_test	流计算 数据采集 数据集成 离线计算 数据挖掘 数据服务	2019-09-02 13:14:08	● 正常	配置项目 修改服务
jhh	2019-08-31 00:25:41	uat_test	流计算 数据采集 数据集成 离线计算 数据挖掘 数据服务	2019-08-31 00:25:41	● 正常	配置项目 修改服务
eewre	2019-08-30 19:53:54	uat_test	流计算 数据采集 数据挖掘	2019-08-30 19:53:54	● 正常	配置项目 修改服务
z830	2019-08-30 17:56:20	uat_test	流计算 数据集成 数据挖掘 数据服务	2019-08-30 20:59:01	● 正常	配置项目 修改服务

创建数据源和topic（以Kafka为例）

新建数据源

最近更新时间: 2019-11-11 07:47:14

为了将数据从指定表中读写，需要按照业务需求创建对应的数据源/库/表。在“数据管理”→“数据源管理”中新建数据源。



serverAddress填写对应Kafka集群地址，可以进行连通性测试，最后点击“确认”。



建行大数据云

控制台首页 产品与服务 帮助与文档

数据管理 > 数据源管理

新增数据源

* 数据源名称: test_kafka

* 数据源环境: 测试环境 生产环境 测试和生产

* 生产环境配置:

* serverAddress: 127.0.0.1:9092

连通性测试

* 是否托管: 是
托管元数据请配置管理员账号,允许平台创建库表并管理数据权限

数据源描述:

上一步 完成

设计态-新建topic

最近更新时间: 2019-11-11 07:47:14

在“数据管理”→“元数据管理”→“库表管理”中新建数据库和数据表（Kafka不需要新建数据库，直接新建topic即可）。

数据管理 > 库表管理

设计态 测试环境 生产环境

关系型数据库
Kafka
COS
ArangoDB
HBase
Elasticsearch

新建Topic 删除 发布到测试

所属项目: 全部 请输入表名称

Topic名称	Topic中文名称	所属项目	创建人	创建时间	版本	表描述	发布状态	发布的数据库	操作
<input type="checkbox"/> test001		liufei_dg_test	ccb_test	2019-08-20 17:23:44	1566293023668		已发布	测试:dgKafka.topicdb_156_23, defaultKafka.topicdb_156_84	编辑 删除
<input type="checkbox"/> source_php_logtime		ss_test_php	ccb_test	2019-08-20 14:26:44	1566287493241		已发布	测试:defaultKafka.topicdb_119_84	编辑 删除

新建topic时需指定topic归属的项目，即此topic可在归属项目下使用。



创建topic

① 配置基本信息

* Topic名称:
仅支持英文、数字、下划线，且不能以数字和下划线开头，最大100字符

Topic中文名:

* 所属项目:

* Topic数据格式: json CSV

* Topic分区数:

生命周期:

Topic描述:

创建topic

① 配置基本信息

② 字段设置

字段名称	字段中文名	字段类型	字段长度	字段密级	字段脱敏	字段描述	操作
<input type="text" value="id"/>	<input type="text"/>	字符串 <input type="text"/>	<input type="text" value="50"/>	<input type="text" value="5"/>	不脱敏 <input type="text"/>	<input type="text"/>	下移 删除
<input type="text" value="name"/>	<input type="text"/>	字符串 <input type="text"/>	<input type="text" value="50"/>	<input type="text" value="5"/>	不脱敏 <input type="text"/>	<input type="text"/>	上移 删除

发布到测试态

最近更新时间: 2019-11-11 07:47:11

完成topic创建动作后，topic信息可在“库表管理的”的发布态查看，点击列表操作中的“发布至测试”，可申请将topic发布到测试环境，并在测试环境使用。

数据管理 > 库表管理

设计态 测试环境 生产环境

关系型数据库 Kafka COS ArangoDB HBase Elasticsearch

新建Topic 删除 发布到测试

所属项目: 全部 请输入表名称

Topic名称	Topic中文名称	所属项目	创建人	创建时间	版本	表描述	发布状态	发布的数据库	操作
s_test_topic		s_test	ccb_test	2019-08-21 10:18:38	1566353918120		待发布		编辑 删除 发布到测试
test001		liufei_dg_test	ccb_test	2019-08-20 17:23:44	1566293023668		已发布	测试:dgKafka.topicdb_156_223, defaultKafka.topicdb_156_84	编辑 删除
course_cho_langtime		cc_test_cho	ccb_test	2019-08-20 14:26:44	1566387403341		已发布	测试:defaultKafka.topic	编辑 删除

发布测试时，需指定此topic归属的数据源。

发布到测试环境

* 数据源类型: kafka

* 数据源名称: defaultKafka

发布说明:

取消 确认

发布测试完毕后，可在测试环境中

筛选查看：

数据管理 > 库表管理

设计态 测试环境 生产环境

我设计的库表 我有权限的表 项目账号的库表 项目有权限的表

发布到生产

数据源种类: 数据源名称: 所属项目: s_test 请输入库名称

Topic名称	库中文名称	所属项目	数据源种类	数据源名称	创建人	创建时间	用途描述	操作
topicdb_160_84		s_test	kafka	defaultKafka	ccb_test	2019-08-21 10:19:41		发布到生产 发布到测试

上一页 1 下一页 每页显示 10行 / 页 共 1 条

topicdb_160_84

打标

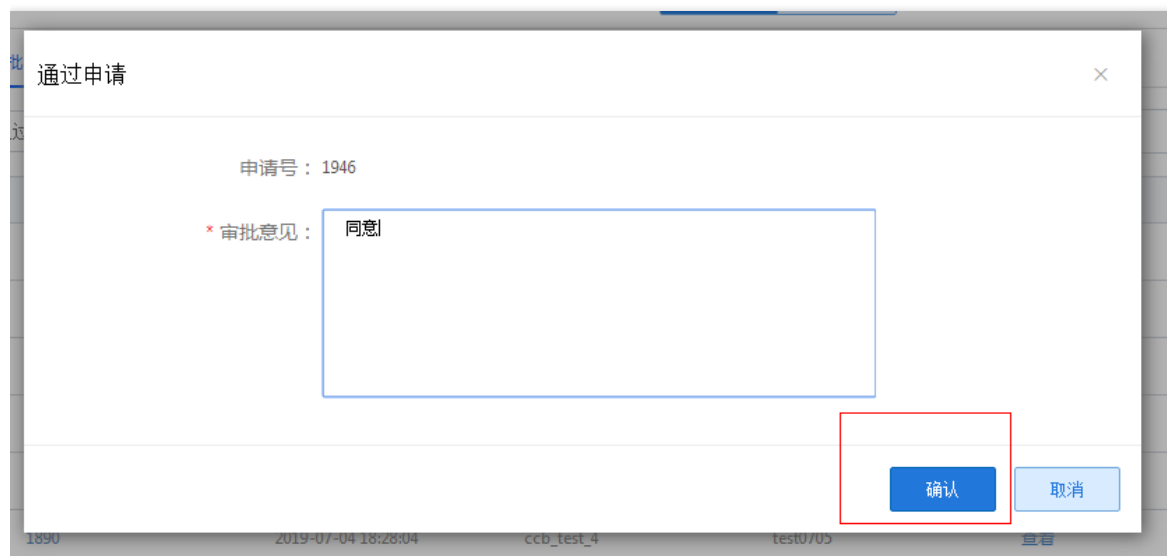
Topic名称	库中文名称	所属项目	数据源种类	数据源名称	创建人	创建时间	用途描述	操作
s_test_topic		s_test	kafka	defaultKafka	ccb_test	2019-08-21 10:18:38		发布到生产 详情

上一页 1 下一页 每页显示 10行 / 页 共 1 条

测试态审批

最近更新时间: 2019-11-11 07:47:11

使用审批账号登陆大数据云平台后，在“数据管理”→“元数据管理”→“发布审批”→“测试环境”→“待我审批”中可以查看待审批的数据库/表，点击审批通过，并在弹出的窗口中，填写审批意见后，点击“确认”，即完成了审批流程。审批后，待发布的数据库/表成功发布至测试环境。



在“测试环境”选项卡下的“我



已审批”选项卡下，可以看到所有已审批通过的数据库/表。

数据管理 > 发布审批

数据空间 测试环境 生产环境

待我审批 我已审批

申请人: 全部 所属项目: 全部 申请状态: 全部 请选择日期

申请号	发布类型	申请时间	申请人	申请项目	流程状态
3838	技术元数据/表级	2019-09-04 16:45:41	uat_test	dts_test_830	查看
3837	技术元数据/表级	2019-09-04 16:44:08	uat_test	dts_test_830	查看
3836	技术元数据/表级	2019-09-04 16:44:06	uat_test	dts_test_830	查看
3835	技术元数据/表级	2019-09-04 16:00:31	uat_test	zeng_test819	查看
3826	技术元数据/表级	2019-09-04 10:11:25	uat_test	test11111	查看
3825	技术元数据/表级	2019-09-04 10:08:27	uat_test	test11111	查看
3824	技术元数据/表级	2019-09-04 10:01:06	uat_test	test11111	查看
3823	技术元数据/表级	2019-09-03 17:06:08	uat_test	atk_proj08	查看

发布到生产态

最近更新时间: 2019-11-11 07:47:11

将topic发布到测试后, 点击“操作”列“显示表”, 页面下方会显示该数据库中所有的表, 其中每个表对应一个topic。同样, 点击“操作”列“发布到生产”, 经过同样的审批流程后, 即可将该topic发布到生产。

数据管理 > 库表管理

设计态 测试环境 生产环境

我设计的库表 我有权限的表 项目账号的库表 项目有权限的表

发布到生产

库名称	库中文名称	所属项目	数据源种类	数据源名称	创建人	创建时间	用途描述	操作
topicdb_160_84		s_test	kafka	defaultKafka	ccb_test	2019-08-21 10:19:41		发布到生产 显示表

上一页 1 下一页 每页显示 10行 / 页 共 1 条

建行大数据云

数据管理

数据源管理

元数据管理

标签管理

业务元数据

库表管理

元数据关联

数据搜索

数据目录

发布审批

发布申请

数据权限管理

控制台首页 产品与服务 帮助与文档 uat_test

数据管理 > 发布审批

数据空间 测试环境 生产环境

待我审批 我已审批

批量通过 批量驳回

申请人: 全部 所属项目: 全部 请选择日期

申请号	发布类型	申请时间	申请人	申请项目	流程状态	操作
3646	数据质量/规则配置	2019-08-30 17:23:36	uat_test	test_lir_0816	查看	通过 驳回
3645	数据质量/规则配置	2019-08-30 17:22:50	uat_test	test_lir_0816	查看	通过 驳回
3644	数据质量/规则配置	2019-08-30 17:20:57	uat_test	ZHH	查看	通过 驳回
3518	数据质量/规则配置	2019-08-29 23:28:50	uat_test	test_lir_0816	查看	通过 驳回
3514	数据质量/规则配置	2019-08-29 22:58:58	uat_test	test_lir_0816	查看	通过 驳回
3511	数据质量/规则配置	2019-08-29 22:33:50	uat_test	test_lir_0816	查看	通过 驳回

审批通过后, 点击“生产环境”, 在“我设计的库表”选项卡下, 按所属项目进行筛选后, 即可看到已发布到生产的库, 点击“显示表”, 即可看到已发布到生产的topic。

数据管理 > 库表管理

设计态 测试环境 生产环境

我设计的库表 我有权限的表 项目账号的库表 项目有权限的表

数据源种类: 数据源名称: 所属项目: s_test 输入库名称: 搜索

库名称	库中文名称	所属项目	数据源种类	数据源名称	创建人	创建时间	用途描述	操作
topicdb_160_84		s_test	kafka	defaultKafka	ccb_test	2019-08-21 10:19:41		显示表

上一页 1 下一页 每页显示 10行 / 页 共 1 条



申请数据权限（跨项目的数据库表）

申请数据权限

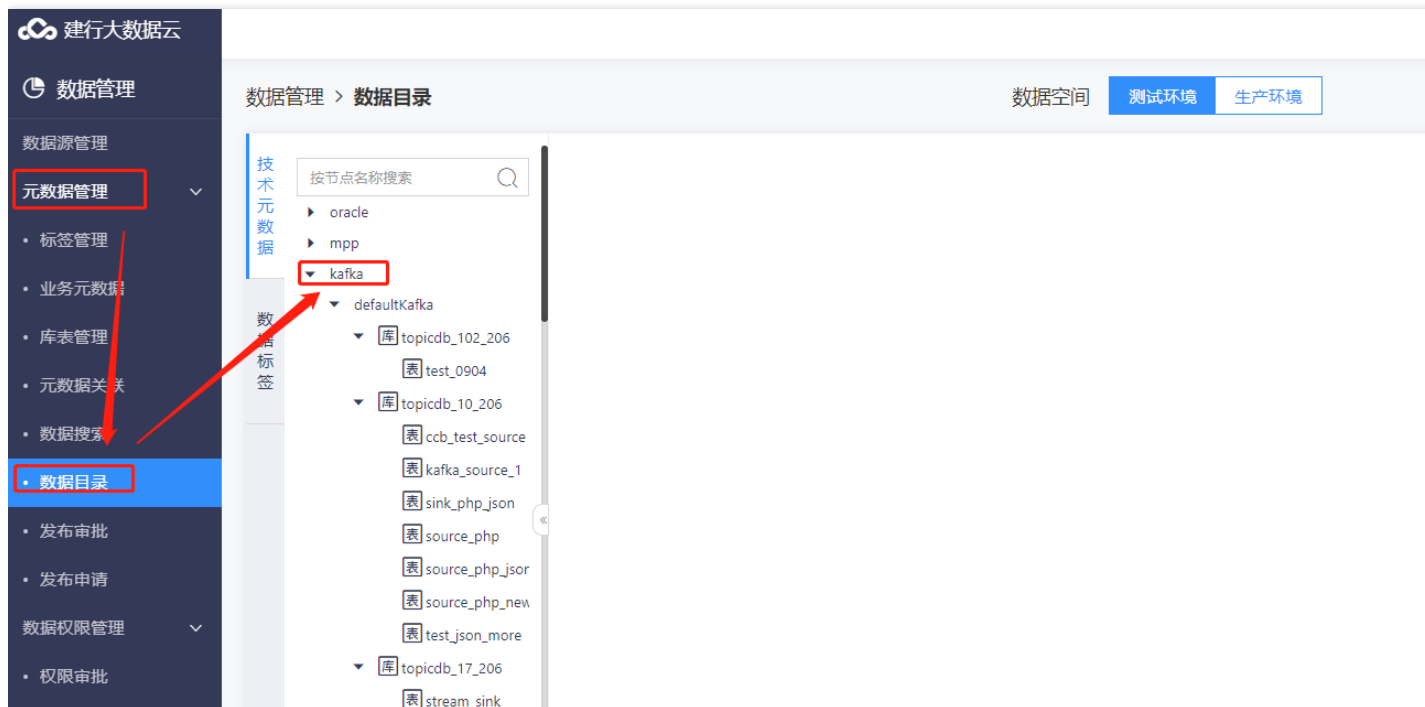
最近更新时间: 2019-11-11 07:48:26

流计算作业如果需要使用到其他项目的数据库表均需要提前通过数据管理组件进行申请和注册。

申请权限

最近更新时间: 2019-11-11 07:47:11

可通过如下步骤申请权限：(1) 除数据管理员角色以外的其他角色的人员（包括：项目管理员，开发人员，运维人员，业务人员等）进入数据管理模块，点击“元数据管理”下面的“数据目录”。



(2) 选中数据源Kafka，搜索表名称，点击对应的表，会展示出表详情。(3) 点击【权限申请】按钮，弹出权限申请的页签，选择申请对象，选择权限，选择申请对象，选择申请项目，点击【确认】按钮，权限申请成功。



说明：1. 申请类型为列权限时，申请对象只能为个人；申请类型为表权限时，申请对象可以为项目，个人或者项目和个人。2. Kafka类型数据源，没有列权限，只有表权限，不需要选择申请类型；Kafka类型数据源的权限选择只有两个选项：insert和select。3. 申请个人和项目的表权限，会同时生成两条申请记录，审批记录也是两条。

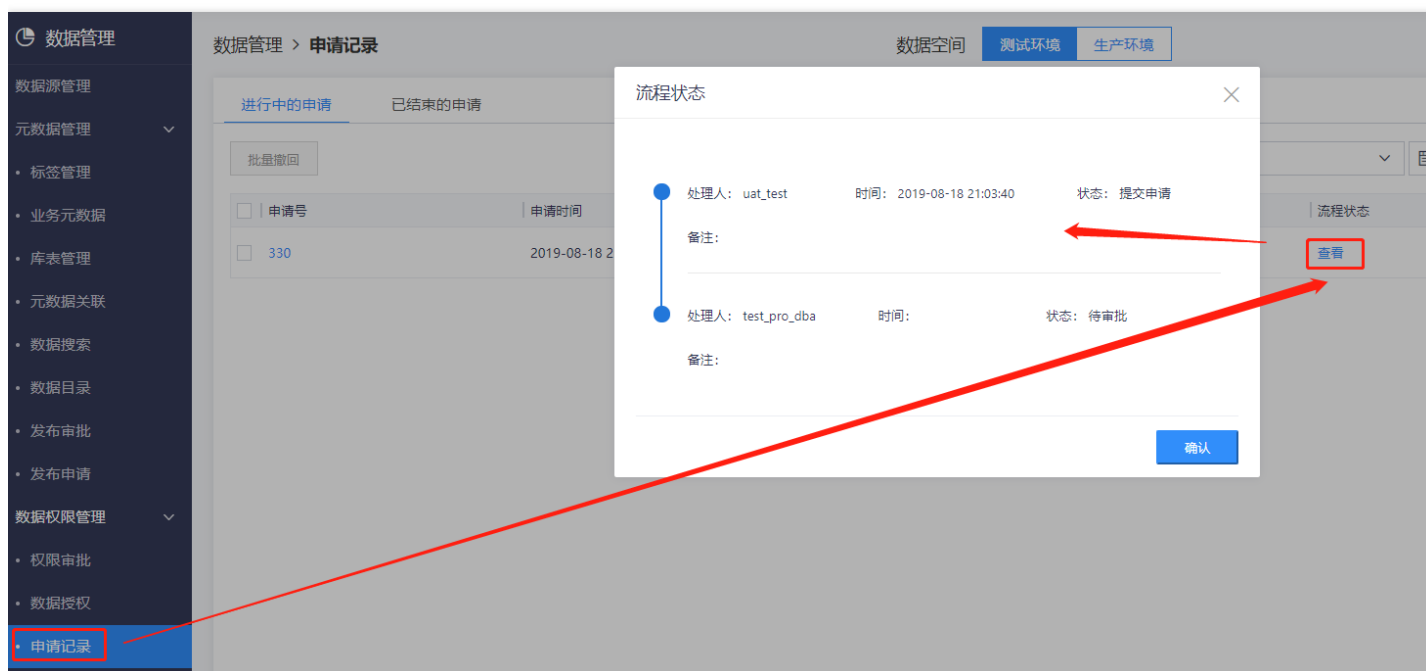
权限审批

最近更新时间: 2019-11-11 07:47:11

可通过如下步骤审批权限：(1) 使用数据管理员账号登录进入数据管理模块，点击“数据权限管理”，展开下拉列表中点击“权限审批”。



(2) 根据申请号选择提交的申请，点击“查看”可查看提交申请的用户，备注中显示提交原因。



(3) 点击“通过”（或者驳回）审批该申请。说明：通过申请，如果申请的时候选择的是项目，该表的项目白名单会新增一条数据，项目管理员对该数据表有对应的权限；如果选择的是个人，用户白名单新增一条数据，申请人对该表有对应的权限。

数据授权

最近更新时间: 2019-11-11 07:47:11

可通过如下步骤进行数据授权：(1) 使用数据管理员或者租户的账号登录，进入数据管理模块，点击“数据权限管理”，展开的下拉框中点击“数据授权”，数据空间显示数据授权列表。

表名	表中文名称	所属项目	数据源种类	数据源名称	创建人	公开状态	黑白名单	操作
table_test_renhang		dts_test_830	hive	defaultHive	uat_test	未公开	查看	公开设置 白名单设置
qqqq		zeng_test819	cos	cos	uat_test	未公开	查看	公开设置 白名单设置
test_table_0904		test11111	mysql	mysql	uat_test	未公开	查看	公开设置 白名单设置
test_0904		test11111	kafka	defaultKafka	uat_test	未公开	查看	公开设置 白名单设置
datatest		atk_proj08	mysql	dj_mysql	uat_test	未公开	查看	公开设置 白名单设置
data_ll		atkproject08	mysql	dj_mysql	uat_test	未公开	查看	公开设置 白名单设置
test1		dsp830	cos	cos	xperia	未公开	查看	公开设置 白名单设置
ccb_test_es_sink		ss_test_main1	elasticsearch	elasticsearch	uat_test	未公开	查看	公开设置 白名单设置
ccb_test_hbase_sink		ss_test_main1	hbase	hbase_online	uat_test	未公开	查看	公开设置 白名单设置
ccb_test_mysql_sink		ss_test_main1	mysql	mysql	uat_test	未公开	查看	公开设置 白名单设置

(2) 选中需要授权的表，点击“白名单设置”按钮，弹出白名单设置标签，选择授权类型，选择权限，选择授权对象，授权项目和人都支持多选，点击“确认”。说明：授权类型为表权限时，授权对象可以为项目，个人或者项目和个人；授权对象为列权限时，授权对象只可以为个人；个人和项目均支持多选和搜索。

白名单设置

- * 授权类型: 表授权 列授权
- * 选择权限: insert delete update select
- * 授权对象: 项目 个人
- * 授权项目: 123qwe atk1 atkproject08
- * 授权个人: liuxiangdong zsl2007 xperia

授权说明:

取消 确认

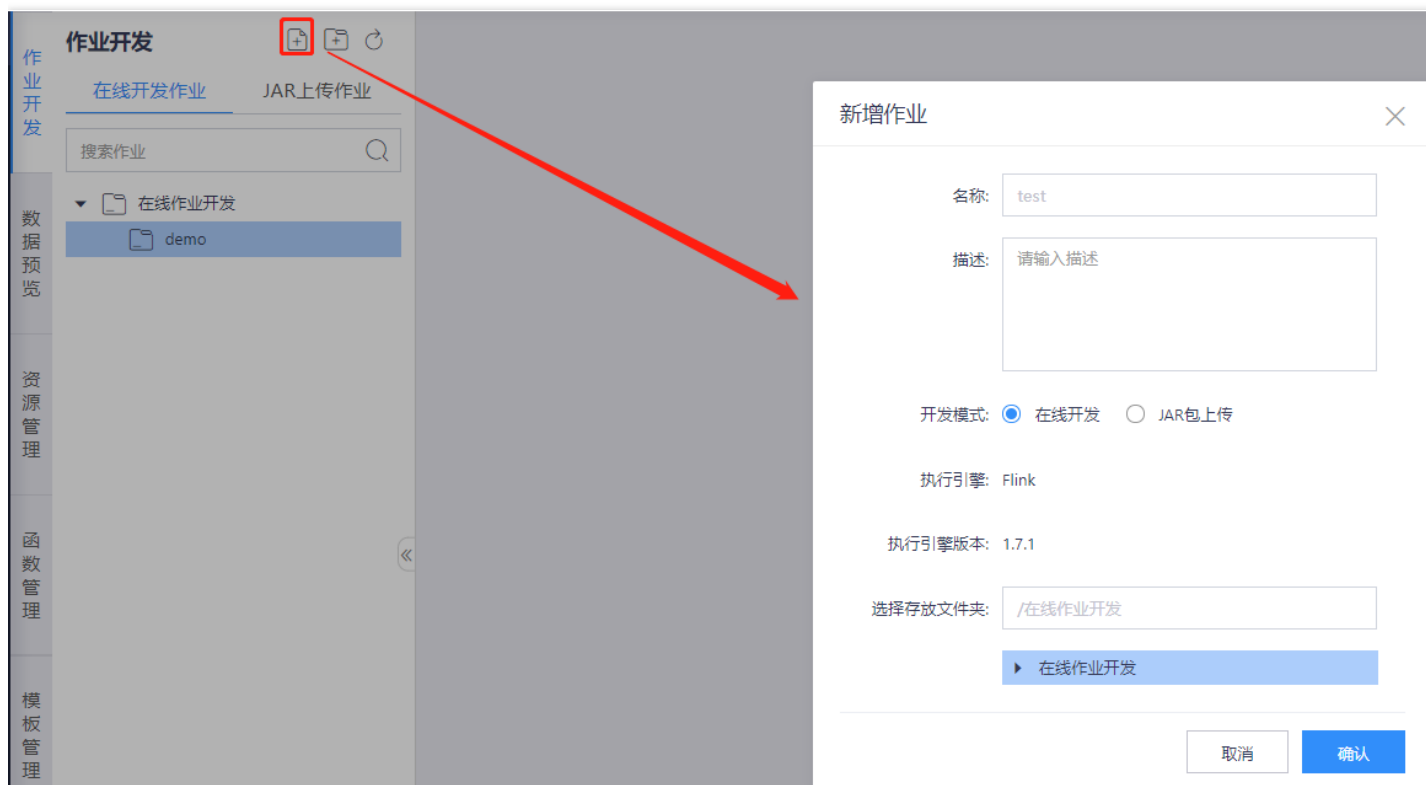
(3) 授权成功后，点击“查看”按钮，可以看到授权的项目白名单和个人白名单。说明：查看项目白名单和用户白名单，点击“新增白名单”按钮，可以给表增加白名单，操作类似白名单设置。

流计算作业开发

新建作业

最近更新时间: 2019-11-11 07:47:11

进入流计算组件, 点击“新建文件夹”和“新建作业”, 创建一个新的流计算作业, 创建完毕后即进入流计算作业的开发IDE界面。左侧为本项目流计算任务树, 右侧为任务组件和画布区。可以通过拖拽算子和连线轻松编辑流计算作业。创建文件夹/作业时, 支持: 在线开发和JAR包上传2种模式, 其中在线开发模式支持Flink引擎; JAR包上传模式, 依赖线下开发并上传至云平台的JAR包, 支持Flink和Spark Streaming两种引擎。



在线开发模式

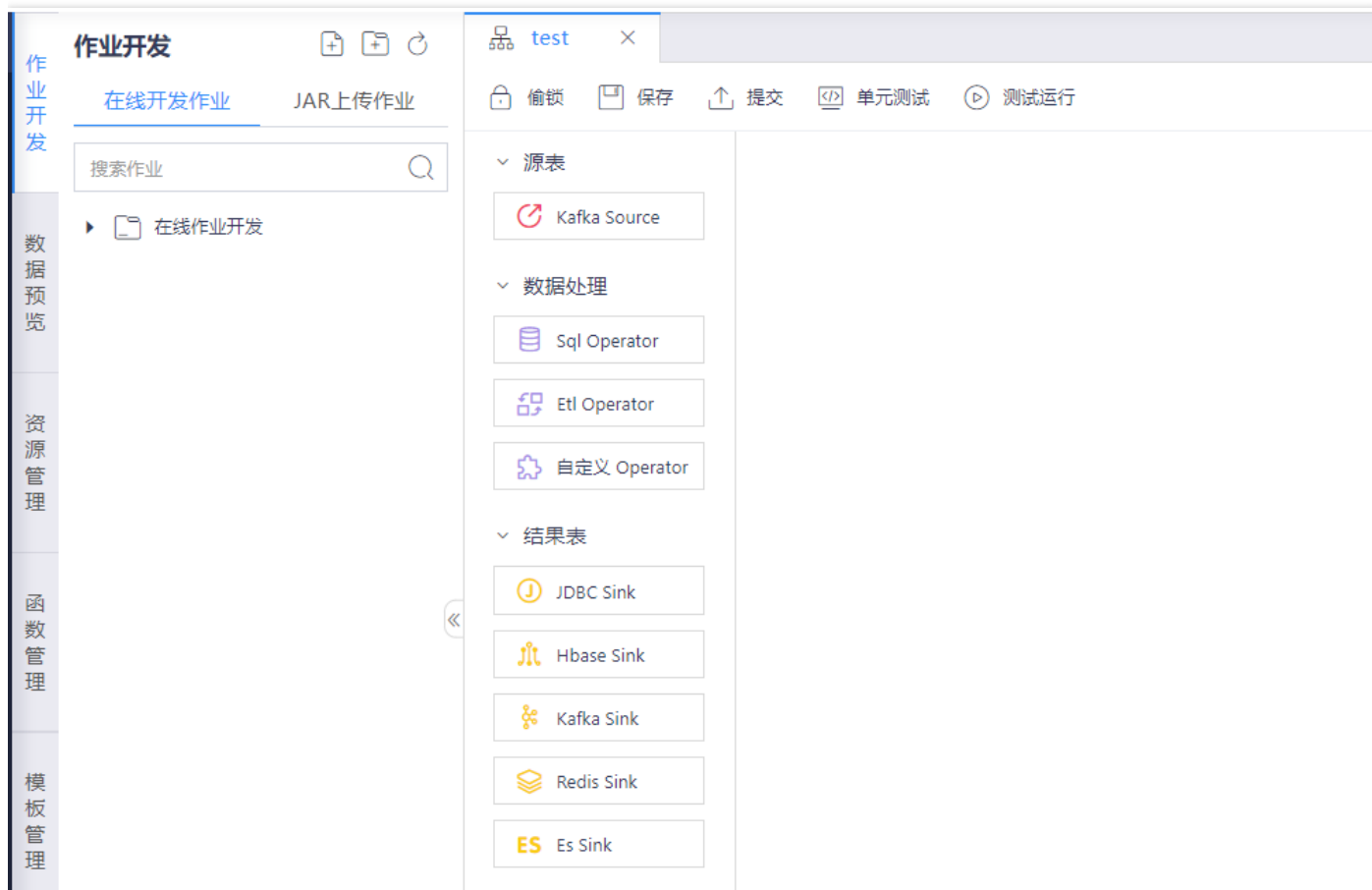
在线开发模式

最近更新时间: 2019-11-26 14:47:12

新建在线开发模式的流计算作业，将进入在线开发的IDE界面，中间的“源表”“数据处理”“结果表”等Operator支持拖拽。拖拽后，将自动弹出算子的参数设置栏，完成参数配置，连接各Operator后，即可生成流作业。流计算支持：

- 双流joining（支持2个Kafka数据源）、
- 维表管理（支持Kafka数据源关联MySQL、Oracle等其他数据源）、
- JAR包依赖（通过上传的JAR包进行ETL等操作）、
- 单元测试（支持线下上传数据包，进行乱序、延时测试）

等功能，灵活支持各种业务场景。





作业开发

在线开发作业 JAR上传作业

test

解锁 保存 提交 单元测试 测试运行

搜索作业

在线作业开发

数据预览

资源管理

函数管理

模板管理

源表

- Kafka Source

数据处理

- Sql Operator
- ETL Operator
- 自定义 Operator

结果表

- JDBC Sink
- Hbase Sink
- Kafka Sink
- Redis Sink
- Es Sink

```
graph TD; K1[KafkaSource_...] --> S[Sql_15675186...]; K2[KafkaSource_...] --> S; S --> J[JdbcSink_156...]
```

JDBC Sink:参数设置

名称: JdbcSink_1567518633803

*数据源类型: 请选择数据源类型

*数据源: 请选择数据源

*目标数据库: 请选择目标数据库

*目标数据表: 请选择目标数据表

batchsize: 1 pieces

operator并行度: 请输入operator并行度

刷新时间(ms): 10000

作业参数配置

最近更新时间: 2019-11-26 14:47:12

新建作业后, 会自动弹出作业参数配置弹窗, 可先对流作业的信息进行配置:

- 时间属性: 支持Event Time和Processing Time 2种, 分别表示流计算应用中, 按照什么时间进行数据处理。
- CheckPoint触发间隔: 即设置检查点的时间间隔。
- 默认并行数: 即流作业的并行数。
- 自定义环境参数: 非必须, 可配置Flink引擎的其他环境参数, 缓存属性等。

作业名: test_002

作业ID: e5960d28445845ac939733166d04e290

作业描述:

开发模式: 在线开发

执行引擎: flink

执行引擎版本: 1.7.1

时间属性: Processing Time Event Time

Checkpoint触发间隔(秒):

默认并行数:

自定义环境参数

序号	参数名	参数值	操作
无数据			
+ 增加一行			

作业配置

版本

维表配置

包依赖

数据预览

最近更新时间: 2019-11-11 07:47:11

数据预览功能可支持提前对source、sink的数据结构和数据进行预览, 辅助开发。

The screenshot displays the 'Data Preview' (数据预览) interface. On the left, a sidebar contains '作业开发' (Job Development) and '数据预览' (Data Preview) sections. The '数据预览' section is active, showing a table of column information:

序号	列名	类型	注释
1	ren1	NUMBER	
2	ren2	STRING	
3	ren3	DECIMAL	
4	ren4	DATE	
5	ren5	TIME	
6	ren6	DATETIM E	

The main workspace shows a workflow diagram with nodes: 'KafkaSource_1564723', 'Custom_1564...', and 'JdbcSink_1564...'. A red arrow points from the table to the 'KafkaSource_1564723' node. The right sidebar shows the 'Kafka Source:参数设置' (Kafka Source: Parameter Settings) panel with fields for name, topic source, topic library, topic name, intermediate table name, operator parallelism, and Kafka consumption position (Latest/Earliest).

添加数据源

最近更新时间: 2019-11-26 14:47:12

在IDE界面拖拽数据源插件“Kafka source”，双击拖拽的Kafka source弹出参数设置框（首次拖拽插件时，自动弹出），配置参数后，点击右上角的“X”进行保存。配置参数说明如下：

- 名称：即插件名称，可修改
- 数据源：可选择有项目下有限的数据源
- topic名称：可选择已选数据源下的topic
- 中间表名称：非必填，为将Kafka中的数据进行结构化后注册表的名称；
- Operator并发度：非必填，如不指定将按照作业默认并行数执行；
- Kafka消费位置：支持Latest（最近时间）、Earliest（最早时间）、Custom（自定义选择开始时间）三种方式，默认Latest time执行。

The screenshot shows the IDE interface with the 'Kafka Source: 参数设置' dialog box open. The dialog contains the following fields and options:

- 名称: KafkaSource_1567577261951
- *topic数据源: defaultKafka
- *topic名称: ss_dim1
- 中间表名称: 请输入中间表名称
- operator并行度: 请输入operator并行度
- *kafka消费位置: Latest Earliest Custom

Below these fields is a section for '自定义环境参数' (Custom Environment Parameters) with a table:

序号	参数名	参数值	操作
无数据			
+ 增加一行			

SQL Operator

最近更新时间: 2019-11-11 07:47:04

在IDE界面拖拽数据源插件“SQL Operator”，编写具体的流计算业务逻辑，SQL Operator不能独立存在，只能在source和sink之间。双击拖拽的“SQL Operator”，即可通过编写Flink SQL进行数据处理，SQL编译框默认提供编写规范，并支持添加预制的模板（具体详见模板管理）。

The screenshot shows the SQL Operator configuration interface. On the left, there is a code editor with the following content:

```
1  /*
2  3  书写规范：
4  1、FlinkSql中的关键字，要使用反引号`。例如：使用count作为字段，正确方式为`count`。
5  2、经过Sql查询出的Schema要和下游算子的Schema匹配，包括字段个数、类型、名称。
6  3、使用窗口函数时，时间类型支持ProcessingTime和EventTime。格式规范为：
7  ProcessingTime时，时间字段使用proctime: TUMBLE(proctime,interval '10' second),
8  EventTime时，时间字段使用rowtime: TUMBLE(rowtime,interval '10' second)。
9  4、当使用多级json时，获取数组内容时，下标从1开始。例如score[1]。
10 */
```

On the right, there is a configuration panel titled "Sql Operator:参数设置" with the following fields:

- 名称: Sql_1567520293420
- 中间表名称: 请输入中间表名称
- *是否使用模板: 是 否
- 模板级别: 请选择模板级别
- 模板名称: 请选择模板名称

A vertical "参数设置" (Parameter Settings) label is on the far right.

The screenshot shows the SQL Operator configuration interface with a SQL query in the code editor:

```
1 select
2 NAME,
3 age,
4 ctime
5 from
6 kafka_table
```

On the right, the configuration panel is the same as in the previous screenshot, but with the following changes:

- 模板级别: 公共模板
- 模板名称: kafkatomysql (highlighted with a red box)

A red arrow points from the "kafkatomysql" template name in the configuration panel to the "kafka_table" in the SQL query. A vertical "参数设置" (Parameter Settings) label is on the far right.

添加目标表

最近更新时间: 2021-09-15 15:59:08

拖拽目标表算子，即可进行添加目标表的操作，流计算服务支持将数据输出至：JDBC（含：MySQL、Oracle二种数据源）、HBase、Kafka、Redis、Elasticsearch七种数据源。

偷锁 保存 提交 单元测试 测试运行

源表

- Kafka Source

数据处理

- Sql Operator
- Etl Operator
- 自定义 Operator

结果表

- JDBC Sink
- Hbase Sink
- Kafka Sink
- Redis Sink
- ES Sink

KafkaSource_...

Sql_15675814...

JdbcSink_156...

HbaseSink_15...

KafkaSink_15...

RedisSink_156...

ES Sink_15675...

偷锁 保存 提交 单元测试 测试运行 下载自定义开发资源包 前往运维

源表

- Kafka Source

数据处理

- Sql Operator
- Etl Operator
- 自定义 Operator

结果表

- JDBC Sink
- Hbase Sink
- Kafka Sink

JdbcSink_...

JDBC Sink:参数设置

*数据源类型: 请选择数据源类型

*数据源: 请选择数据源

*目标数据库: 请选择目标数据库

*目标数据表: 请选择目标数据表

batchsize: 1 pieces

operator并行度: 请输入operator并行度

刷新时间(ms): 10000

JDBC sink、HBase sink、ES sink中都有数据写入批次概念：数据按照批次写入，有两个条件限制：批次大小batchsize和刷新时间flushIntervalms，只要满足其中一个条件，就会触发写数据操作。batchsize表示这一批次的大小，也就是有多少条数，刷新时间表示这一批次等待多长时间写入。

The screenshot shows a web-based configuration interface for a data pipeline. On the left, there is a sidebar with categories: '源表' (Source Tables) containing 'Kafka Source'; '数据处理' (Data Processing) containing 'Sql Operator', 'Etl Operator', and '自定义 Operator' (Custom Operator); and '结果表' (Result Tables) containing 'JDBC Sink', 'Hbase Sink', and 'Kafka Sink'. A 'JdbcSink_15' component is highlighted in the main workspace. A modal window titled 'Redis Sink:参数设置' (Redis Sink: Parameter Settings) is open, displaying the following configuration fields:

- operator并行度: 请输入operator并行度 (operator parallelism: please enter operator parallelism)
- redis操作: 操作类型: 请选择操作类型 (Redis operation: operation type: please select operation type)
- key所需字段: 请输入key所需字段 (key required field: please enter key required field)
- value所需字段: 请输入value所需字段 (value required field: please enter value required field)
- key连接符: 请输入key连接符 (key separator: please enter key separator)
- 过期时间: 请输入过期时间(ms) (expiration time: please enter expiration time in ms)

用户配置Redis sink时，支持多种操作，incr、incrByFloat、set、sadd、zadd、rpush、lpush、hset。key所需字段: key是从上游表中取的某个或多个字段对应的值，如果输入多个字段时用逗号分隔。value所需字段: value的值从上游表中取的某个字段对应的值。key连接符: 拼接key时候的连接符。

连接source、Operator、sink插件

最近更新时间: 2019-11-11 07:47:04

将鼠标悬停在算子上，可浮现连线点，可选择算子的连接顺序，生成作业流。在IDE界面通过拖拽方式将源表Kafka source，SQL Operator，Kafka sink 进行连接，多个作业的时候也可灵活连接。

The screenshot shows the IDE interface with a workflow diagram. On the left, there is a sidebar with categories: 源表 (Source Tables), 数据处理 (Data Processing), and 结果表 (Result Tables). Under 源表, there is a 'Kafka Source' component. Under 数据处理, there are 'Sql Operator', 'Etl Operator', and '自定义 Operator'. Under 结果表, there are 'JDBC Sink', 'Hbase Sink', 'Kafka Sink', and 'Es Sink'. In the workflow diagram, a 'Kafka Source' component is connected to a 'Sql_15647238...' operator, which is then connected to a 'JdbcSink_156...' sink. A 'JDBC Sink:参数设置' dialog box is open on the right, showing configuration options: 名称 (Name): JdbcSink_1564723866553, *数据源类型 (Data Source Type): 请选择数据源类型 (Please select data source type), *数据源 (Data Source): 请选择数据源 (Please select data source), *目标数据库 (Target Database): 请选择目标数据库 (Please select target database), *目标数据表 (Target Table): 请选择目标数据表 (Please select target table), batchsize: 1 pieces, and operator并行度 (Operator Parallelism): 请输入operator并行度 (Please enter operator parallelism).

流计算支持同时对2个Kafka source进行关联处理（双流joining）。

The screenshot shows the IDE interface with a workflow diagram. On the left, there is a sidebar with categories: 源表 (Source Tables), 数据处理 (Data Processing), and 结果表 (Result Tables). Under 源表, there is a 'Kafka Source' component. Under 数据处理, there are 'Sql Operator', 'Etl Operator', and '自定义 Operator'. Under 结果表, there are 'JDBC Sink', 'Hbase Sink', 'Kafka Sink', 'Redis Sink', and 'Es Sink'. In the workflow diagram, two 'KafkaSource_...' components are connected to a 'Sql_15675811...' operator, which is then connected to an 'HbaseSink_15...' sink. This illustrates a double-stream joining operation.

单元测试

最近更新时间: 2019-11-11 07:47:04

单元测试可以快速进行作业测试，支持乱序测试、延迟测试2种方式。读取线下上传的测试数据包中的数据，代替从source中读取数据，并且可以复用画板中source的Schema信息，将测试结果数据输出到开发界面，代替了数据输出到sink。1、单元测试用于测试数据处理算子的可用性；2、单元测试需选择指定数据包，数据包在[资源管理]中上传；3、单元测试使用的数据包将复用画布中拖拽的Kafka source的Schema信息；4、单元测试source算子并行度为1；5、单元测试结果获取会有一定延迟(每3s获取一次，每次上限100条)。

配置单元测试参数

备注:

- 1、单元测试用于测试数据处理算子的可用性;
- 2、单元测试需选择指定数据包。数据包在 [资源管理] 中上传;
- 3、单元测试使用的数据包将复用画布中拖拽的kafka source的schema信息;
- 4、单元测试Source算子并行度为1;
- 5、单元测试结果获取会有一定延迟(每3s获取一次，每次上限100条);

* 选择待测试作业版本: 编辑版本

* 测试数据包: MacFileStreamingTest1Res

* 对应kafka source: 交易明细

测试类型: 延迟数据测试 乱序数据测试

* 单条发送间隔(ms): 1

* 最大延时(ms): 1

* 延时触发比例(%): 1

* 选择资源组: 默认资源分组(CU)

可用/全部资源: 4 / 6

选择资源: 2

进行单元测试

单元测试结果:

版本管理

最近更新时间: 2019-11-11 07:47:04

支持生成多个版本的流作业（点击“提交”按钮生成新版本），点击“版本管理”，可以查看各个版本的信息。

The screenshot displays a workflow management interface. On the left, there is a sidebar with categories: 源表 (Source Tables) containing Kafka Source; 数据处理 (Data Processing) containing Sql Operator, Etl Operator, and 自定义 Operator (Custom Operator); and 结果表 (Result Tables) containing JDBC Sink, Hbase Sink, Kafka Sink, and Redis Sink. The main area shows a workflow diagram with nodes: KafkaSource_..., Etl_15674233..., Sql_15674233..., JdbcSink_156..., and ES EsSink_15674... connected by arrows. On the right, a table lists versions:

序号	选择	版本	提交人	提交时间
1	<input type="checkbox"/>	1	uat_test	2019-09-02 19:54:18
2	<input checked="" type="checkbox"/>	2	uat_test	2019-09-02 20:03:02
3	<input checked="" type="checkbox"/>	3	uat_test	2019-09-02 20:12:29
4	<input type="checkbox"/>	4	uat_test	2019-09-02 20:13:50
5	<input type="checkbox"/>	5	uat_test	2019-09-02 20:19:30

Below the table are buttons for 版本对比 (Version Comparison) and 流程图对比 (Flowchart Comparison). The right sidebar contains options for 作业配置 (Job Configuration), 版本 (Version), 选择函数 (Select Function), and 维表配置 (Dimension Table Configuration).

当多版本时，通过点击“流程图对比”或“版本对比”，可对任意2个版本的流程图和参数进行对比，内容不一致的部分会被标注出来：

作业流版本对比 版本2/版本3



作业版本对比--- 版本2 / 版本3

流计算内容对比

```
2
27 },
28 {
29   "jobName": "HbaseSink_1567423474810",
30   "layout": {"x": 223, "y": 301},
31   "jobId": "1909021924825",
32   "jobType": "HbaseSink",
33   "params": {
34     "commitType": "put",
35     "columnFamily": "cf",
36     "commitColumns": "",
37     "dsId": "239",
38     "parallelism": "",
39     "rowkey": "name",
40     "batchSize": "1",
41     "db": "ss_hbase_tx",
42     "table": "ss_test_create_hbase",
43     "flushIntervalMs": "10000"
44   }
45 },
46 {
47   "jobName": "JdbcSink_1567423473468",
48   "layout": {"x": 86, "y": 294},
49   "jobId": "1909021924532",
50   "jobType": "JdbcSink",
51   "params": {
52     "dsId": "215",
53   }
54 },
55 }
56 }
```

```
3
27 },
28 {
29   "jobName": "JdbcSink_1567423473468",
30   "layout": {"x": 86, "y": 294},
31   "jobId": "1909021924532",
32   "jobType": "JdbcSink",
33   "params": {
34     "dsId": "215",
35   }
36 },
37 {
38   "jobName": "HbaseSink_1567423474810",
39   "layout": {"x": 223, "y": 301},
40   "jobId": "1909021924825",
41   "jobType": "HbaseSink",
42   "params": {
43     "commitType": "put",
44     "columnFamily": "cf",
45     "commitColumns": "",
46     "dsId": "239",
47     "parallelism": "",
48     "rowkey": "name",
49     "batchSize": "1",
50     "db": "ss_hbase_tx",
51     "table": "ss_test_create_hbase",
52     "flushIntervalMs": "10000"
53   }
54 },
55 }
```



自定义开发说明

自定义开发说明

最近更新时间: 2019-11-11 07:46:40

自定义开发相关功能包括ETL、Custom Operator、UDF、维表等，提供JAR包上传和在线代码编写两种实现方式。下面将为每一种功能提供详细说明和并附注案例。



ETL Operator

ETL Operator

最近更新时间: 2019-11-11 07:45:53

ETL功能以算子形式，将Stream流暴露给用户，进行细力度的操作。系统提供JAR/在线开发两种实现方式。

JAR开发

最近更新时间: 2019-11-26 14:47:12

- 代码编写 ETL功能需要继承基类ETLFunction，并重写ETL和tableSchema两个方法。

```
trait ETLFunction extends Serializable {  
  /**  
  * 把row暴露给用户，进行自定义操作  
  *  
  * @param row 任何文本  
  */  
  def ETL(row: Row): Seq[Array[Any]]  
  
  /**  
  * 定义ETL输出的Schema  
  * 使用JSONRowSchemaConverter.convert(tableScheme())解析  
  */  
  def tableScheme(): String  
  
}
```

举例ETLTestSimple为一个ETL简单实现： 1. params为页面设置的自定义参数，以map形式输入，此处只简单输出。 2. with Logging，代码中可直接使用log实例，日志会打印在TaskManager的log信息中。 3. ETL函数，重载父类函数。每一条流消息调用一次本方法。本例中将第一列取出，并拼接一个固定列值。实际使用时以具体业务逻辑替换。 4. tableSchema函数，重载父类函数。因处理逻辑不确定，所以返回Schema需算子自己定义。以标准JSON Schema返回。注意：必须要和ETL函数的返回值对应上，不然后续转换过程和报错。

```
class ETLTestSimple (dac: AppConfig, params: Map[String, Any]) extends ETLFunction with Logging {  
  params.foreach(kv => {  
    log.info(kv._1 + " =>" + kv._2)  
  })  
  override def ETL(row: Row): Seq[Array[Any]] = {  
    log.info("in: " + row.toString)  
    val rESult: Array[Any] = Array(  
      row.getField(0), "fieldAddInETL"  
    )  
    log.info("out:" + rESult.mkString(", "))  
    Seq(rESult)  
  }  
  
  /**  
  * 定义ETL输出的Schema  
  */  
  override def tableScheme(): String = {  
    val Schema = "{ 'type': 'object', 'properties': { 'id': { 'type': 'integer' }, 'name': { 'type': 'string' } } }"  
    Schema  
  }  
}
```

- 代码打包 mvn clean package -DskipTests -Dcheckstyle.skip=true
- 页面配置

1. 在资源管理功能中将打包后的【streaming-flink-test-1.0-SANAPSHOT.jar】上传。

The screenshot shows the 'Resource Management' (资源管理) interface. A modal window titled '资源上传' (Resource Upload) is open, displaying the following fields:

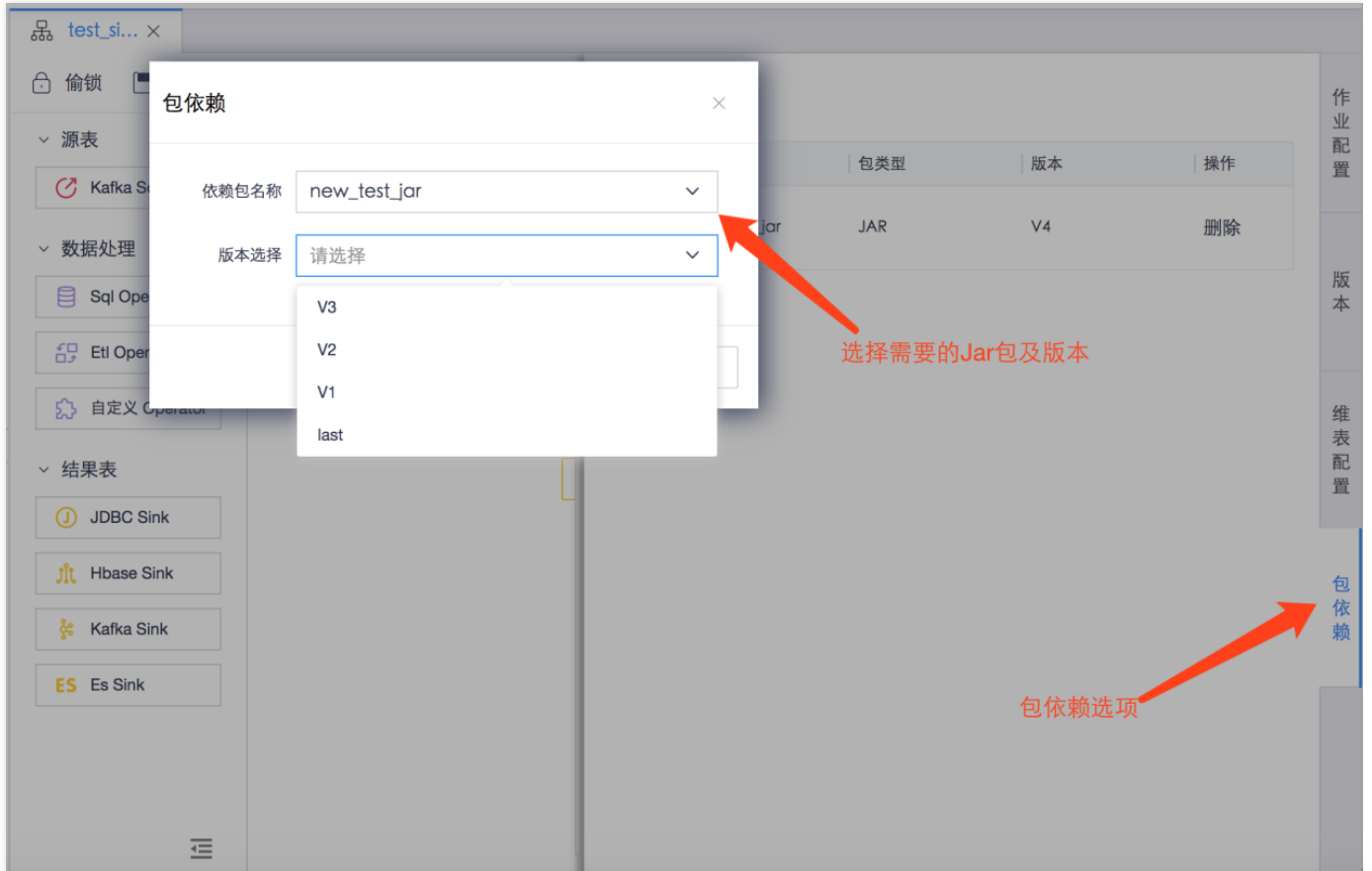
- 操作类型 (Operation Type): 更新 (Update)
- *资源名称 (Resource Name): new_test_jar
- *资源类型 (Resource Type): JAR
- *上传文件 (Upload File): streaming-flink-test-1.0-SAN [选择文件] (Select File)

Buttons for '确认' (Confirm) and '取消' (Cancel) are visible at the bottom of the modal. In the background, a table lists resources:

资源名称	资源类型	更新时间
new_test_jar	JAR	2019-12-2
0712_etl	JAR	2019-10-5
flink_test	JAR	2019-18-C
ddd	TEXT	2019-16:5

An orange arrow points to the '资源管理' (Resource Management) menu item in the left sidebar.

2. 在作业开发界面选择依赖的JAR。



3. 拖拽算子生成graph，并设置ETL Operator参数。其中类名称需为全路径，在执行时从依赖JAR包中反射注入。

源表

- Kafka Source

数据处理

- Sql Operator
- Etl Operator**
- 自定义 Operator

结果表

- JDBC Sink
- Hbase Sink
- Kafka Sink
- ES Es Sink

Etl Operator: 参数设置

名称: Etl_1562848510789

类名称: com.ksyun.dc.streaming.flink.test.ETLWithSide

中间表名称: etl

operator并行度: 1

+ 增加参数

序号	参数名	参数值	操作
1	k1	v1	删除

主类全路径

拖拽算子

自定义参数，可设置多个

在线开发

最近更新时间: 2019-11-11 08:00:19

- **代码编写** 在线代码编写功能在原ETL算子中，增加了“在线开发”选项，并提供Java/Scala两种语法支持。默认会展示最基本的代码框架，如需业务相关模板可在“是否使用模板”功能中选择。

```
1  override def etl(row: Row): Seq[Array[Any]] = {
2      log.info("in: " + row.toString)
3      val result: Array[Any] = Array(
4          row.getField(0), "fieldAddInEtl"
5      )
6      log.info("out:" + result.mkString(", "))
7      Seq(result)
8  }
9
10 override def tableScheme(): String = {
11     val schema = "{ 'type': 'object', 'properties': { 'id'
12     schema
13 }
```

Etl Operator:参数设置

名称: Etl_1565590948394

*开发方式: jar包上传 在线开发 1、选择模式

*中间表名称: tmp

operator并行度: 1

*开发语言: JAVA SCALA

*是否使用模板: 是 否 2、选择支持语法

增加参数 3、选择模板

序号	参数名	参数值	操作
无数据			

5、设置外部参数

双击ETL算子后进入编辑页面。按上图进行操作，其中需注意点如下：

1. 模板需预先在左侧“模板管理”功能中设置，且公共模板整个租户可见。
2. 模板具体参见《Flink自定义开发模板.doc》，只需编写“模板内容”部分，系统会自动填充成“补全后内容”运行。
3. “添加参数”输入的参数于ETL功能中使用，通过params: Java.util.Map[String, Any]直接获取。
4. 如需在ETL中调用维表，请参见《维表与ETL结合方式》。

Custom Operator

最近更新时间: 2019-11-26 14:47:12

Custom Operator功能提供一个自定义的全新算子，可以充当ETL，也可以充当source或者sink。需要自己维护TypeInfo并通过flatMap函数实现业务逻辑，Custom Operator只支持JAR开发模式。

- 代码编写 ETL功能需要继承基类tableOperateProcESSor，并重写innerBuild函数。举例CustomTestSimple为一个Custom Operator简单实现：
- 构造器中定义了四个参数：name：本算子的name，通常设置为在jobGraph中显示名称。chilids：子算子集合 configs: Map格式参数集合 sm: StreamingMate可获取环境变量、维表等全局参数
- with Logging，代码中可直接使用log实例，日志会打印在TaskManager的log信息中。
- innerBuild函数，重载父类函数。将dataStream直接暴露给用户，本例中不做处理直接转给下一个算子。实际使用时以具体业务逻辑替换。
- dataStream.dataType.asInstanceOf[RowTypeInfo]获取输入流的Schema 注意：必须要隐式声明implicit val tpe: TypeInfo[Row]，供下游算子使用

```
class CustomTestSimple(name: String, chilids: mutable.MutableList[ProcESSor[DataStream[Row]]], configs: Map[Any, Any], sm: StreamingMate)
extends tableOperateProcESSor(name, chilids, configs, sm) {

  override def innerBuild(dataStream1: DataStream[Row],dataStream2: DataStream[Row]): DataStream[Row] = {
    configs.foreach(kv => {
      println(kv._1 + " =>" + kv._2)
    })

    val rowType = dataStream1.dataType.asInstanceOf[RowTypeInfo]
    implicit val tpe: TypeInfo[Row] = rowType

    val finalStream = dataStream1.flatMap[Row](new MapMapper(rowType)).setParallelism(1)
    finalStream
  }
}

class MapMapper(rowTypeInfo: RowTypeInfo) extends RichFlatMapFunction[Row, Row] {
  override def flatMap(row: Row, collector: Collector[Row]): Unit = {
    println("CustomTestSimple:" + row)
    collector.collect(row)
  }
}
```

- 代码打包 同ETL
- 页面配置 同ETL

自定义UDF JAR开发

最近更新时间: 2019-11-11 08:00:19

- 代码编写 UDF功能需要继承基类ScalarFunction, 并重写eval方法。

```
import org.apache.flink.table.functions.ScalarFunction
```

```
class UDFTestScala extends ScalarFunction{  
  def eval(s: String): String = {s}  
}
```

- 代码打包

```
mvn clean package -DskipTests -Dcheckstyle.skip=true
```

- 页面配置 在资源管理功能中将打包后的streaming-flink-test-1.0-SANPSHOT.jar上传。

新建函数

开发模式: JAR包上传 在线开发

选择JAR包:

选择JAR包版本:

JAR包描述:

请录入函数信息:

序号	输入函数名	ClassName	函数类型	命令格式	用途	参数说明	返回
1	<input type="text" value="请输入函数名"/>	<input type="text" value="请输入ClassNar"/>	<input type="text" value="请选择函"/>	<input type="text" value="请输入命令格"/>	<input type="text" value="请输入用途"/>	<input type="text" value="请输入参数说"/>	<input type="text" value="请"/>

1、选择模式

2、选择jar包

3、选择jar包版本

4、输入函数信息, 支持批量添加

在线开发

最近更新时间: 2019-11-11 08:00:19

- 代码编写 在左侧菜单栏添加了“函数管理”功能，支持UDF/UDAF的添加和修改:

1、选择【函数管理】菜单

2、Tree中选项UDF/UDAF

3、新增按钮

4、修改功能在函数的版本信息中

双击新增按钮后进入编辑页面。

新建函数

开发模式: JAR包上传 在线开发

函数名:

函数类型:

命令格式:

用途:

参数说明:

返回值:

选择开发语言: JAVA Scala

在线编辑:

```
1 public long eval(String a) {
2     return a == null ? 0 : a.length();
3 }
4 public long eval(String b, String c) {
5     return eval(b) + eval(c);
6 }
```

1、选择在线开发

2、函数名，在SQL中使用

3、函数类型

4、展示用函数属性

5、模板代码

按上图进行操作，其中需要注意点如下：1. 函数名为最终register的函数名称，SQL中使用这个名称。2. 函数属性用以展示，在选择函数时使用。3. 模板具体参见《Flink自定义开发模板.doc》，只需编写“模板内容”部分，系统会自动填充成“补全后内容”运行。

自定义UDAF JAR开发

最近更新时间: 2019-11-11 08:00:19

UDAF功能采用Flink自身UDF语法, 详细说明参见<https://ci.apache.org/projects/Flink/Flink-docs-release-1.7/dev/table/UDFs.html>

- 代码编写 UDAF功能需要继承基类AggregateFunction, 并重写createAccumulator、getValue、accumulate、getResultType四个方法(有些可选方法见官网)。

```
import org.apache.flink.table.functions.AggregateFunction
import java.lang.Long
import java.util.concurrent.atomic.AtomicLong

import org.apache.flink.api.common.typeinfo.{TypeInfo, TypeES}
class UDAFTestScala/*占位符*/ extends AggregateFunction[Long, AtomicLong] {
  //初始化count UDAF的accumulator
  override def createAccumulator: AtomicLong = new AtomicLong(0L)

  //getValue提供了如何通过存放状态的accumulator计算count UDAF的结果的方法
  override def getValue(accumulator: AtomicLong): Long = accumulator.get()

  //accumulate提供了如何根据输入的数据更新count UDAF存放状态的accumulator
  def accumulate(accumulator: AtomicLong, iValue: Long): Unit = {
    accumulator.addAndGet(iValue)
  }

  override def getResultType: TypeInfo[Long] = TypeES.LONG
}
```

- 代码打包

```
mvn clean package -DskipTests -Dcheckstyle.skip=true
```

- 页面配置 在资源管理功能中将打包后的streaming-flink-test-1.0-SANPSHOT.jar上传。

新建函数 ×

开发模式: JAR包上传 在线开发

选择JAR包: 1、选择模式

选择JAR包版本: 2、选择jar包

JAR包描述:

3、选择jar包版本

请录入函数信息: 4、输入函数信息，支持批量添加

序号	输入函数名	ClassName	函数类型	命令格式	用途	参数说明	返回
1	<input type="text" value="请输入函数名"/>	<input type="text" value="请输入ClassNar"/>	<input type="text" value="请选择函"/>	<input type="text" value="请输入命令格"/>	<input type="text" value="请输入用途"/>	<input type="text" value="请输入参数说"/>	<input type="text" value="请"/>

在线开发

最近更新时间: 2019-11-11 08:04:05

在左侧菜单栏添加了“函数管理”功能，支持UDF/UDAF的添加和修改：

1、选择【函数管理】菜单

2、Tree中选项UDF/UDAF

3、新增按钮

4、修改功能在函数的版本信息中

双击新增按钮后进入编辑页面。按上图进行操作，其中需要注意点如下：1. 函数名为最终register的函数名称，SQL中使用这个名称。2. 函数属性用以展示，在选择



函数时使用。3. 模板具体参见《Flink自定义开发模板.doc》，只需编写“模板内容”部分，系统会自动填充成“补全后内容”运行。

新建函数



开发模式: JAR包上传 在线开发

函数名:

函数类型:

命令格式:

用途:

参数说明:

返回值:

选择开发语言: JAVA Scala

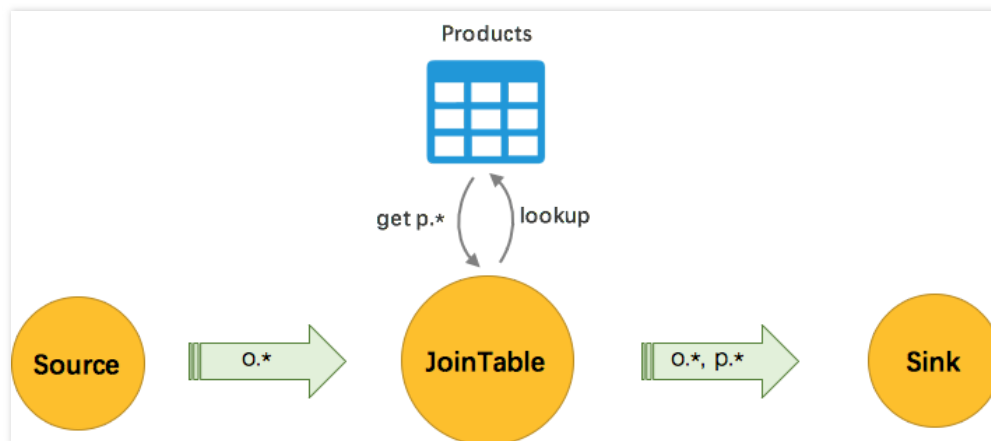
在线编辑:

```
1
```

维表配置 功能说明

最近更新时间: 2019-11-26 14:47:12

当前Flink版本只支持流式数据源，当流数据需要关联外部数据库（如MySQL、Oracle、Redis等），需要采用维表的形式支持。如下图所示，在ETL操作中，每一条



消息会lookup一下外部数据库。

数据库的维表查

询请求，有大量相同 key 的重复请求。如何减少重复请求？本地缓存是常用的方案。本方案目前提供两种缓存方案：LRU 和 ALL。cache = ALL（默认）：全量内存缓存 cacheTTLms:缓存的过期时间(ms) cache = LRU: LRU内存缓存 cachESize: 缓存的条目数量 cacheTTLms:缓存的过期时间(ms)

- All缓存 Async 和 LRU-Cache 能极大提高吞吐率并降低数据库的读压力，但是仍然会有大量的 IO 请求存在，尤其是当 miss key（维表中不存在的 key）很多的时候。如果维表数据不大（通常百万级以内），那么其实可以将整个维表缓存到本地。那么 miss key 永远不会去请求数据库。因为本地缓存就是维表的镜像，缓存中不存在那么远程数据库中也不存在。ALL cache 可以通过 cache='ALL'参数开启，通过cacheTTLms控制缓存的刷新间隔。会为 updater 节点起一个异步线程去同步缓存。在 Job 刚启动时，会先阻塞主数据流，直到缓存数据加载完毕，保证主数据流流过时缓存就已经 ready。在之后的更新缓存的过程中，不会阻塞主数据流。因为异步更新线程会将维表数据加载到临时缓存中，加载完毕后再与主缓存做原子替换。只有替换操作是加了锁的。因为几乎没有 IO 操作，所以使用 cache ALL 的维表性能可以非常高。但是由于内存需要能同时容纳两份维表拷贝，因此需要加大内存的配置。全量内存缓存模式，用户需要实现一个updater，用以定期全量刷新数据；getKey用以生成关联缓存的key。由框架来管理缓存cacheTTLms到期和updater调用。Datasource.select()为具体的查询操作。

```
class MySQLUpdater(config: Map[String, Any])
extends AppConfigUpdater(config) with Logging {
  val MySQLOptions = Map("url"-> config("url").toString,
    "drivername"-> config("drivername").toString,
    "username"-> config("username").toString,
    "password"-> config("password").toString)

  val dbKey = MySQLOptions.get("url").toString
  val SQL= config("SQL").toString
  val rowkeyFields = config("rowkeyFields").toString.split(",")

  Datasource.initIfNeeded(dbKey, MySQLOptions)

  override def update(): Map[String, Any] = {
    Datasource.select(dbKey, SQL, rowkeyFields)
  }

  override def getKey(params: Map[String, Any]): String = {
    Datasource.getKey(params)
  }
}

object MySQLUpdater {
  val key = "MySQLUpdater"
}
```



- LRU缓存 通过 `cache='LRU'`参数可以开启 LRU 缓存优化，会为每个slot节点创建一个 LRU 本地缓存Map。当每个数据进来的时候，先去缓存中查询，如果存在则直接关联输出，减少了一次 IO 请求。如果不存在，再发起数据库查询请求（异步或同步方式，本系统默认全部采用异步方式处理），请求返回的结果会先存入缓存中以备下次查询。为了防止缓存无限制增长，所以使用的是 LRU 缓存，并且可以通过 `cacheSize` 调整缓存的大小。为了定期更新维表数据，可以通过 `cacheTTLms` 调整缓存的失效时间。`cacheTTLms` 是作用于每条缓存数据上的，也就是某条缓存数据在指定 `timeout` 时间内没有被访问，则会从缓存中移除。在使用 LRU 缓存时，如果存在大量的 `invalid key`，或者数据库中不存在的 `key`。由于命中不了缓存，导致缓存的收益较低，仍然会有大量请求打到数据库。因此我们将未命中的 `key` 也加进了缓存，提高了未命中 `key` 和 `invalid key` 情况下的缓存命中率。LRU内存缓存，用户需要实现一个`updater`，用以单条刷新数据，每一次查询都会尝试调用`updater`方法。由框架来管理缓存`cacheTTLms`到期，查询是否命中及`updater`调用。`Datasource.select()`为具体的查询操作。

```
class MySQLQuerier(config: Map[String, Any])
  extends AppConfigQuerier(config) with Logging {
  val MySQLOptions = Map("url" -> config("url").toString,
    "drivename" -> config("drivename").toString,
    "username" -> config("username").toString,
    "password" -> config("password").toString)

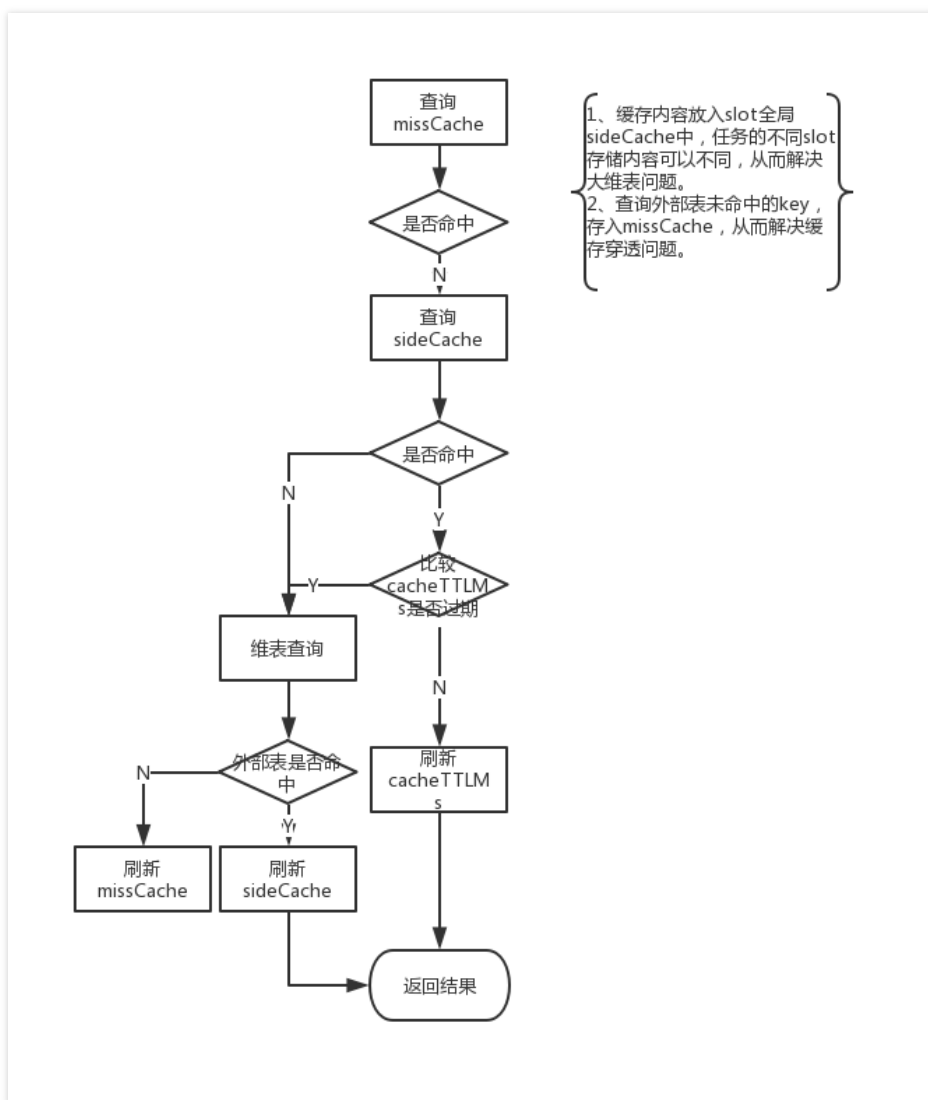
  val dbKey = MySQLOptions.get("url").toString
  val SQL = config("SQL").toString
  val rowkeyFields = config.get("rowkeyFields").getOrElse("").toString.split(",")

  Datasource.initIfNeeded(dbKey, MySQLOptions)

  override def query(params: collection.Map[String, Any]): Any = {
    Datasource.select(dbKey, SQL, params)
  }

  override def getKey(params: collection.Map[String, Any]): String = {
    Datasource.getKey(params)
  }
}

object MySQLQuerier {
  val key = "MySQLQuerier"
}
```



• 执行流程图如下：



与ETL结合方式

最近更新时间: 2019-11-11 08:08:56

本版本维表功能只支持在ETL算子中使用。本文提供MySQL版的维表实现案例，如需其他版本请用户自定义实现。

- 代码编写 `dac.getConfigByKey(MySQLUpdater.key)`为从运行时环境获取缓存对象，其中`MySQLUpdater.key`为对应实现的key。

```
class ETLWithSide (dac: AppConfig, params: Map[String, Any]) extends ETLFunction with Logging {  
  
  /**  
  * 把一个字符串解析成表(多行 多列)  
  *  
  * @param row 任何文本  
  */  
  override def ETL(row: Row): Seq[Array[Any]] = {  
    //ALL 类型返回ALLCACHE  
    val cache1: Cache = dac.getConfigByKey(MySQLUpdater.key).get.asInstanceOf[AllCache]  
    //LRU和NONE 类型返回LRUCache  
    val cache2: Cache = dac.getConfigByKey(MySQLQueryer.key).get.asInstanceOf[LRUCache]  
  
    val result: Array[Any] = Array(  
      row.getField(0),  
      cache1.get(Map("id" -> row.getField(0))).getOrElse(Map()).asInstanceOf[Map[String, Any]].getOrElse("name", ""),  
      cache2.get(Map("id" -> row.getField(0))).getOrElse(Map()).asInstanceOf[Map[String, Any]].getOrElse("name", "")  
    )  
    println("ETLWithSide:" + result.mkString(", "))  
    Seq(result)  
  }  
  
  /**  
  * 定义ETL输出的Schema  
  */  
  override def tableScheme(): String = {  
    val Schema = "{ 'type': 'object', 'properties': { 'id': { 'type': 'integer' }, 'nameFromAllCache': { 'type': 'string' }, 'nameFromLRUCache': { 'type': 'string' } } }"  
    Schema  
  }  
}
```

- 代码打包 同ETL

- 页面配置 其中数据源非必选参数，如使用外部数据源，可完全由自定义参数配置数据库连接。维表可设置多个。

维表数据源配置

维表主类全路径

缓存参数

自定义参数

序号	参数名	参数值	操作
1	rowkeyFields	id	删除

如选择内部数据源，对应的参数对照表见【维表参数说明】



维表参数说明

最近更新时间: 2019-10-28 01:42:05

类型	值	说明	demo
公用参数	type	数据库类型	MySQL
	cachetype	缓存类型	ALL/LRU/NONE
	cacheTimeout	缓存超时时间	1000S
	cacheSize	缓存大小 (LRU 时有效)	1000
	dsId	流实例 ID	
MySQL	drivername	数据库驱动类型	com.MySQL.JDBC.Driver
	url	数据库地址	JDBC:MySQL://127.0.0.1:3306/Flink
	db	库名	testdb
	tablename	表名	MySQL_side_table
	username	用户名	root
	password	密码	1234
Redis	masterNode	主节点	10.77.0.6:6379
	connectType	连接类型	masterSlave
	password	密码	123456
	operations	操作描述	[{"keyFields": "ren1,ren2", "opType": "set", "expire": "", "scoreField": "", "valueFields": "ren1", "keyDelimiter": " "}]

JAR包上传模式

最近更新时间: 2019-11-11 08:08:56

JAR包上传模式支持线下上传JAR包进行流作业处理, JAR包上传模式支持Flink和Spark Streaming两种引擎。

The screenshot shows the '新增作业' (New Job) dialog box. It contains the following fields and options:

- 名称: JAR_test
- 描述: 请输入描述
- 开发模式: 在线开发 JAR包上传
- 执行引擎: Flink Spark
- 执行引擎版本: 1.7.1
- 选择存放文件夹: /JAR上传作业
- Buttons: 取消, 确认

The screenshot shows the main configuration page for the job 'JAR_test'. It includes the following details:

- 作业名称: JAR_test
- 选择JAR包: 请选择JAR包
- JAR包版本: 请选择JAR包版本
- 引擎: flink
- 引擎版本: 1.7.1
- 作业描述: (empty text area)
- 执行命令: 请输入执行命令
- demo展示:

```
flink run -m yarn-cluster -c org.apache.flink.streaming.examples.wordcount.WordCount -yn 1 -ys 1 -p 1 -yjm 1024 -ytm 3072 -yqu default -ynm demo -n -d mainJar param1 param2 param3
```

On the right side, there is a table for '增加包' (Add Package) with columns: 序号, 包名称, 包类型, 版本, 操作. The table currently shows '无数据' (No data).

测试及运行

发布测试

最近更新时间: 2019-11-11 08:24:36

完成作业开发配置后, 点击“保存”按钮保存当前任务, 点击“提交”按钮生成新版本的作业, 点击“运行测试”将作业发布到“运维中心”进行测试运行。

The screenshot shows a job configuration interface. On the left, there is a sidebar with categories: 源表 (Source Tables), 数据处理 (Data Processing), and 结果表 (Result Tables). Under 源表, there is a 'Kafka Source' button. Under 数据处理, there are 'Sql Operator', 'Etl Operator', and '自定义 Operator' (Custom Operator). Under 结果表, there are 'JDBC Sink', 'Hbase Sink', 'Kafka Sink', and 'ES Sink'. The main area displays a workflow diagram with two 'KafkaSource_...' nodes at the top, both pointing to a central 'Sql_15646292...' node, which then points to a 'JdbcSink_156...' node at the bottom. On the right, a 'Kafka Source:参数设置' (Kafka Source: Parameter Settings) panel is open, showing fields for '名称' (Name: KafkaSource_1564629217035), '*topic数据源' (Topic Data Source), '*topic库' (Topic Database), '*topic名称' (Topic Name), '中间表名称' (Intermediate Table Name), 'operator并行度' (Operator Parallelism), and '*kafka消费位置' (Kafka Consumption Position) with radio buttons for 'Latest', 'Earliest', and 'Custom'. Below these is a table for '自定义环境参数' (Custom Environment Parameters) with columns for '序号' (Serial Number), '参数名' (Parameter Name), '参数值' (Parameter Value), and '操作' (Action). The table is currently empty with a '+ 增加一行' (Add a row) button.

This screenshot is similar to the one above, showing the same workflow diagram and configuration panel. However, the '测试运行' (Test Run) button in the top toolbar is highlighted with a red box, indicating the next step in the process. The rest of the interface, including the sidebar and configuration panel, remains the same.

测试运行

最近更新时间: 2019-11-11 08:24:36

提交作业后, 会在“运维中心-测试实例-流计算作业”生成一条流作业数据, 可进行: 【启动】 【终止】 【创建SavePoint】 【查看SavePoint】 【查看Flink UI】 等操作。

- 【启动】: 点击启动, 并选择启动方式(可选择立即启动或SavePoint启动)、配置作业参数后, 可启动作业。
- 【终止】: 点击终止, 可停止作业;

运维中心 > 测试实例

流计算作业

序号	执行ID	作业流名称	项目	作业版本	作业类型	责任人	业务时间	启动时间	停止时间	运行时长	状态	操作
1	8677	kafka_mysql	stream0727	6	在线开发	ccb_test	2019-08-16	2019-08-16 17:52:53	--	0s	运行中	终止 发布生产 创建savepoint 查看savepoint Flink UI 查看日志
2	8614	kafka2mysql805	zb0727	1	在线开发	ccb_test	2019-08-16	2019-08-16 16:37:50	2019-08-16 16:53:05	15m 15s	失败	启动 发布生产 查看savepoint 查看日志

启动设置

启动设置: 立即启动

作业设置:

*checkpoint触发间隔-秒: - 10 +

默认并行数: - 1 +

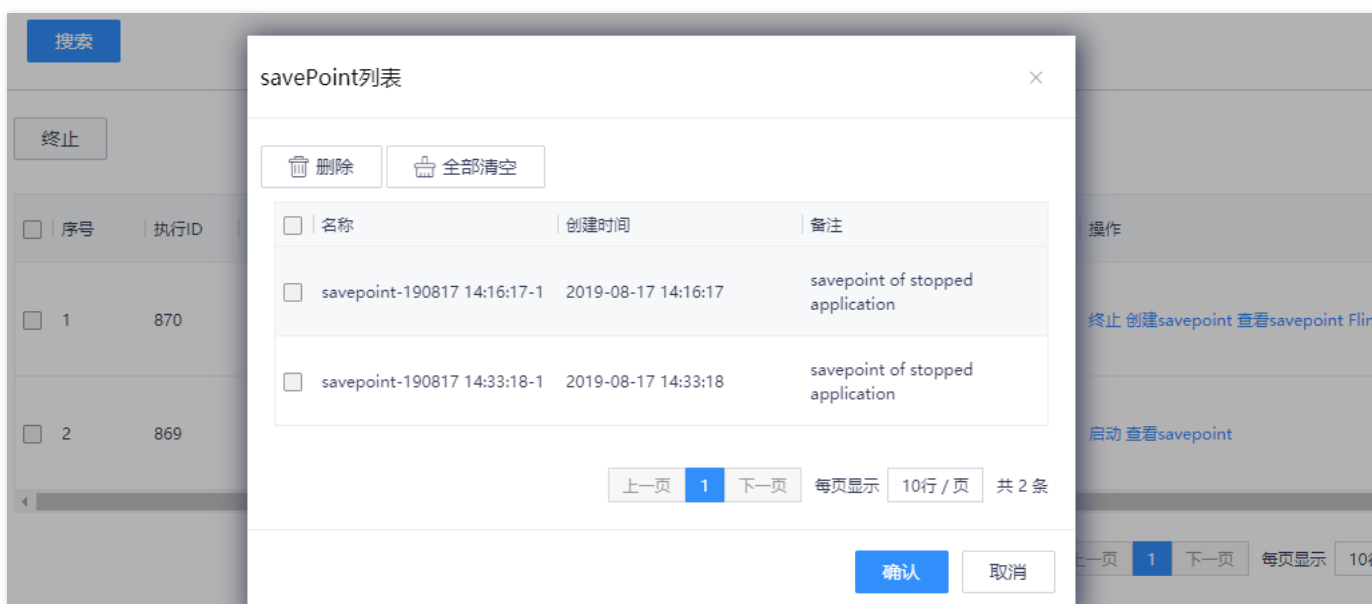
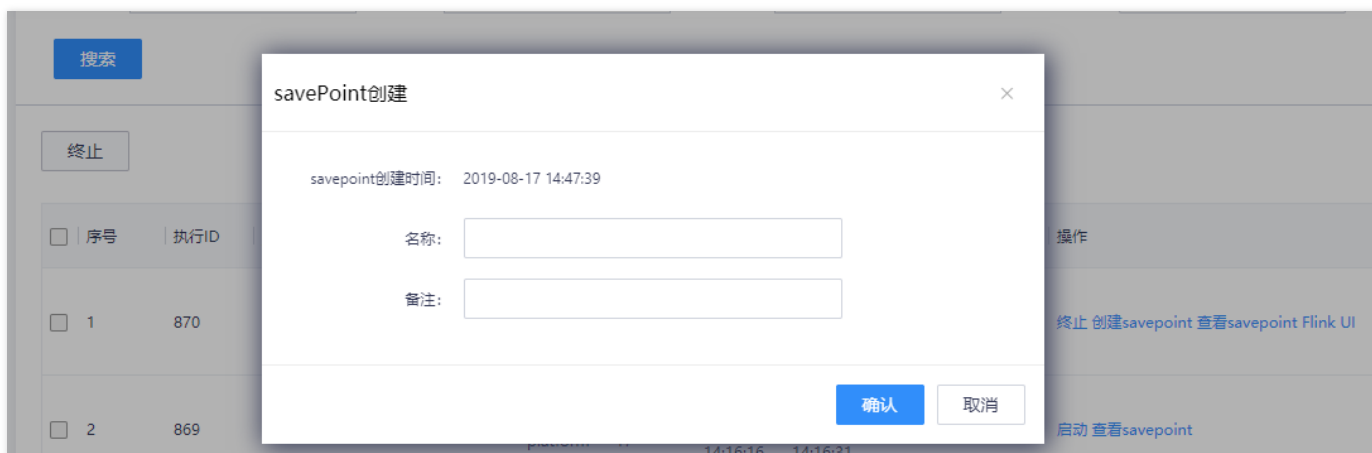
可选参数: 添加参数

确认 取消

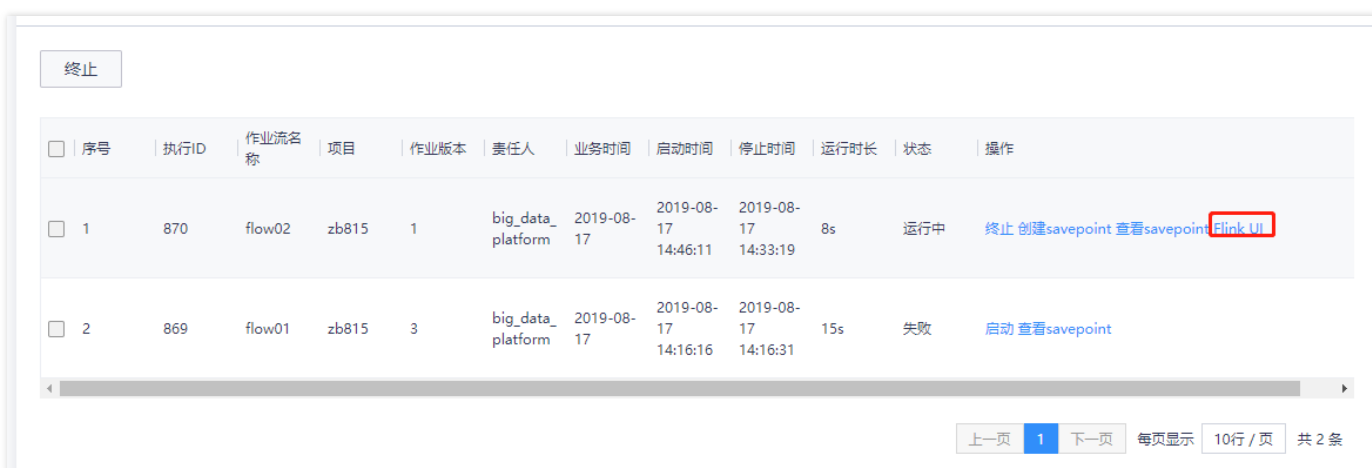
- 【创建SavePoint】: 创建后, 再次启动作业时, 可选择从SavePoint标识的点启动作业;
- 【查看SavePoint】: 可以查看所有的SavePoint列表

终止

序号	执行ID	作业流名称	项目	作业版本	责任人	业务时间	启动时间	停止时间	运行时长	状态	操作
1	870	flow02	zb815	1	big_data_platform	2019-08-17	2019-08-17 14:46:11	2019-08-17 14:33:19	8s	运行中	终止 创建savepoint 查看savepoint Flink UI
2	869	flow01	zb815	3	big_data_platform	2019-08-17	2019-08-17 14:16:16	2019-08-17 14:16:31	15s	失败	启动 查看savepoint



- 【查看Flink UI】备注：查看Flink UI前，需要提前配置部署集群的hosts，否则页面无法正常跳转。





发布生产

最近更新时间: 2019-11-11 08:24:36

提交作业后, 会在【运维中心-测试实例-流计算作业】生成一条流作业数据, 可进行:【启动】【终止】【创建SavePoint】【查看SavePoint】【查看Flink UI】等操作, 对于执行正常的作业, 点击【发布生产】, 可进入发布生产的审批流程。

运维中心 > 测试实例

流计算作业

项目空间: 全部 状态: 全部 责任人: 全部 任务名称: 输入任务名称查询 搜索

终止

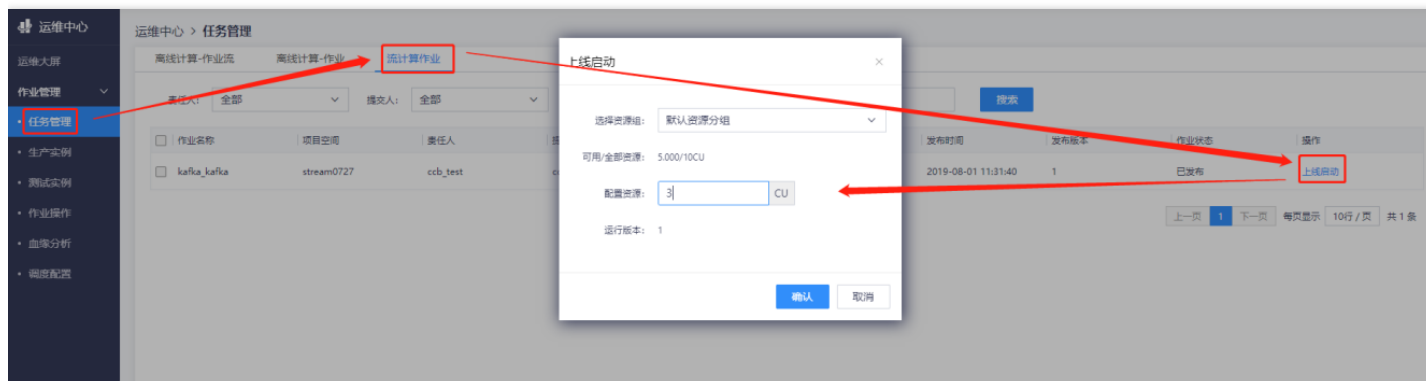
序号	执行ID	作业流名称	项目	作业版本	责任人	业务时间	启动时间	停止时间	运行时长	状态	操作
1	2065	test_001	stream0727	编辑版本	ccb_test	2019-08-01	2019-08-01 11:21:52	2019-08-01 11:21:54	2s	失败	启动 发布生产 查看savepoint

上一页 1 下一页 每页显示 10行/页 共1条

上线启动

最近更新: 2019-11-11 08:23:52

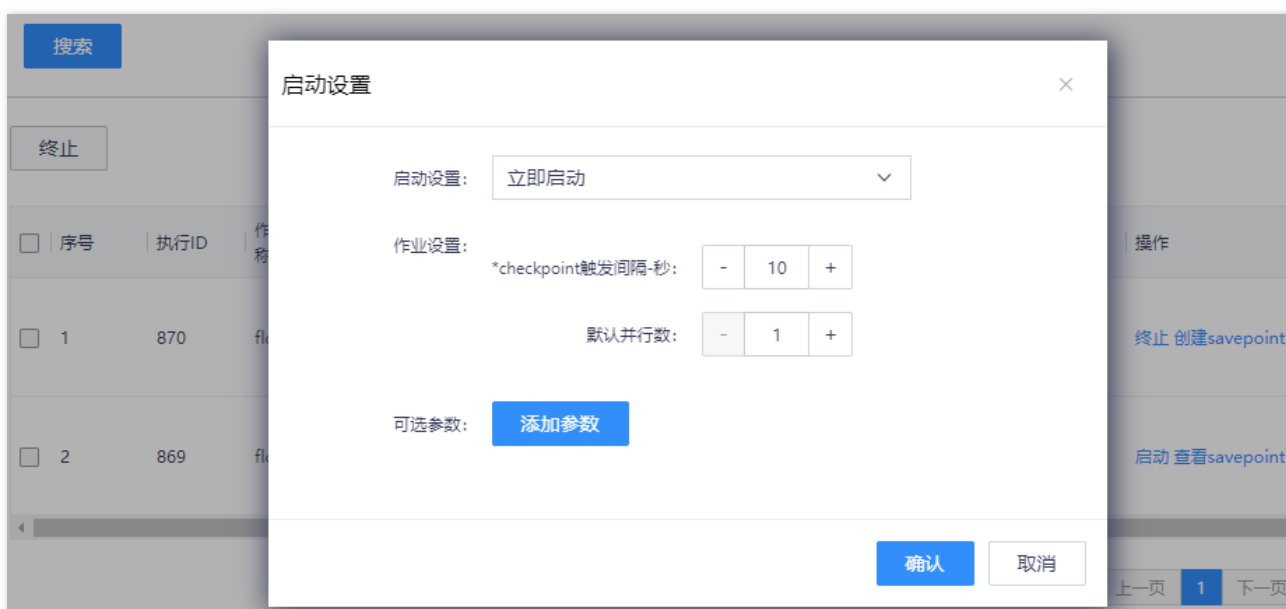
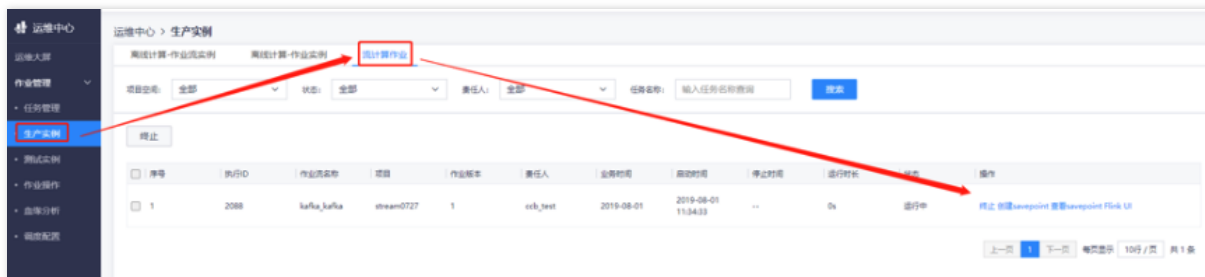
已通过审批的作业流可以在【运维中心-任务管理】页面中看到，并进行相应操作，已拒绝的作业流不可以在【运维中心-任务管理】中看到。点击【上线启动】按钮，需要对待上线的作业分配运行资源。完成资源分配的作业，将提交至【运维中心-生产实例-流计算作业】中，进行调度运行：启动、停止、创建SavePoint、查看SavePoint、查看Flink UI等。



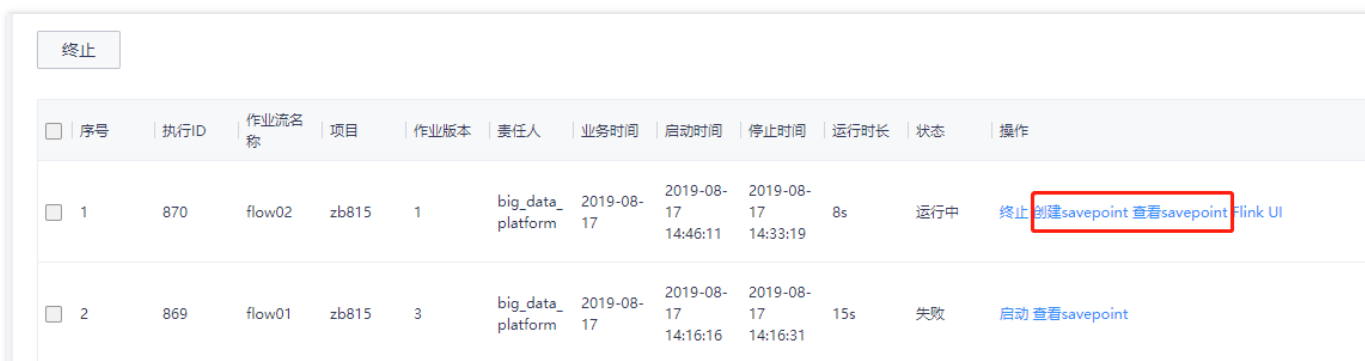
生产运行

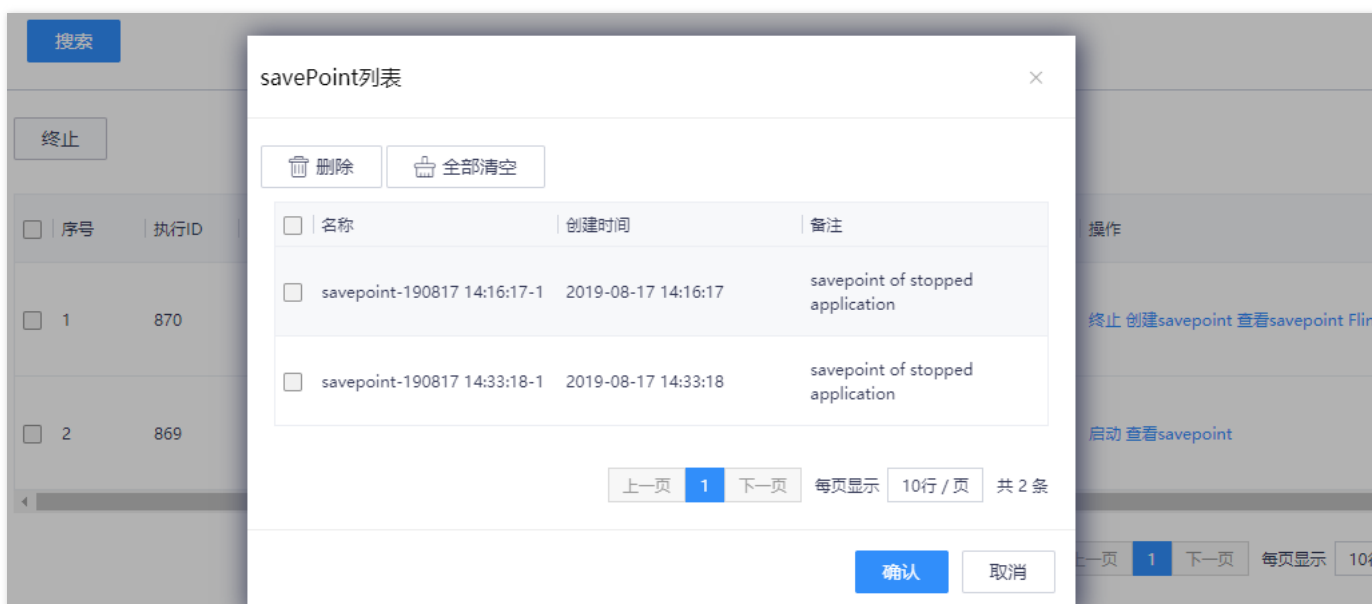
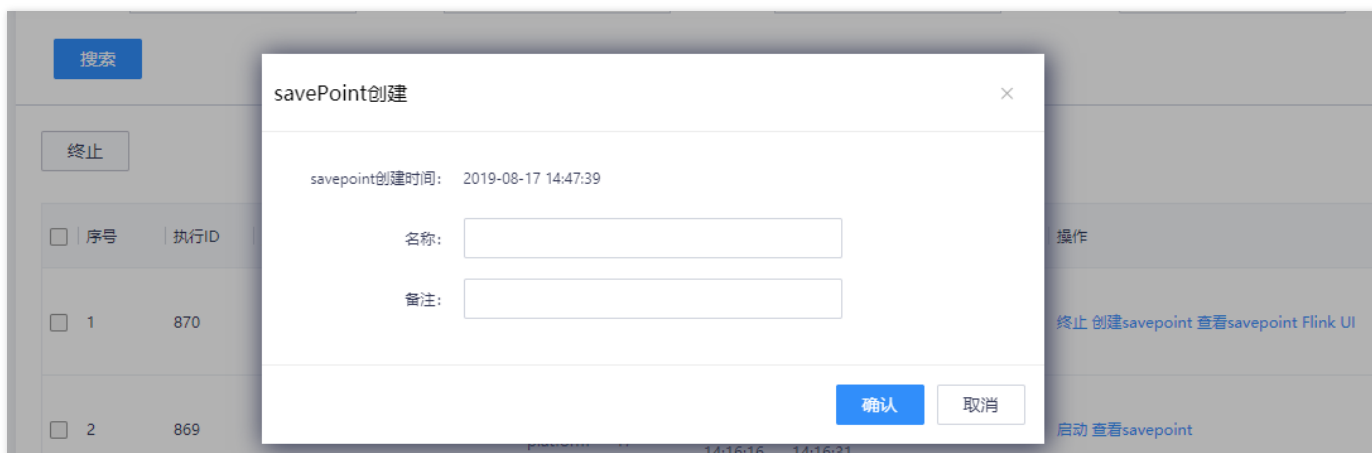
最近更新时间: 2019-11-11 08:23:52

- **【启动】**：选择启动位点，点击按以上配置启动，作业即提交至计算集群运行。流计算作业启动时候您可以指定启动时间。表示从源头数据存储的指定时间点开始读取数据。例如设置启动时间为当前时间或者一个小时之前，也可以立即启动
- **【终止】**：点击终止，可停止作业；

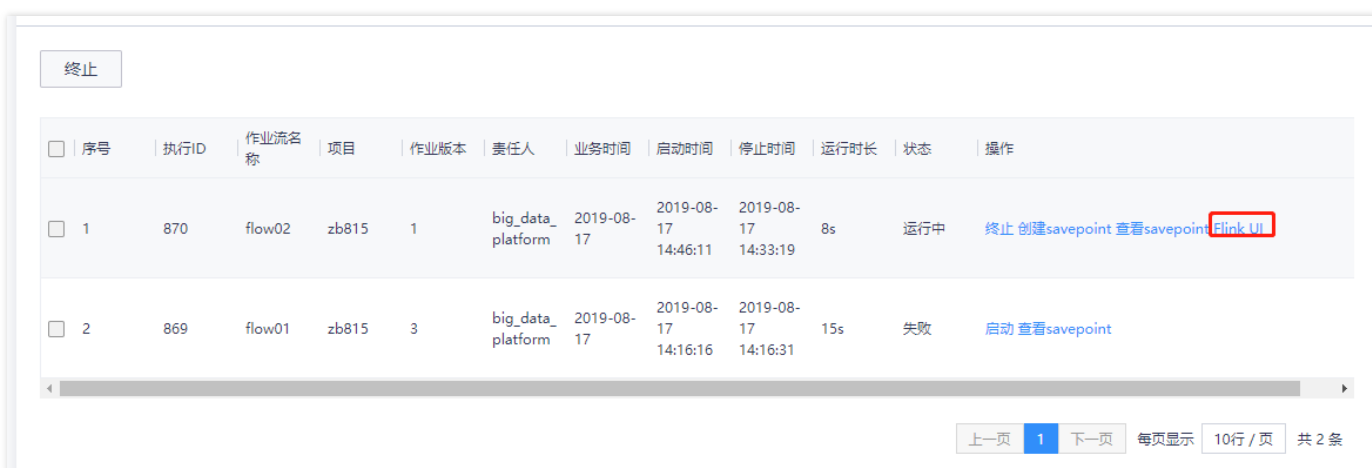


- **【创建SavePoint】**：创建后，再次启动作业时，可选择从SavePoint标识的点启动作业；
- **【查看SavePoint】**：可以查看所有的SavePoint列表。





- 查看Flink UI 备注: 查看Flink UI前, 需要提前配置部署集群的hosts, 否则页面无法正常跳转。





Apache Flink Dashboard

- Overview
- Running Jobs
- Completed Jobs
- Task Managers
- Job Manager
- Submit new Job

Overview

Version: 1.7.1 Commit: 6e9a97b

	1
Task Managers	
	1
Task Slots	
	0
Available Task Slots	

Total Jobs

Running	1
Finished	0
Canceled	0
Failed	0

Running Jobs

Start Time	End Time	Duration	Job Name	Job ID	Tasks	Status
2019-08-17, 14:46:26	2019-08-17, 14:54:58	8m 31s	zb815 flow02	07c8551845f2689d1e9729ea7e51e112	1 0 0 0 0 1 0 0 0 0 0 0	RUNNING



最佳实践

场景实践

背景

最近更新时间: 2019-11-26 14:47:12

随着互联网络技术的发展，直播的概念有了新的拓展和发展，现在更多的人关注网络直播，特别是视频直播生态链更受关注。通过网络信号，在线收看球赛、体育赛事、重大活动和新闻等，让大众有了广阔且自由的选择空间，我们能够真正的随时随地的体验直播的快乐和便捷。

在当前体验为王的时代，任何体验不佳均会导致客户群体的大规模流失。对于网站直播平台而言，需要重点关注用户(主播和粉丝)的使用体验，重点关注的系统指标包括音视频卡顿率、延迟率、丢包情况，同时由于直播平台时效性，平台方更需要实时了解到系统周边问题，并先于用户发现并定位问题和故障。另外，作为平台运营方还应该对于整体网站客户运营情况、视屏爆款情况有及时的跟踪和了解。

下面我们以某个视屏直播平台方为例，讲解下如何使用实时计算进行系统稳定性以及平台运营情况进行实时监控和展现。



业务 业务

最近更新时间: 2019-11-26 14:47:12

为了最大程度活跃用户社群，覆盖更多场景直播频道，最终实现平台方盈利，通常同一个视频直播网站有多个主播，每个主播向一个频道内的用户进行广播，用户可以看到当前频道内的主播视频，并听到其声音，主播可以与频道内的多个用户进行私聊。

主播和粉丝各自持有直播APP软件，该APP软件将每10s向服务器打点，服务器接收到打点信息会输出到本地磁盘，并通过阿里云日志服务采集端将打点数据到日志服务，实时计算通过订阅该日志，实时计算客户端视频播放情况。



业务目标

最近更新时间: 2019-11-26 14:47:12

针对客户端APP的监控, 获取以下指标:

- 房间故障, 故障包括卡顿、丢帧、音视频不同步等。
- 分地域统计数据端到端延迟平均情况。
- 统计实时整体卡顿率 (出现卡顿的在线用户数/在线总用户数*100%, 通过此指标可以 衡量当前卡顿影响的人群范围)。
- 统计人均卡顿次数 (在线卡顿总次数/在线用户数, 通过此指标可以从卡顿频次上衡量 整体的卡顿严重程度)。

数据格式

最近更新时间: 2019-11-11 08:28:13

客户端APP向服务器打点日志格式如下:

字段	含义
ip	用户端ip
agent	端类型
roomid	房间号
userid	用户id
abytES	音频码率
afcnc	音频帧数
adrop	音频丢帧数量
afts	音频时间戳
alat	音频帧端到端延迟
vbytES	视频码率
vcnc	视频帧数
vdrop	视频丢帧数量
vfts	视频时间戳
vlat	视频帧端到端延迟
ublock	上行卡顿次数
dblock	下行卡顿次数
timESamp	打点时间戳
region	地域

日志服务内部是半结构化存储, 上述数据将展现出来格式如下:

```
{
  "ip": "ip",
  "agent": "agent",
  "roomid": "123456789",
  "userid": "123456789",
  "abytES": "123456",
  "afcnc": "34",
  "adrop": "3",
  "afts": "1515922566",
  "alat": "123",
  "vbytES": "123456",
  "vcnc": "34",
  "vdrop": "4",
  "vfts": "1515922566",
  "vlat": "123",
  "ublock": "1",
  "dblock": "2",
  "timESamp": "15151922566",
  "region": "beijing"
}
```



场景

最近更新时间: 2019-10-28 01:42:05

- 房间故障统计 统计房间故障，故障包括卡顿、丢帧、音视频不同步信息，使用10分钟一个窗口进行统计。

```
SELECT
CAST(TUMBLE_START(proctime, INTERVAL '10' MINUTE) as VARCHAR) as app_ts,
roomid,
SUM(ublock) as ublock, --统计10分钟内上行卡顿次数
SUM(dblock) as dblock, --统计10分钟内下行卡顿次数
SUM(adrop) as adrop, --统计10分钟内音频丢包次数
SUM(vdrop) as vdrop, --统计10分钟内视频丢包次数
SUM(alat) as alat, --统计10分钟内音频延迟
SUM(vlat) as vlat, --统计10分钟内视频延迟
FROM
view_app_heartbeat_Stream
GROUPBY
TUMBLE(proctime, INTERVAL '10' MINUTE), roomid
```

- 分地域统计延迟情况 分地域统计数据端到端延迟平均情况，每10分钟统计音频、视频平均延迟情况。

```
SELECT
CAST(TUMBLE_START(proctime, INTERVAL '10' MINUTE) as VARCHAR) as app_ts,
region,
SUM(alat)/COUNT(alat) as alat,
SUM(vlat)/COUNT(vlat) as vlat,
FROM
view_app_heartbeat_Stream
GROUPBY
TUMBLE(proctime, INTERVAL '10' MINUTE), region;
```

- 实时整体卡顿率 统计实时整体卡顿率，即出现卡顿的在线用户数/在线总用户数*100%，通过此指标可以衡量当前卡顿影响的人群范围。

```
SELECT
CAST(TUMBLE_START(proctime, INTERVAL '10' MINUTE) as VARCHAR) as app_ts,
SUM(IF(ublock <> 0 OR dblock <> 0, 1, 0)) / CAST(COUNT(DISTINCT userid) ASDOUBLE) as block_rate,
FROM
view_app_heartbeat_Stream
GROUPBY
TUMBLE(proctime, INTERVAL '10' MINUTE);
```



技术实践

最近更新时间: 2019-10-27 08:04:46

1. Kafka source并发度设置不要超过topic Partitions数量；Kafka Consumer数量不能超过topic Partitions数量，超过的部分会占用slot但不会消费数据。
2. State量大的算子，尽量在前一个算子进行keyBy，将相同key的数据归集到同一个slot。使用Flink Tuples。
3. 当使用类似于groupBy、join或keyBy这些操作时，Flink提供了多种方式以使用户在数据集中选择主键。如果现在使用的是Flink Tuple类型，那么只要简单地指定字段元组的位置，就可以被用作主键。
4. 复用Flink对象，当从用户定义的函数返回数据时，最好使用可变对象。
5. 使用注解功能。由于Flink无法解析和理解代码，所以可以提供一些有利于构建更有效执行计划的重要信息。可以使用以下三个注解：
@ForwardedFields：指定输入值中哪些字段保持不变，哪些字段是用于输出的。
@NotForwardedFields：指定在输出中未保留相同位置的字段。
@ReadFields：指定用来计算结果值的字段。指定的字段应该只在计算中使用，而不仅仅是复制到输出参数中。
6. Select Join Type。
BROADCAST_HASH_SECOND：第二个数据集要小得多
REPARTITION_HASH_FIRST：第一个数据集稍微小一些
REPARTITION_HASH_SECOND：第二个数据集要小一点
REPARTITION_SORT_MERGE：使用排序和合并策略对数据集进行重新分配
OPTIMIZER_CHOOSES：Flink优化器将决定如何join数据集



故障指南

最近更新时间: 2019-10-28 01:42:05

- SQL不规范-带关键字的SQL FlinkSQL的关键字, 要使用反引号`引起来。例如: `select name, sum(1) as count from tmp group by id`, 其中count为关键字, 要改为 `select name, sum(1) as count from tmp group by id`。
- 上下游Schema不匹配 经过SQL查询出的Schema要和下游算子的Schema匹配, 包括字段个数、类型、名称。例如: SQL中查询出name, age, score三个字段: `select name, age, score from tmp` 下游算子MySQL sink表中只有两个字段name、age。这种不匹配的情况将会报错。
- SQL不规范-窗口函数 时间类型支持Processing Time和Event Time。格式规范为: Processing Time时, 时间字段使用proctime: `TUMBLE(proctime, interval '10' second)`, Event Time时, 时间字段使用rowtime: `TUMBLE(rowtime, interval '10' second)`。例如: ProcessingTime时, `select sum(score) as cnt from tmp group by name, TUMBLE(ctime, interval '10' second)` 应该改为 `select sum(score) as cnt from tmp group by name, TUMBLE(proctime, interval '10' second)`
- 使用Kafka中多级JSON报错 获取数组内容时, 下标从1开始。例如 `select score[0] from tmp` 应该改为 `select score[1] from tmp`
- 元数据信息无法获取 1、刷新页面, 重新加载; 2、退出当前用户登录, 重新登录; 3、如果前两步操作仍然无法正常进入数据开发页面, 则可能是服务问题, 可以通过以下步骤检查: (1) 浏览器进入开发者模式; (2) 打开network选项框, 同时点击元信息下拉菜单, 查看 <http://console.bigdata.yun.ccb.com/i/ide/dataManager/ideGateWay> 调用返回值; (3) 如果相应state不是200, 反馈给后台运维人员处理。
- 资源列表无法获取 1、刷新页面, 重新加载; 2、退出当前用户登录, 重新登录; 3、如果前两步操作仍然无法正常进入数据开发页面, 则可能是服务问题, 可以通过以下步骤检查: (1) 浏览器进入开发者模式; (2) 打开network选项框, 同时点击元信息下拉菜单, 查看 <http://console.bigdata.yun.ccb.com/i/ide/dataManager/getResourceGroup> 调用返回值; (3) 如果相应state不是200, 反馈给后台运维人员处理。
- UDF中函数参数与SQL中传入的格式不匹配 在“运维中心-运行日志-YARN日志”功能中直接查看日志, 查找“SQL validation failed. Given parameters of function do not match any signature. Actual: () Expected: ()”异常, 如果有则修改UDF或者业务处理逻辑, 重新调度运行。
- Job无报错长时间无数据输出 流计算开发界面, 检测batchsize、flushInterval参数设置, 可能因为设置的值过大, 还没触发数据写入, 调整成合适值后, 重新调度运行。



常见问题

产品介绍常见问题

最近更新时间: 2021-09-15 15:59:08

Q: 什么是实时流计算?

A: “实时”指实时处理, 计算框架支持按消息时间逐条处理; “流”指数据如水流, 一个接着一个; “计算”指数学运算、数据分析、算法模型执行等。“实时流计算”指实时处理当下正在发生的数据流, 逐条大数据分析或算法运算。

Q: 目前流计算引擎支持哪几种? 在使用上有什么区别?

A: 目前支持Flink、Spark Streaming两种引擎。其中Flink支持上传JAR包、通过IDE拖拽插件的方式进行流计算作业的开发, Spark Streaming只支持上传JAR包的方式。

Q: 目前流计算支持哪几种数据源?

A: source支持: Kafka; sink支持: Kafka、MySQL、Oracle、Elasticsearch、HBase、Redis



产品使用常见问题

最近更新时间: 2021-09-15 15:59:08

Q: 如何快速测试流计算作业是否有问题?

A: 可通过以下步骤进行测试。

- 在IDE开发页面进行在线开发, 开发完毕后, 可点击“单元测试”快速进行作业测试;
- 单元测试使用线下上传的测试数据包充当数据源, 支持对数据进行延迟、乱序测试, 并快速生成结果;
- 单元测试无误后, 可将作业发布到测试环境, 进行试运行。经过数据验证后, 可将该流计算作业发布至生产环境进行正式运行。

Q: 流计算作业怎么更新?

A: 可通过以下步骤进行测试。

- 流计算作业支持多版本, 版本生成动作通过开发IDE界面的“提交”进行触发。
- 作业提交前, 可以在流计算的开发IDE页面进行作业编辑。
- 作业提交后, 为保障已有版本作业的稳定, 以后版本不可再进行编辑操作, 但可以通过生成新版本的方式替换老版本, 实现流作业的更新。

Q: SavePoint有什么作用, 如何创建?

A: SavePoint可以支持程序升级后, 继续从升级前的那个点开始执行计算, 保证数据不中断。SavePoint需要用户手动创建。

Q: 是否支持自定义函数, 自定义函数的作用范围多大?

A: 流计算支持用户自定义函数, 函数统一作用在租户级别。

Q: 流计算是否支持用户自定义扩展插件?

A: 流计算支持Kafka source、SQL Operator、ETL Operator、JDBC sink、HBase sink、ElasticSearch sink、Redis sink、Kafka sink 这8个插件, 可灵活满足用户的各种需求。插件由平台统一维护, 不支持租户自定义扩展。

Q: 流计算有哪些使用限制?

A: Flink SQL不支持TopN、Emit等语法, 可采用ETL实现。Kafka source多流Join, 不可同时支持多个Kafka 版本。