

# UniCloud 集成平台

## 用户手册

紫光云技术有限公司  
[www.unicloud.com](http://www.unicloud.com)

资料版本：5W100-20211130  
产品版本：UniCloud iPaaS (E6101)

© 紫光云技术有限公司 2021 版权所有，保留一切权利。

未经本公司书面许可，任何单位和个人不得擅自摘抄、复制本书内容的部分或全部，并不得以任何形式传播。

对于本手册中出现的其它公司的商标、产品标识及商品名称，由各自权利人拥有。

由于产品版本升级或其他原因，本手册内容有可能变更。紫光云保留在没有任何通知或者提示的情况下对本手册的内容进行修改的权利。本手册仅作为使用指导，紫光云尽全力在本手册中提供准确的信息，但是紫光云并不确保手册内容完全没有错误，本手册中的所有陈述、信息和建议也不构成任何明示或暗示的担保。

# 前言

本手册主要介绍了 UniCloud 集成平台概述、功能介绍、典型配置案例、常见问题解答等内容。前言部分包含如下内容：

- [读者对象](#)
- [本书约定](#)
- [资料意见反馈](#)

## 读者对象

本手册主要适用于如下工程师：

- 网络规划人员
- 现场技术支持与维护人员
- 负责网络配置和维护的网络管理员






## 本书约定

### 1. 图形界面格式约定

格 式	意 义
<>	带尖括号“<>”表示按钮名，如“单击<确定>按钮”。
[ ]	带方括号“[ ]”表示窗口名、菜单名和数据表，如“弹出[新建用户]窗口”。
/	多级菜单用“/”隔开。如[文件/新建/文件夹]多级菜单表示[文件]菜单下的[新建]子菜单下的[文件夹]菜单项。

### 2. 各类标志

本书还采用各种醒目标志来表示在操作过程中应该特别注意的地方，这些标志的意义如下：

 警告	该标志后的注释需给予格外关注，不当的操作可能会对人身造成伤害。
 注意	提醒操作中应注意的事项，不当的操作可能会导致数据丢失或者设备损坏。
 提示	为确保设备配置成功或者正常工作而需要特别关注的操作或信息。
 说明	对操作内容的描述进行必要的补充和说明。
 窍门	配置、操作、或使用设备的技巧、小窍门。

### 3. 端口编号示例约定

本手册中出现的端口编号仅作示例，并不代表设备上实际具有此编号的端口，实际使用中请以设备上存在的端口编号为准。

## 资料意见反馈

如果您在使用过程中发现产品资料的任何问题，可以通过以下方式反馈：

**E-mail: [unicloud-ts@unicloud.com](mailto:unicloud-ts@unicloud.com)**

感谢您的反馈，让我们做得更好！

# 目 录

1 概述	1-1
1.1 简介	1-1
1.2 产品架构	1-1
2 访问数字平台的集成服务	2-1
2.1 首页	2-1
2.2 退出登录	2-2
3 功能介绍	3-1
3.1 工作空间管理	3-1
3.2 数据源管理	3-1
3.3 数据集成	3-1
3.3.1 主要功能	3-1
3.3.2 ETL 任务组件简介	3-3
3.4 服务集成	3-8
3.5 消息集成	3-10
3.6 告警管理	3-11
3.7 资产目录	3-11
3.8 资产市场	3-12
3.9 系统	3-12
3.10 运维	3-13
3.11 个人中心	3-14
4 典型配置案例	4-1
4.1 数据集成	4-1
4.1.1 增量抽取场景-时间戳	4-1
4.1.2 增量抽取场景-自增 ID	4-11
4.1.3 数据转换场景	4-22
4.1.4 实时流场景	4-28
4.1.5 大数据组件场景	4-36
4.1.6 数据清洗场景	4-42
4.1.7 整库迁移	4-50
4.1.8 GPLoad 加载	4-52
4.1.9 REST 抽取	4-58
4.2 服务集成	4-69

4.2.1	将数据库字段共享开放成接口场景-数据 API .....	4-69
4.2.2	接入第三方系统接口场景-通用 API .....	4-76
4.2.3	复杂场景下第三方接口对接-函数 API .....	4-79
4.2.4	画布方式实现多接口编排场景 .....	4-82
4.2.5	资产市场订阅使用接口场景 .....	4-84
4.3	消息集成 .....	4-88
4.3.1	作为消息中间件的生产消费场景 .....	4-88
<b>5</b>	<b>典型应用案例 .....</b>	<b>5-1</b>
5.1	医保云案例 .....	5-1
5.1.1	应用现状 .....	5-1
5.1.2	解决方案 .....	5-1
5.1.3	示例详细流程 .....	5-1
5.2	疫苗接种案例 .....	5-18
5.2.1	需求介绍 .....	5-18
5.2.2	操作流程 .....	5-19
5.2.3	抽取基础数据 .....	5-19
5.2.4	创建数据源 .....	5-22
5.2.5	新建数据表 .....	5-24
5.2.6	构建业务流程 .....	5-32
5.2.7	数据查询 .....	5-43
5.2.8	结果数据发布 .....	5-44
5.2.9	数据最终呈现 .....	5-46
<b>6</b>	<b>常见问题解答 .....</b>	<b>6-1</b>
6.1	HBase、Hive、HDFS、Kafka 等大数据组件开启了 Kerberos 认证,连接这些数据源时如何配置 Kerberos 认证信息 .....	6-1
<b>7</b>	<b>附录 .....</b>	<b>7-1</b>
7.1	业务数据库建表语句示例 .....	7-1
7.2	stg_2_ods_cep_stcdb_mdtrt_d.sql 脚本内容 .....	7-3
7.3	ods_2_dwd_cep_stcdb_mdtrt_d.sql 脚本内容 .....	7-16
7.4	ods、stg 及 dwd 建表语句 .....	7-18

# 1 概述

## 1.1 简介

UniCloud 集成平台是一个全栈式的集成平台，旨在打通应用和数据孤岛，实现异构数据/API/消息/设备集成，提供异构数据集成、应用间通信集成能力、API 接口集成能力及物理设备设备集成能力，助力打造标准统一、融会贯通、资产化、服务化、闭环自优化的智能数据体系、以驱动应用创建，适用于多种常见的企业系统集成场景。

UniCloud 集成平台定义为提供消息集成、服务集成、数据集成的统一集成平台，各组件既能够单独运行也能组合成套件，各集成共享相同的技术底座，用户可根据需要自由配置部署。

## 1.2 产品架构

图1-1 产品架构



通过集成平台可以完成以下功能：

- 工作空间管理：  
工作空间的新增、编辑、修改，删除等。支持工作空间的导入、导出及资产（任务、数据源、API、Topic）的详情查看。
- 数据源管理：

丰富的数据源支持。支持的数据源类型包括 DB2、达梦、Greenplum、HBase、MPP、MySQL、Oracle、PostgreSQL 等，为数据来源提供统一的管理。

- 数据集成

数据集成是一个以调度、监控和管理 ETL 过程为核心功能的应用系统。该系统通过图形化工具，快速灵活地设计与部署，实现数据抽取、转换及加载，并能在设计中设置统一的清洗规则，从而提升数据的质量，能为企业和组织提供一套完备的数据集成解决方案。

- 消息集成

消息集成旨在为集成平台提供可靠、无状态、满足各应用间信息最终一致性的消息集成服务。支持原生的 Kafka 特性，具备原生 Kafka 所有消息处理特性；支持安全的消息传输，通过 sasl 认证、消息存储加密等措施加强网络访问控制；支持消息数据高可靠，支持消息持久化、多副本存储机制。

- 服务集成

服务集成主要用于数据库表字段的开放、第三方接口的代理转发、文件资源的开放等。支持 API 注册、API 测试、API 部署、API 授权、API 编辑、API 删除、API 版本管理等全生命周期管理。可以对 API 的访问进行统计分析，记录访问日志，实现对 API 访问的审计功能。

- 告警管理

告警管理功能可实时检测系统运行时出现的各种告警。告警列表模块可查看系统中的当前告警和历史告警，进行确认或查看解决方案和详情日志。

- 资产目录

用户可对类目进行增、删、改。页面展示选定类目下的资产目录列表，用户可对资产目录可以进行新增、查看详情、编辑、下线、发布、批量发布、批量下线、删除、查看详情等操作，还可以根据名称、共享条件、创建时间等条件对目录进行进行搜索

- 资产市场

用于展示系统中已发布的数据资产、服务资产、AI 资产以及消息资产。用户可根据实际情况，订阅需要的资产进行使用

- 系统管理

提供用户、组织、权限、操作日志、流程审批、系统升级、软件授权等基础服务管理。

- 运维

运维提供了数字平台的运行维护功能，包括服务的管理、资源监控及系统巡检等功能。



## 2 访问数字平台的集成服务



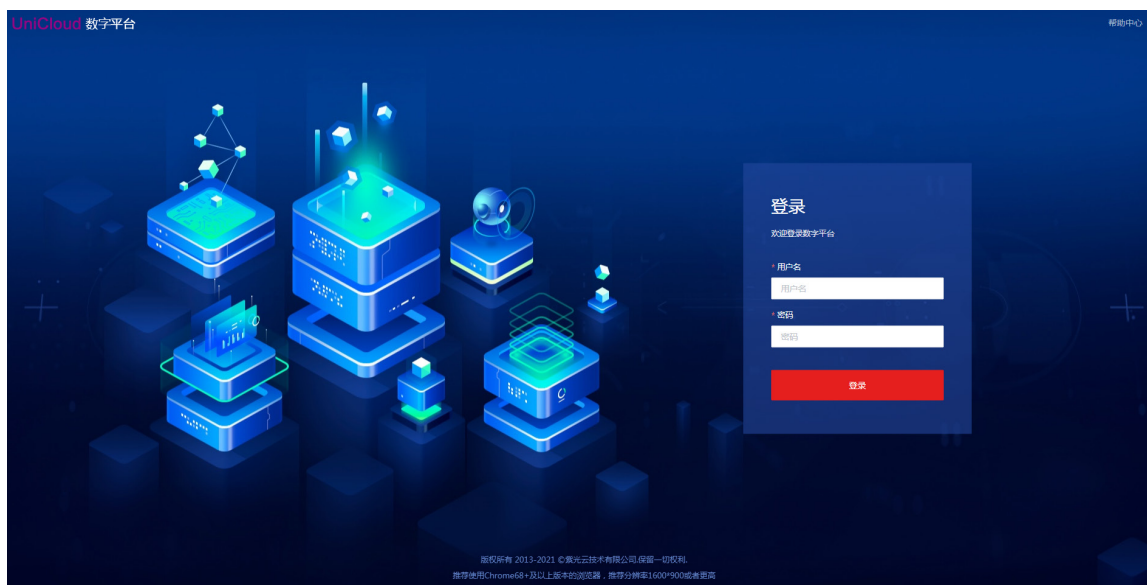
注意

为保证系统安全，登录系统后请及时修改登录密码。

数字平台的系统服务安装成功以后，数字平台的 URL 地址会在安装脚本成功执行完成后显示，格式为：<https://VIP:32015>。

在浏览器中输入地址：<https://VIP:32015>，进入登录页面，如[图 2-1](#)所示。输入正确的用户名和密码，单击<登录>按钮即可登录数字平台。数字平台缺省的超级管理员用户名为 `admin`，缺省密码为 `Passw0rd@_`。如果用户名或密码不正确，系统会弹出相应的错误提示。

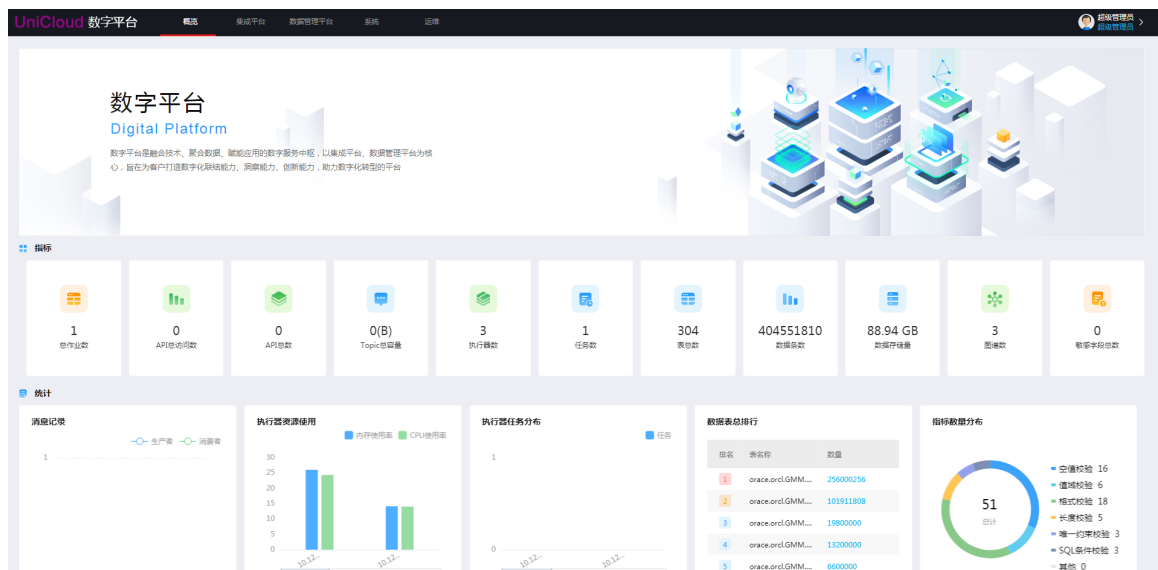
图2-1 登录页面



### 2.1 首页

首页展示了数字平台的统计信息，首页如[图 2-2](#)所示。

图2-2 首页




## 2.2 退出登录

登录成功后，在右上角登录的当前用户下拉菜单中选择[退出]菜单项，即可退出系统。

# 3 功能介绍



说明

- 集成平台中包含联机帮助，单击页面左上角的  帮助按钮，弹出窗口中提供各功能详细的配置说明、操作指导以及注意事项等，可帮助用户更好的使用集成平台。
- 本章节仅对各功能进行概括说明，以便用户快速了解数据集成平台提供的主要功能，关于各功能的详细说明请参见集成平台的联机帮助。

## 3.1 工作空间管理

工作空间是为了让用户更好的将项目相关资源进行统一管理。比如用户可以新建工作空间，然后将同一项目的相关资源创建在一个工作空间下，方便后续查看及操作。组织内的用户创建的工作空间属于该组织。该组织下的用户可以查看该组织下所有的工作空间并进行操作。

## 3.2 数据源管理

不同使用场景下需要连接不同的数据源，在数据源管理页面可新增数据源并管理系统中所有已添加的数据源。集成平台支持丰富的数据源类型，包括：DB2、达梦、Greenplum、HBase、MPP、MySQL、Oracle、PostgreSQL 等。

## 3.3 数据集成

数据集成是一个以设计、部署、调度、监控和管理 ETL 过程为核心功能的应用系统。

操作者借助该平台可以通过流程图式的图形化工具快速、灵活地设计 ETL 过程，并能方便的进行部署、调度及监控等管理活动，真正地提供一体化数据集成开发环境。

功能特性包括：

- 数据集成任务模型全生命周期管理。
- 对作业进行多维度调度、监控，实现作业自动化处理。
- 根据已存在任务生成任务模板，任务模板组成作业模板，通过作业模板快速生成任务及作业，并自动部署作业，使大量重复的工作完成自动化。
- 通过标签标记任务及作业，在有大量任务和作业的场景中可根据标签快速定位作业或任务。

### 3.3.1 主要功能

如表 3-1 所示，数据集成具备的主要功能包括：定制化首页、资源配置、任务管理、作业管理、整库迁移、标签管理、定期清理等。

表3-1 数据集成主要功能说明

功能		说明
定制化首页		监控功能在数据集成首页以定制化首页的形式展示，通过在自定义面板下拉列表中选择，在首页中增加监控面板。监控功能包括：集成统计概述、完成情况趋势、集成节点任务分布、运行监控和重点作业监控
资源配置	执行器管理	执行器是执行任务的实际容器，用于执行服务端发送的作业中的任务
	调度器监控	调度器用于接收服务端下发的作业，并将作业中的环节调度到具体的执行器中运行
	资源收藏	展示同组织内用户收藏的任务组件资源。用户在新建数据集成任务时，可将常用组件添加到收藏夹中，方便组织内的用户使用
任务管理	任务列表	<p>任务列表模块可新增/导入任务并管理系统中所有已创建的任务。数据集成支持任务类型如下：</p> <ul style="list-style-type: none"> <li>• 普通 ETL 任务可对源数据进行抽取和转换，然后将结果加载至目标数据源。数据集成支持丰富的数据源类型和多种数据转换操作，包含的组件详情请参见 <a href="#">3.3.2 ETL 任务组件简介</a></li> <li>• Sqoop 任务支持丰富的数据源类型和 5 种任务类型（数据库到 HDFS、数据库到 HBase、数据库到 Hive、HDFS 到数据库、Hive 到数据库）</li> <li>• Shell 任务支持以 Shell 脚本形式直接运行任务</li> <li>• SQL 任务支持 SQL 脚本形式直接运行任务</li> <li>• Flume 任务支持在指定主机上运行 flume-ng agent 命令执行 flume 任务</li> <li>• Kettle 任务支持连接 Kettle carte 客户端，直接运用 Kettle 里任务保存后得到的.ktr 文件</li> <li>• KettleJob 任务连接 Kettle carte 客户端，直接运用 Kettle 里作业保存后得到的.kjb 文件</li> <li>• CDC 任务基于数据库日志解析，完成数据库数据实时同步</li> <li>• DataX 任务支持连接 DataX 客户端</li> </ul>
	任务模板	通过已创建的任务，将变量部分替换为参数得到任务模板，然后填入参数值可直接生成任务。任务模板模块可新增任务模板并管理系统中所有已创建的任务模板。该功能适用于有大量重复任务，且这些任务的差异集中在部分参数设置值不同的场景中
作业管理	作业列表	作业由若干任务构成，是数据集成系统中运行的最小单位。作业列表模块可新增作业并管理系统中所有已创建的作业。作业支持多种调度类型，且下发作业时根据每个执行器上作业执行的情况以及剩余资源等，系统通过计算会提供执行器的推荐选择星级
	作业模板	作业模板由任务模板构成，通过作业模板可直接快速生成作业中的任务和作业，

功能		说明
		<p>并下发部署作业。作业模板模块可新增作业模板并管理系统中所有已创建的作业模板，同时可查看各作业运行实例的监控历史详情，包括任务运行的统计数据 and 运行速度详情、运行日志等</p> <p>该功能适用于有大量重复作业，且这些作业的差异集中在部分参数设置值不同的场景中</p>
整库迁移		<p>全量迁移模块可直接创建作业，将源数据库中选定的某些表批量复制到目标数据库，节省批量创建任务的时间</p>
标签管理		<p>根据业务特点可以给作业和任务自定义创建标签，以实现快速查找。标签管理模块可新增标签并管理系统中所有已创建的标签</p> <p>该功能主要适用于有大量作业及任务的场景</p>
定期清理		<p>配置数据定期清理策略，可防止系统数据库因存在过多的冗余数据而导致系统变慢。定期清理模块可新增定期清理策略并管理系统中所有已创建的定期清理策略</p>

### 3.3.2 ETL 任务组件简介

针对每个普通 ETL 组件，拖拽至任务设计面板上双击打开弹窗，单击页面左上角的 ? 帮助按钮，可查看各组件的操作步骤、配置说明及使用示例，帮助用户更好的使用组件。

表3-2 普通 ETL 任务数据抽取组件

数据抽取组件	描述
表抽取	表抽取组件用来利用已配置好的数据库连接和SQL语句，从数据库中读取数据
文件抽取	通过文件抽取组件，可以读取单个或多个文本文件、指定读取的文件列表或者用正则表达式表示的目录列表。该组件可支持抽取本地文件或目录、FTP文件或目录、SFTP文件或目录、HTTP文件
生成记录	生成记录组件可根据配置输出指定数量的记录行，缺省为空。可选包括一定数量的静态字段
Kafka抽取	Kafka抽取组件可以根据配置从Kafka消息系统中抽取数据
Kafka流抽取	Kafka流抽取组件从Kafka抽取流数据，并运行子转换，该转换根据消息批量大小或持续时间执行，可近乎实时地处理连续的数据流
HDFS抽取	HDFS抽取组件可对HDFS文件系统中结构化的数据进行抽取
HBase抽取	HBase抽取组件可以将存储在HBase表中的数据抽取到其他类型数据库中
REST抽取	REST抽取组件用来从REST服务中抽取数据
Excel抽取	Excel抽取可以将存储在Excel表中的数据抽取到其他类型数据库中

数据抽取组件	描述
WebService抽取	WebService组件用来从WebService服务中抽取数据
MongoDB抽取	MongoDB抽取组件可以将存储在MongoDB中的数据抽取到其他类型数据库中
Redis抽取	Redis抽取组件可以将数据从Redis数据库抽取出来，目前支持的数据类型包括String、Set、List、Hash、Hashall及ZSet等。该组件不能作为转换任务的开始步骤，需要在其前面有其他输入组件进行值传递
JSON抽取	JSON抽取组件用于从JSON输入源中抽取数据
XML抽取	XML抽取组件使用XPath规范将XML文件的数据解析到流中，支持本地/FTP的XML文件 (*.xml) 和以其他方式读取的XML格式数据的解析
XML输入流	XML输入流组件使用StAX解析器从XML源中读取数据，与XML抽取组件使用的DOM解析器相比，该组件适用于XML文档结构复杂和数据量大的场景
JMS抽取	JMS抽取组件可以从Apache ActiveMQ JMS服务器或IBM MQ中间件消费流数据
MQTTConsumer	MQTTConsumer组件可以从MQTT代理或客户端抽取流数据。MQTTConsumer步骤运行于子转换，该转换根据消息批量大小或持续时间执行，可近乎实时地处理连续的数据流
从结果获取记录	和复制记录到结果 组件一起使用，其可以获取其他ETL任务中复制到结果的数据，用于作业中的数据传递
CSV文件抽取	CSV文件抽取组件提供从定界文件中读取数据的功能
HttpServer	支持根据用户自定义的配置实现http服务器的实时部署，通过请求发布的http服务器接口可以实现用户数据主动推送到数据集成侧，接收到的数据可以通过其它组件进行数据处理

表3-3 普通 ETL 任务数据转换组件

数据转换组件	描述
字符串操作	字符串操作组件可以对输入字段类型为String的字段进行补位、改变大小写、转义等操作
字符串切割	该组件用于切割字符串的一部分。如果指定字段的起始位置超出范围，则返回空白
字段选择	字段选择组件用来选择字段、重命名字段、指定字段的长度或精度
正则表达式	正则表达式组件可以将输入字段的字符串值与正则表达式定义的文本模式匹配，将匹配结果存入新的字段中可以为后续连接的组件使用
对称加密	该组件可以对流中数据进行加密和解密，支持的对称加密算法有：AES/DES/DESeed三种。建议配合生成密钥组件一起使用
JS代码	JS代码组件允许通过JavaScript语言对数据做复杂的运算。右侧JavaScript函数列表包含了该组件支持的JavaScript的函数

数据转换组件	描述
去重记录	去重记录组件从输入流中移除重复的记录，可以输入字段名称直接去重
排序记录	排序记录组件可以利用指定的字段对行按照升序或降序进行排序。当行数超过5000行的时候，将使用临时文件来排序行
增加常量	增加常量组件用于添加常量到流中，用字符串形式指定名称、类型和值。利用选择的数据类型指定转换格式
字符串替换	字符串替换组件允许指定的字符串替换输入流中的原指定字符串，并生成新的输出字段
计算器	计算器组件提供一个功能列表，可以在字段值上运行，操作更加便捷，且运算效率比JavaScript脚本更高
设置字段值	设置字段值组件用于将流中某一字段的值赋值给另一字段
增加序列	增加序列组件可以在流中新增一列，这一列为在某个起始值和增量的基础上的序列。可以生成两种序列，一种为自定义的序列，可以自己设置起始值，增长值及最大值，但每次转换运行的时候序列的值又会重新循环一次；另一种为使用数据库的序列
公式	公式组件可以实现复杂的逻辑判断，或者使用该组件自带的函数进行数据的转换清洗
过滤记录	过滤记录组件允许根据条件和比较符来过滤记录。这个步骤一旦连接到先前的步骤中，就可以通过简单的单击“<field>”、“=”和“<value>”区域来构建条件。该步骤支持2层条件嵌套，如果有复杂判断条件，需要配合JS代码组件一起使用
列拆分为多行	列拆分为多行组件可将输入数据行集中的某个列按照条件拆分为多行，这种条件可以是简单的一个分隔符，也可以指定正则表达式
拆分字段	该组件可根据分隔符信息拆分字段，拆分后的字段将形成新的列
行转列	行转列组件用于将数据某一列中的唯一值转换为输出中的多个列来旋转该数据，即将列值转换为列名，并可以在输出的时候对最后输出中的任何其余列值进行聚合
列转行	该组件可以将输入中的列转换为行
数据检验	数据检验组件通常用于确保传入数据的质量。通过定义一些简单的检验规则来确保传入的数据符合指定的数据规范，例如：数据的范围、长度和类型等，该组件可以同时定义多种检验规则
行比较	行比较组件用于比较两个不同来源的数据，这两个来源的数据分别为旧数据源和新数据源，该步骤将旧数据和新数据按照指定的关键字匹配、比较后进行合并，该组件主要用于数据表的同步更新
分组	分组组件允许通过定义字段实现分组后再执行计算，即可以按照某一个或某几个字段进行分组，同时将其余字段按照某种规则进行计算。例如：计算产品的平均销售额，获取某商品库存的数量等
数值范围	该组件比较一个字段的数字范围，在这个范围的话就输出一个新的列，并且打印出来你定好

数据转换组件	描述
	的结果。如果不在指定的范围则打印出来unknown，或者打印自定义的结果
流连接	流连接组件可以将多个前面步骤的数据以INNER/FULL OUTER的方式合并数据
Java代码	Java代码组件允许使用Java语言编写一个step插件
流查询	流查询组件允许使用来自其他组件的信息查找数据。来自查询组件的数据首先被读取到内存中，然后根据与查询关键字所匹配的字段从查询组件中查找数据
记录集连接	记录集连接组件用于将2个输入步骤的数据执行合并数据的操作
值映射	该组件用于将字符串值从一个值映射到另一个值
数据库查询	该组件可以使用设置的关键字在目标表中查询，并从查询结果中返回指定的字段
添加XML	该组件可以生成XML格式数据

表3-4 普通 ETL 任务数据加载组件

数据加载组件	描述
加载至表	加载至表组件可以将数据加载到数据库表中
加载至Excel表	加载至Excel表组件可以将数据加载到Excel表中
加载至文件	加载至文本文件组件用于将数据加载到文本文件中，并可以通过配置路径信息加载至远程文件中
丢弃数据	丢弃数据组件，表示不对数据做任何操作，常用于在调试ETL任务中测试和丢弃数据
加载至ES	加载至ES组件可以将数据写入到Elasticsearch中
加载至HDFS	加载至HDFS组件可以将数据加载至HDFS集群上的文本文件中
加载至HBase	加载至HBase组件可以将数据加载至HBase数据库中
加载至Kafka	加载至Kafka组件可以将数据发布至Kafka消息系统中
Kafka流加载	Kafka流加载组件允许以近乎实时的方式将消息发布到Kafka服务器
数据删除	数据删除组件可以利用查询关键字在表中搜索行。如果能被找到，行就会被删除
数据库更新	数据库更新组件可以利用查询关键字在表中搜索行。如果字段能被找到，并且没有任何改变，字段不更新；如果字段有改变，字段就会被更新
加载至Redis	加载至Redis组件可以将数据加载到Redis数据库中，目前支持的数据类型包括String、Set、List、Hash及Sorted Set等。该组件常和“字段选择”组件一起使用，即可以在“字段选择”组件中将需要的流字段按顺序先选择出来，然后再给该组件使用
GPLoad	GPLoad组件可以使用Greenplum的外部表并行加载功能进行大规模并行加载数据



数据加载组件	描述
加载至XML	加载至XML组件用于加载数据到XML文件中，可以根据配置路径加载至远程或本地文件中
加载至JSON	加载至JSON组件可以将数据生成JSON块，并选择输出到文件，或者生成包含JSON块的数据流
加载至JMS	加载至JMS(Java Messaging Service)组件几乎可以实时地将消息发布到Apache ActiveMQ JMS服务器或IBM MQ中间件
MQTTProducer	MQTTProducer组件允许以近乎实时的方式将消息发布到MQTT代理
插入更新	插入更新组件可以利用查询关键字在表中搜索行。如果找不到，则插入该行。如果可以找到，并且要更新的字段相同，则不执行任何操作；如果不完全相同，则更新表中的行
Oracle批量加载	Oracle批量加载组件可以将源数据以适当的加载格式写入，然后调用Oracle SQL*Loader将其传输到指定的表
MySQL批量加载	该组件可以将数据批量加载到MySQL数据库中，批量插入的原理是利用了MySQL一个高效导入方法load data infile将文本文件中的数据导入表中，文本文件是由系统后台根据前序步骤的数据流生成的FIFO文件，因为用户数据的不可控性，故后台自动生成的FIFO文件的相关配置需要用户自行配置
加载至MongoDB	加载至MongoDB组件可以将其它地方的数据，加载到MongoDB数据库中
PostgreSQL批量加载	PostgreSQL批量加载组件可以批量加载数据到PostgreSQL数据库中
复制记录到结果	和“从结果获取记录”步骤配合使用，此步骤将上一步骤的数据进行记录，供作业中其它ETL类型的任务使用此数据
数据同步	该组件可以与“行比较”组件结合使用。“行比较”组件使用标志字段来保存比较结果，该组件使用此标志字段对数据库表进行更新、插入及删除

表3-5 普通 ETL 任务其他组件

其他组件	描述
存储过程	存储过程组件可以运行一个数据库存储过程，并获取返回结果
执行SQL	执行SQL组件用于对目标数据库执行SQL脚本，可以在转换初始化的时候执行。执行方式分为执行每一行、作为一个语句执行以及变量替换三种
发送邮件	发送邮件组件用来给指定邮箱发送E-Mail
检查文件存在	检查文件存在组件可以在本地系统或FTP服务器中检验是否存在某个文件，并将布尔标志字段添加到输出字段
FTP下载	FTP下载组件可以从FTP服务器上获取一个或者多个文件

其他组件	描述
HTTP下载	HTTP下载组件可以通过HTTP协议从Web服务器上获取一个文件
文件解压	文件解压组件可以解压文件，但只能解压zip格式的本地文件
表结构复制	表结构复制组件可以获取源数据库中的表结构，并在目标库中创建相同结构的表
从流中获取记录	从流中获取记录组件用于返回另一个任务生成的记录，您可以输入该任务所期望的字段元数据。该组件需作为任务的开始步骤，所在任务可以作为JMS抽取、MQTTConsumer抽取以及Kafka流抽取组件的子转换使用
生成密钥	该组件用于生成一定数量的密钥，生成的密钥可以在对称加密组件中使用
设置变量	设置变量组件可以在作业的维度内设置一个变量值，作业内的其他组件如果使用了“变量替换”功能，那么该组件内的所有变量都会在运行时被替换为变量组件中设置的同名变量值

表3-6 CDC 任务组件

CDC 组件	描述
MySQL CDC	捕获MySQL数据变更的Binlog信息，解析出相应数据并发送到Kafka Topic
Oracle CDC	使用LogMiner解析Oracle归档日志，获取数据库的数据变化信息，解析出相应数据并发送到Kafka Topic
JDBC Producer	从Kafka中获取CDC日志，解析后在目标数据库中进行相应的增删改操作

## 3.4 服务集成

服务集成主要用于数据库表字段的开放、第三方接口的代理转发、文件资源的开放等。支持 API 注册、API 测试、API 部署、API 授权、API 编辑、API 删除、API 版本管理等全生命周期管理。可以对 API 的访问进行统计分析，记录访问日志，实现对 API 访问的审计功能。

表3-7 服务集成功能说明

特性	描述
概览	概览页面可以实时展示当前用户所属组织中开放 API 的今日访问次数、今日访问成功率、部署 API 的数量、授权 API 的数量以及激活工作空间的个数、组织内 API 开放统计、工作空间内 API 开放统计、API 类型开放统计以及 API 来源开放统计
API工厂	API管理功能主要用于API注册、继承注册、API导入、API编辑、API测试、API部署、API撤销部署、API修改记录查询以及API删除等 <ul style="list-style-type: none"> <li>数据 API 设计：数据 API 通过向导式的设计将数据库表里字段开放出去，以标准 <code>resetful</code> 接口的形式对外提供</li> </ul>

特性	描述	
	<ul style="list-style-type: none"> <li>● 通用 API 设计：通用 API 设计用来对已经存在的第三方接口进行代理</li> <li>● 函数 API 设计：函数 API 可以在系统内通过编写 JavaScript 脚本，设计并生成标准 <code>resetful</code> 接口</li> <li>● 继承注册：继承已经创建好的 API，快速生成一个新的 API</li> <li>● API 导入：通过模板快速导入通用 API 和函数 API</li> <li>● API 测试：API 注册后可以通过在线测试查看接口是否可用</li> <li>● API 部署：API 测试通过后可以部署到网关节点，进行授权访问</li> <li>● API 上架：API 部署后可以申请上架，管理员审批后可以展示在资产市场供他人订阅</li> <li>● API 版本管理：用来查看 API 的修改记录，并可以通过载入功能查看对应版本 API 的详情信息，并支持在该历史版本的基础上进行修改</li> </ul>	
服务编排	服务编排页面主要是对现有一些业务接口按照特定的业务执行流进行组织和编排，生成一个统一的对外访问接口，该接口实现了多个业务接口的功能	
认证模板	通过在认证模板页面添加认证模板，注册通用 API 时关联认证模板，解决带认证接口的访问权限问题，自动帮用户维护接口访问所需认证授权信息，使用户更能专注于自己的业务	
API 组管理	API 组管理功能主要用于管理 API 注册时选择的分组。API 组管理页面展示了当前用户所属组织内创建的 API 组，展示内容包括组名称、组内包含的 API 数量、描述、创建人、更新时间等信息，并可对 API 组进行新增、编辑及删除操作	
密码箱管理	密码箱功能主要用于函数 API 设计时对工作空间秘钥的获取	
环境配置	环境配置功能主要用于函数 API 设计时获取环境变量	
API网关	API列表	API 列表下展示了指定工作空间下已部署的的 API 信息。并可以根据实际需要 API 授权/取消授权给指定的工作空间，并可对 API 的访问配置访问规则和熔断规则
	授权的API	授权的 API 页面下展示了当前工作空间下被授权的 API 信息。并可以根据实际需要 API 进行查看详情、测试及取消授权等操作
	监控统计	监控统计页面对应用访问 API 的日志信息进行多维度分析展示，让用户对 API 的使用情况一目了然，包括 API 调用趋势分析和根据客户端 IP、响应状态、响应时间进行分类统计等
	日志搜索	日志搜索页面可以根据搜索条件来查看用户所在组织下的指定的 API 访问详情日志信息
	访问网络	访问网络页面定义访问网络分类，用于对转发节点进行分类标记。访问网络页面可维护平台部署的网络信息，支持多个网络的维护
	节点管理	节点管理页面用于配置可用的跨网转发节点，作为网关来代理发布接口。维护部署

特性		描述
		的转发节点信息，支持不同网络的转发节点维护，可通过不同网络的转发节点实现接口的跨网发布
	规则管理	规则管理分为 IP 规则和超限规则，IP 规则可以限制能访问 API 的 IP 地址或 IP 网段，超限规则可以限制每小时或者每天的访问次数，支持 IP 黑白名单
	熔断管理	熔断管理分为信号量隔离模式和线程隔离模式，可以设置熔断开启的条件，在网关访问第三方接口发生错误时，出发熔断，避免系统雪崩
文件管理	文件列表	文件列表页面展示系统中共享的文件列表，并可对文件进行共享、编辑、删除及授权操作。文件可以申请上架到资产市场，管理审批通过后可以在资产市场展示供他人订阅
	授权的文件	授权的文件页面可以查看已授权给当前工作空间的文件资源详情并下载文件
	文件资源	文件资源页面以列表形式展示当前用户所属的文件资源以及用户所属组织下的文件资源
	存储配额	存储配额管理页面用于对组织内用户的文件空间存储大小进行配置管理
服务分类	行业领域	行业领域页面用于管理服务所属行业领域分类，并可对行业领域进行新增、删除及编辑操作
	API 来源	API 来源页面用于管理 API 来源分类，并可对 API 来源进行新增、删除及编辑操作
订阅管理		订阅管理页面展示了当前组织下发布的服务列表以及订阅信息，可以进行规则配置和熔断配置，可以终止订阅

### 3.5 消息集成

消息集成定义为集成平台中各集成服务之间提供可靠的、可持久化的、高吞吐量的准实时消息管道系统。消息集成使用统一的消息接入机制，标准化的消息通道，具有如下几方面的优势：

- 支持原生的 Kafka 特性：具备原生 Kafka 所有消息处理特性。
- 支持安全的消息传输：通过 sasl 认证，消息存储加密等措施加强网络访问控制。
- 支持消息数据高可靠：支持消息持久化，多副本存储机制。支持节点级扩容与 Topic 重分配。

表3-8 消息集成功能说明

功能		说明
概览		概览页面主要展示了Topic总容量、历史记录，消费者，生产者等情况的统计信息
Topic管理	Topic列表	Topic列表页面展示Topic的相关信息，用户可以对Topic进行新建、编辑、删除、授权、清空数据等操作
	Topic配置	Topic配置页面主要提供对Topic属性配置进行新增、删除、查看配置操作。当用户

功能		说明
		对已创建的Topic有属性配置的诉求时，可以使用该模块进行配置
	重分配	Topic重分配功能用于Topic修改分区所在broker节点位置
消息查询		消息查询主要提供毫秒级、可视化的Kafka集群中的消息查询能力，支持按照分区和生产时间进行过滤
消息转发		主要用于将一个Topic中的数据经过一定规则过滤后写入另外一个Topic。用户可在消息转发页面下配置新的消息转发规则，并对这些消息转发规则进行管理
消费进度		显示组消费进度和活跃Topic的情况
监控	Broker监控	显示整个系统中Brokers性能度量情况
	Kafka监控	显示Kafka运行情况，可根据自己需要选择指标进行展示
	Zookeeper监控	显示Zookeeper监控信息，包括发包数量、收包数量、活跃连接数、排队请求数在各时间节点的数据分布情况
集群信息		集群信息页面展示了Kafka集群和Zookeeper的相关信息

## 3.6 告警管理

表3-9 告警管理功能说明

特性	描述
告警类型	包含任务失败、心跳丢失、调度无效、Kafka集群异常、API接口调用异常、共享存储服务异常等告警类型
告警状态	包含严重告警、重要告警
当前告警	当前告警页面会向后台查询系统的当前告警信息
历史告警	历史告警页面显示系统中的历史告警信息

## 3.7 资产目录

表3-10 资产目录主要特性列表

特性	描述
资产目录管理	资产目录页面左侧以树结构展示系统中的类目，用户可对类目进行增、删、改。页面右侧展示选定类目下的资产目录列表。用户可对资产目录进行新增、编辑、删除、下线、发布、批量发布、批量下线、批量删除、查看详情等操作，还可以根据名称、状态、创建时间等条件

特性	描述
	对目录进行进行搜索
类目管理	类目管理可以对类目进行统一的管理操作。类目管理以树形结构对系统类目进行整体展示，默认分为主题、基础和部门三大模块，可以根据需要进行新增。单击左侧类目树，右侧会相应显示该类目下的资产目录

## 3.8 资产市场

表3-11 资产市场主要特性列表

特性	描述
总览	总览页面展示了资产统计指标，并可以根据资产名称进行全局搜索
数据资产	数据资产页面展示了系统中的数据资产和文件资产
服务资产	服务资产页面展示了系统中的服务资产，服务资产是由集成下的服务集成模块提供的API。包括数据API、通用API、函数API
AI资产	AI资产页面展示了AIOS上架到资产市场的AI应用服务
消息资产	消息资产页面展示了系统中的消息资产。消息资产是由集成下的消息集成提供的Topic，将Topic的生产/消费权限作为资产上架到资产市场，提供给用户

## 3.9 系统

表3-12 系统管理主要特性列表

特性	描述
组织管理	用于对系统内的所有组织进行管理，如新建、编辑或删除组织等。系统管理员可根据用户需求为根组织划分子组织，便于对用户进行分类管理
用户管理	用户即可以登录本系统并使用系统中功能和资源的人。用户管理提供了统一管理系统中用户的功能。主要包括系统未对接认证通和对接认证通两种情况下的配置
角色管理	用于管理本系统中用户的角色。不同角色具有不同的系统资源使用权力。本系统通过角色实现分组式的权限管理，每种角色都有自己的权限集，被赋予该角色的用户即具有权限集所包含的权利。系统内置了基本角色，也支持创建自定义角色
操作日志	是以用户为行为主体，从而产生的行为日志，需要包含时间、操作用户的名称、产生操作行为的IP地址、操作对象、级别、操作结果以及操作描述等主

特性		描述
		要信息，方便后续的审计和统计
系统配置	菜单管理	用于配置数字平台中各功能的布局及是否显示等，方便超级管理员对界面功能展示进行配置
	流程配置	系统管理员通过配置流程模板可以有效规范普通用户对资源的申请和使用，即系统提供了审批流程自定义功能，用户可以根据自身企业的情况自定义审批流程，通过在流程中添加多个审批人员实现多级审批的需求
	安全设置	对系统中的安全策略进行管理，包括密码策略、登录认证策略、SSO认证
	基础设置	基础设置提供了系统消息的配置功能和关联大数据平台的配置功能
	大数据集群配置	大数据集群配置用于配置数据管理平台所使用的大数据集群相关信息
	大数据集群资源	大数据集群资源用于管理系统中各组织的集群资源，包括指定Kerberos用户，可使用的YARN队列和Spark Oozie Launcher队列
软件授权		系统完成安装部署后，默认未完成License注册的系统为试用版本，用户可以使用试用授权来对资源进行申请。试用期到期后，用户需申请正式授权才能继续使用本系统

### 3.10 运维

运维提供了数字平台的运行维护功能，包括服务的管理、资源监控及系统巡检等功能。

表3-13 系统功能特性

特性	描述
服务管理	服务管理用于可视化部署数字平台各服务，如集成或数据管理平台中的各服务功能等，并对各服务进行统一管理
资源监控	<p>资源管理提供主机监控、磁盘管理、主机管理的功能：</p> <ul style="list-style-type: none"> <li>主机监控可以对部署服务的主机进行监控，为用户提供主机当前的运行状态，帮助用户更合理地使用主机资源；</li> <li>磁盘监控可以对部署服务的主机磁盘分区进行监控，为用户提供磁盘分区的实时使用率信息，帮助用户灵活调整数据存储位置，提高磁盘利用效率；</li> </ul>
系统巡检	系统维护支持一键巡检、定时巡检、查看报告、下载报告的功能

## 3.11 个人中心

表3-14 个人中心主要特性列表

特性	描述
个人信息	提供了当前登录用户修改个人信息和密码的功能
我的申请	展示了当前用户提交的申请流程，用户可查看申请流程详情、审批进度，并可根据实际需要撤销流程
待办审批	为用户提供了审批流程管理功能。我的审批页面展示了当前登录用户待办的审批流程，用户可对这些流程进行审批（同意/驳回）、更改责任人、删除等操作
所有申请	所有申请页面展示了当前登录用户相关的所有的申请流程，包括审批中和已完成的申请流程
我的订阅	展示了当前登录用户申请订阅的资产信息
我的收藏	展示了当前登录用户收藏的资产信息，包括：资产名称、资产类型、收藏的时间及描述信息，并提供了相关操作
工作空间	工作空间管理页面展示了用户已创建的工作空间列表，并可对工作空间进行新增、编辑及删除操作等



# 4 典型配置案例

## 4.1 数据集成

### 4.1.1 增量抽取场景-时间戳

目前很多企业都会构建自己的数据仓库，通过整合所有业务系统数据形成企业的数据资产，并对其进行深度分析挖掘数据背后的价值，以支持企业做出一些重要的决策。一般情况下，企业会采用增量数据抽取方式完成业务系统数据向数据仓库中的同步。为了实现数据的增量抽取，在设计数据表时新增一个时间戳字段，数据在入库时，会为时间戳字段赋值为入库时间。

#### 1. 场景描述

XX 交通运输管理公司的业务系统数据库中，维护一张车辆通过卡口信息记录表 `vv_1000`。现公司需要将该表中每天新增的数据定时向数据仓库中的 `vv_1000000` 表做数据同步。公司业务系统使用 PostgreSQL 数据库，数据仓库使用 SeaSQL MPP。

表4-1 源表（`vv_1000`）结构

KKWZBM	HPHM	HPZL	TGSJ	HBSJ	SSSD	RYLX	RKSJ
3116250005	莆P3Z907	2	2020/3/10 15:30:58	26	68	01	2020/3/10 20:10:47

#### 2. 场景分析

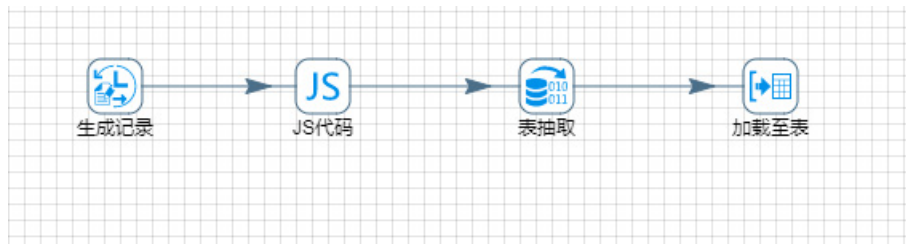
使用数据集成创建普通 ETL 任务，通过作业定时调度实现增量的抽取。

#### 3. ETL 设计方案

数据流向：生成记录 -> JS 代码 -> 表抽取 -> 加载至表

ETL 方案：数据增量是基于当前调度时间，要求每天凌晨 2 点将前一天的数据同步到目标表。比如当前调度时间为 2020/3/10 02:00:00，则增量抽取的数据是 2020/3/9 00:00:00 到 2020/3/9 23:59:59 这段时间入库的数据，这个时间范围可以使用 JS 代码组件依据当前调度时间计算得出，然后将起始时间作为参数传入表抽取步骤，将源表中时间戳（RKSJ 字段值）在指定时间范围内的数据传给“加载至表”步骤。

图4-1 ETL 任务设计图示



#### 4. 示例前置条件

PostgreSQL 数据库中 vv\_1000 表已创建完成。

SeaSQL MPP 数据库中 vv\_1000000 表已创建完成。

#### 5. 示例详细步骤

##### 创建任务

进入[任务管理/任务列表]页面，单击<新增>按钮，新建任务。如图 4-2 所示，创建任务类型为：普通 ETL 任务，完成任务名称、任务描述的填写，单击<跳转任务设计页面>按钮，跳转至任务设计页面。

图4-2 新增任务图示



按照设计方案，拖拽组件步骤、建立连接。各步骤详细配置：

##### (1) 生成记录

该步骤为了配合 JS 代码组件，生成基于当前调度时间的抽取时间范围。

图4-3 生成记录具体配置图示

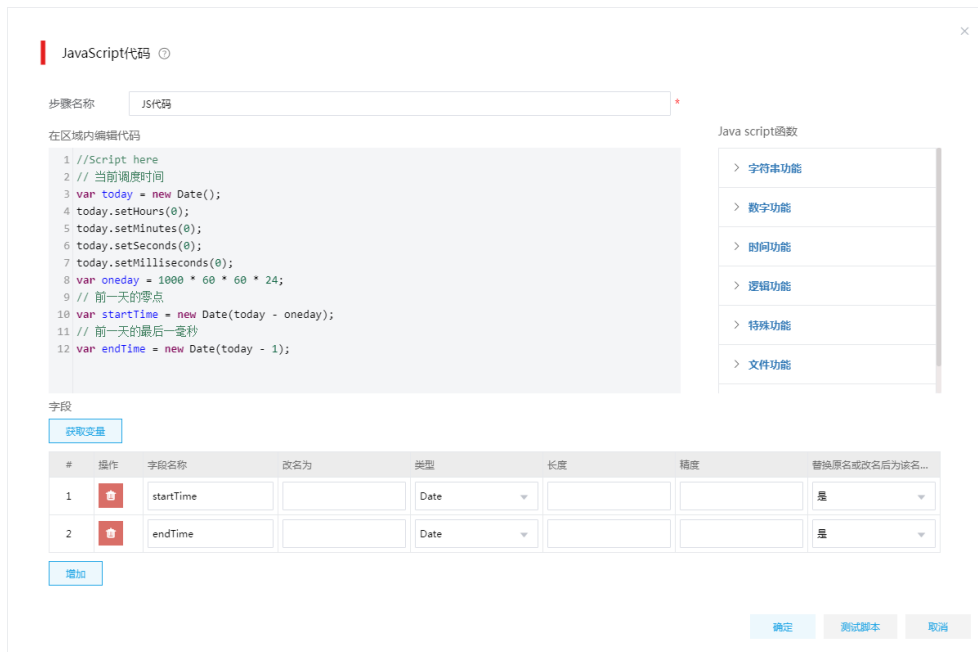


#	操作	名称	类型	格式	值	长度	精度
1	查	startTime	Date				
2	查	endTime	Date				

## (2) JS 代码

JS 代码步骤根据当前调度时间，取得一个入库的时间段。

图4-4 JS 代码具体配置图示



## (3) 表抽取

通过数据表抽取步骤，源表为 PostgreSQL 数据库中名为 vv\_1000 的表。在“从步骤插入数据”字段选择步骤二 JS 代码，SQL 语句如图 4-5 所示，SQL 中可以自动将数据表抽取步骤获取的变量按照顺序替换给 SQL 脚本中的“?”。

图4-5 数据表抽取具体配置图示



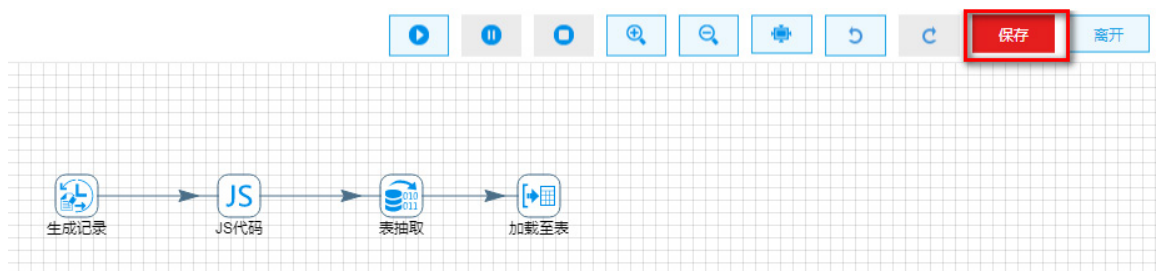
(4) 加载至数据库表

将最终的增量数据加载至 SeaSQL MPP 库中 vv\_1000000 表中。

图4-6 加载至数据库表具体配置图示

(5) 设计好任务后，单击<保存>按钮，添加至任务并退出。

图4-7 任务添加并退出



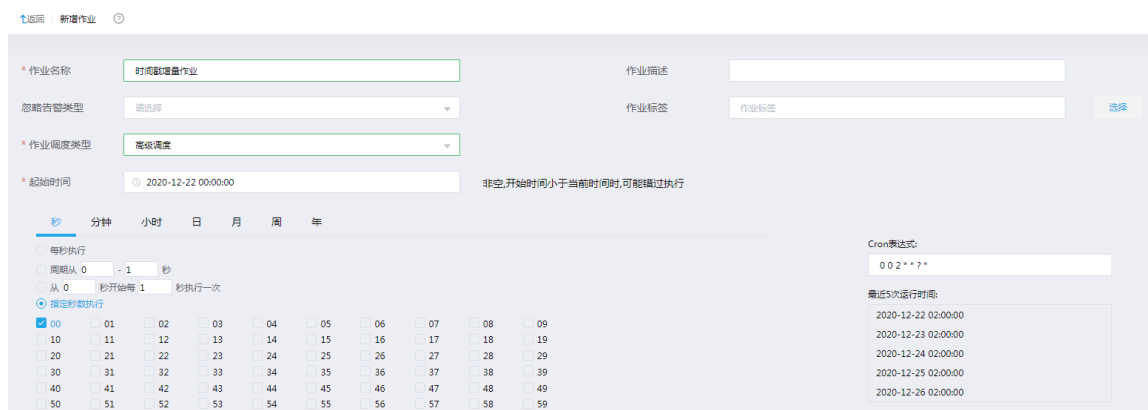
此时在[任务管理/任务列表]页面，可以在任务列表中查看此任务。

**新建作业**

进入[作业管理/作业列表]页面，单击<新增>按钮，新建作业，步骤如下：

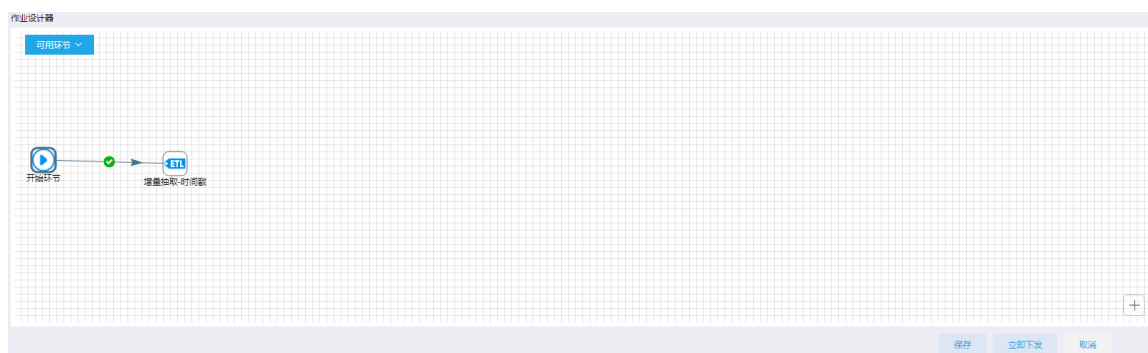
- (1) 如图 4-8 所示，在新增作业中配置成高级调度，实现每天的凌晨两点执行一次作业。
- 调度配置方法：对于 Cron 表达式实现定义时间规则为每天凌晨两点执行的具体配置，即指定每一分钟的第 0 秒，每小时的第 0 分钟，指定小时执行（勾选 02 前复选框）。日、月、周、年都按默认配置，即可将 Cron 表达式定义为每天凌晨两点执行。

图4-8 新建作业



- (2) 进行作业设计，然后进行立即下发或保存即可。

图4-9 作业设计



## 使用任务模板创建增量任务

- (1) 创建任务模板

任务模板是根据任务创建的。首先根据增量抽取任务创建任务模板，进入[任务管理/任务列表]页面，在“增量抽取-时间戳”任务右侧<更多操作>下拉框里，单击<创建模板>按钮，进入任务模板编辑页面。

在左侧填写参数区域增加参数，并分别双击画布上组件将组件内容替换为参数的替换值，参数及替换内如图 4-10 所示。

图4-10 任务模板编辑

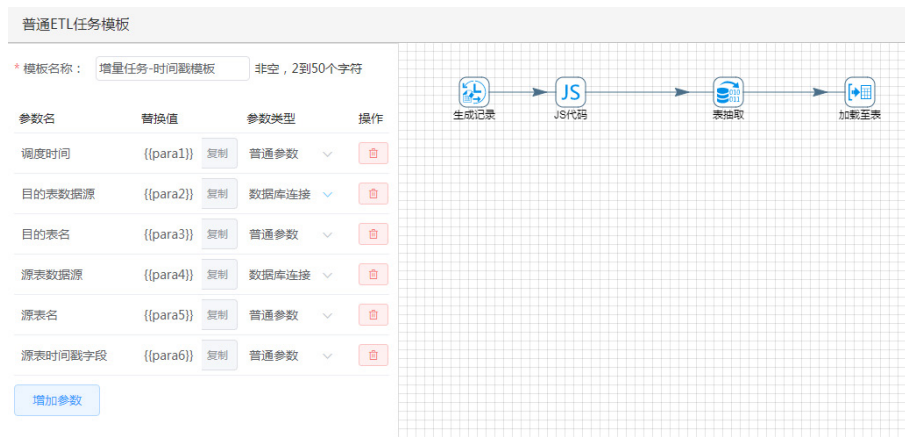


图4-11 JS 代码替换参数

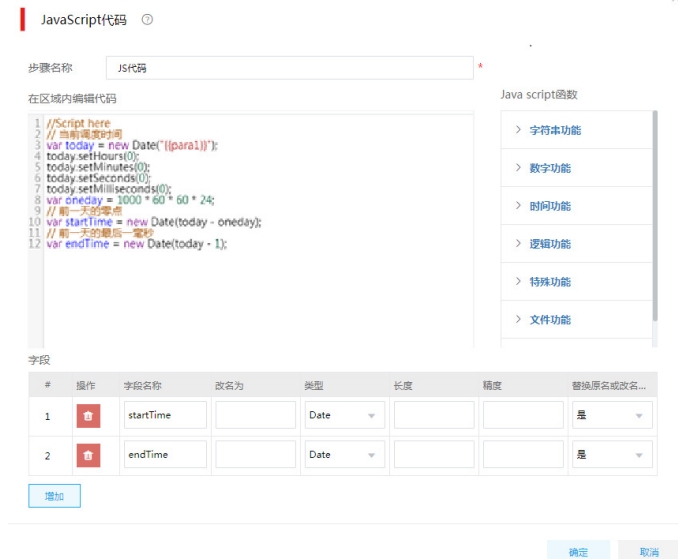


图4-12 表抽取替换参数

**数据表抽取** ⓘ

\* 步骤名称  非空，2到50个字符

数据库连接

SQL

```
1 SELECT *
2 FROM "{{para5}}"
3 WHERE "{{para6}}" >= ? AND "{{para6}}" <= ?
4 ORDER BY "{{para6}}"
```

允许简易转换

替换SQL语句里的变量

从步骤插入数据

执行每一行

记录数量限制

图4-13 加载至表替换参数

**加载至数据表** ⓘ

步骤名称  \*

数据库连接

目标模式

目标表

提交记录数量

清空表

忽略插入错误

指定数据库字段

表分区数据

分区字段

每个月分区数据

每天分区数据

使用批量插入

表名定义在一个字段里

包含表名的字段

最后，单击右上角<保存>按钮即可保存任务模板，系统将会跳转至任务模板列表页面。该任务模板的初始状态为草稿状态，还不能进行任务部署，单击改变模板状态为“使用”即可通过模板部署任务。

## (2) 模板部署任务

进入[任务管理/任务模板]页面，单击模板列表里模板右侧的<部署>按钮，弹出部署任务弹出框，配置如图 4-14 所示。

图4-14 任务模板创建任务参数

部署任务 ①

任务信息

\* 任务名称 增量任务-时间戳2

参数信息

调度时间 2020/03/10 02:00:00

目的表数据源 MPP1-72\_185 选择

目的表名 vw\_1000000

源表数据源 PG-74\_29 选择

源表名 vw\_1000

源表时间戳字段 RKSJ

确定 取消

(3) 然后单击<创建任务>按钮，即可在任务列表中生成新的任务。

图4-15 任务生成

任务名	任务描述	最近修改日期	任务类型	操作
增量任务-时间戳2	由任务模板【增量任务-时...	2020-03-11 13:50:32	普通ETL任务	编辑 运行 更多操作
增量抽取-时间戳	将PostgreSQL数据库中一...	2020-03-11 12:49:10	普通ETL任务	编辑 运行 更多操作

## 使用作业模板创建增量作业

### (1) 创建作业模板

创建作业模板与创建作业的规则相同。不同点：在作业模板设计器中，拖拽的控件不是任务，而是任务模板。这就要求进入[任务管理/任务模板]页面，首先根据“增量抽取-时间戳”任务创建任务模板，在任务模板列表里面修改模板状态为“使用”。

进入[作业管理/作业模板]页面，单击左上角<新增>按钮，在创建作业模板页面进行配置，图 4-16 为创建一个简单调度的类型的作业模板实例。配置完成后，单击<下一步>按钮，跳转至创建作业模板的参数归一页面。



图4-16 编辑简单调度作业模板实例

新增作业模板

---

\* 模板名称  作业调度类型

是否可重复

开始时间  非空,开始时间小于当前时间时,可能错过执行

作业模板设计器

任务模板列表

```
graph LR; A[开始环节] --> B[ ]; B --> C[增量任务-时间戳模板];
```

需要注意的是，在上一步编辑作业链中有几个作业模板单元，参数归一页面就有几个便签块，每个都分别是作业模板中的参数。有些时候任务模板中相同的变量却取了不同的名字，在这里可以替换成相同的名字，即可在作业模板部署成作业时减少填参数的工作量。这里不需修改，接受默认直接单击<保存>按钮，创建作业模板成功。

(2) 作业模板创建作业及任务

进入[作业管理/作业模板]页面，单击模板列表里模板右侧的<部署>按钮，弹出部署作业弹出框，配置如图 4-17 所示。

图4-17 作业模板创建作业参数

**部署作业** ?

**作业信息**

\* 作业名称

作业标签  选择

作业部署后是否下发  是  否

**参数信息**

源表名

源表数据源  选择

源表时间戳字段

目的表名

目的表数据源  选择

调度时间

**作业环节**

增量任务-时间...

确定 取消

然后单击<创建作业>按钮，即可在作业列表自动生成由模板创建的作业，在任务列表自动生成作业中的任务。

图4-18 由作业模板生成的作业

作业名	描述信息	最近修改日期	调度类型	执行器	状态	操作
增量作业-时间戳模板测试	由作业模板【增量...	2020-03-11 14:09:16	简单调度	10.121.72.77_37317	等待调度	<a href="#">编辑</a> <a href="#">停止</a> <a href="#">更多操作</a> <span>▼</span>

图4-19 由作业模板生成的任务

任务名	任务描述	最近修改日期	任务类型	操作
增量抽取-时间戳环节	由任务模板【增量任务-时间...	2020-03-11 14:09:16	普通ETL任务	<a href="#">编辑</a> <a href="#">运行</a> <a href="#">更多操作</a> <span>▼</span>

## 4.1.2 增量抽取场景-自增 ID

为了实现数据的增量抽取，除了增加时间戳字段外，还可以为数据表设置自增 ID，当一条数据入库时，无需传入 ID 字段，数据库自动获取当前表中最大 ID，然后将其加 1 作为当前记录的 ID。

### 1. 场景描述

XX 交通运输管理公司的业务系统数据库中，维护一张车辆通过卡口信息记录表 `kk_1000`。现公司需要将该表中新增的数据向数据仓库中的 `kk_1000000` 表做数据同步。公司业务系统使用 PostgreSQL 数据库，数据仓库使用 SeaSQL MPP。

表4-2 源表（kk\_1000）结构

ID	KKWZBM	HPHM	HPZL	TGSJ	HBSJ	SSSD	RYLX	CRSJ
1	3116250005	莆P3Z907	2	2014/4/8 23:36:58	26	68	01	2016/6/6 10:10:47

### 2. 场景分析

使用数据集成创建普通 ETL 任务，通过作业定时调度实现增量的抽取。

### 3. ETL 设计方案

数据流向：表抽取-> 表抽取 -> 加载至表

ETL 方案：[表 4-2](#) 中有一列名为 ID 的自增序列，判断需要被加载的表 ID 字段的最大数值，作为参数传入第二个数据抽取步骤，与数据库表抽取 2 中被抽取的表 ID 进行对比，数据库表抽取 2 中 ID 比传入参数大的数据传给“加载至表”步骤。

图4-20 ETL 任务设计图示



### 4. 示例前置条件

PostgreSQL 数据库中 `kk_1000` 表已创建完成，创建序列 `seq_id`，设置 ID 字段的默认值为：`next('seq_id')`

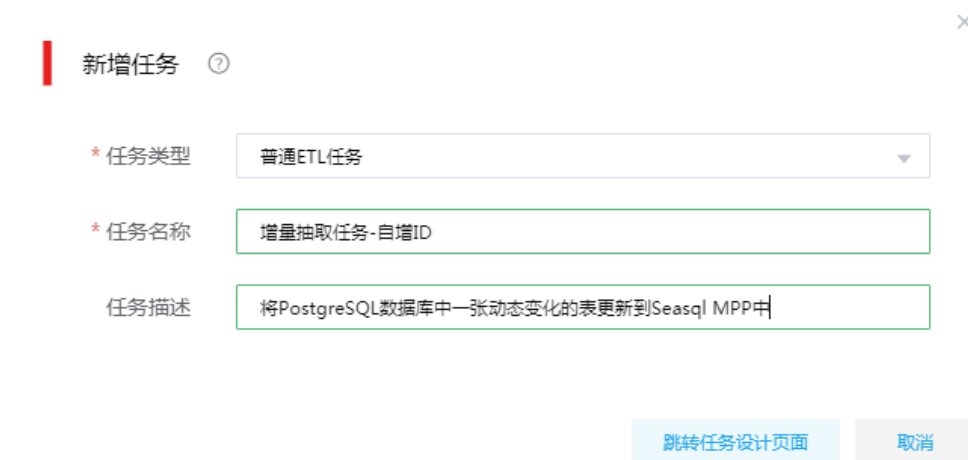
SeaSQL MPP 数据库中 `kk_1000000` 表已创建完成。

### 5. 示例详细步骤

#### 创建任务

进入[任务管理/任务列表]页面，单击<新增>按钮，新建任务。如[图 4-21](#)所示，创建任务类型为：普通 ETL 任务，完成任务名称、任务描述的填写，单击<跳转任务设计页面>按钮，跳转至任务设计页面。

图4-21 新增任务图示



新增任务 ⓘ

\* 任务类型 普通ETL任务

\* 任务名称 增量抽取任务-自增ID

任务描述 将PostgreSQL数据库中一张动态变化的表更新到Seasql MPP中

跳转任务设计页面 取消

按照设计方案，拖拽组件步骤、建立连接。各步骤详细配置：

(1) 数据表抽取

最终被插入数据的目标表为 kk\_1000000。在这一步骤中通过 SQL 语句查询出 kk\_1000000 表中的最大 id 作为参数传输至下一环节。

图4-22 数据表抽取具体配置图示



数据表抽取 ⓘ

步骤名称 表抽取 \* 非空，2到50个字符

数据库连接 MPP1-72\_185 选择 清除缓存 查询语句

SQL

```
1 SELECT MAX("ID")
2 FROM *kk_1000000*
3
```

允许简易转换

替换SQL语句里的变量

从步骤插入数据 指定步骤名

执行每一行

记录数量限制 0

确定 预览 取消

(2) 数据表抽取 2

通过数据表抽取步骤，源表为 PostgreSQL 数据库中名为 kk\_1000 的表。在“从步骤插入数据”字段选择步骤一数据表抽取，SQL 语句如[图 4-23](#)所示，SQL 中可以自动将数据表抽取步骤获取的变量按照顺序替换给 SQL 脚本中的“?”。

图4-23 数据表抽取 2 具体配置图示

数据表抽取 ①

步骤名称  \* 非空，2到50个字符

数据库连接

SQL

```
1 SELECT *
2 FROM "kk_1000"
3 WHERE "ID" > ?
4
```

允许简易转换

替换SQL语句里的变量

从步骤插入数据

执行每一行

记录数量限制

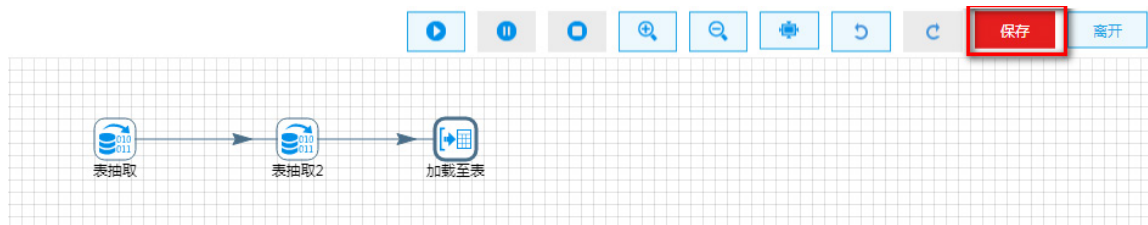
### (3) 加载至数据库表

将最终的增量数据加载至 SeaSQL MPP 库中 kk\_1000000 表中。

图4-24 加载至数据库表具体配置图示

(4) 设计好任务后，单击<保存>按钮，添加至任务并退出。

图4-25 任务添加并退出



(5) 此时在[任务管理/任务列表]页面，可以在任务列表中查看此任务。

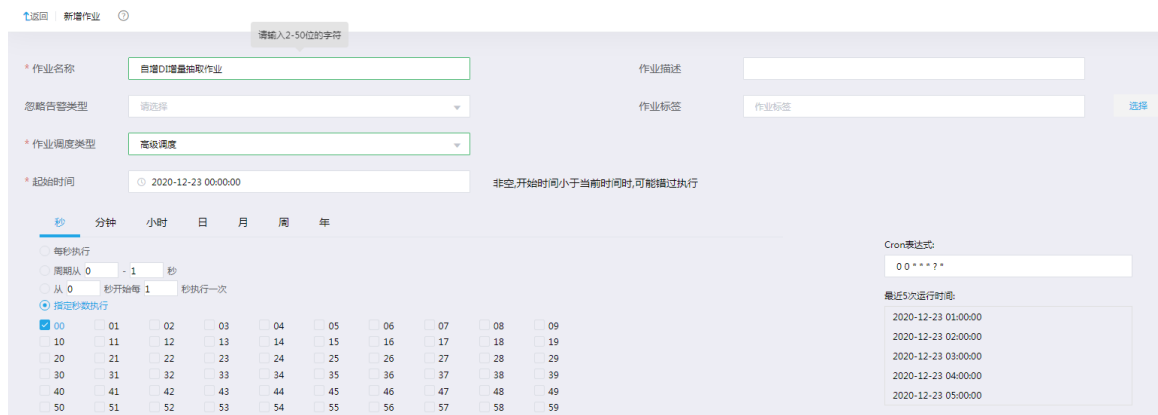
### 新建作业

进入[作业管理/作业列表]页面，单击<新增>按钮，新建作业，步骤如下：

(1) 如图 4-26 所示，在新增作业中配置成高级调度，实现每小时执行一次作业。

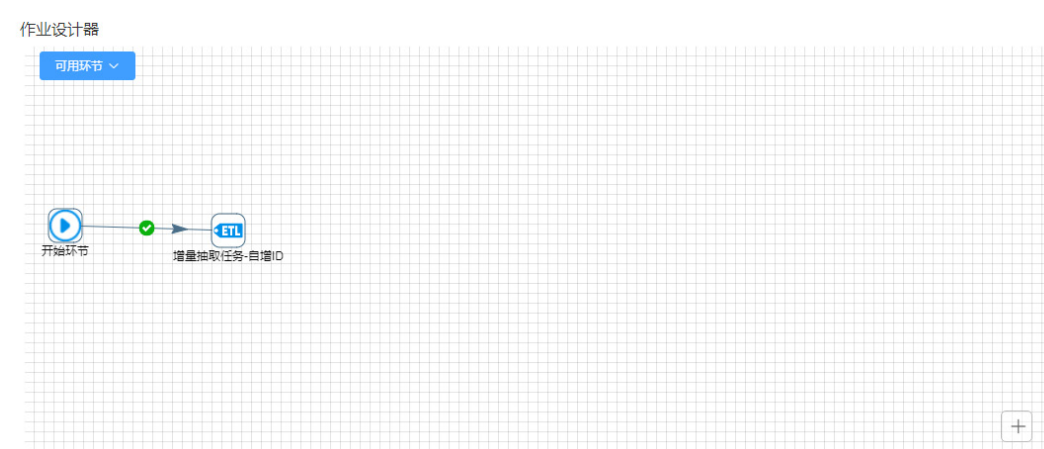
调度配置方法：对于 Cron 表达式实现定义时间规则为每小时执行的具体配置，即指定每一分钟的第 0 秒，每小时的第 0 分钟，每小时执行。日，月，周，年都按默认配置，即可将 Cron 表达式定义为每小时执行。

图4-26 新建作业



(2) 进行作业设计，然后进行立即下发或保存即可。

图4-27 作业设计



### 使用任务模板创建增量任务

#### (1) 创建任务模板

任务模板是根据任务创建的。首先根据增量抽取任务创建任务模板，进入[任务管理/任务列表]页面，在“增量抽取-自增 ID”任务右侧<更多操作>下拉框里，单击<创建模板>按钮，进入任务模板编辑页面。

在左侧填写参数区域增加参数，并分别双击画布上组件将组件内容替换为参数的替换值，参数及替换内如[图 4-28](#)所示。

图4-28 任务模板编辑

普通ETL任务模板

\* 模板名称:  非空, 2到50个字符

参数名	替换值	参数类型	操作
目的表数据源	{{para1}}	复制 数据库连接	
目的表名	{{para2}}	复制 普通参数	
源表数据源	{{para3}}	复制 数据库连接	
源表名	{{para4}}	复制 普通参数	
目的增量字段	{{para5}}	复制 普通参数	
源表增量字段	{{para6}}	复制 普通参数	

图4-29 表抽取替换参数

**数据表抽取**

步骤名称:  \* 非空, 2到50个字符

数据库连接:

SQL

```

1 SELECT MAX("{{para5}}")
2 FROM "{{para1}}"
3

```

允许简易转换

替换SQL语句里的变量

从步骤插入数据

执行每一行

记录数量限制



图4-30 表抽取 2 替换参数

**数据表抽取** ⓘ

步骤名称  \* 非空，2到50个字符

数据库连接

SQL

```
1 SELECT *  
2 FROM "{{para4}}"  
3 WHERE "{{para6}}" > ?  
4
```

允许简易转换

替换SQL语句里的变量

从步骤插入数据

执行每一行

记录数量限制

图4-31 加载至表替换参数

加载至数据表 ?

步骤名称  \* 非空, 2到50个字符

数据库连接

目标模式

目标表

提交记录数量

清空表

忽略插入错误

指定数据库字段

主选项 数据库字段

表分区数据

分区字段

每个月分区数据

每天分区数据

使用批量插入

表名定义在一个字段里

确定 取消

最后单击右上角<保存>按钮即可保存任务模板，系统将会跳转至任务模板列表页面。该任务模板的初始状态为草稿状态，还不能进行任务部署，单击改变模板状态为“使用”即可通过模板部署任务。

## (2) 模板部署任务

进入[任务管理/任务模板]页面，单击模板列表里模板右侧的<部署>按钮，弹出部署任务弹出框，配置如图 4-32 所示。

图4-32 任务模板创建任务参数

部署任务 ?

任务信息

\* 任务名称

参数信息

目的表数据源  选择

目的表名

源表数据源  选择

源表名

目的增量字段

源表增量字段

确定 取消

然后单击<创建任务>按钮，即可在任务列表中生成新的任务。

图4-33 任务生成

任务名	任务描述	最近修改日期	任务类型	操作
增量任务-自增ID2	由任务模板【增量任务-自增L...	2020-03-16 01:09:33	普通ETL任务	<a href="#">编辑</a> <a href="#">运行</a> <a href="#">更多操作</a> <span>▼</span>
增量抽取任务-自增ID	将PostgreSQL数据库中一张...	2020-03-16 00:43:44	普通ETL任务	<a href="#">编辑</a> <a href="#">运行</a> <a href="#">更多操作</a> <span>▼</span>

## 使用作业模板创建增量作业

### (1) 创建作业模板

创建作业模板与创建作业的规则相同。不同点：在作业模板设计器中，拖拽的控件不是任务，而是任务模板。这就要求进入[任务管理/任务模板]页面，首先根据“增量抽取-时间戳”任务创建任务模板，在任务模板列表里面修改模板状态为“使用”。

进入[作业管理/作业模板]页面，单击左上角<新增>按钮，在创建作业模板页面进行配置，[图 4-34](#)为创建一个简单调度的类型的作业模板实例。配置完成后，单击<下一步>按钮，跳转至创建作业模板的参数归一页面。

图4-34 编辑简单调度作业模板实例

新增作业模板

\* 模板名称 增量作业-自增ID模板 作业调度类型 简单调度

是否可重复 是 重复次数 -1 -1即无限重复

重复间隔 1 时 0 分 0 秒 若间隔太短,可能导致调度问题,建议不小于5秒

开始时间 2020-03-16 00:00:00 非空,开始时间小于当前时间时,可能错过执行

结束时间 2020-03-31 00:00:00 调度次数和结束时间同时存在,以两者最近时间为准

作业模板设计器

任务模板列表

开始环节 → 增量任务-自增ID模板

需要注意的是,在上一步编辑作业链中有几个作业模板单元,参数归一页面就有几个便签块,每个都分别是作业模板中的参数。有些时候任务模板中相同的变量却取了不同的名字,在这里可以替换成相同的名字,即可在作业模板部署成作业时减少填参数的工作量。这里不需修改,接受默认直接单击<保存>按钮,创建作业模板成功。

(2) 作业模板创建作业及任务

进入[作业管理/作业模板]页面,单击模板列表里模板右侧的<部署>按钮,弹出部署作业弹出框,配置如图 4-35 所示。

图4-35 作业模板创建作业参数

**部署作业** ?

**作业信息**

\* 作业名称

作业标签  选择

作业部署后是否下发  是  否

**参数信息**

源表名

源表增量字段

源表数据源  选择

目的增量字段

目的表名

目的表数据源  选择

**作业环节**

增量任务-自增...

确定 取消

然后单击<创建作业>按钮，即可在作业列表自动生成由模板创建的作业，在任务列表自动生成作业中的任务。

图4-36 由作业模板生成的作业

作业名	描述信息	最近修改日期	调度类型	执行器	状态	操作
增量作业-自增ID模板测试	由作业模板【增量...	2020-03-16 01:24:58	简单调度	10.121.72.77_37317	等待调度	<a href="#">编辑</a> <a href="#">停止</a> <a href="#">更多操作</a>

图4-37 由作业模板生成的任务

任务名	任务描述	最近修改日期	任务类型	操作
增量抽取-自增ID环节	由任务模板【增量任务-自增ID...	2020-03-16 01:24:58	普通ETL任务	<a href="#">编辑</a> <a href="#">运行</a> <a href="#">更多操作</a>

### 4.1.3 数据转换场景

半结构与结构化之间的数据转换，支持半结构化（xml/JSON）数据转换为结构化数据，也支持结构化数据转换为半结构化数据，比如：JSON 数据采集并解析写入关系型数据库，或者关系型数据库抽取转为 JSON 格式文件等。

#### 1. 场景描述

XX 公司对接客户 REST 接口，数据交换格式为 JSON，XX 公司需将对接的数据解析并写入关系型数据库。

- 接口示例：
  - REST: `http://10.121.56.134:6688/getinfo/1`
  - 请求方式: GET
- 采集 JSON 数据示例：

```
[  
  {"sid": "100001", "sname": "fox"},  
  {"sid": "100002", "sname": "tom"},  
  {"sid": "100003", "sname": "fix"},  
  {"sid": "100004", "sname": "bug"},  
  {"sid": "100005", "sname": "dog"}  
]
```

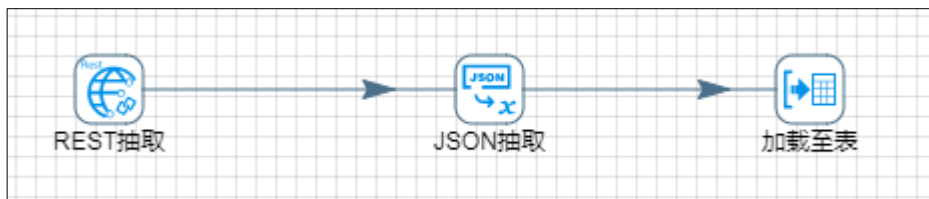
#### 2. 场景分析

数据集成已支持 REST 接口数据采集，可采用 JSON 抽取组件实现数据转换（将半结构化数据转换为结构化数据），并通过加载至表将数据写入数据库。

#### 3. ETL 设计方案

- 数据流向：REST 抽取组件—>JSON 抽取组件—>加载至表组件。ETL 图示如[图 4-38](#)。

图4-38 ETL 图示



- ETL 方案：REST 抽取采集 REST 接口数据并输出，JSON 抽取将输入的 JSON 格式数据解析为字段然后输出，加载至表将数据加载至数据库。

#### 4. 示例前置条件

- (1) 模拟一个 REST 接口，接口返回数据参考 JSON 数据示例。
- (2) MySQL 数据源，名称：mysql-gzh，目标表：test\_stu，SQL 脚本：

```
CREATE TABLE `test_stu` (  
  `sid` varchar(18) DEFAULT NULL,
```

```
`sname` varchar(50) DEFAULT NULL  
) ENGINE=InnoDB DEFAULT CHARSET=utf8
```

## 5. 示例详细步骤

### (1) 第一步：REST 抽取组件配置

“通用”页配置如[图 4-39](#)。

- URL: `http://10.121.56.134:6688/getinfo/1`
- HTTP 方法: 下拉选中 GET

图4-39 “通用页”图示

REST抽取 ⓘ

步骤名称  \*

通用	Query参数	Body参数	Matrix参数	HTTP头部	HTTP认证	SSL	输出字段
----	---------	--------	----------	--------	--------	-----	------

URL

从字段中获取URL

URL字段

HTTP方法

从字段中获取HTTP方法

HTTP方法字段

内容类型

“输出字段”配置如[图 4-40](#)。

- 响应参数字段名称: 此字段接收接口返回数据。

图4-40 “输出字段” 页图示

REST抽取 ?

步骤名称  \*

通用	Query参数	Body参数	Matrix参数	HTTP头部	HTTP认证	SSL	输出字段
----	---------	--------	----------	--------	--------	-----	------

响应参数字段名称

响应状态码字段名称

响应时间字段名称

响应头部字段名称

(2) 第二步：JSON 抽取组件配置

“文件” 页配置如[图 4-41](#)。

- 从前面的步骤获取源：勾选。
- 前一步骤名：下拉选中“REST 抽取”。
- 保存源的字段：下拉选择“result”，注意需要与上一步骤“REST 抽取”输出字段的“响应参数字段名”保持一致。



图4-41 “文件”页图示

JSON抽取 ?

步骤名称  \*

**文件** | 内容 | 字段 | 其它输出字段

本地文件或目录

正则表达式

正则表达式(排除)

选中的文件

#	操作	文件/目录	通配符	通配符(排除)	要求	包含子目录
	从前面的步骤获取源	<input checked="" type="checkbox"/>				
	前一步骤名	<input type="text" value="REST抽取"/>				
	保存源的字段	<input type="text" value="result"/>				
	源是一个文件名	<input type="checkbox"/>				
	源是一个URL	<input type="checkbox"/>				
	从结果中移除源字段	<input type="checkbox"/>				

“字段”页配置如[图 4-42](#)。

- 字段名称：配置示例 JSON 的字段名称。
- 路径：参考 JSONPath 填写，填写对应要读取的 JSONPath。
- 类型：下拉选中 String。

图4-42 “字段”页图示

JSON抽取 ⓘ

步骤名称  \*

文件 内容 字段 其它输出字段

获取字段

#	操作	名称	路径	类型	格式	长度
1		<input type="text" value="sid"/>	<input type="text" value="\$.*.sid"/>	String	<input type="text"/>	<input type="text"/>
2		<input type="text" value="sname"/>	<input type="text" value="\$.*.sname"/>	String	<input type="text"/>	<input type="text"/>

(3) 第三步：加载至表组件配置

- 数据库连接：单击“选择”，然后选中“mysql-gzh”（选择要写入的数据库）。
- 目标表：单击“选择”，然后选中“test\_stu”（选择要写入的表）
- 指定数据库字段：勾选。

图4-43 “主选项” 图示

加载至数据表 ?

\* 步骤名称  非空，2到50个字符

数据库连接

目标模式

目标表

提交记录数量

清空表

忽略插入错误

指定数据库字段

表分区数据

分区字段

每个月分区数据

每天分区数据

使用批量插入

表名定义在一个字段里

包含表名的字段

存储表名字段

“数据库字段” 图示：

- 通过获取字段，配置“表字段”和“流字段”的映射关系。“表字段”是 test\_stu 表的字段名称，“流字段”是数据流中的字段。

图4-44 数据库字段图示

指定数据库字段

主选项 数据库字段

获取字段 输入字段映射

#	操作	表字段	流字段
1		sid	sid
2		sname	sname

增加

确定 SQL 取消

#### 4.1.4 实时流场景

与消息中间件通信，支持向 Kafka、ActiveMQ、IBM MQ 等消息队列发送数据，也支持从这些消息队列中抽取数据，比如：订阅 ActiveMQ 中某个主题，实时拉取主题中发布的数据，并经过解析处理后写入关系型数据库。

##### 1. 场景描述

XX 物联网平台接入很多传感器设备，这些设备使用 MQTT 协议与物联网平台通信，传递数据为 JSON 文档，主要内容是各设备上报的自身状态信息和各类传感器采集到的环境数据，平台使用 ActiveMQ 作为 MQTT 服务端来接收这些数据。现要求实时抽取消息队列中的数据，经过简单转换处理之后写入关系型数据库。

传感器上报 JSON 数据示例：

```
{
  "deviceid": "100001",
  "devicetype": "2",
  "status": "ok",
  "data": "25,42",
  "data": "2020/03/11 16:00:00"
}
```

##### 2. 场景分析

数据集成目前已经支持从消息队列中实时采集数据，可使用 MQTTConsumer 组件实时抽取数据，通过设置组件的批量参数，实现定时/定量触发子转换。子转换中对 MQTTConsumer 抽取过来的一批数据使用 JSON 抽取组件进行解析，并通过加载至表将数据写入数据库。

##### 3. ETL 设计方案

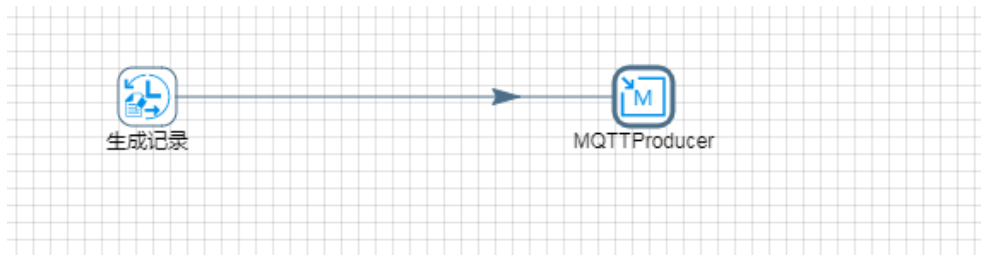
- 数据流向：MQTTConsumer 组件—>从流中获取记录—>JSON 抽取组件—>加载至表组件。

- **ETL 方案：**这里需要设计两个 ETL 任务，我们分别称之为主任务、子任务。主任务中主要使用 MQTTConsumer 组件，实时采集 ActiveMQ 中某个主题的数据；子任务中负责解析入库。原则是：先创建子任务，再创建主任务，然后将子任务配置到主任务的 MQTTConsumer 中。
- **运行逻辑：**子任务设计完成后，不能单独执行，只作为 MQTTConsumer 组件的一项配置。运行主任务时，根据设置的批量参数，周期性触发配置的子任务，同时将一批数据传进子任务中，然后 JSON 抽取组件将输入的 JSON 格式数据解析为字段，最后加载至表将数据加载至数据库。  
 注意事项：子任务必须以“从流中获取记录”组件为开始节点，该组件专门用于接收从主任务中 MQTTConsumer 组件传递进来的数据。

#### 4. 示例前置条件

为了便于实时采集任务设计调试，需要模拟传感器上报数据到 ActiveMQ 的过程。

图4-45 “模拟上报”任务设计图示



使用生成记录组件，不停地生成一个 JSON 字符串，JSON 字符串具体内容参见场景描述中的数据示例。

图4-46 生成记录具体配置图示

生成记录 ①

\* 步骤名称  非空，2到50个字符

限制

从不停止生成

间隔毫秒数(延迟)

当前行时间字段名称

以前行时间字段名称

字段

	操作	名称	类型	格式	值	长度	精度
1		message	String		{"deviceid": "1000C"}		

MQTTProducer 组件使用 MQTT 协议将生成的 JSON 串不断发送到 ActiveMQ 的“device\_info”主题中。

图4-47 MQTTProducer 基本设置

MQTT Producer ⓘ

非空，2到50个字符

步骤名称  \*

设置
  安全
  选项

连接

客户端ID

主题  指定主题  从字段获取主题名

主题名

服务质量

消息字段

图4-48 MQTTProducer 安全设置

MQTT Producer ⓘ

步骤名称  \*

设置
  安全
  选项

身份认证

用户名

密码

使用安全协议

	操作	名称	值
4	<input type="button" value="删除"/>	ssl.keyStore	E:\key\client1.ks
5	<input type="button" value="删除"/>	ssl.keyStorePassword	passwd
6	<input type="button" value="删除"/>	ssl.keyStoreProvider	
7	<input type="button" value="删除"/>	ssl.keyStoreType	JKS
8	<input type="button" value="删除"/>	ssl.protocol	TLS
9	<input type="button" value="删除"/>	ssl.trustManager	
10	<input type="button" value="删除"/>	ssl.trustStore	E:\key\client1.ts
11	<input type="button" value="删除"/>	ssl.trustStorePassword	passwd
12	<input type="button" value="删除"/>	ssl.trustStoreProvider	

MQTTProducer 的安全配置包含两部分内容：身份认证、SSL 认证，各配置项的值取决于 ActiveMQ 中 broker 的设置。具体而言：

- 若 broker 设置为不允许匿名访问，则需要在身份认证部分填写创建 broker 时设置的用户名、密码。
- 若 broker 开启了 SSL 认证，这里需要勾选“使用安全协议”，然后如实填写图 4-48 所示的几个参数。

完成上述配置工作后，保存并运行“模拟上报”任务。此时，不断发送数据到 device\_info 主题中。

## 5. 示例详细步骤

准备工作完成后。现在按照设计的实时采集方案来创建任务。

图4-49 子任务 ETL 图示

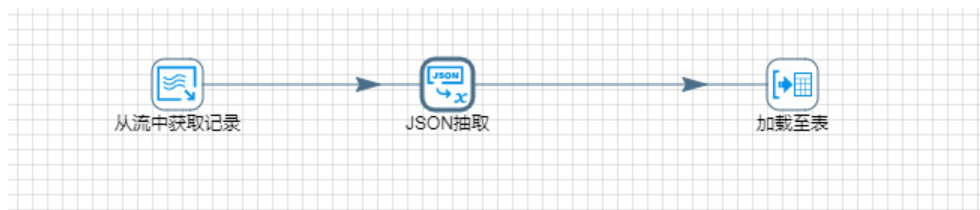
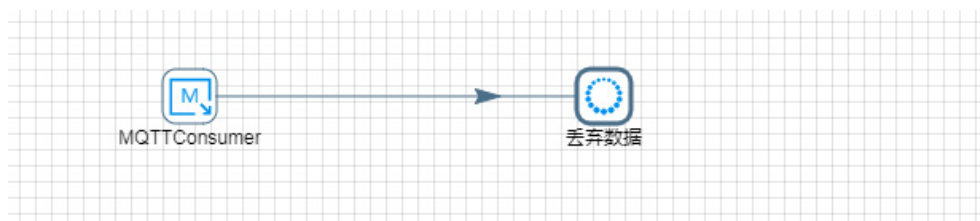


图4-50 主任务 ETL 图示



子任务中各组件配置说明：

图4-51 从流中获取记录具体配置图示

从流中获取记录 ⓘ

步骤名称  \*

字段

#	操作	字段名称	类型	长度	精度
1		<input type="text" value="Message"/>	<input type="text" value="String"/>	<input type="text"/>	<input type="text"/>

这里只需要增加一个 Message 字段，该字段名称对应主任务中 MQTTConsumer 组件的输出字段。如果后面步骤需要用到消息所属主题，这里还可以再增加一个 Topic 字段。

图4-52 JSON 抽取配置图示-文件页签

JSON抽取 ⓘ

非空，2到50个字符

步骤名称: JSON抽取 \*

文件 | 内容 | 字段 | 其它输出字段

本地文件或目录: /opt/file.json [增加] [浏览]

正则表达式: [ ]

正则表达式(排除): [ ]

选中的文件

#	操作	文件/目录	通配符	通配符 (排除)	要求	包含子目录
从前面的步骤获取源 <input checked="" type="checkbox"/>						
前一步骤名		从流中获取记录				
保存源的字段		Message				
源是一个文件名		<input type="checkbox"/>				
源是一个URL		<input type="checkbox"/>				
从结果中移除源字段		<input checked="" type="checkbox"/>				

[确定] [预览] [显示文件名] [取消]

图4-53 JSON 抽取配置图示-字段页签

JSON抽取 ⓘ

\* 步骤名称: JSON抽取 非空，2到50个字符

文件 | 内容 | 字段 | 其它输出字段

获取字段

	操作	名称	路径	类型	格式	长度	精度
1	<input type="checkbox"/>	deviceid	\$.deviceid	String			
2	<input type="checkbox"/>	devicetype	\$.devicetype	String			
3	<input type="checkbox"/>	health	\$.health	String			
4	<input type="checkbox"/>	data	\$.data	String			
5	<input type="checkbox"/>	date	\$.date	Date	yyyy/MM/dd HH:mm:ss		

[增加]

[确定] [预览] [显示文件名] [取消]

JSON 抽取组件，选择从前面步骤获取数据，这里的前一步骤也就是“从流中获取记录”组件。在字段页签下需要手动添加字段，这就需要预先了解采集的 JSON 数据的结构，重点关注 **date** 字段的配置，因为 **date** 字段在 JSON 文档中是一个字符串，比如：“2020/03/11 00:00:00”，所以这里选择使用 **Date** 类型去解析这个 JSON 属性情况下，需要指定与其匹配的格式。



图4-54 加载至表具体配置图示

加载至数据表 ?

步骤名称  \*

数据库连接  选择

目标模式

目标表  选择

提交记录数量

清空表

忽略插入错误

指定数据库字段

主选项 数据库字段

表分区数据

分区字段

每个月分区数据

每天分区数据

确定 SQL 取消

完成上述配置工作后，保存任务并退出设计器。在任务列表中。  
主任务中各组件配置说明：

图4-55 MQTTConsumer 基本设置

MQTTConsumer ⓘ

\* 步骤名称  非空，2到50个字符

转换

配置

连接

主题

	操作	名称
1	<input type="button" value="删除"/>	device_info

服务质量  ▾

单击<浏览>按钮，打开任务列表，选择之前创建的子任务，配置连接地址，添加一个主题。这里的连接地址和主题使用与模拟上报数据任务一样的配置。

图4-56 MQTTConsumer 安全设置

MQTTConsumer ?

\* 步骤名称  非空，2到50个字符

转换

**配置** 安全 批量 字段 结果字段 选项

身份认证

用户名

密码

使用安全协议

	操作	名称	值
4	<input type="button" value="🗑️"/>	ssl.keyStore	E:\key\client1.ks
5	<input type="button" value="🗑️"/>	ssl.keyStorePassword	passwd
6	<input type="button" value="🗑️"/>	ssl.keyStoreProvider	
7	<input type="button" value="🗑️"/>	ssl.keyStoreType	JKS
8	<input type="button" value="🗑️"/>	ssl.protocol	TLS
9	<input type="button" value="🗑️"/>	ssl.trustManager	
10	<input type="button" value="🗑️"/>	ssl.trustStore	E:\key\client1.ts
11	<input type="button" value="🗑️"/>	ssl.trustStorePassword	passwd

安全页签下的配置与 MQTTProducer 保持一致，填写说明参见模拟上报数据任务。

图4-57 MQTTConsumer 批量设置

MQTTConsumer ?

\* 步骤名称  非空，2到50个字符

转换

**配置** 安全 批量 字段 结果字段 选项

持续时间  必填，区间[0,2^53-1]

记录数量  必填，区间[0,2^53-1]

批量页签下的两个配置项：持续时间、记录数量是触发子任务的两个条件。

- 持续时间：每隔多长时间触发一次子任务，单位毫秒。
- 记录数量：每采集到多少条数据触发一次子任务。

两个条件任意一个满足就会触发子任务，每次触发子任务，就会重新计时、计数。

至此，实时流数据采集任务设计全部结束。同时运行“模拟上报”任务、主任务，观察数据库中目标表数据量变化，发现有数据被周期性地写入。

#### 4.1.5 大数据组件场景

对于 HDFS、Hive、HBase 等大数据组件，数据集成也做了适配。通过简单的拖拉拽，既可实现大数据组件中数据与关系数据库或者文本数据的相互转换。

##### 1. 场景描述

XX 公司存在多个孤立的业务系统，数据分散在各个业务数据库中。这样就导致了存在重复数据，并且数据核心价值没有很好被利用。现公司要求，将分散的数据统一集中，作为数据仓库。该公司数据主要分散在文本文件、关系型数据库（例如 MySQL）以及部分消息系统中，数据仓库建设以 Hive 为主，分层存储；ODS 与源数据保持同步，采用 Hive 外部表存储；DW 存储经过沉淀积累下来的数据，使用 Hive 内部表实现。

数据流向：各类数据源--> ODS 层—> DW 层。

##### 2. 场景分析

分析业务需求，针对 XX 公司的数据量，基本分为两种采集方式，全量采集和增量采集。以 100W 表中存量数据为分界，大于 100W 数据采用增量采集方式，小于则采用全量采集方式。

现仅以 dig(数据库名称 dig，是 MySQL)数据库为例，全量采集 behavior 表中。

注：增量采集思路，可参考其他增量采集思想，此处不做重复介绍。

##### 3. ETL 设计方案

全量采集作业全流程：开始环节—>采集 behavior 数据到 ODS 层—>同步 ODS 数据到 DW。

图4-58 设计方案图示



图中，“mysql 数据采集”是一个任务，将 MySQL 中数据抽取到 ODS 层指定的 hdfs 路径下，“加载至 hive\_sql”也是一个任务，通过 insert overwrite 将数据同步到 DW 内部表中。

##### 4. 示例前置条件

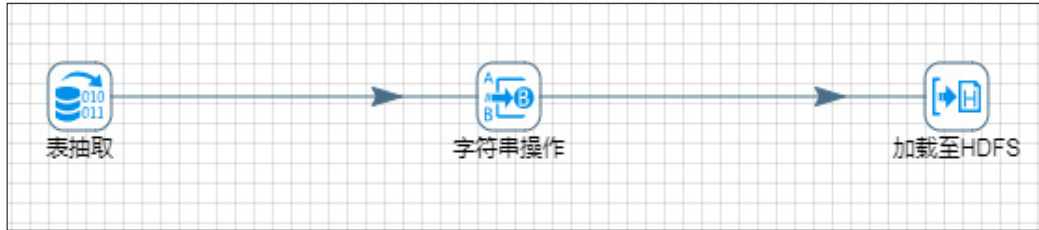
- ODS 层 ods\_behavior 表已创建完成，location 路径 “/usr/hdfs/gzh/behavior/”。
- DW 层 dw\_behavior 表已创建完成。
- MySQL 数据源 dig，采集表 behavior，正常有数据。

## 5. 示例详细步骤

### 任务一：MySQL 数据采集

ETL 流程：MySQL 表数据抽取—> 去除空格和换行字符 --> 数据加载至 HDFS 中。

图4-59 MySQL 数据采集图示



(1) 第一步：表抽取配置，数据源选择 dig，查询 SQL 如图 4-60。

图4-60 数据表抽取图示

数据表抽取 ⓘ

\* 步骤名称  非空，2到50个字符

数据库连接

SQL

```
1 SELECT
2   "uname"
3   ,"cdate"
4   ,"ip"
5   ,"destination"
6   ,"remark"
7   ,"tenant_id"
8 FROM "di_behavior"
9
```

允许简易转换

替换SQL语句里的变量

从步骤插入数据

执行每一行

记录数量限制

(2) 第二步：字符串操作，获取字段并批量设置去除空格，清洗数据两端多余空格。“去除空格”批量设置 both，清洗数据两端空格。

图4-61 字符串操作图示

字符串操作 ⓘ

\* 步骤名称  非空，2到50个字符

字段

#	操作	输入字段	输出字段	去除空格	小写/大写	补位方式	补位字符	补位长度	首字母大写
1	<input type="checkbox"/>	uname		both	none	none			none
2	<input type="checkbox"/>	ip		both	none	none			none
3	<input type="checkbox"/>	destinatic		both	none	none			none
4	<input type="checkbox"/>	remark ▼		both	none	none			none
5	<input type="checkbox"/>	tenant_id		both	none	none			none

- (3) 第三步：加载至 HDFS 文件，浏览配置 HDFS 数据源及 HDFS 路径。
- Folder/File: /usr/hdfs/gzh/behiver/data。
  - 指定字段：勾选。
  - 创建父目录：勾选。

图4-62 加载至 HDFS 文件图示

加载至HDFS文件 ⓘ

\* 步骤名称  非空，2到50个字符

文件 内容 字段

HDFS连接

Hadoop集群

Folder/File

登陆用户

指定写入HDFS用户组

指定字段

创建父目录

从字段中获取文件名

文件名字段

扩展名

文件名里包含日期

文件名里包含时间

(4) 加载至 HDFS 文件“内容”页签配置：

- 分隔符：使用竖线“|”，如图。
- 格式：下拉选择 **Unix**。
- 编码方式：下拉选择 **UTF-8**。

图4-63 HDFS 文件“内容”页签配置图示

加载至HDFS文件 ?

非空，2到50个字符

步骤名称  \*

文件 | **内容** | 字段

追加方式

文件不存时自动创建

分隔符  插入Tab

封闭符

强制在字段周围加封闭符

头部

尾部

格式

编码方式

确定 取消

(5) “字段”页签配置：获取字段，并去除不需要的字段即可。

图4-64 “字段”页签配置图示

加载至HDFS文件 ? ×

步骤名称  \*

文件 | 内容 | **字段**

获取字段

#	操作	名称	类型	格式	长度	精度	货币	小数	分组
1	<span>🗑️</span>	<input type="text" value="uname"/>	String	请选择	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>
2	<span>🗑️</span>	<input type="text" value="ip"/>	String	请选择	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>
3	<span>🗑️</span>	<input type="text" value="remark"/>	Integer	请选择	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>
4	<span>🗑️</span>	<input type="text" value="tenant_id"/>	String	请选择	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>
5	<span>🗑️</span>	<input type="text" value="destination"/>	String	请选择	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>	<input type="text"/>

增加

确定 取消



## 任务二：ODS 数据同步到 DW

执行 SQL，将数据从 Hive 外部表加载至内部表中。

图4-65 执行 SQL 图示



“执行 SQL”配置：

- 数据库连接：选择 Hive2-1。
- SQL 语句：采用 insert overwrite 方式同步数据，如[图 4-66](#)所示。

图4-66 执行 SQL 配置图示

执行SQL ⓘ

\* 步骤名称  非空，2到50个字符

数据库连接

SQL语句（“;”分割语句，“?”替代参数）

```
1 insert overwrite table dw.dw_behiver (select * from ods.ods_behiver);
```

为每一行执行SQL

作为一个语句执行

变量替换

参数：

操作	作为参数的字段
----	---------

## 4.1.6 数据清洗场景

### 1. 场景描述

数据集成支持文本、数据库、大数据组件等多种数据源，在数据集成过程中，经常会需要将数据中不符合规则的数据进行清洗或转换，格式归一后存储到目标仓库中。

### 2. 场景分析

XX 公司需要将数据库中一张表数据，按照规定的格式加载至远程文件中。该表是一个车辆信息表，要求去除车牌号不符合格式的数据，另外将人员类型（RYLX）为空的默认设置为 2。

图4-67 示例数据源（test 表）

	id	kkwzbm	hphm	hpzl	rylx	tgsh
1	1	311,001	豫A11111	1	2	1990-01-01 12:12:12
2	2	311,002	xyz	1	2	2010-01-01 12:12:12
3	3	311,001	豫A11111	1	[NULL]	1990-01-01 12:12:12

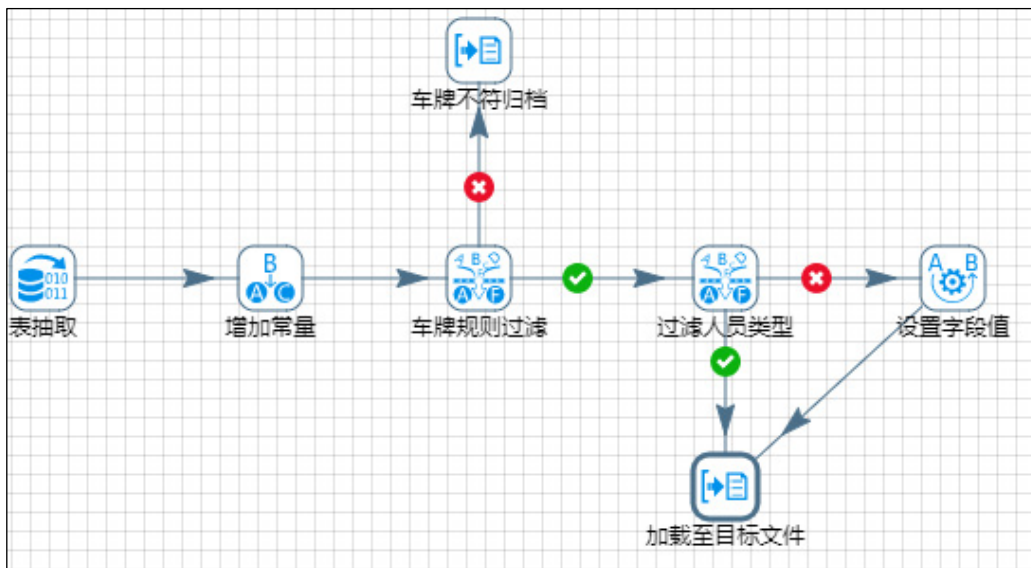
测试数据源说明：

- hphm: 是车牌号字段，其中可能出现未知或不符合车牌号码的车牌号。要求将此字段数据进行清洗，去除不符合规则的纪录。
- rylx: 是人员类型字段，可能为空，需要对此字段进行处理，将空值转换为 2。

### 3. ETL 设计方案

使用表抽取组件进行数据抽取，通过过滤记录组件中正则表达式将 hphm 中不符合要求的数据过滤掉，加载至文本文件归档。正确的数据再次过滤 rylx 为空的数据，为空的数据传递设置字段值，将空值赋值为 2，最后将据全部加载至步骤加载至目标文件。

图4-68 任务图示



#### 4. 示例前置条件

无

#### 5. 示例详细步骤

(1) 第一步：表抽取配置，配置要采集的数据源及 SQL 语句。

图4-69 表抽取配置图示

数据表抽取 ⓘ

\* 步骤名称  非空，2到50个字符

数据库连接

SQL

```
1 SELECT
2 "id"
3 , "kkwzbm"
4 , "hphm"
5 , "hpzl"
6 , "rylx"
7 , "tgsj"
8 FROM "test"
9
```

允许简易转换

替换SQL语句里的变量

从步骤插入数据

执行每一行

记录数量限制

(2) 第二步：增加常量配置，配置 rylx\_default 的默认值为 2。

图4-70 增加常量配置图示

### 增加常量 ?

步骤名称  \*

字段

#	操作	名称	类型	格式
1		<input type="text" value="rylx_defalut"/>	<input type="text" value="Integer"/>	<input type="text" value="#"/>

(3) 第三步：车牌规则过滤，使用过滤记录组件，配置过滤车牌的规则。用于测试用的车牌的正则表达式为： $^[\u4e00-\u9fa5]{1}[A-Z]{1}[A-Z_0-9]{5}$$ 。

图4-71 车牌规则过滤图示

### 过滤记录 ?

\* 步骤名称  非空，2到50个字符

发送true数据给步骤：

发送false数据给步骤：

条件：


表达式1

(4) 第四步：车牌不符归档配置，使用加载至文件。

(5) “文件”页签配置：

- 本地文件或目录：指定归档文件，“D:\opt\errorfile”。

图4-72 “文件” 页签配置图示

加载至文件 

步骤名称  \*

**文件** 内容 字段

本地文件或目录

创建父目录

启动时不创建文件

从字段中获取文件名

文件名字段

扩展名

文件名里包含日期

文件名里包含时间

指定日期时间格式

时间日期格式

结果中添加文件名

- (6) “字段” 页签配置：
- 注意：需移除 `rylx_default` 字段。

图4-73 “字段”页签配置图示

加载至文件

\* 步骤名称  非空，2到50个字符

文件 内容 字段

获取字段

	操作	名称	类型	格式	长度	精度	货币	小数	分组
1		id	Integer	####0;-#				.	,
2		kkwzbm	Integer	####0;-#				.	,
3		hphm	String					.	,
4		hpzl	Integer	####0;-#				.	,
5		rylx	Integer	####0;-#				.	,
6		tgsl	Timestamp	yyyy-MM-				.	,

增加

确定 显示文件名 取消

- (7) 第五步：过滤人员类型，使用过滤记录，过滤出 `rylx` 为空的记录。
- 发送 `true` 数据给步骤：指定校验通过的数据流向，此处为“加载至目标文件”。
  - 发送 `false` 数据给步骤：指定校验未通过的数据流向，此处为“设置字段值”。
  - 条件：配置表达式 1 的内容，如 [图 4-74](#) 所示。

图4-74 过滤记录图示

**过滤记录** ⓘ

\* 步骤名称  非空，2到50个字符

发送true数据给步骤：

发送false数据给步骤：

条件：

表达式1

- (8) 第六步：设置字段值，将 rylx 的值设置默认值（流中 rylx\_default 的值）。
- 输入字段名称：rylx
  - 替换字段名称：rylx\_default

图4-75 设置字段值图示

**设置字段值** ⓘ

步骤名称  \*

字段

#	操作	输入字段名称	替换字段名称
1	<input type="button" value="🗑"/>	<input type="text" value="rylx"/>	<input type="text" value="rylx_defalut"/>

- (9) 第七步：加载至目标文件，使用加载至文件组件。
- 本地文件或目录：指定目标文件路径。

图4-76 加载至文件图示

加载至文件 ⓘ

步骤名称  \*

文件	内容	字段
本地文件或目录	<input type="text" value="D:\opt\succesfile"/>	<input type="button" value="浏览"/>
创建父目录	<input checked="" type="checkbox"/>	
启动时不创建文件	<input type="checkbox"/>	
从字段中获取文件名	<input type="checkbox"/>	
文件名字段	<input type="text" value="请选择"/>	
扩展名	<input type="text" value="txt"/>	
文件名里包含日期	<input type="checkbox"/>	
文件名里包含时间	<input type="checkbox"/>	
指定日期时间格式	<input type="checkbox"/>	
时间日期格式	<input type="text" value="请选择"/>	
结果中添加文件名	<input checked="" type="checkbox"/>	

- (10) “内容”页签配置：
- 格式：下拉选择 Unix 格式。
  - 编码：下拉选择 UTF-8 格式。



图4-77 “内容” 页签配置

加载至文件 ⓘ

步骤名称  \*

文件 内容 字段

追加方式

分隔符  插入Tab

封闭符

强制在字段周围加封闭符

禁用封闭符修复

头部

尾部

格式

压缩

编码

快速数据存储 (无格式)

按行分拆文件


添加文件结束行

确定 显示文件名 取消

(11) “字段” 页签配置:






注意: 删除字段 `rylx_default` 字段。

图4-78 “字段”页签配置

加载至文件 

\* 步骤名称  非空，2到50个字符

文件 内容 字段

操作	名称	类型	格式	长度	精度	货币	小数	分组
1 <input type="button" value="删除"/>	id	Integer 	####0;-#				.	,
2 <input type="button" value="删除"/>	kkwzbm	Integer 	####0;-#				.	,
3 <input type="button" value="删除"/>	hphm	String 					.	,
4 <input type="button" value="删除"/>	hpzl	Integer 	####0;-#				.	,
5 <input type="button" value="删除"/>	rylx	Integer 	####0;-#				.	,
6 <input type="button" value="删除"/>	tgsl	Timestamp	yyyy-MM-				.	,

## 4.1.7 整库迁移

在做数据迁移时，会遇到数据表较多的情况，虽然使用数据集成模板的批量部署可完成数据迁移。但最优的选择是“整库迁移”，它提供了自动建表、自动采集的能力，可通过一次简单的配置，即可完成多个表的迁移工作。

### 1. 场景描述

XX 公司 MySQL 数据库中存在有历史数据，大约 20+张表，需要将数据迁移到 postgres 数据库中。

### 2. 场景分析

根据 XX 公司的需求，我们可通过数据集成的全量抽取方式完成数据迁移，但需要重复创建 20+的任务。而使用整库迁移，可通过一次创建即可完成多表的全量抽取。

### 3. ETL 设计方案

使用数据集成的整库迁移功能完成。

### 4. 示例前置条件

源库：MySQL 数据库。

目标库：postgres 数据库。

### 5. 示例详细步骤

(1) 第一步：选择[整库迁移]，进入整库迁移页面，单击“创建作业”，弹出配置框。

图4-79 创建整库迁移作业

创建整库迁移作业 ①

\* 作业名称

作业标签  [选择](#)

是否立即下载  是  否

是否清空表  是  否

\* 每个任务复制表数

源数据库  [选择](#)

要复制的表 [选择](#)

#	操作	源表模式	源表名	目标模式	新表名
<a href="#">增加</a>					

目标数据库  [选择](#)

[确定](#) [取消](#)

(2) 第二步：配置要迁移的表及目标库。

- 作业名称：MySQL-2-PG。
- 是否立即下载：勾选“否”。
- 是否清空表：勾选“是”（建议勾选）。
- 每个任务复制表数：6，此处我有 27 张表，且数据量不大。
- 要复制的表：通过单击“选择”，勾选要迁移的源表。
- 目标模式：选择 mysql\_bak。
- 目标数据库：选择“mysql\_bak\_pg”。

图4-80 配置图示

\* 作业名称

作业标签

是否立即下载  是  否

是否清空表  是  否 会勾选所有表输出组件的清空表选项

\* 每个任务复制表数

源数据库

要复制的表

#	源表模式	源表名	目标模式 <input type="text" value="mysql_bak"/>	新表名	操作
1		NewTable	mysql_bak	NewTable	<input type="button" value="删除"/>
2		aaa	mysql_bak	aaa	<input type="button" value="删除"/>
3		abc01	mysql_bak	abc01	<input type="button" value="删除"/>
4		bjg	mysql_bak	bjg	<input type="button" value="删除"/>
5		cdgx	mysql_bak	cdgx	<input type="button" value="删除"/>
6		g001	mysql_bak	g001	<input type="button" value="删除"/>

目标数据库

(3) 第三步：单击创建作业，完成作业创建。在作业列表，下发执行即可。

#### 4.1.8 GPLoad 加载

GPLoad 加载组件适用于 GreenPlum、SeaSQL MPP 及 Generic JDBC 数据库，在大数据量、字段较多场景下，相较于加载至表组件，提供了更好的加载性能。

## 1. 场景描述

XX 集团业务系统数据库使用 PostgreSQL，其中有一张物料信息记录表 mm\_1000，字段数近 100 个，现有数据 100 万条，每天物料信息的变化会导致表中记录的修改和增加。该集团现需要定期将业务系统数据库表 mm\_1000 中的数据，同步到数据仓库 SeaSQL MPP 中的 mm\_1000000。

## 2. 场景分析

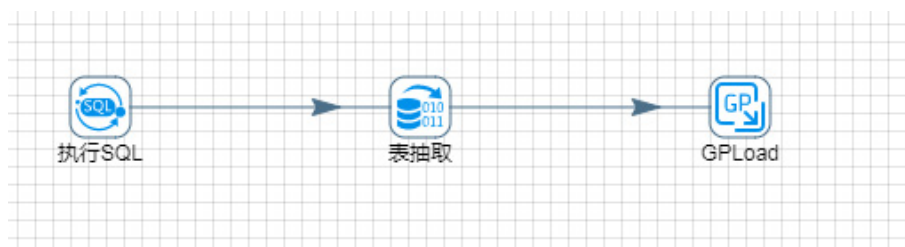
考虑到数据量较大、表字段数较多，这里可以使用数据集成的 GPLoad 批量加载组件。GPLoad 组件运行时先将源表数据写入到落地文件，然后经由 greenplum-loaders 插件提供的文件分发功能，将落地文件发送给 SeaSQL MPP 数据库的各个计算节点，然后在数据库端通过落地文件建立外部表，最后 SeaSQL MPP 各个计算节点在主节点调度下，并行地从外部表查询数据插入到目标表中。

## 3. ETL 设计方案

数据流向：执行 SQL—>表抽取—>加载至表。

ETL 方案：考虑每次同步数据时，源表数据部分发生了修改，同时又有新增数据，首先通过执行 SQL 组件，清空目标表 mm\_1000000，然后使用表抽取组件对源表 mm\_1000 进行全量抽取，最后使用 GPLoad 加载至目标表 mm\_1000000。

图4-81 ETL 任务设计图示



## 4. 示例前置条件

PostgreSQL 数据库中 mm\_1000 表已创建完成，数据量不少于 10 万条。

SeaSQL MPP 数据库中 mm\_1000000 表已创建完成。

## 5. 示例详细步骤

进入[任务管理/任务列表]页面，单击<新增>按钮，新建任务。如[图 4-82](#)所示，创建任务类型为：普通 ETL 任务，完成任务名称、任务描述的填写，单击<跳转任务设计页面>按钮，跳转至任务设计页面。

图4-82 新增任务图示

**新增任务** ?

\* 任务类型

\* 任务名称

任务描述

[跳转任务设计页面](#) [取消](#)

按照设计方案，拖拽组件步骤、建立连接。各步骤详细配置：

(1) 执行 SQL

最终被插入数据的目标表为 mm\_1000000。在这一步骤中通过 SQL 语句清空 mm\_1000000 表中所有数据。

图4-83 执行SQL 具体配置图示

执行SQL ⓘ

步骤名称  \* 非空，2到50个字符

数据库连接

SQL语句（“;”分割语句，“?”替代参数）

```
1 truncate table mm_1000000
```

为每一行执行SQL

作为一个语句执行

变量替换

参数：

<input type="button" value="获取字段"/>		
	操作	作为参数的字段
<input type="button" value="增加"/>		

## (2) 数据表抽取

通过数据表抽取步骤，抽取 PostgreSQL 数据库中表 mm\_1000 的所有数据。

**【注意】**：目标数据库 SeaSQL MPP 是基于 PostgreSQL 实现的分布式数据库，因为这里的业务库是 PostgreSQL，所以简单的使用如[图 4-84](#)所示的查询语句，若业务库使用的是其他类型数据库，比如：mysql、oracle 等，则需要考虑源表中数据可能包含 PostgreSQL 无法处理的字符，此时查询语句需要适当调整。

图4-84 数据表抽取具体配置图示

数据表抽取 ?

步骤名称  \* 非空, 2到50个字符

数据库连接

SQL

```
1 SELECT *
2 FROM "mm_1000"
3
```

允许简易转换

替换SQL语句里的变量

从步骤插入数据

执行每一行

记录数量限制

### (3) Gpload 加载

将从源表中抽取出的数据，使用 Gpload 组件批量加载至 SeaSQL MPP 库的 mm\_1000000 表中。



图4-85 Gpload 组件-基本配置

**Gpload** ⓘ

步骤名称:  \* 非空, 2到50个字符

数据库连接:  选择

目标模式:

目标表:  选择

加载方式:

使用管道:

使用后删除cfg/dat文件:

字段 **GP配置**

加载行为:  更新条件:

获取字段 字段映射

	操作	表字段	流字段	是否匹配	是否更新
1	<input type="checkbox"/>	id	id	否	否
2	<input type="checkbox"/>	name	name	否	否
3	<input type="checkbox"/>	entr date	entr date	否	否

确定 SQL 取消

图4-86 Gpload 组件-GP 配置

**Gpload** ⓘ

使用后删除cfg/dat文件:

字段 **GP配置**

gpload路径:

控制文件:

日志文件:

Data文件:

记录格式错误:

源文件格式:

null as:

quote(引用字符):

逃逸字符:

头部行:

最大错误行:

最大行长度:

编码:

分隔符:

确定 SQL 取消

(4) Gpload 组件配置大概分两部分：

- 配置目标表相关信息，包括数据库、目标表、字段等。
- GP 配置，该部分信息用于指导生成落地文件，为 greenplum-loaders 插件运行做准备。主要包括以下几项配置：
  - 控制文件：生成的控制文件全路径，要求文件的父目录必须存在。
  - 日志文件：生成的插件运行日志文件全路径，要求文件的父目录必须存在。
  - Data 文件：生成的数据文件的全路径，要求文件的父目录必须存在，Data 文件中会被写入源表抽取出的全部数据。
  - 源文件格式：将源表数据写入 Data 文件时使用的文件格式，可选 TEXT/CSV。
  - quote(引用字符)、逃逸字符、分隔符参数的含义详见组件帮助信息，要求分割符不能与 quote(引用字符)、逃逸字符相同。建议从下拉菜单中选择“非可打印字符”，避免与源表数据冲突。
  - 最大错误行：设置最大允许出错行数，仅适用于格式问题导致的某行数据加载错误，若出错行数不超过该参数值，这些错误数据会被忽略，其他数据正常入库，若错误数超过该参数，所有数据加载失败。
  - 最大行长度：Data 文件中一行数据的最大字节数，根据源数据一行记录的实际大小设置。

(5) 任务设计完成后，单击<运行>，等待数据同步完成。

## 4.1.9 REST 抽取

越来越多的数据服务提供商，通过 REST API 为客户提供数据服务，比如天气、位置、公交、火车、交通违章、快递等数据定制查询服务等。API 接口常用的数据交换格式有 XML、JSON 等。

### 1. 场景描述

XX 公司对接客户 REST 接口，该接口提供了健康人员注册信息查询服务，数据交换格式为 JSON，XX 公司需将对接的数据解析并写入关系型数据库。

接口基本信息：

- 接口地址：http://10.121.57.38:8089/person
- 请求方式：GET
- 请求参数：
  - appKey
    - 参数描述：header，授权码；
    - 参数类型：varchar
    - 参数位置：header
    - 是否必填：是
  - pageNumber.cludove
    - 参数描述：请求页数，页数从 0 开始；
    - 参数类型：varchar
    - 参数位置：query
    - 是否必填：是

- recordNumber.cludove
  - 参数描述: 请求每页显示条数, 最大 1000 条, 不传此参数默认为每页请求 200 条信息
  - 参数类型: varchar
  - 参数位置: query
  - 是否必填: 是
- LAST\_MODIFY\_TIME
  - 参数描述: 时间段查询参数, 以天为单位步进查询, 必填, 例如要查询 2020 年 3 月 19 日的的数据, 则参数值为 2020-03-19 00:00:00,2020-03-19 23:59:59 ,以小值在前,大值在后,中间以英文逗号分隔的形式传参, 因为参数值中有空格等特殊字符, 需要将这个参数的值进行两次 url 编码后, 再用来查询
  - 参数类型: varchar
  - 参数位置: query
  - 是否必填: 是

- 请求示例:

例如: 查询时间段 2020 年 3 月 19 日这天的数据 (从第 0 页开始, 每页 1000 条记录), LAST\_MODIFY\_TIME 的参数值实际为 2020-03-19 00:00:00,2020-03-19 23:59:59 , 下面的值是经过了两次 URL 编码的:

```
pageNumber.cludove=0&recordNumber.cludove=1000&LAST_MODIFY_TIME=2020-03-19%2b00%253a00%253a00%252c2020-03-19%2b23%253a59%253a59
```

- 采集 JSON 数据示例:

```
{
  "records": [
    {
      "ID_CARD": "411678199905092345" ,
      "NAME": "张三",
      "POSITION": "郑州市高新区金梭社区云都会小区南门",
      "TEMPERATURE": "36.6",
      "LAST_UPDATE_TIME": "2020-03-23 17:35:26"
    },
    {...},
    ...,
    {...}
  ]
}
```

## 2. 场景分析

数据集成已支持 REST 接口数据采集, 可采用数据集成的 JSON 抽取组件实现数据转换 (将半结构化数据转换为结构化数据), 并通过加载至表将数据写入数据库。

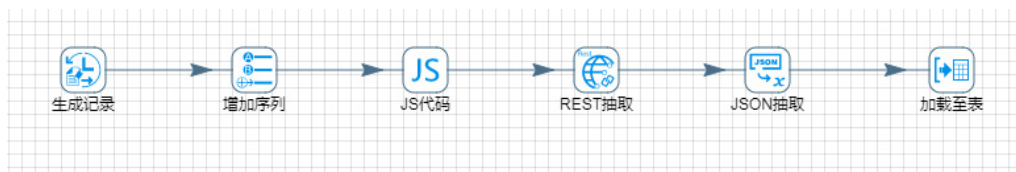
## 3. ETL 方案设计

数据流向: 生成记录 -> 增加序列 -> JS 代码 -> REST 抽取 -> JSON 抽取 -> 加载至表

ETL 方案: 要求每次运行任务时, 将当前运行时间前一天更新的数据同步到目标表。比如当前调度时间为 2020/3/20 09:00:00, 则 REST 抽取组件的 Query 参数中的 LAST\_MODIFY\_TIME 应该设置为 2020-03-19 00:00:00,2020-03-19 23:59:59 , 这个字符串格式的“时间范围”可以使用 JS

代码组件依据当前调度时间计算得出；另外 REST 抽取还需要传入 `pageNumber.cludove`（请求页数）、`recordNumber.cludove`（每页记录数）、`appKey`（API 调用授权码），其中请求页数应该是递增的，可以使用增加序列组件产生，每页记录数和授权码是固定的，使用生成记录组件生成即可。

图4-87 ETL 任务设计图示



#### 4. 示例前置条件

PostgreSQL 数据库中 `hh_1000` 表已创建完成。

#### 5. 示例详细步骤

进入[任务管理/任务列表]页面，单击<新增>按钮，新建任务。如[图 4-88](#)所示，创建任务类型为：普通 ETL 任务，完成任务名称、任务描述的填写，单击<跳转任务设计页面>按钮，跳转至任务设计页面。

图4-88 新增任务图示

### 新增任务 ?

\* 任务类型

\* 任务名称

任务描述

设计方案，拖拽组件步骤、建立连接。各步骤详细配置：

##### (1) 生成记录

该步骤为了配合 JS 代码组件，生成基于当前调度时间的 `LAST_MODIFY_TIME` 的参数值，同时为了生成 REST 抽取组件需要用到的 `recordNumber.cludove`（每页记录数）、`appKey`（API 调用授权码）两个参数的值。

图4-89 生成记录具体配置图示

生成记录 ?

步骤名称  \*

限制

从不停止生成

间隔毫秒数 (延迟)

当前行时间字段名称

以前行时间字段名称

字段

#	操作	名称	类型	格式	值	长度	精度
1		appKey	String		123456		
2		recordNumber.cl...	String		10000		
3		updateTime	String		2020-0319 00:00:00		

(2) 增加序列

增加序列步骤主要是生成递增的 `pageNumber.cludove` (请求页数) 参数。

图4-90 增加序列具体配置图示

**增加序列** ?

步骤名称  \*

字段名称

自定义序列

起始值

增长值

最大值

根据数据库Sequence生成序列

数据库连接  [选择](#)

目标序列  [浏览](#)

[确定](#) [取消](#)

(3) JS 代码

JS 代码步骤根据当前调度时间，计算拼接出 LAST\_MODIFY\_TIME 参数的值。

图4-91 JS 代码具体配置图示

JavaScript代码 ⓘ

步骤名称  \* 非空，2到50个字符

在区域内编辑代码：

```

1 //Script here
2 var now = new Date();
3 var oneday = 1000 * 60 * 24;
4 var before = new Date(now - oneday);
5 var y = year(before);
6 var m = month(before);
7 var d = before.getDate() - 1;
8 var time = y + "-" + m + "-" + d;
9 var searchTime = time + "%2b00%253a00%253a00%252c" + time +
  "%2b23%253a59%253a59";

```

JavaScript函数：

- > 字符串功能
- > 数字功能
- > 时间功能
- > 逻辑功能
- > 特殊功能
- > 文件功能

字段：

操作	字段名称	改名为	类型	长度	精度	替换原名或改名后为该名称的值
<input type="button" value="删除"/>	searchTime	updateTime	String			<input type="checkbox"/>

#### (4) REST 抽取

通过 REST 抽取步骤，抽取接口数据。

图4-92 REST 抽取-通用

REST抽取 ⓘ

\* 非空，2到50个字符

通用 Query参数 Body参数 Matrix参数 HTTP头部 HTTP认证 SSL 输出字段

URL

从字段中获取URL

URL字段

HTTP方法

从字段中获取HTTP方法

HTTP方法字段

内容类型

按照待抽取的 REST 接口定义，在组件的通用页签下，需要配置 URL、HTTP 方法

URL: http://10.121.57.38:8089/person

HTTP 方法: GET

图4-93 REST 抽取-Query 参数

REST抽取 ?

步骤名称: REST抽取 \*

通用 Query参数 Body参数 Matrix参数 HTTP头部 HTTP认证 SSL 输出字段

获取字段

#	操作	字段名称	请求参数名称
1	🗑️	pageNumber.cludove	pageNumber.cludove
2	🗑️	recordNumber.cludove	recordNumber.cludove
3	🗑️	LAST_MODIFY_TIME	LAST_MODIFY_TIME

增加

确定 取消

按照待抽取的 REST 接口定义，在组件的 Query 参数下，需要添加三个参数：  
pageNumber.cludove、recordNumber.cludove、LAST\_MODIFY\_TIME

图4-94 REST 抽取-HTTP 头部

REST抽取 ?

步骤名称: REST抽取 \*

通用 Query参数 Body参数 Matrix参数 HTTP头部 HTTP认证 SSL 输出字段

获取字段

#	操作	字段名称	头部字段名称
1	🗑️	appKey	appKey

增加

确定 取消

按照待抽取的 REST 接口定义，在组件的 HTTP 头部页签中，需要添加 appKey

#### (5) JSON 抽取

选择从 REST 抽取步骤中接收数据（JSON 字符串），然后解析出各个字段的值，并将解析结果发送给加载至表组件。



图4-95 JSON 抽取-文件页签

JSON抽取 ? 非空，2到50个字符

步骤名称  \*

**文件** | 内容 | 字段 | 其它输出字段

本地文件或目录  增加 浏览

正则表达式

正则表达式(排除)

选中的文件

#	操作	文件/目录	通配符	通配符 (排除)	要求	包含子目录

从前面的步骤获取源

前一步骤名

保存源的字段

源是一个文件名

源是一个URL

从结果中移除源字段

确定 预览 显示文件名 取消

在文件页签下，勾选从前面的步骤获取源，选择前一步骤为“REST 抽取”，选择保存源的字段为“result”，该字段名称在 REST 抽取组件的“输出字段”页签下配置，可以手动修改。

图4-96 JSON 抽取-字段页签

JSON抽取 ⓘ

步骤名称  \* 非空，2到50个字符

[文件](#) [内容](#) [字段](#) [其它输出字段](#)

[获取字段](#)

	操作	名称	路径	类型	格式	长度
1		ID_CARD	\$.records[*].ID_CARD	String	▼	
2		NAME	\$.records[*].NAME	String	▼	
3		POSITION	\$.records[*].POSITION	String	▼	
4		TEMPERATURE	\$.records[*].TEMPERATL	String	▼	
5		LAST_UPDATE_TIME	\$.records[*].LAST_UPDA	String	▼	

[增加](#)

[确定](#) [预览](#) [显示文件名](#) [取消](#)

在字段页签下，添加 REST 请求返回结果中包含的各个字段，根据接口返回的 JSON 格式数据示例，我们这里添加 5 个待解析字段，每个字段需要填写字段名称、路径、类型。其中字段名称与示例数据保持一致，与最终加载到数据库表的列名对应，路径则是表名该字段在 JSON 字符串中的层级定位，遵循标准的 JsonPath 定义，类型则统一使用 String，后面若需要用到其他类型，也可以结合字段选择等转换组件做处理。

(6) 加载至数据库表

将 JSON 解析出来的数据加载至 PostgreSQL 数据库中 hh\_1000 表中。

图4-97 加载至数据库表具体配置图示

加载至数据库表 ?

步骤名称  \* 非空, 2到50个字符

数据库连接

目标模式

目标表

提交记录数量

清空表

忽略插入错误

指定数据库字段

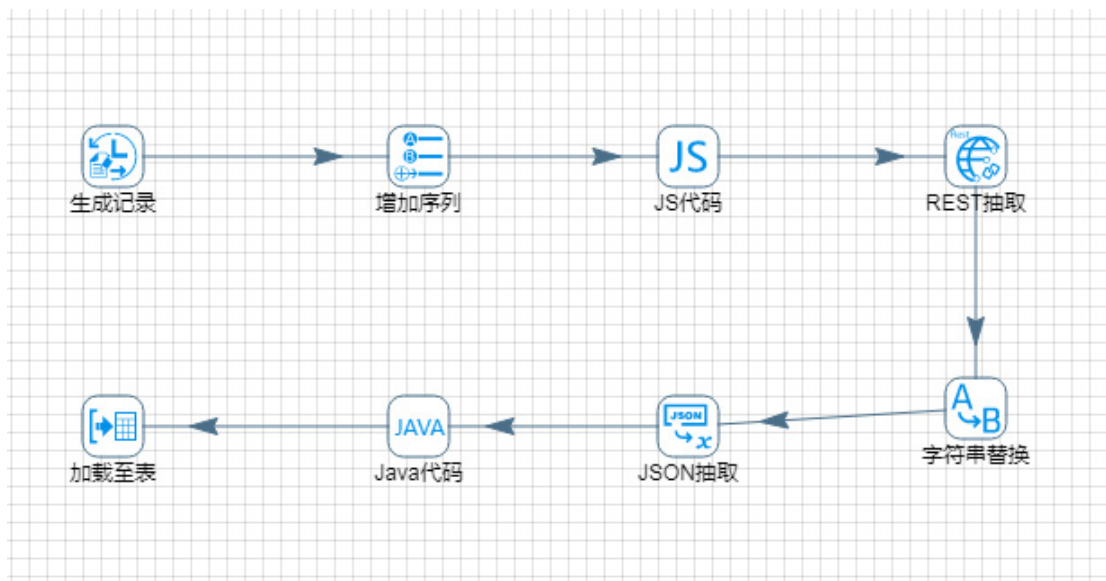
	操作	表字段	流字段
1	<input type="button" value="删除"/>	ID_CARD	ID_CARD
2	<input type="button" value="删除"/>	NAME	NAME
3	<input type="button" value="删除"/>	POSITION	POSITION
4	<input type="button" value="删除"/>	TEMPERATURE	TEMPERATURE
5	<input type="button" value="删除"/>	LAST_UPDATE_TIME	LAST_UPDATE_TIME

(7) 设计好任务后, 单击<运行>按钮, 根据我们之前的设置会发送 100 次 REST 请求, 每次请求抽取 1000 条数据。如果前一天更新的数据总数少于 10 万条, 则当前设置即可满足抽取前一天修改的全部数据, 如果大于 10 万条, 我们还需调整生成记录组件的生成数量、以及增加序列中的最大值。

## 6. 补充场景

根据现场接口变化, REST 接口返回的 JSON 数据中每个“”前都会有“\”字符, 需要替换掉才能进行 JSON 解析; 同时, 某些隐私字段会进行加密处理, 比如: 证件号码、住址以及电话号码等。因此, 根据以上场景需求, 需要进行设计补充, 如图 4-98 所示。

图4-98 补充场景设计图



(1) “字符串替换”组件

该组件用于替换 REST 接口返回的 JSON 数据中多余的符号“\\”，使返回数据符合 JSON 解析格式。图 4-99 是该组件具体的配置信息

图4-99 字符串替换组件配置

字符串替换 ?

步骤名称  \*

获取字段

#	操作	输入流字段	输出流字段	使用正则表达式	被替换的值	使用...替换
1	替换	result		否	\\\\*	*

增加

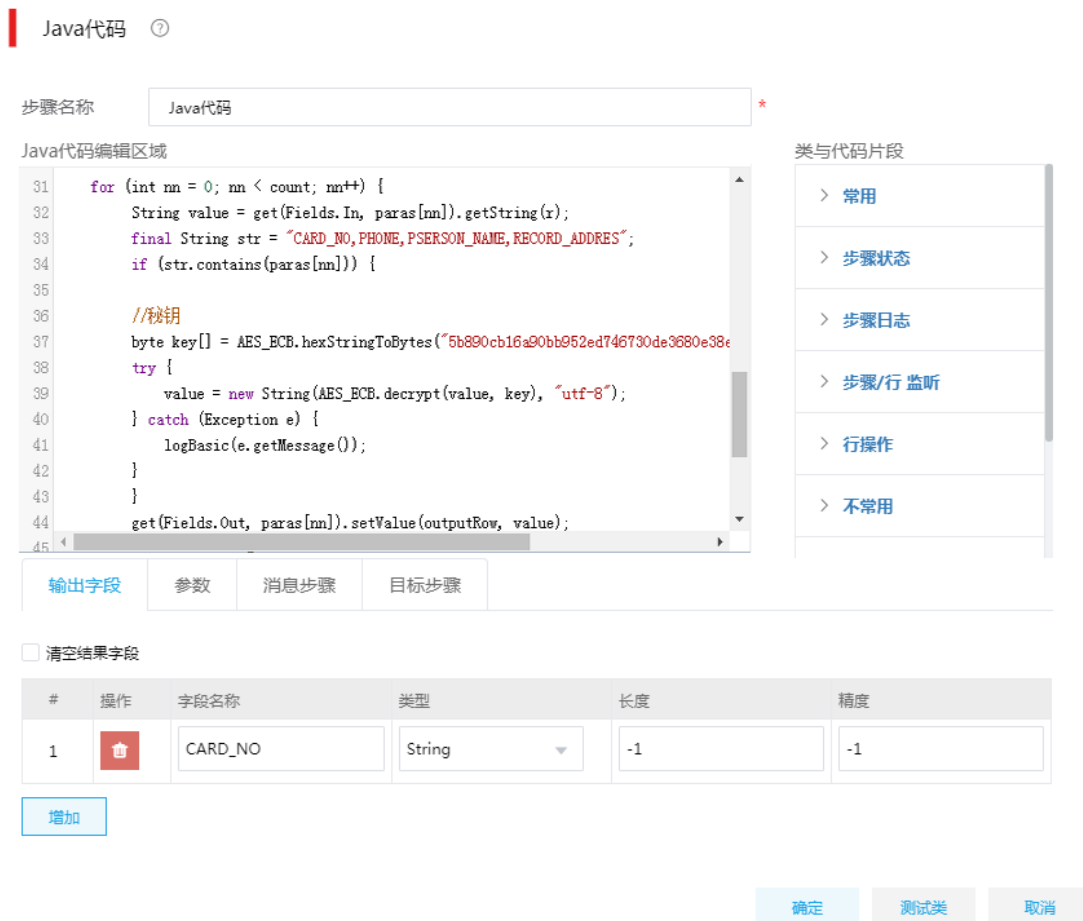
确定 取消

(2) “Java 代码”组件

Java 代码组件主要是对“JSON 抽取”组件输出的字段中，需要进行解密的字段进行解密操作；实现方式可以有多种方式，可以通过引入 REST 接口方提供的加解密工具包，也可以自己通过写代码实现加解密工具包（一般都是需要第三方工具包导入 DI 包路径中，重启服务即可直接使用）。

图 4-100 是“java 代码”组件的具体配置页面，我们只需要按照组件的要求规范在“Java 代码编辑区域”根据实际需求编写 java 代码实现对字段的处理即可。

图4-100 Java 代码组件配置页面



## 4.2 服务集成

### 4.2.1 将数据库字段共享开放成接口场景-数据 API

#### 1. 场景描述

数据所有者希望将数据库里的数据提供给使用者，但是为了安全考虑，又不想直接将数据库的用户名和密码提供给对方，这时可以通过调用接口的方式获取数据库里的数据。

#### 2. 场景分析

集成平台的数据源管理可以添加数据源，服务集成可以选择需要开放的数据源的库表字段，以 restful 接口形式对外提供，返回数据格式为 JSON。

#### 3. 示例详细步骤

(1) [数据源管理]页面，单击<新增>按钮，新增需要对外开放的数据源。

图4-101 新增数据源

新增数据源

\* 数据源名称 非空，2到50个字符

\* 数据源类型

\* 数据源范围  内部数据源  外部数据源

描述信息 0/512

属性列表

#	操作	属性名称	属性值
---	----	------	-----

增加

提交 取消

(2) [服务集成/API 工厂/API 管理]页面，单击<API 注册>按钮，选择注册类型为“数据 API”。

图4-102 选择注册类型

API设计-类型选择

请选择注册的API类型

通用API

函数API

数据API

(3) 进入数据 API 设计页面，填写数据 API 相关信息。

图4-103 配置数据 API 基本属性

The screenshot displays the 'Basic Properties' configuration interface for a data API. At the top, there is a navigation bar with three steps: 'Basic Properties' (selected), 'Parameter Configuration', and 'Preview Completion'. Below this, the 'Basic Properties' section contains the following fields:

- \* API名称: 数据API测试1214
- \* 工作空间: defaultWorkspace
- \* 行业领域: 智慧园区
- \* API来源: 数据能力
- 传输加密:
- \* API组: test\_group
- \* 描述: test
- \* 图标: Includes a '本地上传图片' button and a gallery of default icons.
- \* 版本号: v1.0.0
- \* 版本描述: test

A red '下一步' (Next Step) button is located at the bottom right of the configuration area.

- (4) 基本属性配置完成后，单击<下一步>，进入参数配置页面。配置接口的请求路径，选择刚才配置的数据源及数据源下的数据库、数据表和字段，生成 SQL 查询语句。

图4-104 参数配置

The screenshot shows the '参数配置' (Parameter Configuration) step in an API design tool. It includes fields for request path, request/response formats, data source, database instance, and table. A table for '关联条件' (Association Conditions) is currently empty. An SQL query is provided, and a '生成参数' (Generate Parameters) button is visible at the bottom.

请求路径: /api/v1/testdata

请求参数格式: JSON

返回参数格式: JSON

数据源: pgdata

数据库(实例): cp\_dxedu

数据模式: public

数据表: ability\_base

关联条件: 暂无数据

SQL语句: SELECT "ability\_base"."id","ability\_base"."ability\_name","ability\_base"."ability\_desc","ability\_base"."ability\_status" FROM "public"."ability\_base" WHERE "ability\_base"."id"=:id

生成参数

(5) 单击<生成参数>，生成接口的输入和输出参数。

图4-105 生成接口的输入和输出参数

This screenshot shows the '生成参数' (Generate Parameters) step. It displays the SQL query from the previous step and lists the generated input and output parameters in tables.

SQL语句: SELECT "ability\_base"."id","ability\_base"."ability\_name","ability\_base"."ability\_desc","ability\_base"."ability\_status" FROM "public"."ability\_base" WHERE "ability\_base"."id"=:id

生成参数

输入参数

字段别名	数据类型	字段描述	操作
id	string		删除

输出参数

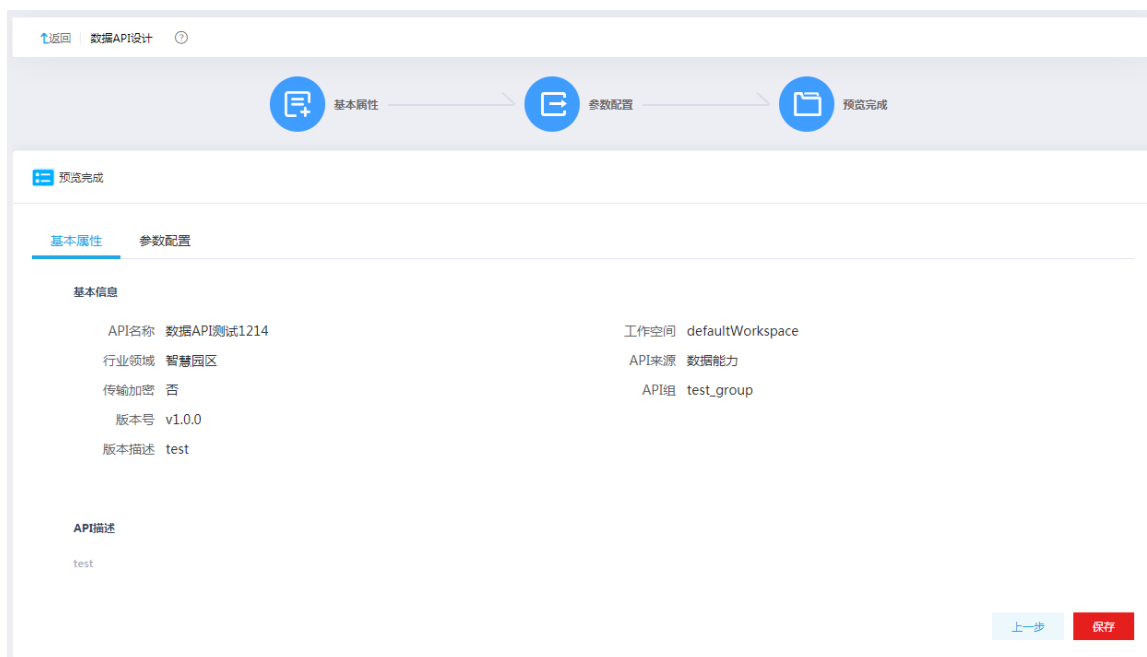
字段别名	数据类型	字段描述	操作
id	string		删除
ability_name	string		删除
ability_desc	string		删除
...	.	.	.

上一步 下一步



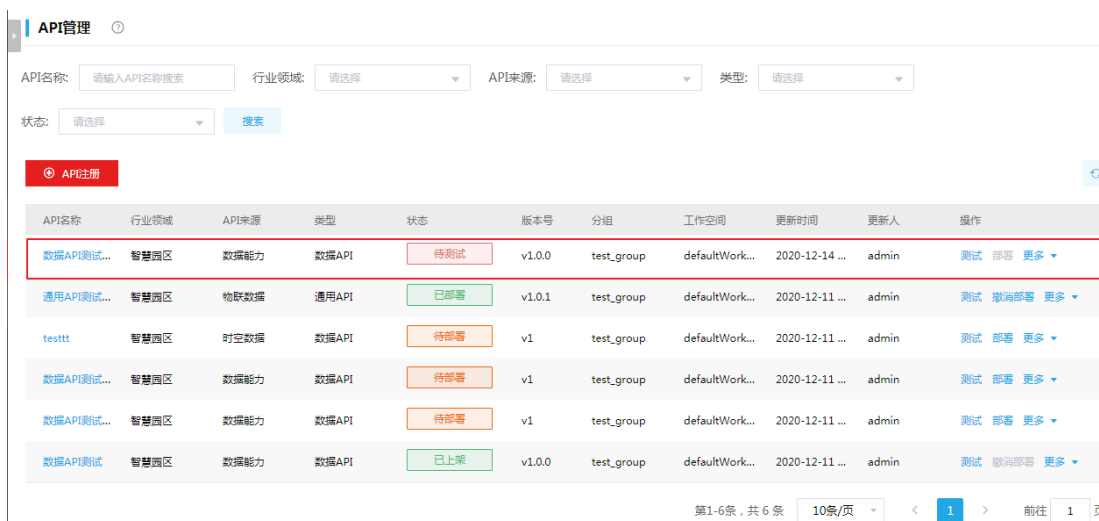
(6) 单击<下一步>, 查看配置的数据 API 的整体信息。

图4-106 查看预览信息



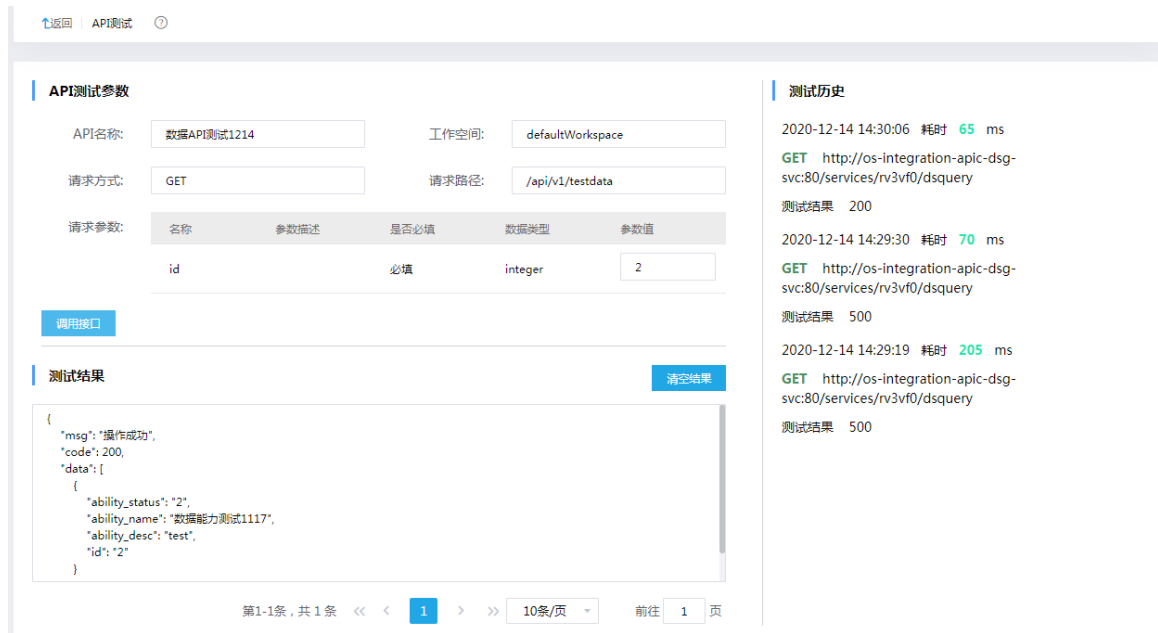
(7) 查看预览信息无误后, 单击<保存>按钮, 返回 API 管理页面。在 API 管理页面可以看到新增的 API。

图4-107 查看 API



(8) 单击<测试>按钮, 进入 API 测试页面, 填写输入输出参数, 单击<调用接口>, 即可查看测试结果。

图4-108 测试 API



(9) API 测试通过后，返回[API 管理]页面，单击<API 部署>，选择部署的网关节点，将 API 部署到网关。

图4-109 API 部署



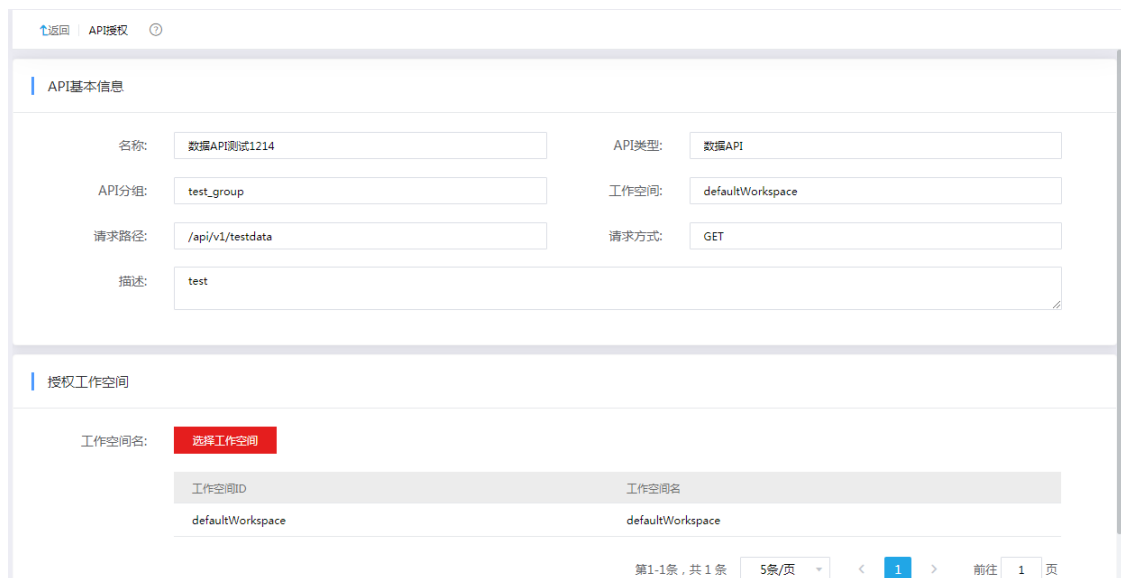
(10) 部署完成后，在[服务集成/API 网关/API 列表]中可以查看部署到网关的 API。

图4-110 查看部署到网关的 API



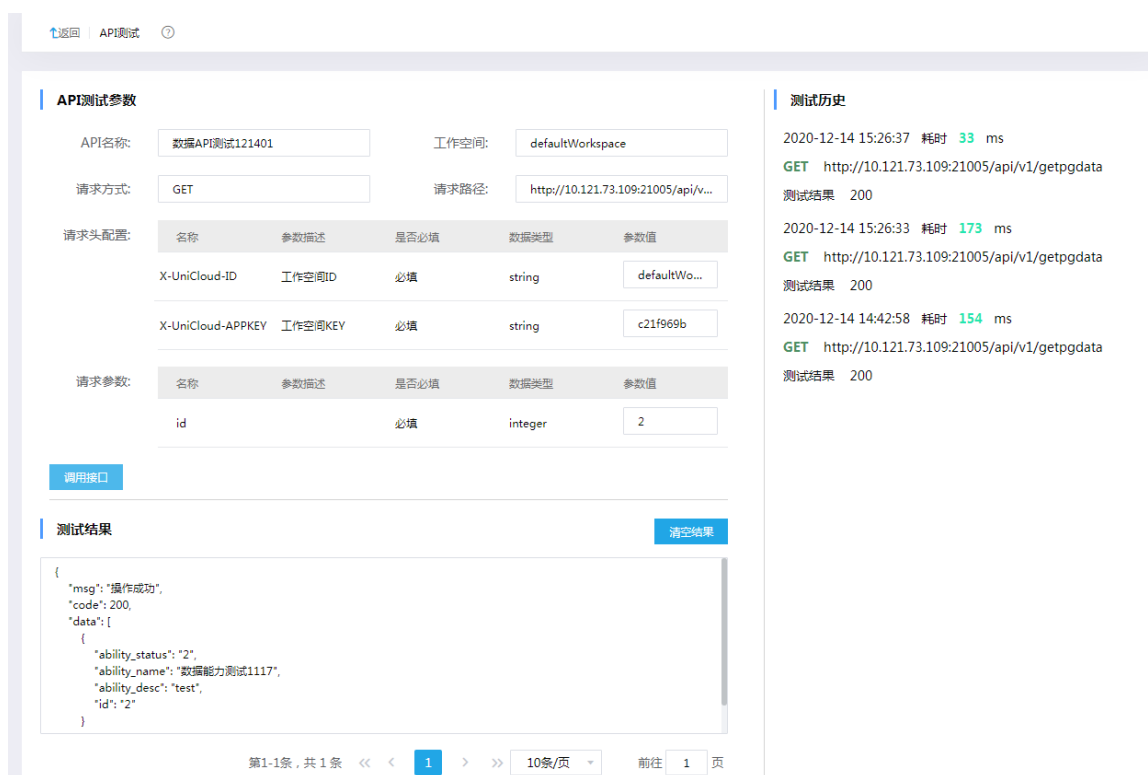
(11) [API 列表]页面，单击<授权>按钮，进入 API 授权页面，将 API 授权给需要的工作空间。

图4-111 API 授权



(12) [授权的 API]页面，可以查看授权给当前工作空间的 API 列表，单击<测试>按钮，进行接口测试。

图4-112 API 测试



(13) 测试页面的请求路径即网关代理后的地址，配合请求头、输入参数，可以在外部进行接口调用。

## 4.2.2 接入第三方系统接口场景-通用 API

### 1. 场景描述

存在多个业务系统，如系统 A、B、C ……，各个业务系统间存在接口调用情况，这时首先需要打通各个系统之间的网络连接，其次一个系统调用其他业务系统时，需要和各个系统进行认证，错综复杂，这时就需要一个公共的平台将所有系统的业务接口统一管理，对外提供统一的出口，每个业务系统只需要和接口管理平台来统一认证对接即可，简化服务间的对接工作。

### 2. 场景分析

服务集成可以通过通用 API 接入第三方业务系统的接口，实现接口的注册、代理和转发。

### 3. 示例前置条件

第三方接口对应的接口文档，如果是带动态认证的第三方接口，需要提前添加认证模板。[服务集成 /API 工厂/认证模板]页面，单击<新增>按钮，添加认证模板。

图4-113 配置认证模板

返回 | 新增模板

\* 模板名称: 认证模板测试

\* 工作空间: defaultWorkspace

\* URL: https://10.121.64.233:443/api/sys/oapi/v1/double\_factor/lo...

\* 请求方式: POST

\* 过期时间(s): 1000

\* 描述: 认证测试模板

**请求头配置** + 新增

参数名称	数据类型	参数值	描述	操作
暂无数据				

**入参配置** + 新增

参数名称	参数类型	数据类型	参数值	描述	操作
password	Body参数	string	UGFzc3cwcmlRAXw==		编辑 删除
domain	Body参数	string	default		编辑 删除
username	Body参数	string	admin		编辑 删除

**出参模板** + 选择出参

参数名称	参数值	数据类型	描述	传递类型	操作
x-auth-token	{res.token}	string		请求头	编辑

#### 4. 示例详细步骤

- (1) [服务集成/API 工厂/API 管理]页面，单击<API 注册>按钮，选择注册类型为“通用 API”，进入通用 API 设计页面，配置通用 API 基本属性。

图4-114 配置基本属性

返回 | 通用API设计

基本属性 → 调用信息 → 预览完成

**基本属性**

\* API名称: 通用API测试1211

\* 工作空间: defaultWorkspace

\* 行业领域: 智慧园区

\* API来源: 物联数据

传输加密:

\* API组: test\_group

\* 描述: test

\* 图标: 本地上传图片

\* 版本号: v1.0.2

\* 版本描述: test

下一步

- (2) 配置完基本属性后，单击<下一步>按钮，配置接口的调用信息。配置接口的请求路径，接入类型选择外部 API，路由地址填写第三方完整的 url 地址，配置请求方式，如果接口带认证，则认证模板关联提前添加好的模板，配置请求头信息、输入输出参数信息，单击<下一步>。

图4-115 配置调用信息

通用API设计

\* 请求路径: /api/v1/getusers

\* 接入类型:  外部API  内部API

\* 路由地址: https://10.121.64.233:443/api/sys/oapi/v1/roles

\* 请求方式: GET

认证模板: 认证模板测试

\* 请求参数格式: JSON

\* 返回参数格式: JSON

**请求头配置** 新增

参数名称	数据类型	是否必填	默认值	描述	操作
暂无数据					

**入参配置** 新增

参数名称	参数类型	数据类型	是否必填	默认值	描述	操作
暂无数据						

**出参配置** 新增

参数名称	数据类型	描述	操作
code	string		<a href="#">↶</a> <a href="#">🗑</a>
res	object		<a href="#">+</a> <a href="#">↶</a> <a href="#">🗑</a>
msg	string		<a href="#">↶</a> <a href="#">🗑</a>

- (3) 单击<下一步>，查看配置的通用 API 的整体信息。

图4-116 API 预览

通用API设计

基本属性 → 调用信息 → 预览完成

预览完成

**基本属性** | 调用信息

**基本信息**

API名称: 通用API测试1211	工作空间: defaultWorkspace
行业领域: 智慧园区	API来源: 物联数据
传输加密: 否	API组: test_group
版本号: v1.0.2	
版本描述: test	

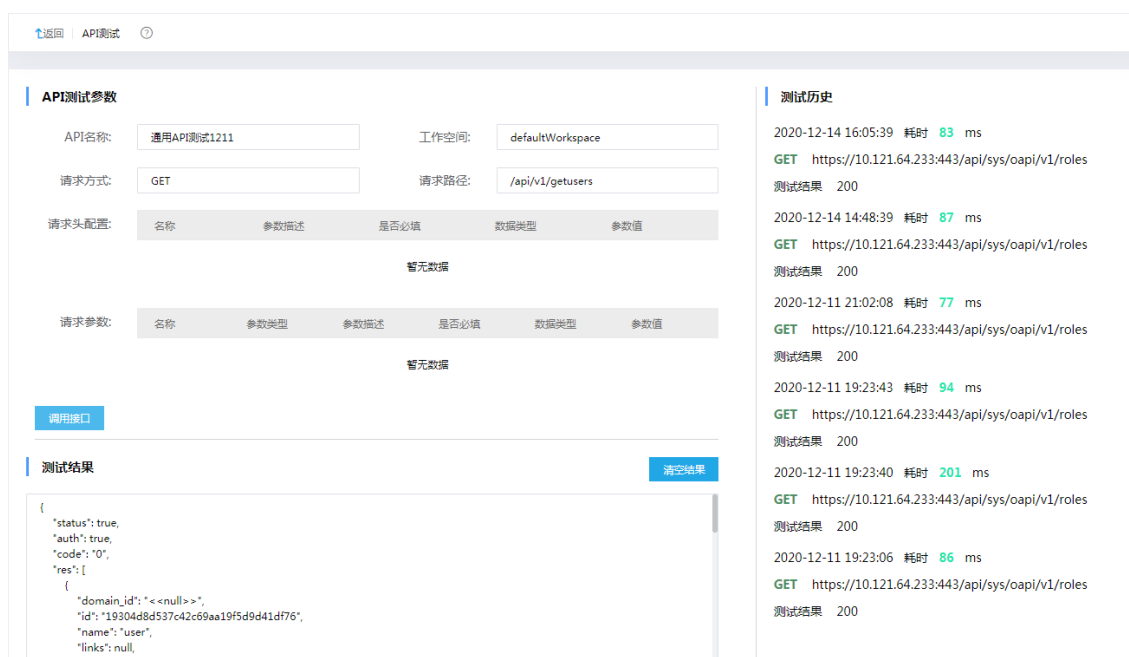
**API描述**

test

上一步 保存

- (4) 单击<保存>按钮，返回 API 列表页面。在 API 列表页面可以看到新增的 API。单击<测试>按钮，可进行接口测试。

图4-117 API 测试



- (5) 接下来的部署、授权操作，可参考 [4.2.1 3. \(9\)](#)中数据 API 的操作步骤。

## 4.2.3 复杂场景下第三方接口对接-函数 API

### 1. 场景描述

第三方接口简单的代理转发可以通过通用 API 来实现，如果涉及参数的转换或者多个接口的编排调用，通用 API 则无法实现，这时需要一种自定义编排的方法来注册接口。

### 2. 场景分析

函数 API 可以通过编写函数脚本来实现复杂场景下的参数转化和多接口调用编排，利用内部实现的工具类，通过 JS 脚本可以实现灵活的编排调用。

### 3. 示例前置条件

在[服务集成/API 工厂/密码箱管理]页面及[服务集成/API 工厂/环境配置]页面分别添加编写函数 API 时需要用到的环境变量和密码箱。

### 4. 示例详细步骤

- (1) [服务集成/API 工厂/API 管理]页面，单击<API 注册>按钮，选择注册类型为“函数 API”，进入函数 API 设计页面，配置函数 API 基本属性。

图4-118 配置函数 API 基本属性

↑返回 | 函数API设计

基本属性 → 参数配置 → 服务脚本 → 预览完成

**基本属性**

\* API名称: 数据变化消息分发api

\* 行业领域: 智慧园区

传输加密:

\* 描述: 数据变化消息分发api

\* 图标:  可选择以下默认图片:

\* 版本号: v1.4.3

\* 版本描述: 111

(2) 配置完基本属性后，单击<下一步>按钮，进入函数 API 的参数配置页面。

图4-119 参数配置

↑返回 | 函数API设计

基本属性 → 参数配置 → 服务脚本 → 预览完成

**参数配置**

\* 请求路径: /io/io.test.oasis.iot/device/distribute

\* 请求方式: POST

\* 请求参数格式: JSON

\* 返回参数格式: JSON

**请求头配置**

参数名称	数据类型	是否必填	默认值	描述	操作
暂无数据					

**入参配置**

参数名称	参数类型	数据类型	是否必填	默认值	描述	操作
body	Body参数	object	必填			+ ↺ 🗑

**出参配置**

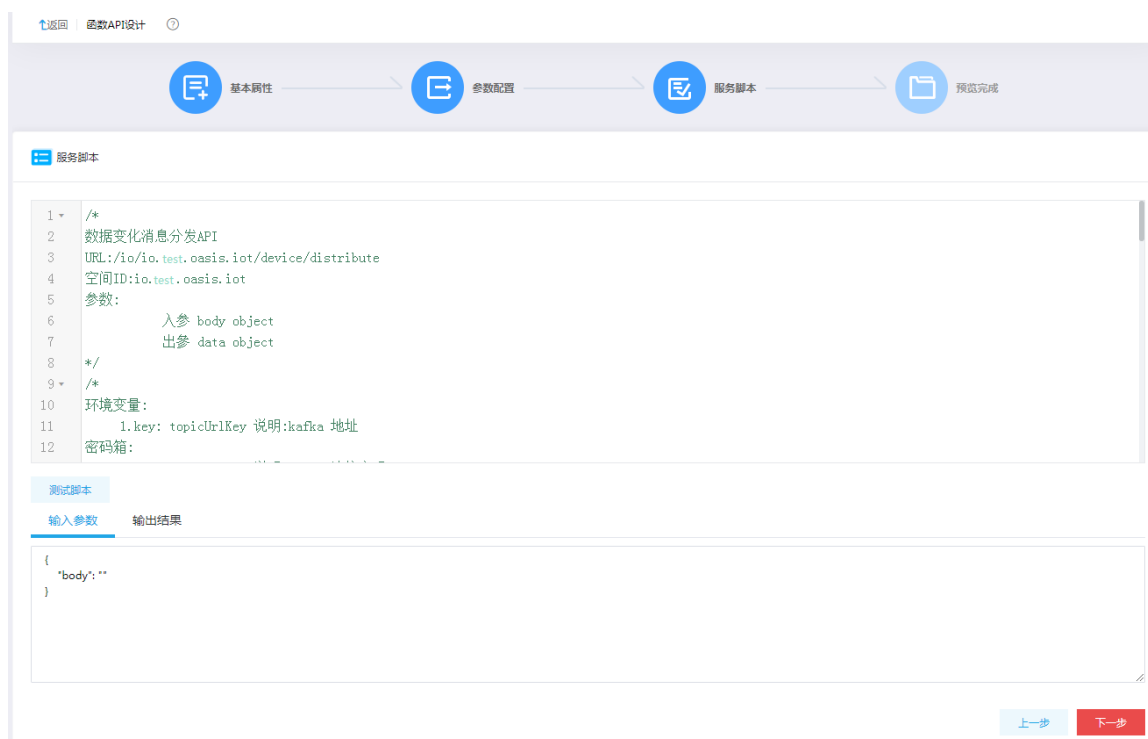
参数名称	数据类型	描述	操作
data	object		+ ↺ 🗑

**错误码**



(3) 单击<下一步>, 编写 JS 脚本

图4-120 编写 JS 脚本。



(4) 单击<下一步>, 查看配置的函数 API 的整体信息。

图4-121 预览信息



(5) 单击<保存>按钮, 返回 API 列表页面。在 API 列表页面可以看到新增的 API。接下来的测试、部署、授权操作, 可参考 [4.2.1 3. \(8\)](#)中数据 API 的操作步骤。

## 4.2.4 画布方式实现多接口编排场景

### 1. 场景描述

对于一个业务接口，如果想对接口的输出参数进行过滤，或者将多个接口的执行编排在一起，函数 API 可以实现这种场景，但是需要编写 JS 脚本，对人员要求较高，此时需要一种简便快速的方法来实现接口的编排。

### 2. 场景分析

服务集成的服务编排可以通过画布的方式，在页面利用拖拉拽来实现输出参数的过滤和多个接口的编排。

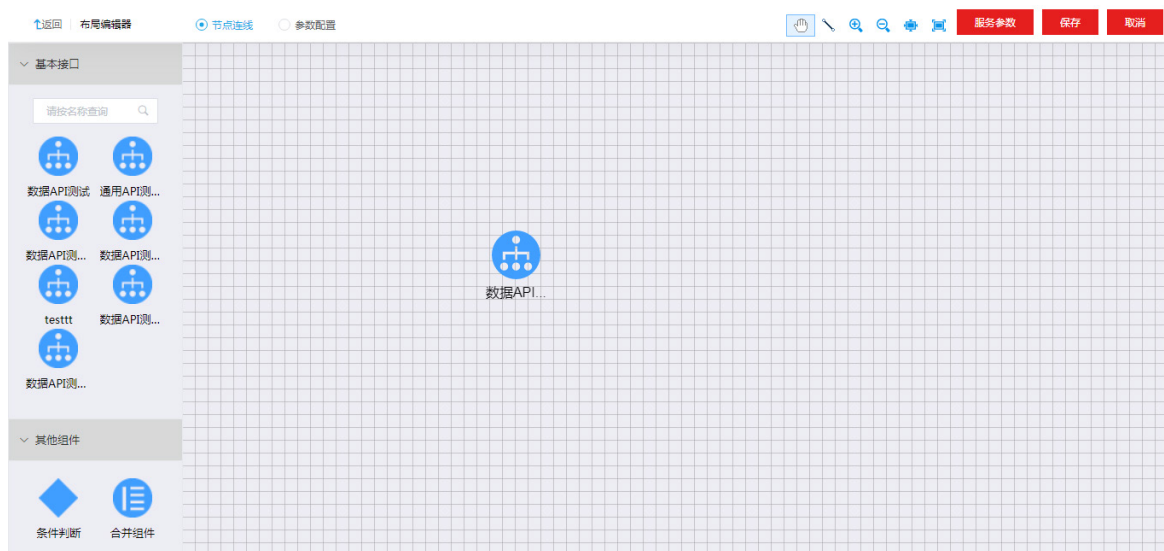
### 3. 示例前置条件

在[API 工厂/API 管理]页面提前注册需要进行服务编排的 API 并测试通过。

### 4. 示例详细步骤

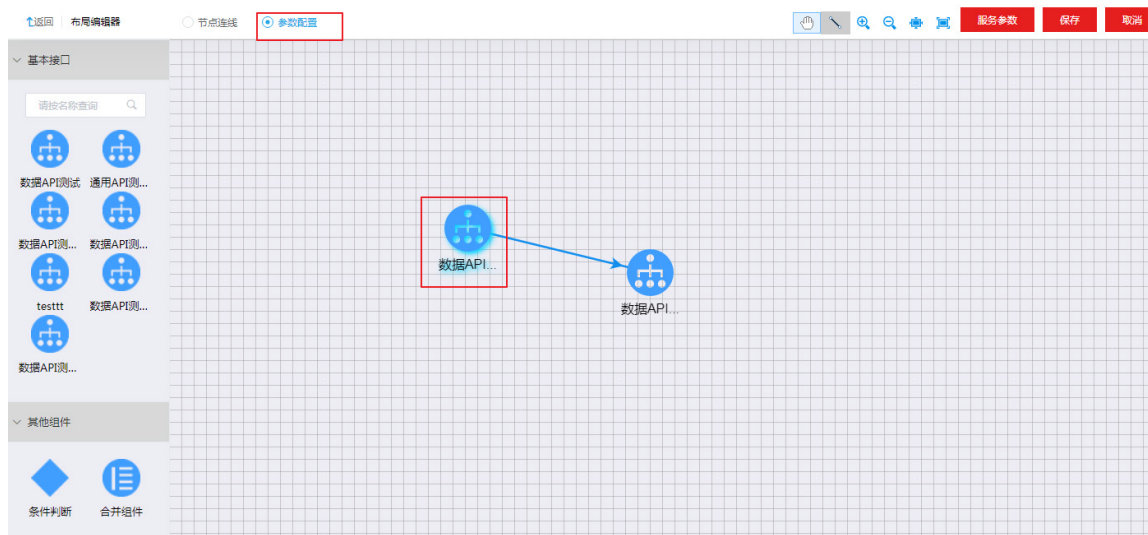
(1) [服务集成/API 工厂/服务编排]页面，单击<新增>按钮，填写基本信息，跳转到设计页面。

图4-122 服务编排设计页面



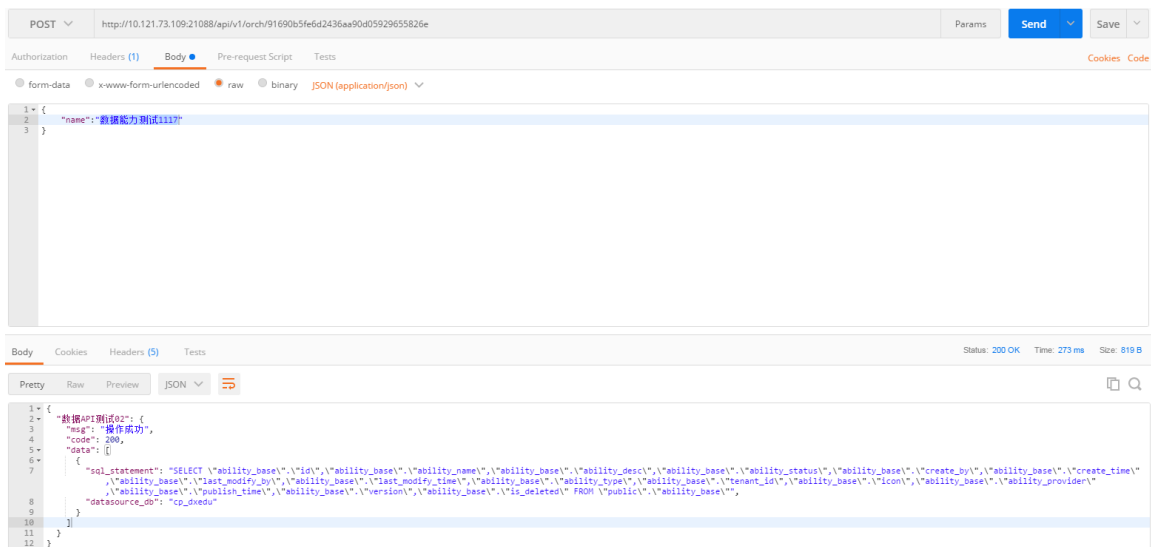
(2) 在设计页面，将需要进行服务编排的 API 拖拽到画布中央，对两个接口进行编排，一个接口的输出作为另一个接口的输入。选择页面上方的参数配置，右键单击所选的 API，弹出框中选择<参数配置>按钮进行输入参数配置。

图4-123 参数配置



- (3) 右键单击第二个 API，弹出参数配置对话框，配置第一个接口的出参为第二个接口的入参。入参配置完成后，单击页面右上角<服务参数>，配置服务出参。
- (4) 另外服务编排支持“条件判断”和“合并组件”操作。
  - 合并组件：目前仅支持数据 API 的结果合并，合并组件主要是对主键和外键的配置（也就是关联字段的配置）。合并组件连接两个父节点后，鼠标右键单击合并组件，弹出窗中选择“合并字段”可进行参数配置。
  - 判断组件：判断组件主要是对任务执行流程的控制。判断组件连接一个父节点后，鼠标右键单击判断组件，弹出窗中选择“选择字段”进入选择字段页面，单击<参数选择>按钮进行判断字段的选择。然后鼠标右键单击判断组件和子节点之间的连线，可在连线上添加判断条件。
- (5) 服务编排完成后，服务编排完成后，用户可以对编排后的服务进行测试，确保编排正确性。在[服务集成/API 工厂/服务编排]页面，单击对应服务操作栏的<详情>按钮。在详情页面，可以看到服务的 url，通过服务编排设计的服务的请求类型均为 POST，参数均在 body 或者 header 中传递。将服务 url 和请求头参数、输入参数填写在 Postman 测试工具中进行接口的测试。

图4-124 通过第三方工具类进行接口测试。



## 4.2.5 资产市场订阅使用接口场景

### 1. 场景描述

集成平台系统内存在多个组织，组织 A 发布的 API 通过主动授权时只能授权给组织 A 内的工作空间使用，组织 B 如果想要使用组织 A 发布的 API，需要主动去订阅使用。

### 2. 场景分析

服务集成发布的 API，部署到网关后，可以单击操作列的<上架>按钮，申请 API 上架，经过管理员审批后，完成上架。上架后的 API 可以在资产市场的公共资产展示，平台内的所有用户都可以看见，并且可以进行收藏和订阅。

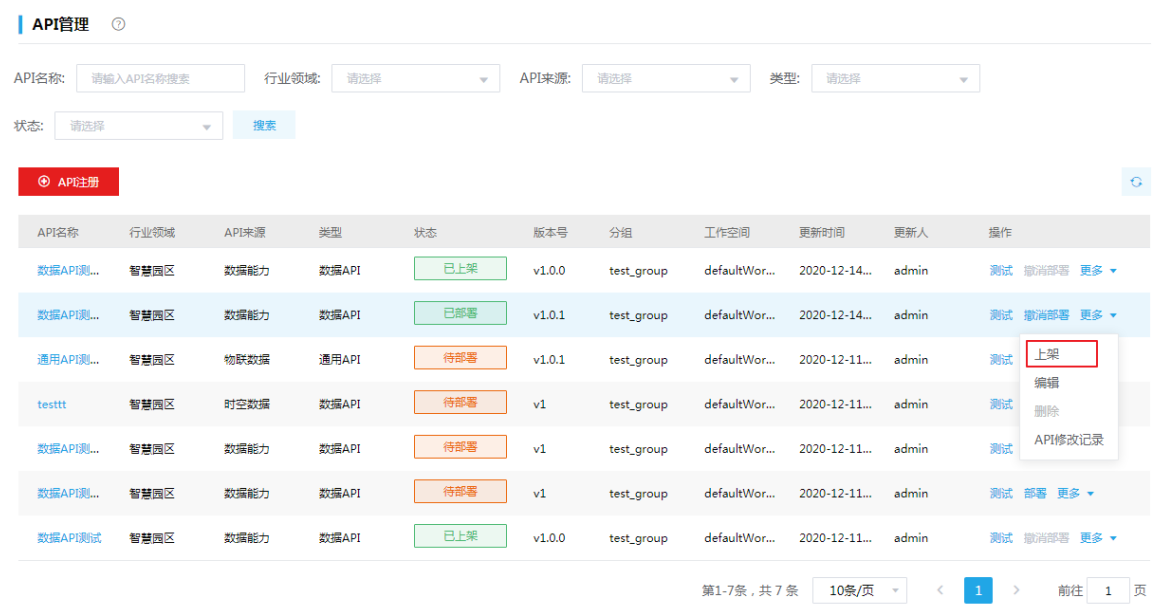
### 3. 示例前置条件

完成 API 的部署及上架。注意，API 上架后，需要经过管理员审批，管理员审批通过后 API 完成上架，在[资产市场]页面展示。

### 4. 示例详细步骤

- (1) [API 工厂/API 管理]页面，将已经部署好的 API 进行上架。API 列表中，单击<更多>按钮，在下拉框中选择上架，系统会自动发出上架审批流程，API 状态变为“上架审批中”。

图4-125 API 上架



(2) 管理员审批通过后，API 状态变为“已上架”，用户可在[资产市场/服务]页面查看已上架的资产。

图4-126 公共资产页面



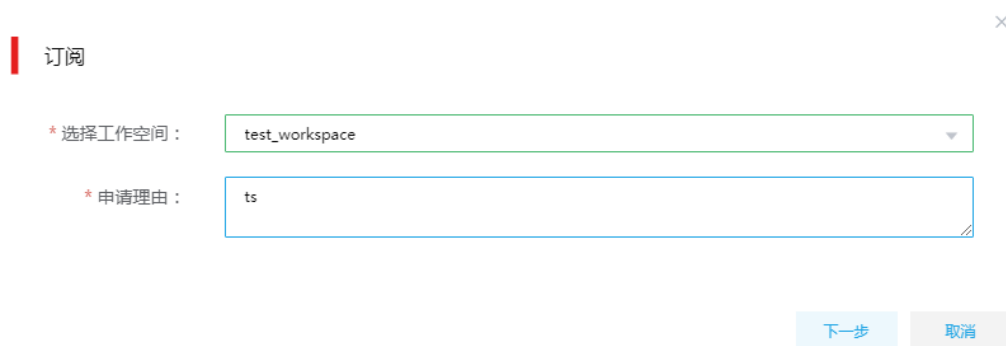
(3) 单击资产，可查看资产详情。API 详情页面，可以查看 API 的基本属性和调用信息。

图4-127 API 详情



- (4) 单击<订阅>按钮，弹出订阅对话框，选择工作空间，填写申请理由，进行订阅申请。系统会发出一个订阅流程给管理员审批。

图4-128 订阅申请



- (5) 管理员登录集成平台系统后，可在[个人中心/待办审批]页面查看到待审批的流程。

图4-129 待办审批

名称	状态	审批结果	审批级别	申请人	责任人	申请时间	操作
API订网_数据API测试1...	处理中	--	一级审批	admin	组织管理员	2020-12-14 17:12:42	<a href="#">更改责任人</a> <a href="#">删除</a>
API上架_数据API测试1...	处理中	--	一级审批	admin	组织管理员	2020-12-14 17:06:47	<a href="#">更改责任人</a> <a href="#">删除</a>
API上架_数据API测试1...	已完成	同意	一级审批	admin	admin	2020-12-14 15:00:05	<a href="#">删除</a>
API上架_三项电表数据...	处理中	--	一级审批	admin	组织管理员	2020-12-14 11:37:11	<a href="#">更改责任人</a> <a href="#">删除</a>
API订网_数据API测试	已完成	同意	一级审批	admin	admin	2020-12-11 21:20:10	<a href="#">删除</a>
API上架_数据API测试	已完成	同意	一级审批	admin	admin	2020-12-11 21:18:53	<a href="#">删除</a>
API上架_文件测试	已完成	同意	一级审批	admin	admin	2020-12-11 21:18:41	<a href="#">删除</a>

第1-7条, 共 7 条 << < 1 ∨ /1 >> 10条/页

(6) 管理员审批通过后，订阅者在[个人中心/我的订阅]页面可以查看审批通过的资产。

图4-130 已审批资产列表

我的订阅

名称:  类型:  [搜索](#)

[已审批](#) [待审批](#)

名称	类型	描述	工作空间	订阅时间	操作
数据API测试	数据API	test	耀州IoT	2020-12-11 21:20:39	<a href="#">详情</a> <a href="#">测试</a> <a href="#">取消订阅</a>

第1-1条, 共 1 条 << < 1 > >> 10条/页

(7) 在我的订阅页面的列表中，单击待查看详情资产对应操作列中的<详情>按钮，进入资产详情页面。未订阅前，用户可以查看资产相关介绍信息；订阅后用户可以查看资产具体的使用信息，用户可通过这些信息使用服务资产。

图4-131 资产详情

基本属性 [调用信息](#) [请求示例](#) [关联目录](#)

**接口信息**

请求路径: /test0326

请求方式: POST

请求参数格式: JSON

返回参数格式: JSON

**请求头说明**

参数名称	数据类型	是否必填	参数说明
X-UniCloud-Id	string	必填	工作空间ID
X-UniCloud-APPKEY	string	必填	工作空间KEY
cmp-token	string	必填	参数描述
auth-token	string	必填	参数描述

**输入参数说明**

参数名称	数据类型	是否必填	参数说明
param1	string	必填	参数描述
param2	string	必填	参数描述
person2	string	必填	参数描述
person1	object	必填	

## 4.3 消息集成

### 4.3.1 作为消息中间件的生产消费场景

消息集成系统作为消息中间件使用，支持 Kafka 客户端进行消息的生产消费，这是消息集成使用最多的场景，以该场景为例，进行实操案例的描述。

#### 1. 场景描述

用户生产、消费消息到消息集成的 Kafka 集群中。

#### 2. 操作示例

(1) [Topic 管理/Topic 列表]页面，单击<新建>按钮，弹出新建窗口，创建 Topic。

图4-132 创建 Topic



The screenshot shows a '创建Topic' (Create Topic) dialog box with the following fields and controls:

- \* Topic名称**: Text input field.
- \* Topic别名**: Text input field.
- \* 权限**: Dropdown menu with '生产+消费' selected.
- \* 老化时间(小时)**: Spin box with '72'.
- \* 分区数**: Spin box with '1'.
- \* 副本数**: Spin box with '1'.
- \* 同步复制**: Toggle switch (off).
- \* 同步落盘**: Toggle switch (off).
- 附件**: '点击上传' button and help icon.
- 描述**: Text area.
- 新建** and **取消** buttons at the bottom.

(2) [Topic 管理/Topic 列表]页面，单击 Topic 列表中的<权限>按钮，弹出 Topic 权限配置窗口，可为指定工作空间赋予该 Topic 的读写权限。



图4-133 配置 Topic 权限



- (3) 为工作空间分配 Topic 权限后，前往工作空间管理页面，单击<密钥管理>按钮，查询相应工作空间的 ID 和密钥。

图4-134 查看密钥



- (4) 以 ID 和密钥分别作为 Kafka 客户端用户的用户和密码，用户可以自行构造 Kafka 生产、消费客户端。

# 5 典型应用案例

## 5.1 医保云案例

### 5.1.1 应用现状

在全国医保的建设要求中，省级医疗保障局将根据国家医疗保障局信息化建设指导意见，坚持医疗保障信息化建设“一盘棋”的原则，依托省级医疗保障平台与国家医疗保障平台之间的协作联通。建设医保云面临如下挑战：

- 时效：数据中台需实现在医保决策分析系统中，支撑实时监控和 T+1 两种类型医保指标的统计和显示
- 数据库响应：需满足高并发、低延迟、实时计算要求，为典型的互联网 OLTP 场景
- 数据质量：需按照 4000+规范要求 进行适配和代码转变，以保证提交上级的数据符合质量要求
- 复杂业务场景：数据库需满足海量数据场景下交易型事务处理

### 5.1.2 解决方案

使用数字平台构建数据中台，从各医疗保障相关数据源中采集或接入数据，并由数据管控和数据开发功能对数据进行处理，然后通过 API 接口和数据库接口进行数据服务集成，实现对上层数据应用提供数据服务。

其中，数据处理部分通过实时计算从实时数据流提取出高价值密度的结构化数据，通过离线计算从结构化、半结构化和非结构化基础数据中提取出高价值密度的结构化数据，并通过数据仓库的分层建模整理数据，最终实现数据整合，统一标准。

医保数据处理全流程：

- (1) 历史数据迁移：集成平台创建 DataX 任务，将原医保平台的历史数据迁移至新的医保云平台。
  - Oracle->生产库
- (2) 建仓：在大数据平台上建立统一的 Hive 数据仓库对所有生产数据进行管理，最大化发挥数据的价值。
  - 生产库->Hive（历史全量）
  - 生产库->Hive（T+1）
- (3) 交换：生产数据导入 Hive 数仓后，通过在集成平台中创建 ETL 任务，定期从仓库中抽取数据，向省交换库做数据同步。
  - Hive->省交换库（T+1）
  - 生产库->省交换库（5 分钟定时）

### 5.1.3 示例详细流程

#### 1. 步骤一：历史数据迁移

原医保平台的历史数据需要迁移至新的医保云平台，一般情况下的数据流向为：Oracle -> DRDS。

集成平台支持创建 DataX 任务，适用于关系型数据库间大宗数据的迁移。

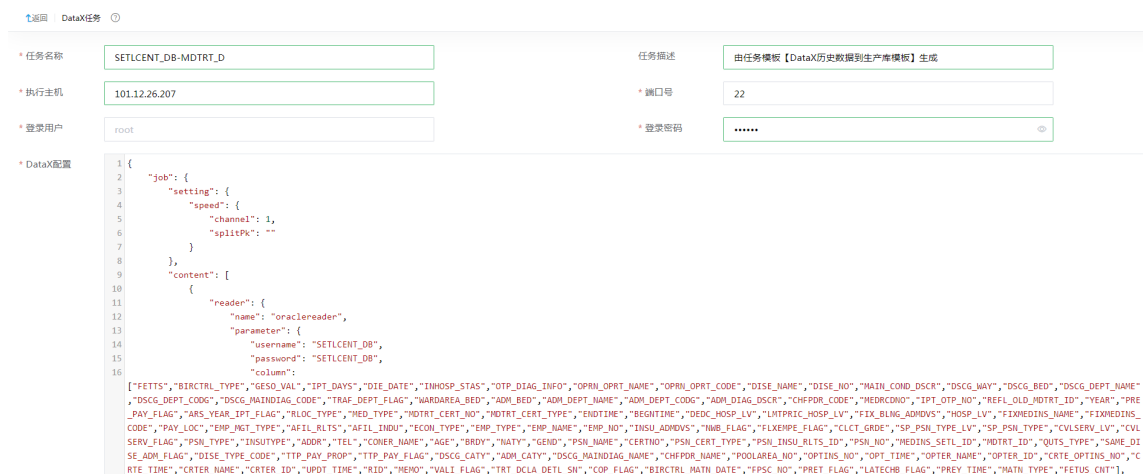
(1) 在[集成平台/数据集成/任务管理/任务列表]页面单击<新增>按钮，弹出新增任务窗口。

图5-1 新增任务



(2) 新增任务窗口，选择 DataX 任务类型。选择任务类型后，单击<跳转任务设计页面>可跳转至任务设计页面，用户可根据实际需要进行任务设计。

图5-2 DataX 任务



DataX 任务配置参数说明如下，用户可根据实际情况进行配置。

- 任务名称：任务名称要求在整个系统中唯一（不区分工作空间）。

- 执行主机：安装了 DataX 客户端的节点 IP。
- 端口号：缺省端口为 22，此处端口指利用 SSH 通道连接执行主机的端口号。
- 登录用户：缺省为 root，此处不可更改。
- 登录密码：该节点的 root 用户密码。
- DataX 配置：根据需要在 DataX 配置内容框中填写 DataX 执行的 JSON 内容。

DataX 配置示例（仅供参考）：

```

{
  "job": {
    "setting": {
      "speed": {
        "channel": 1,
        "splitPk": ""
      }
    },
    "content": [
      {
        "reader": {
          "name": "oraclereader",
          "parameter": {
            "username": "SETLCENT_DB",
            "password": "SETLCENT_DB",
            "column":
["FETTS","BIRCTRL_TYPE","GESO_VAL","IPT_DAYS","DIE_DATE","INHOSP_STAS","OTP_DIAG_INFO","
OPRN_OPRT_NAME","OPRN_OPRT_CODE","DISE_NAME","DISE_NO","MAIN_COND_DSCR","DSCG_W
AY","DSCG_BED","DSCG_DEPT_NAME","DSCG_DEPT_CODG","DSCG_MAINDIAG_CODE","TRAF_DEPT
_FLAG","WARDAREA_BED","ADM_BED","ADM_DEPT_NAME","ADM_DEPT_CODG","ADM_DIAG_DSCR","
CHFPDR_CODE","MEDRCDNO","IPT_OTP_NO","REFL_OLD_MDTRT_ID","YEAR","PRE_PAY_FLAG","AR
S_YEAR_IPT_FLAG","RLOC_TYPE","MED_TYPE","MDTRT_CERT_NO","MDTRT_CERT_TYPE","ENDTIME
","BEGNTIME","DEDC_HOSP_LV","LMTPRIC_HOSP_LV","FIX_BLNG_ADMDVS","HOSP_LV","FIXMEDINS
_NAME","FIXMEDINS_CODE","PAY_LOC","EMP_MGT_TYPE","AFIL_RLTS","AFIL_INDU","ECON_TYPE","
EMP_TYPE","EMP_NAME","EMP_NO","INSU_ADMDVS","NWB_FLAG","FLXEMPE_FLAG","CLCT_GRDE","
SP_PSN_TYPE_LV","SP_PSN_TYPE","CVLSERV_LV","CVLSERV_FLAG","PSN_TYPE","INSUTYPE","ADD
R","TEL","CONER_NAME","AGE","BRDY","NATY","GEND","PSN_NAME","CERTNO","PSN_CERT_TYPE","P
SN_INSU_RLTS_ID","PSN_NO","MEDINS_SETL_ID","MDTRT_ID","QUTS_TYPE","SAME_DISE_ADM_FLA
G","DISE_TYPE_CODE","TTP_PAY_PROP","TTP_PAY_FLAG","DSCG_CATY","ADM_CATY","DSCG_MAIN
DIAG_NAME","CHFPDR_NAME","POOLAREA_NO","OPTINS_NO","OPT_TIME","OPTER_NAME","OPTER_
ID","CRTE_OPTINS_NO","CRTE_TIME","CRTER_NAME","CRTER_ID","UPDT_TIME","RID","MEMO","VALI
_FLAG","TRT_DCLA_DETL_SN","COP_FLAG","BIRCTRL_MATN_DATE","FPSC_NO","PRET_FLAG","LATEC
HB_FLAG","PREY_TIME","MATN_TYPE","FETUS_CNT"],
            "splitPk": "",
            "connection": [
              {
                "table": [
                  "SETLCENT_DB.MDTRT_D"
                ],
                "jdbcUrl": [
                  "jdbc:oracle:thin:@//101.12.54.35:1521/orcl"
                ]
              }
            ]
          }
        ]
      }
    ]
  }
}

```

## 2. 步骤二：建仓

医保系统包含有核心经办、基金监管在内的多个业务子系统，数据分散在各个业务数据库中。需要建立统一的数据仓库对所有生产数据进行管理，最大化发挥数据的价值。

集成平台支持创建 Sqoop 任务，适用于关系型数据库和 HDFS、Hive、HBase 等大数据组件之间的数据迁移。同时，集成平台的 ETL 任务也对 HDFS、Hive、HBase 等做了适配，通过创建 ETL 任务，进行一些简单的拖拉拽和配置，既可实现大数据组件中数据与关系型数据库或者文本数据的相互转换。

### 目标

建仓是数字平台承载的核心业务之一，分为初始全量数据导入和 T+1 增量导入。通过数字平台创建 Hive 数仓（分层），并准确、准时将生产数据抽取并写入 Hive 数仓。

### 示例前置条件

已提前部署好大数据平台，并在大数据平台上安装好了 Hive。

### 建库建表

数据仓库使用 Hive，数据分层管理。主要包括：STG、ODS、DWD、ADS 等。

- STG 层：临时数据层，存放原始数据，直接加载原始数据，数据保持原貌不做处理。数据根据 updt\_time 字段按天分区（同一条数据在 STG 层可能会有多条记录，分布在不同的分区中）。
- ODS 层：结构和粒度与原始表保持一致，对 STG 层数据进行简单清洗（去除空值、脏数据、超过极限范围的数据）。数据根据 crte\_time 字段按年或月分区（同一条数据在 ODS 层只有一条，经过行级质检、表级质检）。
- DWD 层：数据来源于 ODS，根据子系统进一步分区管理，以 ODS 层为基础伴随业务需要进行轻度汇总。
- ADS 层：数据集市层，为各种统计报表提供数据，供上层应用使用。

#### (1) 建库

登录任意 DE 集群节点，连接到 Hive，若集群开启了 Kerberos 认证，需要使用集群超级用户进行认证。使用如下命令创建几层数据库（仅为示例，用户根据实际情况进行创建）：

```
create database stg_prd;
create database ods_prd;
create database dwd_prd;
```

- (2) 登录任意 DE 集群节点，连接到 Hive（或者通过数据管理平台进行可视化建表）。使用如下语句创建 Hive 表（仅为示例，用户根据实际情况创建需要的表，本应用案例中创建了 stg\_prd.stg\_cep\_stcdb\_mdtrt\_d、ods\_prd.ods\_cep\_stcdb\_mdtrt\_d、dwd\_prd.dwd\_dgn\_mdtrt\_d。如下为 ods 的建表语句，ods、stg 及 dwd 的建表语句具体内容可参见 [7.4 ods、stg 及 dwd 建表语句](#)）：

---

```

CREATE TABLE `ods_prd.ods_cep_stcdb_mdtrt_d` (
  `mdtrt_id` string COMMENT '就诊ID',
  `medins_setl_id` string COMMENT '医药机构结算ID',
  `psn_no` string COMMENT '人员编号',
  `psn_insu_rlts_id` string COMMENT '人员参保关系ID',
  -----用户可根据实际需要自由扩展字段-----
  `emp_type` string COMMENT '单位类型',
  `dty_flag` string COMMENT '是否脏数据(0否1是)',
  `local_dty_flag` string COMMENT '本地规则-脏数据标识(0否1是)',
  `exch_updt_time` timestamp COMMENT '入仓时间')
COMMENT "
PARTITIONED BY (
  `dt` string,
  `region` string);

```

---

## 全量迁移

为了支撑医保云上层的监管、决策类子系统的正常运行，需要将近几年乃至全量的医保生产数据一次性导入到 Hive 仓库中。

关键点：部分表的历史数据较多，覆盖的时间范围较大，STG 层要求基于数据的更新时间按天分区，所以在执行历史数据初始导入时涉及的分区数会很多，需要重点关注是否存在数据倾斜，并相应调整运行时参数。如：某地医保项目在做历史数据导入时，单表 1.4 亿数据，以 updt\_time 字段按天分区后发现位于某一天的数据达到了 5000w+，Sqoop 单个 map 无法在默认的 SQL 执行超时（1 小时）之前完成这一天的数据抽取，因此在运行前需要优化一部分参数，包括数据库服务端调大超时时间（DRDS 中参数为 sqlTimeout）、Sqoop 命令的 jdbc url 添加自定义参数。如：socketTimeout=3600000。

数据流向：生产 -> STG 层 -> ODS 层 -> DWD 层 -> ADS 层。整体操作思路如下：

- (1) 执行 Sqoop 命令从生产库抽取全量数据，导入 STG。
- (2) 执行 Spark SQL 对 STG 数据进行质检，数据分流至 ODS 和脏库。
- (3) 执行 Spark SQL 将 ODS 数据导入 DWD。

具体步骤如下：

- (1) 从生产库抽取全量数据，导入 STG 层。使用集成平台创建 Shell 任务，执行预先准备好的数据同步脚本。
  - a. 在[集成平台/数据集成/任务管理/任务列表]页面单击<新增>按钮，弹出新增任务窗口。

图5-3 新增任务



- b. 新增任务窗口，选择 **Shell** 任务类型，选择任务类型后，单击<跳转任务设计页面>可跳转至任务设计页面，用户可根据实际需要进行任务设计。

图5-4 Shell 任务类型



- c. **Shell** 任务配置参数说明如下，用户可根据实际情况进行配置。
- 任务名称：任务名称要求在整个系统中唯一（不区分工作空间）。
  - 执行主机：任意一台能与本机网络互通并且可执行 **Shell** 脚本的节点 IP。
  - 端口号：缺省端口为 **22**，此处端口指利用 **SSH** 通道连接执行主机的端口号。
  - 登录用户/密码：**Shell** 主机上任一用户名/密码。
  - **Shell** 脚本内容：在 **Shell** 脚本内容框中输入脚本内容。



d. Shell 脚本中 `property.properties` 用于定义一些脚本执行参数，如：数据库连接信息、公共脚本路径等，示例如下：

---

```
#生产库连接信息
setlcent_db_src_db_url=jdbc:mysql://101.12.60.51:3323
setlcent_db_src_db_user=test
setlcent_db_src_db_pwd=passwd
#stg层的数据库名
target_stg_db=stg_prd
#各类预置脚本的存放目录
run_job_shell_path=/home/yibao/sjzt/shell
#封装了通过beeline客户端执行hive sql的命令
run_job_shell_name=run_job.sh
#封装了通过spark客户端执行spark sql的命令
run_spark_shell_name=stg_ods_spark_job.sh
#封装了通过spark客户端执行spark sql的命令，申请的资源较多
run_big_spark_shell_name=stg_ods_big_spark_job.sh
#封装了通过beeline客户端执行hive sql的命令
run_beeline_shell_name=run_beeline_job.sh
#存放将生产数据同步至stg的sqoop脚本路径
src_stg_sqoop_path=/home/yibao/sjzt/sqoop
#stg同步至ods的sql文件目录名
stg_ods_sql_path=stg_ods_sql
#ods同步至dwd的sql文件目录名
ods_dwd_sql_path=ods_dwd_sql
#所有入仓相关脚本、配置文件根目录
path=/home/yibao/sjzt
```

---

e. `src_2_stg_cep_stcdb_mdtrt_d_his.sh` 内容示例如下：

---

```

#!/bin/sh

v_time=`date "+%Y-%m-%d %H:%M:%S"`
v_date=`date -d "$v_time" +%Y%m%d`
v_date_ago_1=`date -d "$v_date -1 day" +%Y-%m-%d`
etl_date=${v_date_ago_1}

source /home/yibao/sjzt/property/property.properties

url=${setlcent_db_src_db_url}

username=${setlcent_db_src_db_user}

password=${setlcent_db_src_db_pwd}

hive_db=${target_stg_db}

sql="alter table ${hive_db}.stg_cep_stcdb_mdtrt_d drop if exists partition (dt='${etl_date}');"

query_sql="select
mdtrt_id,medins_setl_id,psn_no,psn_insu_rlts_id,psn_cert_type,certno,psn_name,gend,naty,brdy,age,coner_n
ame,tel,addr,insutype,psn_type,cvlserv_flag,cvlserv_lv,sp_psn_type,sp_psn_type_lv,clct_grde,flxempe_flag,nw
b_flag,insu_admdvs,emp_no,emp_name,emp_type,econ_type,afil_indu,afil_rlts,emp_mgt_type,pay_loc,fixmedi
ns_code,fixmedins_name,hosp_lv,fix_blng_admdvs,lmtpric_hosp_lv,dedc_hosp_lv,begntime,endtime,mdtrt_cer
t_type,mdtrt_cert_no,med_type,rloc_type,ars_year_ipt_flag,pre_pay_flag,year,refl_old_mdtrt_id,ipt_opt_no,med
rcdno,chfpdr_code,adm_diag_dscr,adm_dept_codg,adm_dept_name,adm_bed,wardarea_bed,traf_dept_flag,d
scg_maindiag_code,dscg_dept_codg,dscg_dept_name,dscg_bed,dscg_way,main_cond_dscr,dise_no,dise_na
me,oprn_oprt_code,oprn_oprt_name,otp_diag_info,inhosp_stas,die_date,ipt_days,geso_val,birctrl_type,fetts,fet
us_cnt,matn_type,prey_time,latechb_flag,pret_flag,fpsc_no,birctrl_matn_date,cop_flag,trt_dcla_detl_sn,valid_fla
g,memo,rid,updt_time,crter_id,crter_name,crte_time,crte_optins_no,opter_id,opter_name,opt_time,optins_no,p
oolarea_no,chfpdr_name,dscg_maindiag_name,adm_caty,dscg_caty,ttp_pay_flag,ttp_pay_prop,dise_type_cod
e,same_dise_adm_flag,quts_type,'cep' as subsys_codg_src,'0' dty_flag from mdtrt_d where updt_time
<'${v_date}' and \${CONDITIONS}"

beeline -e "${sql}" -e "!exit"

sqoop import \
--connect ${url}/SETLCENT_DB \
--username ${username} \
--password ${password} \
--hcatalog-database ${hive_db} \
--hcatalog-table stg_cep_stcdb_mdtrt_d \
--hcatalog-partition-keys dt \
--hcatalog-partition-values ${etl_date} \
--query "${query_sql}" --split-by updt_time -m 12
--boundary-query "select min(updt_time), max(updt_time) from mdtrt_d where updt_time <'${v_date}' "

```

---

## (2) STG 到 ODS（质检、分流）

对 STG 层数据进行行级质检、表级质检，并根据质检结果将数据分流至 ODS 和脏数据库。使用数字平台创建 Shell 任务，执行预先准备好的质检和数据分流脚本。其中 `stg_2_ods_cep_stcdb_mdtrt_d.sql` 内容可参见 [7.2 stg\\_2\\_ods\\_cep\\_stcdb\\_mdtrt\\_d.sql 脚本内容](#)。

图5-5 STG 到 ODS



## (3) ODS 到 DWD

基于 ODS 层数据进行数据重分布，按照子系统进一步分区管理。使用数字平台创建 Shell 任务，执行脚本。`ods_2_dwd_cep_stcdb_mdtrt_d.sql` 脚本内容可参见 [7.3 ods\\_2\\_dwd\\_cep\\_stcdb\\_mdtrt\\_d.sql 脚本内容](#)。

图5-6 ODS 到 DWD



## (4) DWD 到 ADS

该部分内容与相应的医保子系统业务强相关，数字平台主要提供作业管理、运行和调度。执行的脚本无需关注。

### T+1 入仓

一般情况下，全量入仓只需要在系统初始化时做一次。后续的持续建仓过程需要通过在数字平台创建周期调度的作业，定期执行预先准备好的 T+1 入仓脚本来完成。

增量流程与全量入仓类似，只是在相应的环节中执行的脚本、SQL 语句稍微有一些区别。

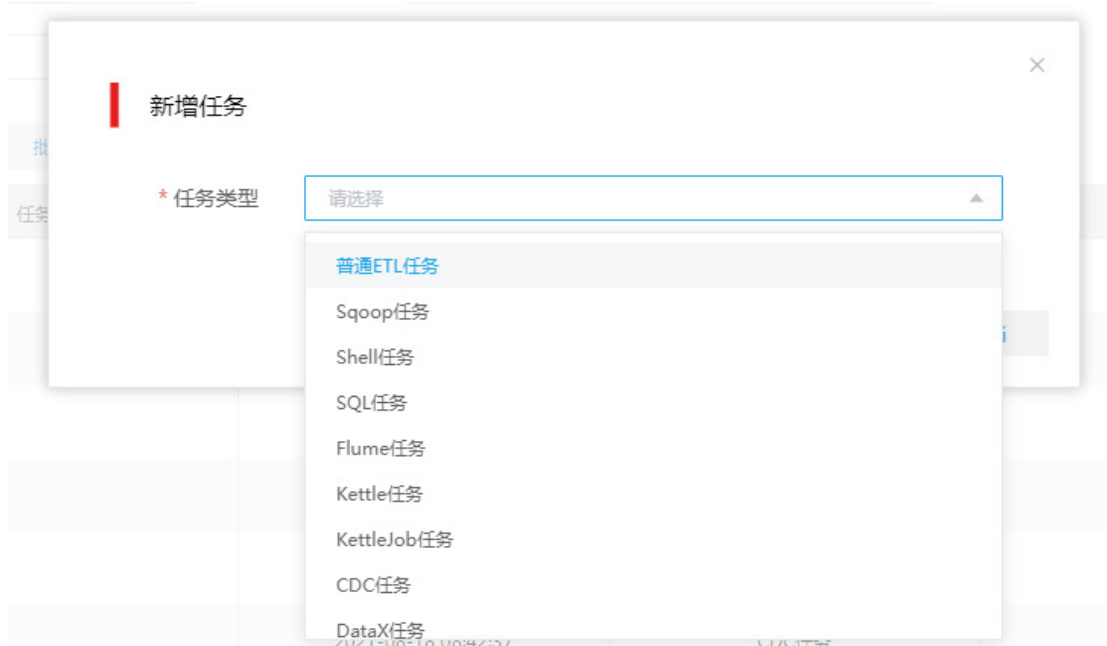
### 3. 步骤三：交换

生产数据导入 Hive 数仓后，需要定期从仓库中抽取数据，向省交换库做数据同步。

#### T+1 交换

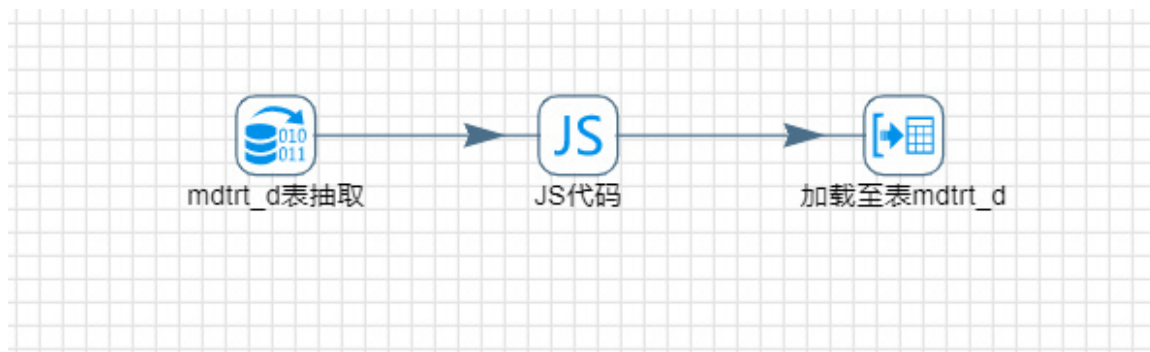
(1) 在[集成平台/数据集成/任务管理/任务列表]页面单击<新增>按钮，弹出新增任务窗口。

图5-7 新增任务



(2) 新增任务窗口，选择普通 ETL 任务类型。选择任务类型后，单击<跳转任务设计页面>可跳转至任务设计页面。用户可根据实际需要进行 ETL 任务设计，每天将增量数据同步至省交换库。

图5-8 ETL 任务设计



(3) 其中，“mdtrt\_d 抽取”步骤负责从 Hive 的 ods 库中抽取 mdtrt\_d 表的 T+1 增量数据。

图5-9 mdtrt\_d 抽取

数据表抽取 ?

步骤名称  \*

数据库连接

SQL

```
1 select b.* from ods_prd.ods_cep_stcdb_mdtrt_d b
2 where substr(b.updt_time,1,10) = date_sub(current_date(), 1)
3 and substr(b.dt, 1, 4) in (
4     select substr(a.crte_time,1,4) from
5     stg_prd.stg_cep_stcdb_mdtrt_d a
6     where a.dt = replace(date_sub(current_date(),1),'-', '')
7 )
```

将时间转换为字符串

时间格式

允许简易转换

替换SQL语句里的变量

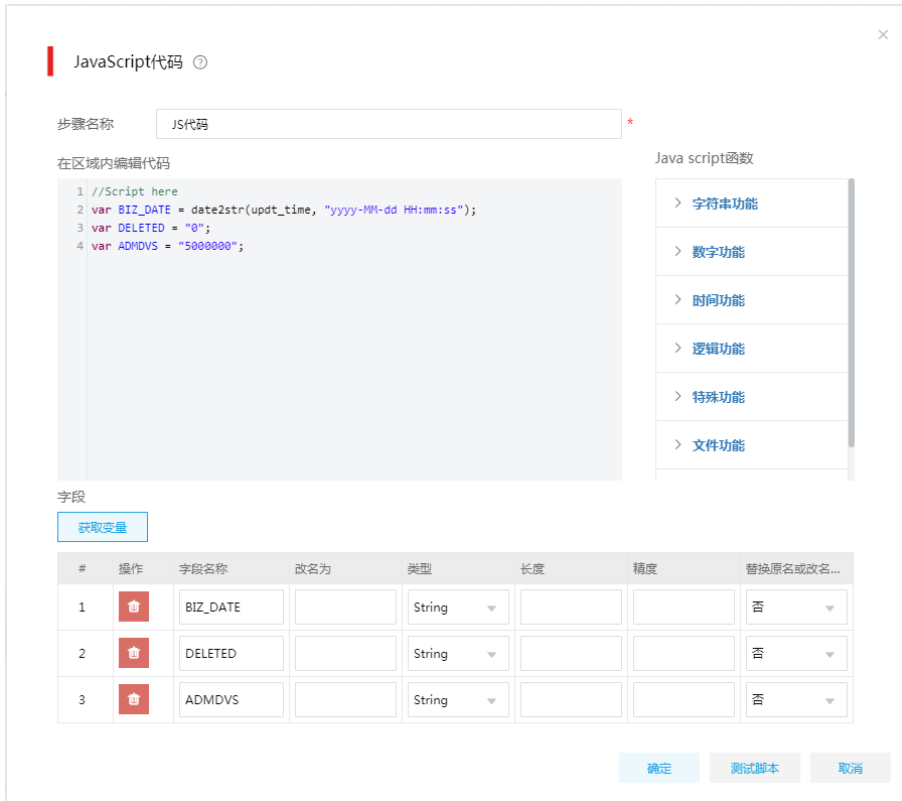
从步骤插入数据

执行每一行

记录数量限制

(4) “JS 代码” 步骤用于添加几个交换库专有字段，如：业务日期、医保区划等。

图5-10 JS 代码



(5) “加载至表 mdtrt\_d” 步骤用于向交换库加载数据，需要配置正确的字段映射。

图5-11 加载至数据表

加载至数据表 ①
✕

步骤名称  \*

数据库连接  选择 清除缓存

目标模式

目标表  选择

提交记录数量

清空表

忽略插入错误

指定数据库字段

主选项

数据库字段

获取字段

输入字段映射

#	操作	表字段	流字段
1	<span style="color: red;">✕</span>	<input style="width: 100%;" type="text" value="MDTRT_ID"/>	<input style="width: 100%;" type="text" value="mdtrt_id"/>
2	<span style="color: red;">✕</span>	<input style="width: 100%;" type="text" value="MEDINS_SETL_ID"/>	<input style="width: 100%;" type="text" value="medins_setl_id"/>
3	<span style="color: red;">✕</span>	<input style="width: 100%;" type="text" value="PSN_NO"/>	<input style="width: 100%;" type="text" value="psn_no"/>
4	<span style="color: red;">✕</span>	<input style="width: 100%;" type="text" value="PSN_INSU_RLTS_ID"/>	<input style="width: 100%;" type="text" value="psn_insu_rlts_id"/>

确定

SQL

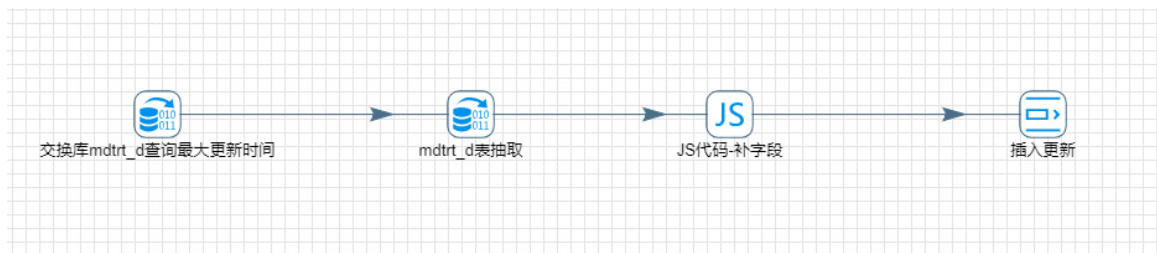
取消

### 准实时交换

生产库中部分表数据需要准实时同步至交换库，数据流向：生产库 -> 交换库。

(1) 通过集成平台创建 ETL 任务，每隔 5 分钟将增量数据同步至省交换库。

图5-12 准实时交换 ETL 任务



(2) 其中，“交换库 mdtrt\_d 查询最大更新时间”步骤用于从省交换库中抽取 mdtrt\_d 表的最新更新时间。

图5-13 抽取最大更新时间

数据表抽取 ?

步骤名称: 交换库mdtrt\_d查询最大更新时间 \*

数据库连接: DRDS\_EXC\_RT [选择] [清除缓存] [查询语句]

SQL

```
1 SELECT
2   ifnull(max(updt_time), subdate(now(), interval 1 year)) as
   max_updt_time
3 FROM SETLCENT_DB_EXC_RT.mdtrt_d
4
```

将时间转换为字符串

时间格式: yyyyMMddHHmmss

允许简易转换

替换SQL语句里的变量

从步骤插入数据: 指定步骤名

执行每一行

记录数量限制: 0

[确定] [预览] [取消]

- (3) “mdtrt\_d 表抽取”步骤使用上一步骤查询的最大时间，从生产库中查询增量数据。“从步骤插入数据”配置为“交换库 mdtrt\_d 查询最大更新时间”。



图5-14 查询增量数据

数据表抽取 ?

步骤名称  \*

数据库连接

SQL

```
101 , `dscg_caty`  
102 , `ttp_pay_flag`  
103 , `ttp_pay_prop`  
104 , `dise_type_code`  
105 , `same_dise_adm_flag`  
106 , `quts_type`  
107 FROM SETLCENT_DB.mdrtrt_d  
108 where updt_time >= ?  
109
```

将时间转换为字符串

时间格式  ▼

允许简易转换

替换SQL语句里的变量

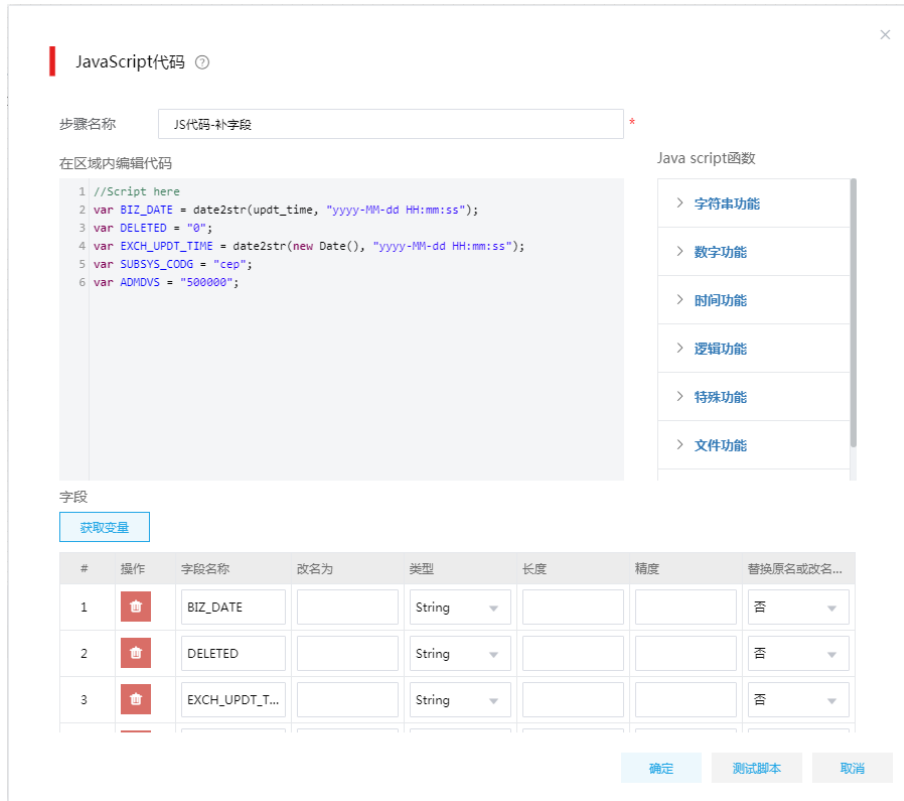
从步骤插入数据

执行每一行

记录数量限制

- (4) “JS 代码-补字段”步骤用于添加几个交换库专有字段，如：医保区划、删除标记、交换时间、子系统代码等。

图5-15 添加交换库专有字段



- (5) “插入更新”步骤用于向交换库做数据同步，这里配置查询字段为表的主键（联合主键），更新字段页签中添加所有表字段。为了避免丢数据，在从生产库中查询增量数据时，还包含了上次同步的部分数据（上次同步的最大更新时间对应的数据），因此这里使用了插入更新组件，来处理重复的数据。

图5-16 插入更新

插入更新 ②

步骤名称: 插入更新 \*

数据库连接: DRDS\_EXC\_RT 选择 清除缓存

目标模式: SETLCENT\_DB\_EXC\_RT

目标表: mdtrt\_d 选择

提交记录数量: 0

不执行任何更新

查询字段 更新字段

获取字段

#	操作	表中字段	比较符	流中字段1	流中字段2
1	🗑️	rid	=	rid	
2	🗑️	biz_date	=	BIZ_DATE	
3	🗑️	subsys_codg	=	SUBSYS_CO...	

增加

确定 取消

## 5.2 疫苗接种案例

### 5.2.1 需求介绍

在疫情防控中，接种疫苗作为重要的一环，是控制病毒传染扩散的重要手段。随着疫苗的推广，接种人员越来越多，为准确迅速地掌握疫苗接种情况，需要对登记的人员信息、人员类别信息、辖区信息、人员接种信息等进行汇总计算，并将结果提供给大屏展示。

为提供大屏展示所需的数据，需在数据库建立业务数据表（创建业务数据表的参考 SQL 语句请参见 [7.1 业务数据库建表语句示例](#)），记录原始数据，并对这些数据进行处理和计算，得出如下 4 个指标：

- 行业接种统计数据
- 社区街道接种统计数据
- 各年龄段接种统计数据
- 一针接种至今各个时间段接种人数统计

## 5.2.2 操作流程

主要步骤及说明如下：

- (1) 本例中的数据来源于业务库中的原始数据，而为了保证这些原始数据信息不受影响，需要通过 DI 将业务库中的原始数据（存量和增量）抽取至数字平台的 ODS（Operation Data Store）层中，作为基础数据。
- (2) 将 ODS 层中保存的基础数据，在数据管理平台中注册为数据源，以便后续步骤中使用。
- (3) 在数据管理平台中，创建对应基础数据的表，表的结构需要与基础数据存储表的结构一致，使系统能够正确读取识别基础数据。此外，还需创建存放清洗后数据的表和统计结果的表，以便在后续步骤中使用。
- (4) 在数据管理平台中，对基础数据进行清洗处理，并将清洗后的数据存放至专用的数据表中。
- (5) 在数据管理平台中，对清洗后的数据进行计算等处理，得出统计结果数据，并存放至预先准备好的表中。
- (6) 对统计结果数据进行查询验证，无误后即可通过集成平台的服务集成进行发布和授权。第三方应用可以调用数据结果用于大屏展示等。

## 5.2.3 抽取基础数据

基础数据的抽取需要通过集成平台的数据集成服务完成，涉及在集成平台-数据集成中执行，完成对基础数据的全量抽取及增量抽取。简要说明如下：

- 注册 MySQL 数据源（用户业务库，存储基础数据），注册 HDFS 数据源（承载从用户业务库中抽取的基础数据）。
- 创建 ETL 任务将用户业务库（MySQL）中的数据抽取至 HDFS 中的 Hive 数据文件中（ODS 层数据）。
- 创建 DI 作业运行 ETL 任务。

### 1. 注册数据源

注册数据源需要组织管理员级别的用户账号。业务数据库为 MySQL，ODS 数据源为 Hive 数据源（Hive 数据存储于 HDFS 中），因此我们这里直接使用 HDFS 数据源（写入 HDFS 比 JDBC 形式写入 Hive 要快）。

#### 注册 MySQL 数据源

如下图所示，在集成平台的数据源管理页面中，新建一个 MySQL 类型的数据源。数据源中输入数据来源的 MySQL 信息，如 IP 地址，用户名和密码等信息。

图5-17 注册 MySQL 数据源

新增数据源

\* 数据源名称: 非空, 2到50个字符

\* 数据源类型: MySQL

\* 驱动: com.mysql.jdbc.Driver

\* IP地址或域名: 127.0.0.1

\* 端口号: 3306

\* 数据库名: 数据库名

\* 用户名: 用户名

\* 密码: 密码

是否采集元数据:  是  否

描述信息: 0/512

属性列表

#	操作	属性名称	属性值
---	----	------	-----

提交 测试连接 取消

新增 MySQL 时，部分参数说明如下：

- 驱动：缺省填入，不可修改。
- IP 地址或域名：必填，目标数据库所在的 IP 地址或域名。
- 端口号：必填，MySQL 数据库使用的端口号，缺省为 3306。
- 数据库名：必填，待连接的已存在的数据库名称。
- 用户名：必填，能够访问对应数据库的用户名。
- 密码：必填，用户名对应的登录密码。
- 是否采集元数据：必填，如果勾选，则会在添加数据源成功后，自动在资产中心中创建相应的元数据采集任务并执行；否则不创建采集任务。
- 描述信息：选填，自定义的描述信息。
- 属性列表：选填，数据源的扩展属性，关于详细属性请参考官方文档。

配置完成后，单击<测试连接>按钮，可检查所填写的信息是否无误，如果测试通过，即证明数据库信息可用。注册完数据以后，即可在任务中使用该数据源信息。

### 注册 HDFS 数据源

为了支持输入数据到 Hive 表中，我们在数据源管理页面注册 HDFS 数据源。使用写入 HDFS 的方式而不是直接使用 Hive 数据源，是因为通过 Hive 的 JDBC 形式写入，性能比较差。没有写入 HDFS 方便，写入 HDFS 相当于写入 Hive 的外部表。

## 2. 建立 ETL 任务

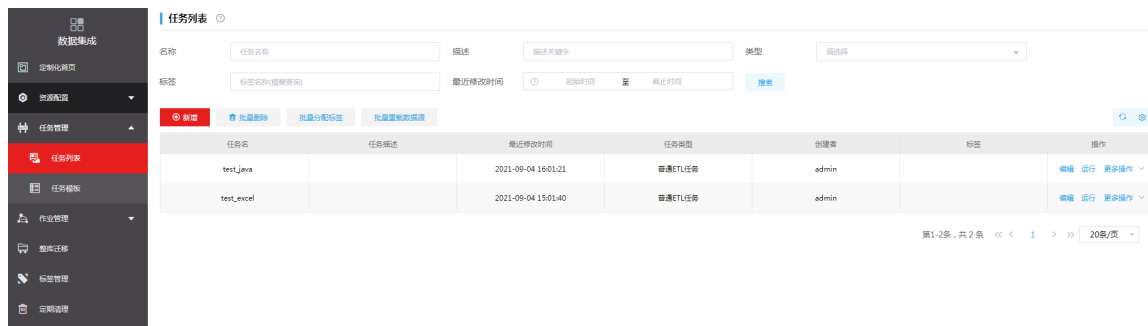
新建数据集成 ETL 任务，抽取表数据至 HDFS（即 Hive 数据文件）。

### 数据集成任务管理

如下图所示，进入[集成平台/数据集成/任务管理/任务列表]页面，单击<新增>按钮，新增 ETL 任务，用于抽取疫苗接种数据、人员信息、行业分类信息等数据，需要新建如下四个 ETL 任务。

- 建立抽取人员信息数据的 ETL 任务。
- 建立抽取人员接种信息 ETL 任务。
- 建立抽取人员分类信息 ETL 任务。
- 建立抽取街道信息 ETL 任务。

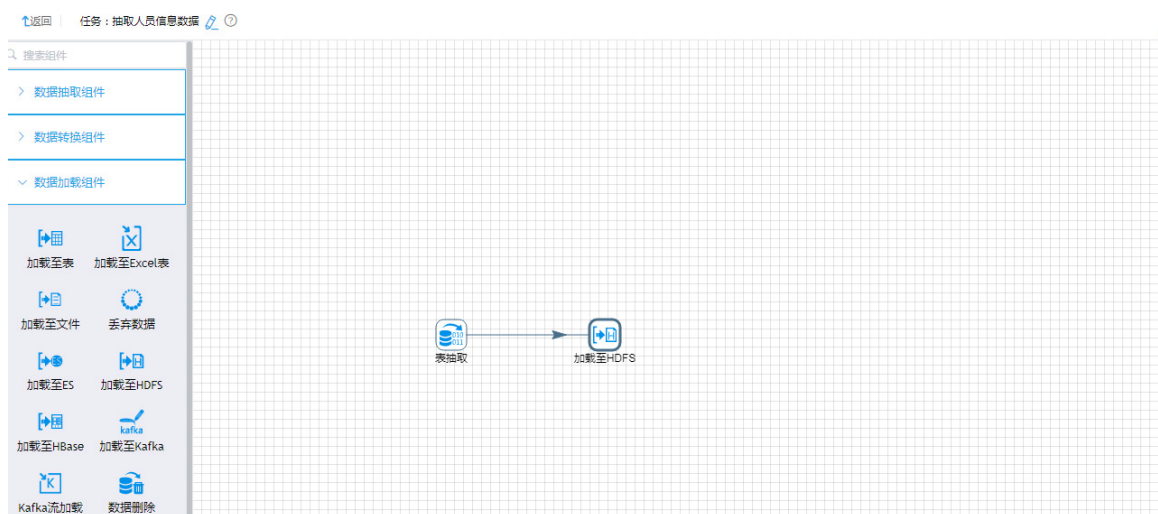
图5-18 任务管理页面



每个任务创建的步骤相同，如下面所述步骤：

- (1) 新建 ETL 类型任务，选取数据抽取组件中的表抽取组件，然后选取数据加载组件列表中的 HDFS 组件。
- (2) 在表抽取组件中输入检索 SQL 语句，从 MySQL 数据库中抽取数据，然后在数据加载组件中设置 HDFS 路径等参数。
- (3) 建立数据抽取组件和数据加载组件之间的联系。

图5-19 编排 ETL 任务



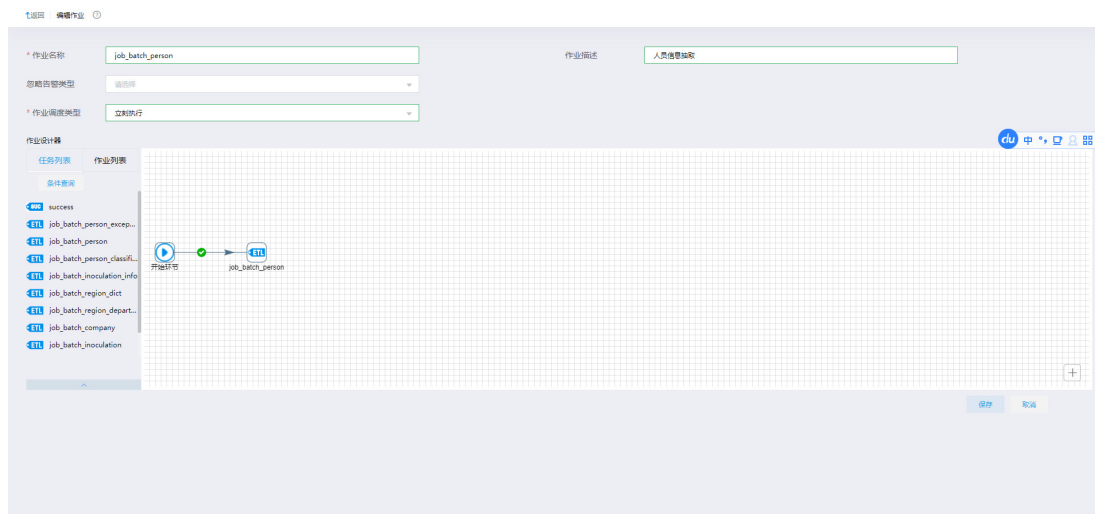
### 3. 建立 DI 作业

新建立即执行的 DI 作业，在每个作业中分别使用上一步骤中新建的 ETL 任务。

#### 建立抽取人员信息数据的 DI 作业

在数据集成模块下的[作业管理/作业列表]页面，单击<新增>按钮，进入新增作业页面，用户可根据实际需要新增作业。输入作业名称，然后在作业调度类型中选择立即执行。

图5-20 新建作业



画布中，单击<可用环节>按钮，在任务列表中选择上一步骤中新建的 ETL 任务，然后和开始节点建立连接。

按照上述步骤，分别建立抽取人员接种信息的 DI 作业、抽取人员分类信息的 DI 作业、抽取街道信息的 DI 作业。

## 5.2.4 创建数据源

在数据管理中需创建 Hive 数据源和 Greenplum 数据源，其中：

- **Hive 数据源：**将前序步骤中，集成平台的 HDFS 数据源中承载数据的 Hive 数据文件，作为 ODS 层数据源。
- **Greenplum 数据源：**用于存放进行计算处理后的结果数据，作为 DWS 层数据源。

### 1. 创建 Hive 数据源

(1) 在[数据管理平台/数据源管理]模块中，单击右上角<创建>按钮，进行数据源的创建操作，如图 5-21 所示。

图5-21 数据源配置页面



(2) 选择 Hive 数据源，并配置参数，如图 5-22 所示。其中：

- Kerberos 用户等信息可以在数据管理平台所使用的大数据平台管理页面中查看。
  - hive principal 参数格式为：hive/IP 地址对应节点的主机名@集群名称大写.COM。
  - krb5.conf 和 Keytab 文件为 Kerberos 认证文件，需要从大数据平台中的集群管理页面下载。
- (3) 在大数据平台获得 Kerberos 的相关信息和认证文件后，返回数据管理平台的新建数据源页面，如图 5-22 所示，填写各个 Kerberos 参数并上传所需的认证文件即可。

图5-22 新增 Hive 数据源

- (4) 填写完毕注册数据源所需要的信息之后，可以单击<测试连接>按钮，测试数据源连通性。
- (5) 提示“连接测试成功”信息，单击<确定>按钮，执行注册数据源。之后即可在数据源列表中看到注册成功的数据源概要信息。

## 2. 创建 Greenplum 数据源

- (1) 在[数据管理平台/数据源管理]模块中，单击右上角<创建>按钮，进行数据源的创建操作，如图 5-23 所示。

图5-23 数据源配置页面

- (2) 选择 Greenplum 数据源，并配置参数，如图 5-24 所示。



图5-24 新增 Greenplum 数据源

↑返回 | 新建数据源 ?

* 数据源名称	vaccination_data_display	24/50
* 数据源类型	Greenplum	?
* 驱动	org.postgresql.Driver	
* IP地址或域名	10.190.31.22	
* 端口号	5434	
* 数据库名	local	
* 用户名	admin	
* 密码	.....	

- (3) 填写完毕注册数据源所需要的信息之后，可以单击<测试连接>按钮，测试数据源连通性。
- (4) 提示“连接测试成功”信息，单击<确定>按钮，执行注册数据源。之后即可在数据源列表中看到注册成功的数据源概要信息。

## 5.2.5 新建数据表

本例中，需要根据 Hive 数据源中的基础数据新建对应的数据表，并提前创建针对人员接种信息数据的清洗表以及后续存储数据处理结果的各结果表。

### 1. 新建基础信息表

为使系统能够正确识别 Hive 数据源中的基础数据，并检测到数据源的表结构，方便后续业务流程中作业的 SQL 处理，需要对 Hive 数据源中的基础信息数据表在[数据管理平台/数据开发]的表管理中新建数据表，这些表为 ODS 层的表。

- (1) 在[数据管理平台/数据开发]模块中，选择左侧导航树中的[表管理]菜单项，进入表管理页面。
- (2) 在页面右上角选择组织，本例中选择“根组织”。
- (3) 单击左上角的<新建>按钮，进入新建表页面。
- (4) 选择 Hive 数据源类型，并选择 [5.2.4 创建数据源](#)中创建的数据源。
- (5) 配置表名等基本属性参数和物理模型设计参数。其中，表名根据实际情况配置，本例中为“ods\_d\_inter\_person\_inoculation\_d”（人员接种信息）；物理模型设计的“外部表”参数


需设为  状态，并指定 Hive 数据源中数据文件存放的 HDFS 路径。

图5-25 基本属性配置

基本属性

* 表名	ods_d_inter_person_inoculation_d	32/100
中文表名	人员接种信息	6/100
主题	请选择主题	+ ↻
标签	请选择标签	+ ↻
描述(可选)		0/200

图5-26 物理模型设计配置

物理模型设计

分层	请选择分层	+ ↻
外部表	<input checked="" type="checkbox"/>	
分区字段	<input type="text"/> 请选择	+ ⊕
存储方式	TEXTFILE	
hdfs路径	/city/ods/dos_d_inter_person_inoculation_d	42/200
字段分隔符	,	1/10
元素间分隔符	-	1/10
kv分隔符	:	1/10

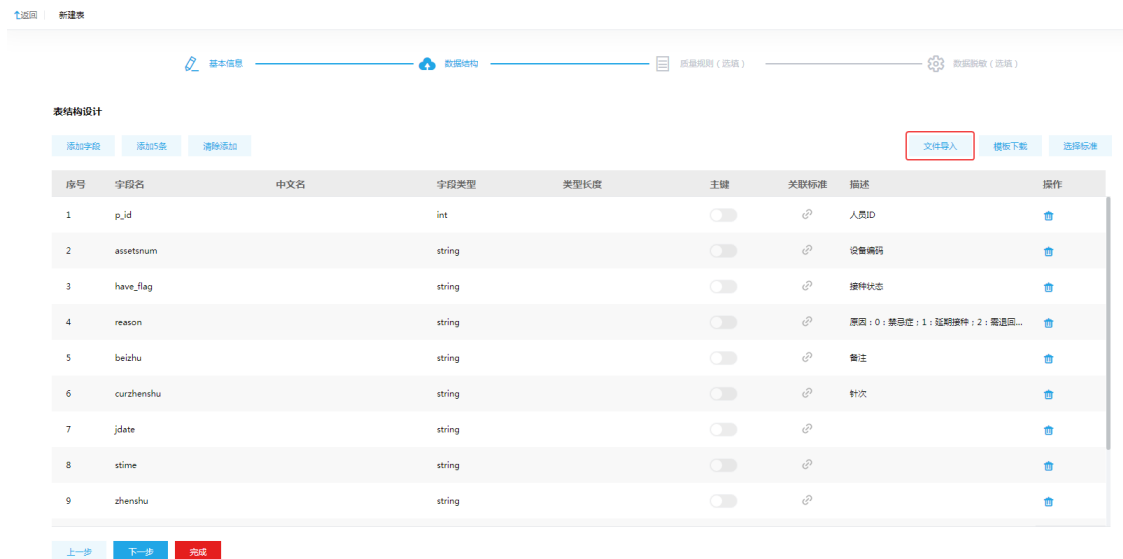
- (6) 单击<下一步>按钮，进入数据结构配置页面。
- (7) 在数据结构配置页面中，可通过<模板下载>按钮下载模板，然后填入字段等信息，再通过<文件导入>按钮进行导入。人员接种信息表示例字段如表 5-1 所示，通过文件导入后如图 5-27 所示。

表5-1 人员接种信息表示例字段信息

字段名称	字段类型	描述
p_id	int	人员ID

字段名称	字段类型	描述
assetsnum	string	设备编码
have_flag	string	接种状态
reason	string	原因：0：禁忌症；1：延期接种；2：需退回上级重新分配；3：不符合接种人群范围
beizhu	string	备注
curzhenshu	string	针次
jdate	string	-
stime	string	-
zhenshu	string	-
yimiao	string	疫苗种类
jinjizheng	string	禁忌症
created	string	接种时间

图5-27 表结构设计



- (8) 单击<完成>按钮，表新建完成。
- (9) 重复步骤(3)-步骤(8)，依次新建人员信息表（字段如表 5-2 所示）、辖区（街道）字典表（字段如表 5-3 所示）、人员分类字典表（字段如表 5-4 所示）。

表5-2 人员信息表示例字段信息

字段名称	字段类型	描述
id	int	序号
name	string	姓名
sex	string	性别
age	int	年龄
mobile	string	手机号
cardno	string	身份证号
classification_id	int	人员分类一级
content	string	备注
company_id	int	单位ID
region_id	int	辖区ID对应属地的摸底工作部门
declare_department_id	int	申报部门ID
created_time	string	创建时间
modified_time	string	修改时间
uuid	string	UUID
addr	string	现住址
streetId	string	街道
provinceid	string	省ID
cityid	string	市ID
disctrictid	string	区域ID
flag	string	是否本市住户
area	string	小区名字
subclass_id	int	人群分类二级

表5-3 辖区（街道）字典表示例字段信息

字段名称	字段类型	描述
id	int	序号


字段名称	字段类型	描述
region	string	辖区
streetId	string	街道
userid	string	User_id
category	string	有没有二级分类（1表示有）
sort	int	顺序编号

表5-4 人员分类字典表示例字段信息

字段名称	字段类型	描述
id	int	序号
classsfication	string	人员分类
created_time	string	创建时间
modified_time	string	修改时间
category	string	级别

## 2. 新建人员接种信息数据清洗表

在基础的人员接种信息表中，可能存在错误或不完整的数据，为保证后续的数据处理可以正常进行，需要对基础信息表中的人员接种信息表进行清洗处理。人员接种信息数据清洗表即用于存放清洗后的人员接种信息数据，需要在数据清洗操作执行前，先新建该表。该表结构与基础信息表中的人员接种信息表相同。该表为 DWB 层的表。

- (1) 在[数据管理平台/数据开发]模块中，选择左侧导航树中的[表管理]菜单项，进入表管理页面。
- (2) 在页面右上角选择组织，本例中选择“根组织”。
- (3) 单击左上角的<新建>按钮，进入新建表页面。
- (4) 选择 Hive 数据源类型，并选择 [5.2.4 创建数据源](#)中创建的的数据源。
- (5) 配置表名等基本属性参数和物理模型设计参数。其中，表名根据实际情况配置，本例中为“dwb\_filtered\_person\_inoculation\_d”（人员接种信息数据清洗表）；物理模型设计的“外部表”参数需设为  状态，将该表设置为内部表，以便于管理和使用。
- (6) 单击<下一步>按钮，进入数据结构配置页面。
- (7) 在数据结构配置页面中，表的字段信息与人员接种信息表示例字段一致，如[表 5-1](#)所示，可通过文件导入，导入后如[图 5-27](#)所示。
- (8) 单击<完成>按钮，表新建完成。

### 3. 新建结果表

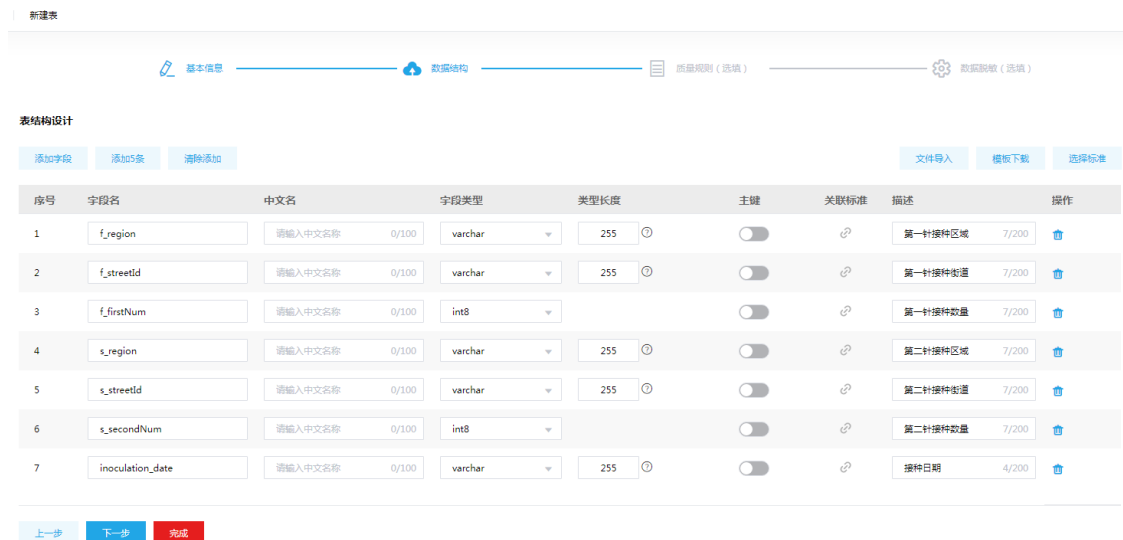
为便于存储数据处理后的结果数据，需要先新建各结果数据表，这些表为 DWS 层的表。

- (1) 在[数据管理平台/数据开发]模块中，选择左侧导航树中的[表管理]菜单项，进入表管理页面。
- (2) 在页面右上角选择组织，本例中选择“根组织”。
- (3) 单击左上角的<新建>按钮，进入新建表页面。
- (4) 选择 Greenplum 数据源类型，并选择 [5.2.4 创建数据源](#)中创建的数据源。
- (5) 配置表名等基本属性参数和物理模型设计参数。其中，表名根据实际情况配置，本例中为“dws\_region\_inoculation\_day\_statistics”（社区街道疫苗接种按天统计结果表）；物理模型设计部分，存储模式使用默认值“row”，模式选择 public。
- (6) 单击<下一步>按钮，进入数据结构配置页面。
- (7) 在数据结构配置页面中，表的字段信息与人员接种信息表示例字段一致，如所示，可通过文件导入后如所示。

表5-5 社区街道疫苗接种按天统计结果表字段信息

字段名称	字段类型	描述
f_region	varchar(255)	第一针接种区域
f_streetId	varchar(255)	第一针接种街道
f_firstNum	int8	第一针接种数量
s_region	varchar(255)	第二针接种区域
s_streetId	varchar(255)	第二针接种街道
s_secondNum	int8	第二针接种数量
inoculation_date	varchar(255)	接种日期

图5-28 表结构设计



- (8) 单击<完成>按钮，表新建完成。
- (9) 重复步骤(3)-步骤(8)，依次新建社区街道疫苗接种全量统计结果表（字段信息如表 5-6 所示）、行业疫苗接种按天统计结果表（字段信息如表 5-7 所示）、行业疫苗接种全量统计结果表（字段信息如表 5-8 所示）、各年龄段疫苗接种按天统计结果表（字段信息如表 5-9 所示）、各年龄段疫苗接种全量统计结果表（字段信息如表 5-10 所示）、第一针接种至今各个时间间隔人数统计结果表（字段信息如表 5-11 所示）。

表5-6 社区街道疫苗接种全量统计结果表字段信息

字段名称	字段类型	描述
f_region	varchar(255)	第一针接种区域
f_streetId	varchar(255)	第一针接种街道
first_total_num	int8	第一针接种数量
s_region	varchar(255)	第二针接种区域
s_streetId	varchar(255)	第二针接种街道
second_total_num	int8	第二针接种数量

表5-7 行业疫苗接种按天统计结果表字段信息

字段名称	字段类型	描述
f_classification	varchar(255)	第一针行业分类
firstNum	int8	第一针接种数量
s_classification	varchar(255)	第二针行业分类

字段名称	字段类型	描述
secondNum	int8	第二针接种数量
inoculation_date	varchar(255)	接种日期

表5-8 行业疫苗接种全量统计结果表字段信息

字段名称	字段类型	描述
f_classification	varchar(255)	第一针行业分类
firstNum	int8	第一针接种数量
s_classification	varchar(255)	第二针行业分类
secondNum	int8	第二针接种数量

表5-9 各年龄段疫苗接种按天统计结果表字段信息

字段名称	字段类型	描述
f_ageRange	varchar(255)	第一针年龄段
firstNum	int8	第一针接种数量
s_ageRange	varchar(255)	第二针年龄段
secondNum	int8	第二针接种数量
inoculation_date	varchar(255)	接种日期

表5-10 各年龄段疫苗接种全量统计结果表字段信息

字段名称	字段类型	描述
f_ageRange	varchar(255)	第一针年龄段
first_total_num	int8	第一针接种数量
s_ageRange	varchar(255)	第二针年龄段
second_total_num	int8	第二针接种数量

表5-11 第一针接种至今各个时间间隔人数统计结果表字段信息

字段名称	字段类型	描述
interval_period	varchar(255)	时间间隔



字段名称	字段类型	描述
person_num	int8	人数

## 5.2.6 构建业务流程

准备工作完成后，即可开始构建业务流程，包括创建业务流程，并在业务流程画布中增加数据清洗作业和各类数据计算作业。

### 1. 创建业务流程

- (1) 在[数据管理平台/数据开发]模块中，选择左侧导航树中的[调度中心]菜单项，进入调度中心页面。
- (2) 在页面右上角选择组织，本例中选择“根组织”。
- (3) 单击左上角的<新建>按钮，进入新建业务流程。
- (4) 输入业务流程名称和描述信息，本例中名称为“疫苗接种数据统计”。
- (5) 单击<确定>按钮，业务流程创建成功，页面进入该业务流程的画布编辑页签。
- (6) 将左侧的作业组件拖入画布中，生成业务流程中的作业节点。双击该作业节点，可在弹窗中配置节点参数。

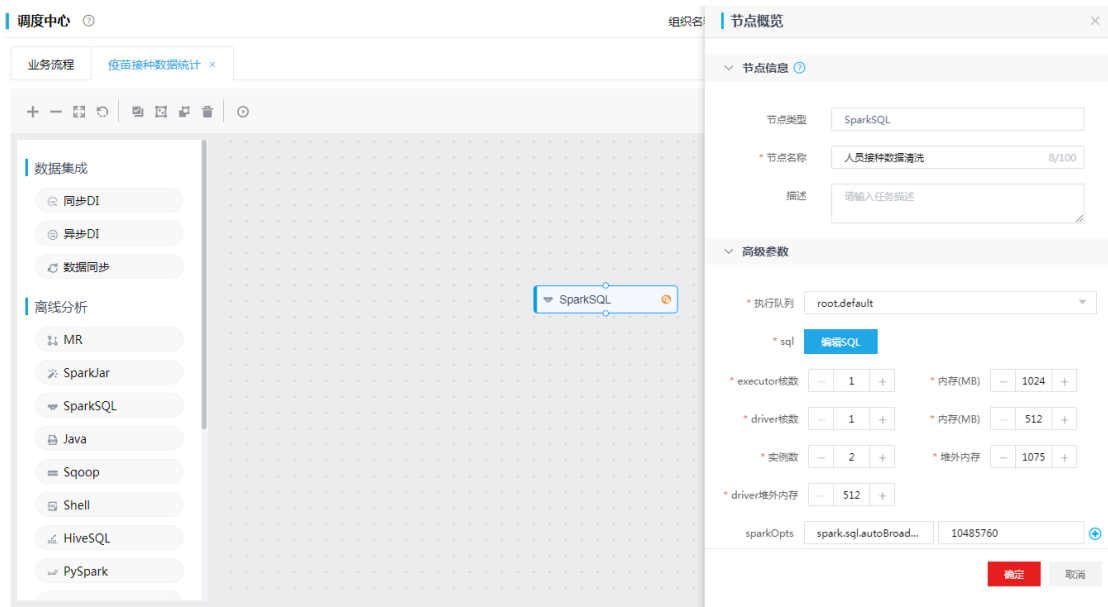
本例中需要增加人员接种数据清洗作业和各数据计算作业。

### 2. 添加数据清洗作业

业务流程创建后，需要在业务流程的画布中增加人员接种数据清洗作业。

- (1) 在业务流程的画布编辑页签中，选择左侧离线分析下的 **SparkSQL** 组件，并拖入画布中。
- (2) 双击画布中的 **SparkSQL** 作业节点，弹出作业节点参数编辑窗口。
- (3) 本例中，配置节点名称为“人员接种数据清洗”，选择执行队列为缺省队列。

图5-29 配置节点参数



(4) 单击<编辑 SQL>按钮，在弹出框中编写 SQL 语句，示例如下：

```
insert into
  default.dwb_filtered_person_inoculation_d
select
  *
from
  default.ods_d_inter_person_inoculation_d
where
  p_id is not null
  and haveflag is not null
  and curzhenshu is not null
  and created is not null;
```

图5-30 配置 SQL 语句



(5) 编写完成，并通过语法校验后，单击<确定>按钮，保存 SQL 语句。

(6) 单击<确定>按钮，数据清洗作业节点配置完成。

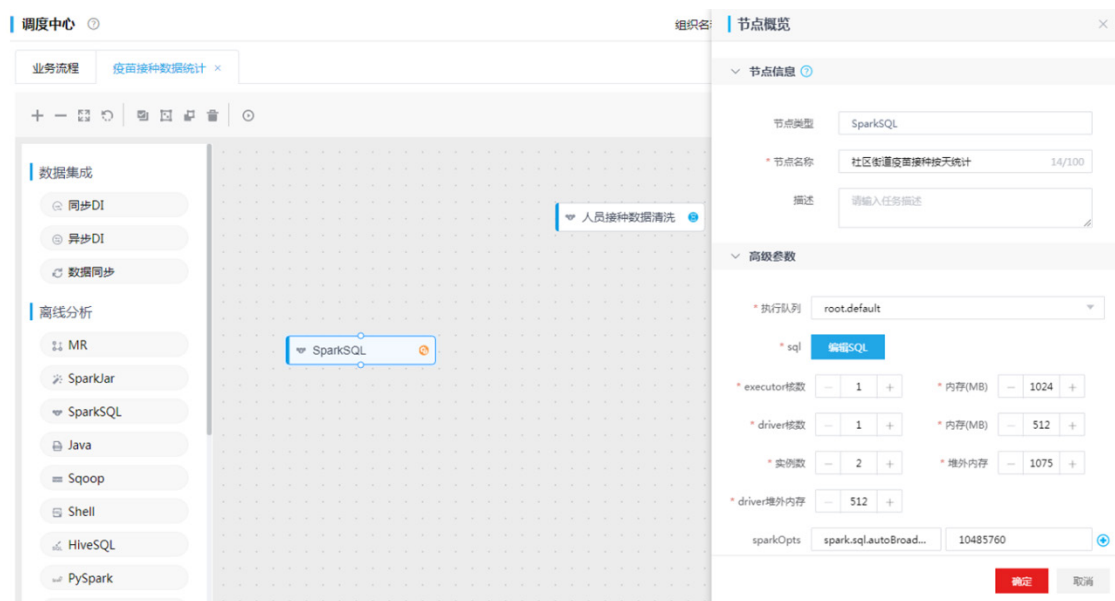
### 3. 添加数据计算作业

(1) 在业务流程的画布编辑页签中，选择左侧离线分析下的 SparkSQL 组件，并拖入画布中。

(2) 双击画布中的 SparkSQL 作业节点，弹出作业节点参数编辑窗口。

(3) 本例中，配置节点名称为“社区街道疫苗接种按天统计”，选择执行队列为缺省队列。

图5-31 配置节点参数



(4) 单击<编辑 SQL>按钮，在弹出框中编写 SQL 语句，示例如下：

```
select
  f.region as f_region,
  f.streetId as f_streetId,
  f.firstNum as f_firstNum,
```

```

s.region as s_region,
s.streetId as s_streetId,
s.secondNum as s_secondNum,
f.f_inoculation_date as inoculation_date
from
(
  select
    case
      when rd.region is not null then rd.region
      else '无区域'
    end as region,
    case
      when p.streetId is not null then p.streetId
      else '无街道'
    end as streetId,
    pi.f_inoculation_date,
    count(pi.p_id) as firstNum
  from
    (select *, substring_index(created, ' ', 1) as f_inoculation_date from
    default.dwb_filtered_person_inoculation_d) pi
    left join ods_d_inter_person_d p on pi.p_id = p.id
    left join ods_d_inter_region_dict_d rd on rd.id = p.region_id
  where
    pi.haveflag = 'true'
    and curzhenshu = '0'
  group by
    rd.region,
    p.streetId,
    pi.f_inoculation_date
) f full
join (
  select
    case
      when rd.region is not null then rd.region
      else '无区域'
    end as region,
    case
      when p.streetId is not null then p.streetId

```

```

        else '无街道'
    end as streetId,
    pi.s_inoculation_date,
    count(pi.p_id) as secondNum
from
    (select *, substring_index(created, ' ', 1) as s_inoculation_date from
default.dwb_filtered_person_inoculation_d) pi
    left join ods_d_inter_person_d p on pi.p_id = p.id
    left join ods_d_inter_region_dict_d rd on rd.id = p.region_id
where
    pi.haveflag = 'true'
    and curzhenshu = '1'
group by
    rd.region,
    p.streetId,
    pi.s_inoculation_date
) s on f.region = s.region
    and f.streetId = s.streetId and f.f_inoculation_date = s.s_inoculation_date;

```

- (5) 编写完成，并通过语法校验后，单击<确定>按钮，保存 SQL 语句。
- (6) 单击<确定>按钮，社区街道疫苗接种按天统计结果节点配置完成。
- (7) 依次增加其他数据计算作业，各作业使用的 SQL 语句如所示。

表5-12 各数据计算作业使用的 SQL 语句

作业	SQL 语句
社区街道疫苗接种 全量统计	<pre> select f.f_region as f_region, f.f_streetId as f_streetId, f.first_toal_num as first_total_num, s.s_region as s_region, s.s_streetId as s_streetId, s.second_toal_num as second_total_num from (select f_region, f_streetId, sum(f_firstNum) as first_toal_num from default.dws_region_inoculation_day_statistics group by f_region, f_streetId) as f full join (select s_region, s_streetid, sum(s_secondNum) as second_toal_num from default.dws_region_inoculation_day_statistics group by s_region, s_streetid) as s on f.f_region = s.s_region and f.f_streetId = s.s_streetId ; </pre>
行业疫苗接种按天 统计	<pre> select     f.classification as f_classification, </pre>

作业	SQL 语句
	<pre> f.firstNum as firstNum, s.classification as s_classification, s.secondNum as secondNum, f.f_inoculation_date as inoculation_date from(   select     case       when pc.classification is not null then pc.classification       else '未分类人员' end as classification,     pi.f_inoculation_date as f_inoculation_date,     count(pi.p_id) as firstNum   from     (select *, substring_index(created, ' ', 1) as f_inoculation_date from default.dwb_filtered_person_inoculation_d) pi     left join ods_d_inter_person_d p on pi.p_id = p.id     left join ods_d_inter_person_classification_d pc on pc.id = p.classification_id   where     pi.haveflag = 'true'     and pi.curzhenshu = '0'   group by     pc.classification,     pi.f_inoculation_date ) f full join (   select     case       when pc.classification is not null then pc.classification       else '未分类人员' end as classification,     pi.s_inoculation_date as s_inoculation_date,     count(pi.p_id) as secondNum   from     (select *, substring_index(created, ' ', 1) as s_inoculation_date from default.dwb_filtered_person_inoculation_d) pi </pre>

作业	SQL 语句
	<pre> left join ods_d_inter_person_d p on pi.p_id = p.id left join ods_d_inter_person_classification_d pc on pc.id = p.classification_id where pi.haveflag = 'true' and pi.curzhenshu = '1' group by pc.classification, pi.s_inoculation_date ) s on f.classification = s.classification and f.f_inoculation_date = s.s_inoculation_date; </pre>
行业疫苗接种全量统计	<pre> select f.f_classification as f_classification, f.first_toal_num as first_toal_num, s.s_classification as s_classification, s.second_toal_num as second_toal_num from ( select f_classification, sum(firstNum) as first_toal_num from default.dws_classification_inoculation_day_statistics group by f_classification ) as f full join ( select s_classification, sum(secondNum) as second_toal_num from default.dws_classification_inoculation_day_statistics group by </pre>

作业	SQL 语句
	<pre>s_classification ) as s on f.f_classification = s.s_classification;</pre>
<p>各年龄段疫苗接种 按天统计</p>	<pre>select   f.age_range as f_ageRange,   f.firstNum as firstNum,   s.age_range as s_ageRange,   s.secondNum as secondNum,   f.f_inoculation_date as inoculation_date from   (     select       p.age_range as age_range,       pi.f_inoculation_date as f_inoculation_date,       count(pi.p_id) as firstNum     from       (         select           *,           substring_index(created, ' ', 1) as f_inoculation_date         from           default.dwb_filtered_person_inoculation_d       ) pi     left join (       select         case           when age &lt;= 18 then '18岁以下'           when age &lt;= 59             and age &gt; 18 then '18岁至59岁'           when age &gt; 59 then '大于59岁'           else '未找到年龄信息'         end as age_range,         id</pre>



作业	SQL 语句
	<pre> from ods_d_inter_person_d ) p on pi.p_id = p.id where pi.haveflag = 'true' and pi.curzhenshu = '0' group by p.age_range, pi.f_inoculation_date ) f full join ( select p.age_range as age_range, pi.s_inoculation_date as s_inoculation_date, count(pi.p_id) as secondNum from ( select *, substring_index(created, ' ', 1) as s_inoculation_date from default.dwb_filtered_person_inoculation_d ) pi left join ( select case when age &lt;= 18 then '18岁以下' when age &lt;= 59 and age &gt; 18 then '18岁至59岁' when age &gt; 59 then '大于59岁' else '未找到年龄信息' end as age_range, </pre>

作业	SQL 语句
	<pre> id from ods_d_inter_person_d ) p on pi.p_id = p.id where pi.haveflag = 'true' and pi.curzhenshu = '1' group by p.age_range, pi.s_inoculation_date ) s on s.age_range = f.age_range and f.f_inoculation_date = s.s_inoculation_date; </pre>
各年龄段疫苗接种 全量统计	<pre> select f.f_ageRange as f_ageRange, f.first_total_num as first_total_num, s.s_ageRange as s_ageRange, s.second_total_num as second_total_num from ( select f_ageRange, sum(firstNum) as first_total_num from default.dws_age_range_inoculation_day_statistics group by f_ageRange ) as f full join ( select s_ageRange, sum(secondNum) as second_total_num from default.dws_age_range_inoculation_day_statistics group by s_ageRange ) as s on f.f_ageRange = s.s_ageRange; </pre>

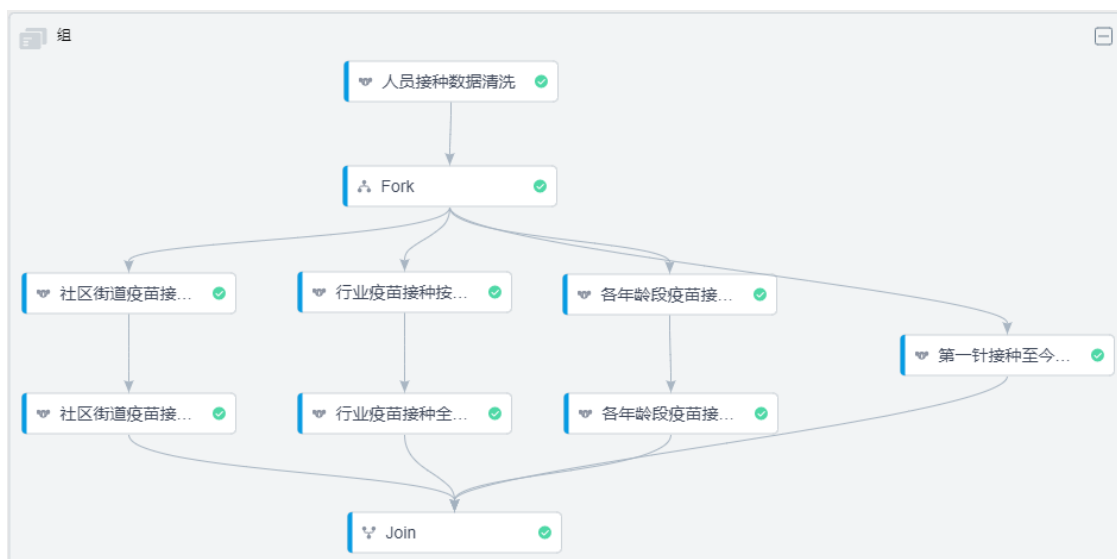
作业	SQL 语句
第一针接种至今各个时间间隔人数统计	<pre> select   case     when datediff(       date_format(CURRENT_DATE, 'yyyy-MM-dd'),       replace(substring_index(created, ' ', 1), '/', '-')     ) &lt;= 21 then '三周以内'     when datediff(       date_format(CURRENT_DATE, 'yyyy-MM-dd'),       replace(substring_index(created, ' ', 1), '/', '-')     ) &gt; 21     and datediff(       date_format(CURRENT_DATE, 'yyyy-MM-dd'),       replace(substring_index(created, ' ', 1), '/', '-')     ) &lt;= 56 then '三周至八周'     else '超过八周'   end as interval_period,   count(p_id) as person_num from   default.dwb_filtered_person_inoculation_d where curzhenshu = 0 and haveflag = 'true' group by   interval_period; </pre>


#### 4. 构建完成作业并运行

各作业创建完成后,需要通过控制节点下的 **Fork** 组件和 **Join** 组件进行连接,构建完整的业务流程。

- (1) 在业务流程的画布编辑页签中,选择左侧控制节点下的 **Fork** 组件和 **Join** 组件,并拖入画布中。
- (2) 依次连接各作业,连接结果如 [图 5-32](#) 所示。

图5-32 关联作业



- (3) 连接完成后，即可单击画布上方的  按钮，运行该业务流程。针对业务流程中的各作业，可配置调度策略，使作业定期自动运行。

## 5.2.7 数据查询

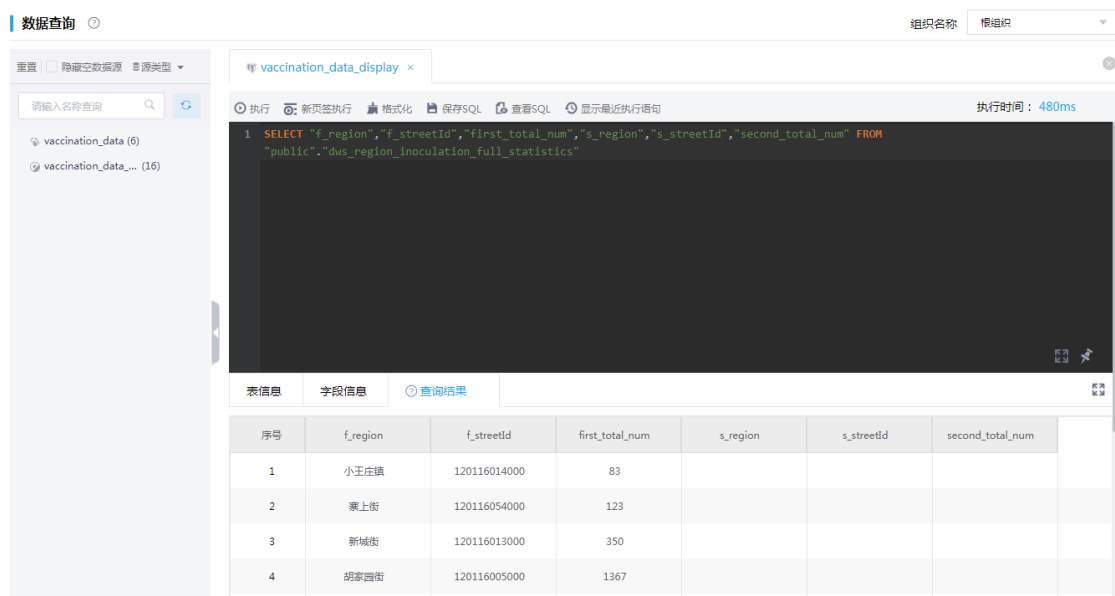
当业务流程运行完成后，统计结果数据会存入预先创建的统计结果表中。

数据管理平台中提供了数据查询功能，可查询统计结果数据。

- (1) 在[数据管理平台/数据开发]模块中，选择左侧导航树中的[数据查询]菜单项，进入数据查询页面。
- (2) 在左侧目录中选择数据源，右侧出现该数据源的数据查询页签。
- (3) 在输入框中输入 SQL 查询语句，查看各统计结果表中的数据。示例语句如下：

```
SELECT "f_region","f_streetId","first_total_num","s_region","s_streetId","second_total_num"
FROM "public"."dws_region_inoculation_full_statistics"
```
- (4) 单击<执行>按钮，执行该查询语句，下方的查询结果页签中，会展示该统计结果表中的数据。

图5-33 查询结果



## 5.2.8 结果数据发布

通过数据计算得出的统计数据存入了统计结果表中，数字平台支持以表为单位，在集成平台的服务集成功能中，将 Greenplum 数据源中的统计结果表发布，并授权给特定的工作空间，以便于第三方应用通过 URL 获取数据。

### 1. API 注册

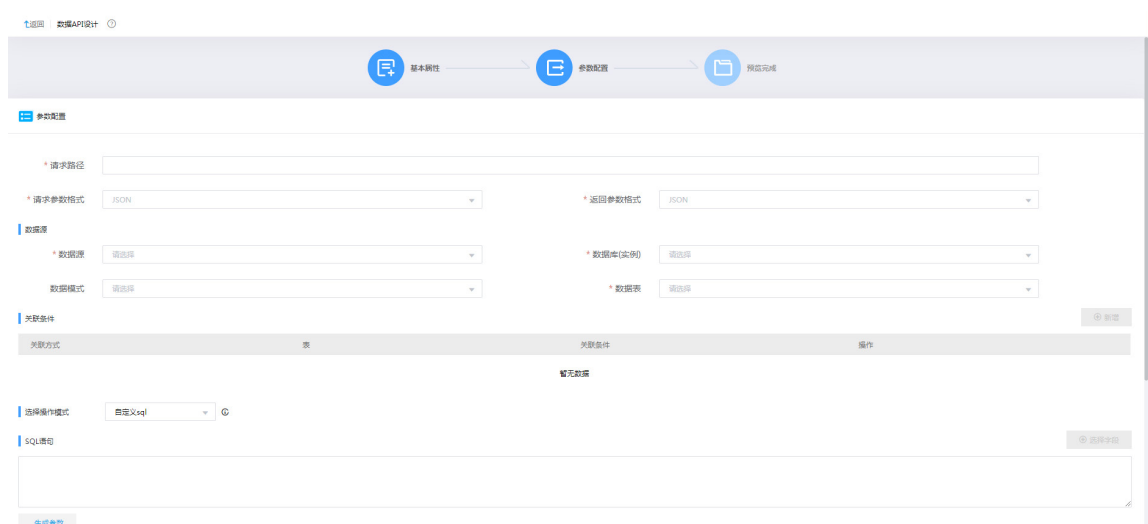
在服务集成模块下的[API 工厂/API 列表]页面，单击<API 注册>按钮，选择注册“数据 API”类型。

图5-34 API 注册



选择 [5.2.5 3. 新建结果表](#) 中的结果表（注意：注册的一个 API 只能发布一张表）作为发布数据对象。配置数据 API 基本属性、选择需要发布的数据表。

图5-35 数据 API 设计



注册完成后，单击新生成的 API 右侧按钮<测试>，对接口进行测试，如下图所示，测试接口是否可用。

图5-36 API 测试



完成测试后，当前 API 状态即为待部署状态。单击右侧的<部署>按钮，即可在弹窗页面中配置部署节点。部署完成后，即可进行 API 授权操作。

## 2. API 授权

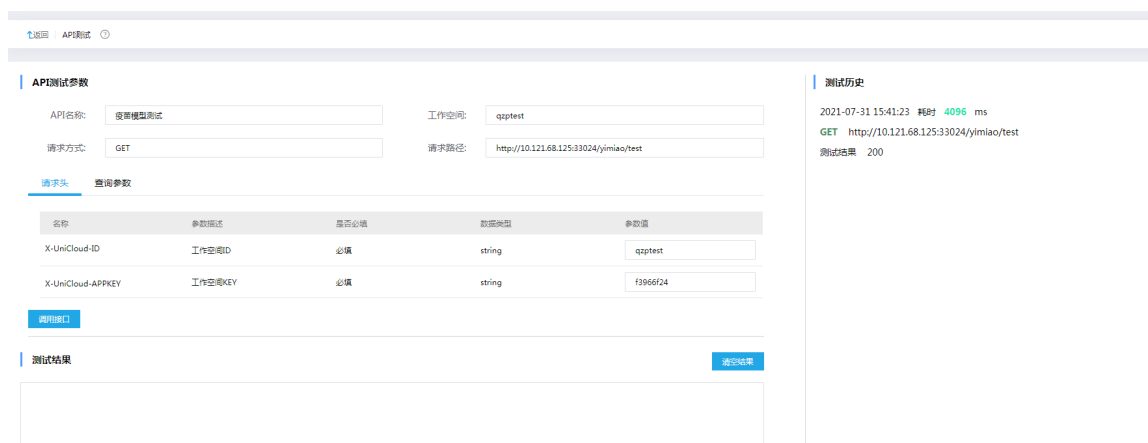
服务集成模块下的[API 网关/API 列表]页面，在列表中选择上一步注册的 API，然后单击右侧的<授权>按钮，进入 API 授权页面，配置需要授权的工作空间。授权完毕以后，在 API 授权页面下方会出现已授权的工作空间，单击操作列中的<测试>按钮。

图5-37 已授权工作空间



单击<测试>按钮后，进入 API 测试页面，API 测试页面完整显示了包括 IP 地址、端口号在内的完整访问路径，第三方应用只需访问此 URL 即可获取数据，无需在请求中携带任何用户信息。

图5-38 API 测试页面



## 5.2.9 数据最终呈现

该案例中的统计数据发布后，支持通过第三方调用展示，如[图 5-39](#)和[图 5-40](#)所示。

图5-39 疫苗接种情况展示（一）



图5-40 疫苗接种情况展示（二）





# 6 常见问题解答

## 6.1 HBase、Hive、HDFS、Kafka等大数据组件开启了Kerberos认证，连接这些数据源时如何配置Kerberos认证信息

- (1) 登录大数据集群任一节点。
- (2) 拷贝 krb5.conf 文件到本地，krb5.conf 文件路径：/etc/krb5.conf。
- (3) 拷贝 keytab 文件到本地，keytab 文件路径：/etc/security/keytabs/xxx.service.keytab。
- (4) 获取 kerberos 用户名。以 Kafka 为例：执行 **klist -kt /etc/security/keytabs/kafka.service.keytab** 命令，查询出的 Principal(kafka/kflt1.hde.com@KFLTZC.COM)即为 Kerberos 用户名，如[图 6-1](#)。

图6-1 获取 Kerberos 用户名

```
[root@kflt1 config]# klist -kt /etc/security/keytabs/kafka.service.keytab
Keytab name: FILE:/etc/security/keytabs/kafka.service.keytab
KVNO Timestamp Principal
-----
2 11/19/2020 00:13:09 kafka/kflt1.hde.com@KFLTZC.COM
2 11/19/2020 00:13:09 kafka/kflt1.hde.com@KFLTZC.COM
2 11/19/2020 00:13:09 kafka/kflt1.hde.com@KFLTZC.COM
2 11/19/2020 00:13:09 kafka/kflt1.hde.com@KFLTZC.COM
2 11/19/2020 00:13:09 kafka/kflt1.hde.com@KFLTZC.COM
```

- (5) 配置 Kerberos 信息时，填写 Kerberos 用户名，上传 krb5.conf 和 keytab 文件，如[图 6-2](#)。

图6-2 配置 Kerberos 认证相关信息



新增数据源 ?

所属集群/共享资源

\* bootstrap服务IP 127.0.0.1

\* bootstrap端口 6667

\* kafka版本 10

Kerberos认证

\* 登录用户 kafka/kft1.hde.com@KFLTZC.COM

\* krb5.conf路径 krb5.conf

\* keytab文件路径 kafka.service.keytab

\* 数据源范围  内部数据源  外部数据源

描述信息 0/512

提交 测试连接 取消

(6) 配置完成后，测试连接并提交即可完成数据源的新增。

# 7 附录

## 7.1 业务数据库建表语句示例

在业务数据库中，可通过 SQL 语句创建记录原始数据的表，本节提供了参考示例。

### 1. 创建人员信息表的 SQL 语句示例

```
CREATE TABLE 'person' (  
  'id' int(11) NOT NULL AUTO_INCREMENT COMMENT '序号',  
  'name' varchar(255) DEFAULT NULL COMMENT '姓名',  
  'sex' varchar(255) DEFAULT NULL COMMENT '性别',  
  'age' int(11) DEFAULT NULL COMMENT '年龄',  
  'mobile' varchar(255) DEFAULT NULL COMMENT '手机号',  
  'cardno' varchar(20) DEFAULT NULL COMMENT '身份证号',  
  'classification_id' int(11) DEFAULT NULL COMMENT '人员分类(一级)',  
  'content' varchar(255) DEFAULT NULL COMMENT '备注',  
  'company_id' int(11) DEFAULT NULL COMMENT '单位名称',  
  'region_id' int(11) DEFAULT NULL COMMENT '辖区 id，对应属地的摸底工作部门',  
  'declare_department_id' int(11) DEFAULT NULL COMMENT '申报部门 id',  
  'created_time' datetime DEFAULT NULL COMMENT '创建时间',  
  'modified_time' varchar(255) DEFAULT NULL COMMENT '修改时间',  
  'uuid' varchar(50) DEFAULT NULL COMMENT 'UUID',  
  'addr' varchar(500) DEFAULT NULL COMMENT '现住址',  
  'streetId' varchar(20) DEFAULT NULL COMMENT '街道',  
  'provinceId' varchar(20) DEFAULT NULL COMMENT '省 ID',  
  'cityId' varchar(20) DEFAULT NULL COMMENT '市 ID',  
  'districtId' varchar(20) DEFAULT NULL COMMENT '区域 ID',  
  'flag' varchar(10) DEFAULT NULL COMMENT '是否本市住户',  
  'area' varchar(255) DEFAULT NULL COMMENT '小区名字',  
  'subclass_id' int(11) DEFAULT NULL COMMENT '人群分类（二级）',  
  PRIMARY KEY ('id') USING BTREE,  
  UNIQUE KEY 'class_index' ('id','classification_id') USING BTREE,  
  KEY 'compant_index' ('company_id') USING BTREE,  
  KEY 'region_id_index' ('region_id') USING BTREE  
) ENGINE=InnoDB AUTO_INCREMENT=1649514 DEFAULT CHARSET=utf8 ROW_FORMAT=DYNAMIC  
COMMENT='人员信息'
```

## 2. 创建人员分类字典表的 SQL 语句示例

```
CREATE TABLE 'person_classification' (  
  'id' int(11) NOT NULL COMMENT 'id',  
  'classification' text COMMENT '人员分类',  
  'created_time' datetime DEFAULT NULL COMMENT '创建时间',  
  'modified_time' datetime DEFAULT NULL COMMENT '修改时间',  
  'category' varchar(255) DEFAULT NULL COMMENT '级别',  
  PRIMARY KEY ('id') USING BTREE  
) ENGINE=InnoDB DEFAULT CHARSET=utf8 ROW_FORMAT=DYNAMIC COMMENT='人员分类字典表'
```

## 3. 创建辖区字典表的 SQL 语句示例

```
CREATE TABLE 'region_dict' (  
  'id' int(11) NOT NULL AUTO_INCREMENT COMMENT '序号',  
  'region' varchar(255) DEFAULT NULL COMMENT '辖区',  
  'userid' varchar(50) DEFAULT NULL COMMENT 'userid',  
  'category' varchar(255) DEFAULT NULL COMMENT '有没有二级分类(1 是有)',  
  'sort' int(11) DEFAULT NULL COMMENT '顺序编号',  
  PRIMARY KEY ('id') USING BTREE,  
  UNIQUE KEY 'region' ('region') USING HASH COMMENT 'region 索引'  
) ENGINE=InnoDB AUTO_INCREMENT=27 DEFAULT CHARSET=utf8 ROW_FORMAT=DYNAMIC COMMENT='辖区字典表'
```

## 4. 创建人员接种信息表的 SQL 语句示例

```
CREATE TABLE 'person_inoculation' (  
  'p_id' int(11) NOT NULL COMMENT 'ID',  
  'assetsNum' varchar(100) DEFAULT NULL COMMENT '设备编码',  
  'haveflag' varchar(10) DEFAULT NULL COMMENT '接种状态 true/false',  
  'reason' varchar(800) DEFAULT NULL COMMENT '状态: 0:禁忌症,1:延期接种,2:需退回上级重新分配,3:不符合接种人群范围',  
  'beizhu' text COMMENT '备注',  
  'curzhenshu' varchar(20) DEFAULT NULL COMMENT '针次',  
  'jdate' varchar(20) DEFAULT NULL,  
  'stime' varchar(20) DEFAULT NULL,  
  'zhenshu' varchar(20) DEFAULT NULL,  
  'yimiao' varchar(10) DEFAULT NULL COMMENT '疫苗种类 0:北京生物,1:北京科兴,2:武汉生物,3:康希诺',  
  'jinjizheng' varchar(255) DEFAULT NULL COMMENT '禁忌症',  
  'created' datetime DEFAULT NULL COMMENT '接种时间',
```

```
KEY 'class_index' ('p_id') USING BTREE
) ENGINE=InnoDB DEFAULT CHARSET=utf8 ROW_FORMAT=DYNAMIC
```

## 7.2 stg\_2\_ods\_cep\_stcdb\_mdtrt\_d.sql脚本内容

```
set hive.exec.max.dynamic.partitions=100000;
set hive.exec.max.dynamic.partitions.pernode=100000;
set hive.exec.dynamic.partition.mode=nonstrict;

use ods_prd;
CACHE TABLE STG_CEP_STCDB_MDTRT_D_cached as
with ETLDATETEMP as (
select
    date_format(CRTE_TIME,'yyyy-MM') as ETL_DATE
    ,case when POOLAREA_NO is null or POOLAREA_NO = '' or length(POOLAREA_NO) <> 6 or
substr(POOLAREA_NO,1,2) <> '50' then '500000' else POOLAREA_NO end as REGION
from stg_prd.STG_CEP_STCDB_MDTRT_D
where DT ='${etl_date}'
group by
    date_format(CRTE_TIME,'yyyy-MM')
    ,case when POOLAREA_NO is null or POOLAREA_NO = '' or length(POOLAREA_NO) <> 6 or
substr(POOLAREA_NO,1,2) <> '50' then '500000' else POOLAREA_NO end
)
,STG_QUALITY_CONTROLL as (
SELECT
    A.MDTRT_ID,
    A.MEDINS_SETL_ID,
    A.PSN_NO,
    A.PSN_INSU_RLTS_ID,
    A.PSN_CERT_TYPE,
    A.CERTNO,
    A.PSN_NAME,
    A.GEND,
    A.NATY,
    A.BRDY,
    A.AGE,
    A.CONER_NAME,
    A.TEL,
    A.ADDR,
```

A.INSUTYPE,  
A.PSN\_TYPE,  
A.CVLSERV\_FLAG,  
A.CVLSERV\_LV,  
A.SP\_PSN\_TYPE,  
A.SP\_PSN\_TYPE\_LV,  
A.CLCT\_GRDE,  
A.FLXEMPE\_FLAG,  
A.NWB\_FLAG,  
A.INSU\_ADMDVS,  
A.EMP\_NO,  
A.EMP\_NAME,  
A.EMP\_TYPE,  
A.ECON\_TYPE,  
A.AFIL\_INDU,  
A.AFIL\_RLTS,  
A.EMP\_MGT\_TYPE,  
A.PAY\_LOC,  
A.FIXMEDINS\_CODE,  
A.FIXMEDINS\_NAME,  
A.HOSP\_LV,  
A.FIX\_BLNG\_ADMDVS,  
A.LMTPRIC\_HOSP\_LV,  
A.DEDC\_HOSP\_LV,  
A.BEGNTIME,  
A.ENDTIME,  
A.MDTRT\_CERT\_TYPE,  
A.MDTRT\_CERT\_NO,  
A.MED\_TYPE,  
A.RLOC\_TYPE,  
A.ARS\_YEAR\_IPT\_FLAG,  
A.PRE\_PAY\_FLAG,  
A.YEAR,  
A.REFL\_OLD\_MDTRT\_ID,  
A.IPT\_OTP\_NO,  
A.MEDRCDNO,  
A.CHFPDR\_CODE,  
A.ADM\_DIAG\_DSCR,

A.ADM\_DEPT\_CODG,  
A.ADM\_DEPT\_NAME,  
A.ADM\_BED,  
A.WARDAREA\_BED,  
A.TRAF\_DEPT\_FLAG,  
A.DSCG\_MAINDIAG\_CODE,  
A.DSCG\_DEPT\_CODG,  
A.DSCG\_DEPT\_NAME,  
A.DSCG\_BED,  
A.DSCG\_WAY,  
A.MAIN\_COND\_DSCR,  
A.DISE\_NO,  
A.DISE\_NAME,  
A.OPRN\_OPRT\_CODE,  
A.OPRN\_OPRT\_NAME,  
A.OTP\_DIAG\_INFO,  
A.INHOSP\_STAS,  
A.DIE\_DATE,  
A.IPT\_DAYS,  
A.GESO\_VAL,  
A.BIRCTRL\_TYPE,  
A.FETTS,  
A.FETUS\_CNT,  
A.MATN\_TYPE,  
A.PREY\_TIME,  
A.LATECHB\_FLAG,  
A.PRET\_FLAG,  
A.FPSC\_NO,  
A.BIRCTRL\_MATN\_DATE,  
A.COP\_FLAG,  
A.TRT\_DCLA\_DETL\_SN,  
A.VALI\_FLAG,  
A.MEMO,  
A.RID,  
A.UPDT\_TIME,  
A.CRTER\_ID,  
A.CRTER\_NAME,  
A.CRTE\_TIME,

```

A.CRTE_OPTINS_NO,
A.OPTER_ID,
A.OPTER_NAME,
A.OPT_TIME,
A.OPTINS_NO,
A.POOLAREA_NO,
A.CHPDR_NAME,
A.DSCG_MAINDIAG_NAME,
A.ADM_CATY,
A.DSCG_CATY,
A.TTP_PAY_FLAG,
A.TTP_PAY_PROP,
A.DISE_TYPE_CODE,
A.SAME_DISE_ADM_FLAG,
A.QUTS_TYPE,
A.SUBSYS_CODG_SRC,
CONCAT_WS(chr(1)
    ,CHECK_COLUMN_LOGIC(CAST( CASE WHEN (A.BEGNTIME) IS NULL THEN " ELSE
UNIX_TIMESTAMP(A.BEGNTIME) END AS STRING), '<=', CAST( CASE WHEN (A.ENDTIME) IS NULL THEN "
ELSE UNIX_TIMESTAMP(A.ENDTIME) END AS STRING), 'DATETIME', CONCAT( A.BEGNTIME,'|',A.ENDTIME),
'S08202102260432', '开始时间不能大于结束时间', '行级', '1')
    ,CHECK_COLUMN_DICT( CASE WHEN A.AFIL_INDU IS NULL THEN " ELSE A.AFIL_INDU END ,
'AFIL_INDU', 'T', 'SJZD202103011542', '就诊信息表(MDTRT_D)中所属行业(AFIL_INDU)数据值应在数据
字典【所属行业 AFIL_INDU】范围内', '行级', '1')
    ,CHECK_COLUMN_NULL( CASE WHEN A.CRTE_TIME IS NULL THEN " ELSE A.CRTE_TIME END ,
'L0000000001', '创建时间不能为空', '行级', '3')
    ,CHECK_COLUMN_NULL( CASE WHEN A.RID IS NULL THEN " ELSE A.RID END , 'L0000000002', '
数据唯一编号不能为空', '行级', '3')
) AS DQ_RESULT
FROM stg_prd.STG_CEP_STCDB_MDTRT_D A
WHERE A.DT = '${etl_date}'
)
-- 分流
select
    A.MDTRT_ID,
    A.MEDINS_SETL_ID,
    A.PSN_NO,
    A.PSN_INSU_RLTS_ID,
    A.PSN_CERT_TYPE,

```



A.CERTNO,  
A.PSN\_NAME,  
A.GEND,  
A.NATY,  
A.BRDY,  
A.AGE,  
A.CONER\_NAME,  
A.TEL,  
A.ADDR,  
A.INSUTYPE,  
A.PSN\_TYPE,  
A.CVLSERV\_FLAG,  
A.CVLSERV\_LV,  
A.SP\_PSN\_TYPE,  
A.SP\_PSN\_TYPE\_LV,  
A.CLCT\_GRDE,  
A.FLXEMPE\_FLAG,  
A.NWB\_FLAG,  
A.INSU\_ADMDVS,  
A.EMP\_NO,  
A.EMP\_NAME,  
A.EMP\_TYPE,  
A.ECON\_TYPE,  
A.AFIL\_INDU,  
A.AFIL\_RLTS,  
A.EMP\_MGT\_TYPE,  
A.PAY\_LOC,  
A.FIXMEDINS\_CODE,  
A.FIXMEDINS\_NAME,  
A.HOSP\_LV,  
A.FIX\_BLNG\_ADMDVS,  
A.LMTPRIC\_HOSP\_LV,  
A.DEDC\_HOSP\_LV,  
A.BEGNTIME,  
A.ENDTIME,  
A.MDTRT\_CERT\_TYPE,  
A.MDTRT\_CERT\_NO,  
A.MED\_TYPE,

A.RLOC\_TYPE,  
A.ARS\_YEAR\_IPT\_FLAG,  
A.PRE\_PAY\_FLAG,  
A.YEAR,  
A.REFL\_OLD\_MDRTRT\_ID,  
A.IPT\_OTP\_NO,  
A.MEDRCDNO,  
A.CHFPDR\_CODE,  
A.ADM\_DIAG\_DSCR,  
A.ADM\_DEPT\_CODG,  
A.ADM\_DEPT\_NAME,  
A.ADM\_BED,  
A.WARDAREA\_BED,  
A.TRAF\_DEPT\_FLAG,  
A.DSCG\_MAINDIAG\_CODE,  
A.DSCG\_DEPT\_CODG,  
A.DSCG\_DEPT\_NAME,  
A.DSCG\_BED,  
A.DSCG\_WAY,  
A.MAIN\_COND\_DSCR,  
A.DISE\_NO,  
A.DISE\_NAME,  
A.OPRN\_OPRT\_CODE,  
A.OPRN\_OPRT\_NAME,  
A.OTP\_DIAG\_INFO,  
A.INHOSP\_STAS,  
A.DIE\_DATE,  
A.IPT\_DAYS,  
A.GESO\_VAL,  
A.BIRCTRL\_TYPE,  
A.FETTS,  
A.FETUS\_CNT,  
A.MATN\_TYPE,  
A.PREY\_TIME,  
A.LATECHB\_FLAG,  
A.PRET\_FLAG,  
A.FPSC\_NO,  
A.BIRCTRL\_MATN\_DATE,

```

A.COP_FLAG,
A.TRRT_DCLA_DETL_SN,
A.VALI_FLAG,
A.MEMO,
A.RID,
A.UPDT_TIME,
A.CRTER_ID,
A.CRTER_NAME,
A.CRTE_TIME,
A.CRTE_OPTINS_NO,
A.OPTER_ID,
A.OPTER_NAME,
A.OPT_TIME,
A.OPTINS_NO,
A.POOLAREA_NO,
A.CHFPDR_NAME,
A.DSCG_MAINDIAG_NAME,
A.ADM_CATY,
A.DSCG_CATY,
A.TTP_PAY_FLAG,
A.TTP_PAY_PROP,
A.DISE_TYPE_CODE,
A.SAME_DISE_ADM_FLAG,
A.QUTS_TYPE,
A.SUBSYS_CODG_SRC,
A.DTY_FLAG,
A.LOCAL_DTY_FLAG,
A.EXCH_UPDT_TIME,
" AS DQ_RESULT,
A.DT as ETL_DATE,
A.REGION as REGION
FROM ods_prd.ODS_CEP_STCDB_MDRTRT_D A
INNER JOIN ETLDATETEMP B
    ON A.DT = B.ETL_DATE
    AND A.REGION = B.REGION
LEFT JOIN STG_QUALITY_CONTROLL C
    ON A.MDRTRT_ID = C.MDRTRT_ID AND A.PSN_NO = C.PSN_NO
WHERE C.MDRTRT_ID is null AND C.PSN_NO is null

```

union ALL

select

MDTRT\_ID,  
MEDINS\_SETL\_ID,  
PSN\_NO,  
PSN\_INSU\_RLTS\_ID,  
PSN\_CERT\_TYPE,  
CERTNO,  
PSN\_NAME,  
GEND,  
NATY,  
BRDY,  
AGE,  
CONER\_NAME,  
TEL,  
ADDR,  
INSUTYPE,  
PSN\_TYPE,  
CVLSERV\_FLAG,  
CVLSERV\_LV,  
SP\_PSN\_TYPE,  
SP\_PSN\_TYPE\_LV,  
CLCT\_GRDE,  
FLXEMPE\_FLAG,  
NWB\_FLAG,  
INSU\_ADMDVS,  
EMP\_NO,  
EMP\_NAME,  
EMP\_TYPE,  
ECON\_TYPE,  
AFIL\_INDU,  
AFIL\_RLTS,  
EMP\_MGT\_TYPE,  
PAY\_LOC,  
FIXMEDINS\_CODE,  
FIXMEDINS\_NAME,  
HOSP\_LV,  
FIX\_BLNG\_ADMDVS,

LMTPRIC\_HOSP\_LV,  
DEDC\_HOSP\_LV,  
BEGNTIME,  
ENDTIME,  
MDTRT\_CERT\_TYPE,  
MDTRT\_CERT\_NO,  
MED\_TYPE,  
RLOC\_TYPE,  
ARS\_YEAR\_IPT\_FLAG,  
PRE\_PAY\_FLAG,  
YEAR,  
REFL\_OLD\_MDTRT\_ID,  
IPT\_OTP\_NO,  
MEDRCDNO,  
CHFPDR\_CODE,  
ADM\_DIAG\_DSCR,  
ADM\_DEPT\_CODG,  
ADM\_DEPT\_NAME,  
ADM\_BED,  
WARDAREA\_BED,  
TRAF\_DEPT\_FLAG,  
DSCG\_MAINDIAG\_CODE,  
DSCG\_DEPT\_CODG,  
DSCG\_DEPT\_NAME,  
DSCG\_BED,  
DSCG\_WAY,  
MAIN\_COND\_DSCR,  
DISE\_NO,  
DISE\_NAME,  
OPRN\_OPRT\_CODE,  
OPRN\_OPRT\_NAME,  
OTP\_DIAG\_INFO,  
INHOSP\_STAS,  
DIE\_DATE,  
IPT\_DAYS,  
GESO\_VAL,  
BIRCTRL\_TYPE,  
FETTS,

FETUS\_CNT,  
 MATN\_TYPE,  
 PREY\_TIME,  
 LATECHB\_FLAG,  
 PRET\_FLAG,  
 FPSC\_NO,  
 BIRCTRL\_MATN\_DATE,  
 COP\_FLAG,  
 TRT\_DCLA\_DETL\_SN,  
 VALI\_FLAG,  
 MEMO,  
 RID,  
 UPDT\_TIME,  
 CRTER\_ID,  
 CRTER\_NAME,  
 CRTE\_TIME,  
 CRTE\_OPTINS\_NO,  
 OPTER\_ID,  
 OPTER\_NAME,  
 OPT\_TIME,  
 OPTINS\_NO,  
 POOLAREA\_NO,  
 CHFPDR\_NAME,  
 DSCG\_MAINDIAG\_NAME,  
 ADM\_CATY,  
 DSCG\_CATY,  
 TTP\_PAY\_FLAG,  
 TTP\_PAY\_PROP,  
 DISE\_TYPE\_CODE,  
 SAME\_DISE\_ADM\_FLAG,  
 QUTS\_TYPE,  
 SUBSYS\_CODG\_SRC,  
 CASE WHEN INSTR(DQ\_RESULT, '"natResult":"FAIL"') > 0 THEN '1' ELSE '0' END DTY\_FLAG,  
 CASE WHEN INSTR(DQ\_RESULT, '"localResult":"FAIL"') > 0 THEN '1' ELSE '0' END LOCAL\_DTY\_FLAG,  
 concat(date\_add(current\_timestamp(),-1),' ',date\_format(current\_timestamp(),'HH:mm:ss')) AS  
 EXCH\_UPDT\_TIME,  
 DQ\_RESULT,  
 date\_format(CRTE\_TIME, 'yyyy-MM') AS ETL\_DATE,

```
        case when POOLAREA_NO is null or POOLAREA_NO = '' or length(POOLAREA_NO) <> 6 or
substr(POOLAREA_NO,1,2) <> '50' then '500000' else POOLAREA_NO end AS REGION
FROM STG_QUALITY_CONTROLL;
```

-- 写 ODS 表

```
REFRESH TABLE ods_prd.ODS_CEP_STCDB_MDTRT_D;
FROM STG_CEP_STCDB_MDTRT_D_cached
INSERT OVERWRITE TABLE ods_prd.ODS_CEP_STCDB_MDTRT_D PARTITION(DT,REGION)
SELECT
    MDTRT_ID,
    MEDINS_SETL_ID,
    PSN_NO,
    PSN_INSU_RLTS_ID,
    PSN_CERT_TYPE,
    CERTNO,
    PSN_NAME,
    GEND,
    NATY,
    BRDY,
    AGE,
    CONER_NAME,
    TEL,
    ADDR,
    INSUTYPE,
    PSN_TYPE,
    CVLSERV_FLAG,
    CVLSERV_LV,
    SP_PSN_TYPE,
    SP_PSN_TYPE_LV,
    CLCT_GRDE,
    FLXEMPE_FLAG,
    NWB_FLAG,
    INSU_ADMDVS,
    EMP_NO,
    EMP_NAME,
    EMP_TYPE,
    ECON_TYPE,
    AFIL_INDU,
```

AFIL\_RLTS,  
EMP\_MGT\_TYPE,  
PAY\_LOC,  
FIXMEDINS\_CODE,  
FIXMEDINS\_NAME,  
HOSP\_LV,  
FIX\_BLNG\_ADMDVS,  
LMTPRIC\_HOSP\_LV,  
DEDC\_HOSP\_LV,  
BEGNTIME,  
ENDTIME,  
MDTRT\_CERT\_TYPE,  
MDTRT\_CERT\_NO,  
MED\_TYPE,  
RLOC\_TYPE,  
ARS\_YEAR\_IPT\_FLAG,  
PRE\_PAY\_FLAG,  
YEAR,  
REFL\_OLD\_MDTRT\_ID,  
IPT\_OTP\_NO,  
MEDRCDNO,  
CHFPDR\_CODE,  
ADM\_DIAG\_DSCR,  
ADM\_DEPT\_CODG,  
ADM\_DEPT\_NAME,  
ADM\_BED,  
WARDAREA\_BED,  
TRAF\_DEPT\_FLAG,  
DSCG\_MAINDIAG\_CODE,  
DSCG\_DEPT\_CODG,  
DSCG\_DEPT\_NAME,  
DSCG\_BED,  
DSCG\_WAY,  
MAIN\_COND\_DSCR,  
DISE\_NO,  
DISE\_NAME,  
OPRN\_OPRT\_CODE,  
OPRN\_OPRT\_NAME,



OTP\_DIAG\_INFO,  
INHOSP\_STAS,  
DIE\_DATE,  
IPT\_DAYS,  
GESO\_VAL,  
BIRCTRL\_TYPE,  
FETTS,  
FETUS\_CNT,  
MATN\_TYPE,  
PREY\_TIME,  
LATECHB\_FLAG,  
PRET\_FLAG,  
FPSC\_NO,  
BIRCTRL\_MATN\_DATE,  
COP\_FLAG,  
TRT\_DCLA\_DETL\_SN,  
VALI\_FLAG,  
MEMO,  
RID,  
UPDT\_TIME,  
CRTER\_ID,  
CRTER\_NAME,  
CRTE\_TIME,  
CRTE\_OPTINS\_NO,  
OPTER\_ID,  
OPTER\_NAME,  
OPT\_TIME,  
OPTINS\_NO,  
POOLAREA\_NO,  
CHFPDR\_NAME,  
DSCG\_MAINDIAG\_NAME,  
ADM\_CATY,  
DSCG\_CATY,  
TTP\_PAY\_FLAG,  
TTP\_PAY\_PROP,  
DISE\_TYPE\_CODE,  
SAME\_DISE\_ADM\_FLAG,  
QUTS\_TYPE,

```

SUBSYS_CODG_SRC,
DTY_FLAG,
LOCAL_DTY_FLAG,
EXCH_UPDT_TIME,
ETL_DATE,
REGION
-- 写脏数据日志
INSERT OVERWRITE TABLE ods_prd.ODS_DTY_DETAIL_X PARTITION(DT,STG_TABLE_NAME)
SELECT
    case when POOLAREA_NO is null or POOLAREA_NO = '' or length(POOLAREA_NO) <> 6 or
substr(POOLAREA_NO,1,2) <> '50' then '500000' else POOLAREA_NO end AS REGION
    , 'MDTRT_D' AS TABLE_NAME
    , 'CEP' AS SUBSYS_CODG
    , GET_JSON_OBJECT(dq_result_t, '$.checkLv') AS VIO_DQ_LVL
    , '${etl_date}' BIZ_DATE
    , EXCH_UPDT_TIME
    , RID AS RID
    , CRTE_TIME
    , GET_JSON_OBJECT(dq_result_t, '$.rawData') RAW_DATA
    , GET_JSON_OBJECT(dq_result_t, '$.ruleCode') VIO_DQ_CODE
    , GET_JSON_OBJECT(dq_result_t, '$.checkRule') MEMO
    , '${etl_date}' AS ETL_DATE
    , 'STG_CEP_STCDB_MDTRT_D' AS STG_TABLE_NAME
lateral view explode(split(dq_result,chr(1))) num AS dq_result_t
WHERE (DTY_FLAG = '1' or LOCAL_DTY_FLAG = '1')
    and dq_result_t is not null and dq_result_t <> ''
;

DROP TABLE STG_CEP_STCDB_MDTRT_D_cached;
REFRESH TABLE ods_prd.ODS_CEP_STCDB_MDTRT_D;
REFRESH TABLE ods_prd.ODS_DTY_DETAIL_X;

```

### 7.3 ods\_2\_dwd\_cep\_stcdb\_mdtrt\_d.sql脚本内容

```

set hive.optimize.sort.dynamic.partition=true;
set hive.exec.dynamic.partition = true;
set hive.exec.dynamic.partition.mode = nonstrict;
set hive.exec.max.dynamic.partitions=100000;
set hive.exec.max.dynamic.partitions.pernode=100000;

```

```
set hive.compute.query.using.stats=false;
```

```
use dwd_prd;
```

```
create temporary table if not exists tmp_yibao_ods_dwd as
```

```
select
```

```
date_format(CRTE_TIME,'yyyy-MM') as cdate
```

```
,case when POOLAREA_NO is null or POOLAREA_NO = '' or length(POOLAREA_NO) <> 6 or
```

```
substr(POOLAREA_NO,1,2) <> '50' then '500000' else POOLAREA_NO end as area
```

```
from stg_prd.stg_cep_stcdb_mdtrt_d
```

```
where dt='${etl_date}'
```

```
group by date_format(CRTE_TIME,'yyyy-MM'),case when POOLAREA_NO is null or POOLAREA_NO = '' or
```

```
length(POOLAREA_NO) <> 6 or substr(POOLAREA_NO,1,2) <> '50' then '500000' else POOLAREA_NO
```

```
end ;
```

```
insert overwrite table dwd_prd.dwd_dgn_mdtrt_d partition(dt,region,subs_code)
```

```
select
```

```
mdtrt_id, medins_setl_id, psn_no, psn_insu_rlts_id, psn_cert_type, certno, psn_name, gend, naty, brdy,
```

```
age, coner_name, tel, addr, insutype, psn_type, cvlserv_flag, cvlserv_lv, sp_psn_type, sp_psn_type_lv,
```

```
clct_grde, flxempe_flag, nwb_flag, insu_admdvts, emp_no, emp_name, emp_type, econ_type, afil_indu,
```

```
afil_rlts, emp_mgt_type, pay_loc, fixmedins_code, fixmedins_name, hosp_lv, fix_blng_admdvts,
```

```
lmtpric_hosp_lv, dedc_hosp_lv, begntime, endtime, mdtrt_cert_type, mdtrt_cert_no, med_type,
```

```
rloc_type, ars_year_ipt_flag, pre_pay_flag, year, refl_old_mdtrt_id, ipt_otp_no, medrcdno, chfpdr_code,
```

```
adm_diag_dscr, adm_dept_codg, adm_dept_name, adm_bed, wardarea_bed, traf_dept_flag,
```

```
dscg_maindiag_code, dscg_dept_codg, dscg_dept_name, dscg_bed, dscg_way, main_cond_dscr, dise_no,
```

```
dise_name, oprn_oprt_code, oprn_oprt_name, otp_diag_info, inhosp_stas, die_date, ipt_days, geso_val,
```

```
birctrl_type, fetts, fetus_cnt, matn_type, prey_time, latechb_flag, pret_flag, fpssc_no, birctrl_matn_date,
```

```
cop_flag, trt_dcla_dctl_sn, vali_flag, memo, rid, updt_time, crter_id, crter_name, crte_time,
```

```
crte_optins_no, opter_id, opter_name, opt_time, optins_no, poolarea_no, chfpdr_name,
```

```
dscg_maindiag_name, adm_caty, dscg_caty, ttp_pay_flag, ttp_pay_prop, dise_type_code,
```

```
same_dise_adm_flag, quts_type, subsys_codg_src, dty_flag
```

```
,concat(date_add(current_timestamp(),-1),' ',date_format(current_timestamp(),'HH:mm:ss')) as
```

```
exch_updt_time
```

```
,date_format(CRTE_TIME,'yyyy-MM') as dt
```

```
,region
```

```
, 'cep_stcdb' as subs_code
```

```
from ods_prd.ods_cep_stcdb_mdtrt_d odstb join tmp_yibao_ods_dwd on
```

```
odstb.dt=tmp_yibao_ods_dwd.cdate and odstb.region=tmp_yibao_ods_dwd.area
```

```
where local_dty_flag='0'
```

```
cluster BY dt,mdtrt_id,psn_no ;
```

## 7.4 ods、stg及dwd建表语句

### 1. stg\_prd.stg\_cep\_stcdb\_mdtrt\_d

```
CREATE TABLE `stg_prd.stg_cep_stcdb_mdtrt_d`(  
  `mdtrt_id` string COMMENT '就诊 ID',  
  `medins_setl_id` string COMMENT '医药机构结算 ID',  
  `psn_no` string COMMENT '人员编号',  
  `psn_insu_rlts_id` string COMMENT '人员参保关系 ID',  
  `psn_cert_type` string COMMENT '人员证件类型',  
  `certno` string COMMENT '证件号码',  
  `psn_name` string COMMENT '人员姓名',  
  `gend` string COMMENT '性别',  
  `naty` string COMMENT '民族',  
  `brdy` date COMMENT '出生日期',  
  `age` string COMMENT '年龄',  
  `coner_name` string COMMENT '联系人姓名',  
  `tel` string COMMENT '联系电话',  
  `addr` string COMMENT '联系地址',  
  `insutype` string COMMENT '险种类型',  
  `psn_type` string COMMENT '人员类别',  
  `cvlserv_flag` string COMMENT '公务员标志',  
  `cvlserv_lv` string COMMENT '公务员等级',  
  `sp_psn_type` string COMMENT '特殊人员类型',  
  `sp_psn_type_lv` string COMMENT '特殊人员类型等级',  
  `clct_grde` string COMMENT '缴费档次',  
  `flxempe_flag` string COMMENT '灵活就业标志',  
  `nwb_flag` string COMMENT '新生儿标志',  
  `insu_admdvs` string COMMENT '参保所属医保区划',  
  `emp_no` string COMMENT '单位编号',  
  `emp_name` string COMMENT '单位名称',  
  `emp_type` string COMMENT '单位类型',  
  `econ_type` string COMMENT '经济类型',  
  `afil_indu` string COMMENT '所属行业',  
  `afil_rlts` string COMMENT '隶属关系',  
  `emp_mgt_type` string COMMENT '单位管理类型',  
  `pay_loc` string COMMENT '支付地点类别',  
  `fixmedins_code` string COMMENT '定点医药机构编号',  
  `fixmedins_name` string COMMENT '定点医药机构名称',
```

`hosp\_lv` string COMMENT '医院等级',  
 `fix\_blng\_admdvs` string COMMENT '定点归属医保区划',  
 `lmtpric\_hosp\_lv` string COMMENT '限价医院等级',  
 `dedc\_hosp\_lv` string COMMENT '起付线医院等级',  
 `begntime` string COMMENT '开始时间',  
 `endtime` string COMMENT '结束时间',  
 `mdtrt\_cert\_type` string COMMENT '就诊凭证类型',  
 `mdtrt\_cert\_no` string COMMENT '就诊凭证编号',  
 `med\_type` string COMMENT '医疗类别',  
 `rloc\_type` string COMMENT '异地安置类别',  
 `ars\_year\_ipt\_flag` string COMMENT '跨年度住院标志',  
 `pre\_pay\_flag` string COMMENT '先行支付标志',  
 `year` string COMMENT '年度',  
 `refl\_old\_mdtrt\_id` string COMMENT '转诊前就诊 ID',  
 `ipt\_otp\_no` string COMMENT '住院/门诊号',  
 `medrcdno` string COMMENT '病历号',  
 `chfptr\_code` string COMMENT '主治医师代码',  
 `adm\_diag\_dscr` string COMMENT '入院诊断描述',  
 `adm\_dept\_codg` string COMMENT '入院科室编码',  
 `adm\_dept\_name` string COMMENT '入院科室名称',  
 `adm\_bed` string COMMENT '入院床位',  
 `wardarea\_bed` string COMMENT '病区床位',  
 `traf\_dept\_flag` string COMMENT '转科室标志',  
 `dscg\_maindiag\_code` string COMMENT '住院主诊断代码',  
 `dscg\_dept\_codg` string COMMENT '出院科室编码',  
 `dscg\_dept\_name` string COMMENT '出院科室名称',  
 `dscg\_bed` string COMMENT '出院床位',  
 `dscg\_way` string COMMENT '离院方式',  
 `main\_cond\_dscr` string COMMENT '主要病情描述',  
 `dise\_no` string COMMENT '病种编号',  
 `dise\_name` string COMMENT '病种名称',  
 `oprn\_oprt\_code` string COMMENT '手术操作代码',  
 `oprn\_oprt\_name` string COMMENT '手术操作名称',  
 `otp\_diag\_info` string COMMENT '门诊诊断信息',  
 `inhosp\_stas` string COMMENT '在院状态',  
 `die\_date` date COMMENT '死亡日期',  
 `ipt\_days` string COMMENT '住院天数',  
 `geso\_val` string COMMENT '孕周数',

```

`birctrl_type` string COMMENT '计划生育手术类别',
`fetts` string COMMENT '胎次',
`fetus_cnt` string COMMENT '胎儿数',
`matn_type` string COMMENT '生育类别',
`prey_time` string COMMENT '妊娠时间',
`latechb_flag` string COMMENT '晚育标志',
`pret_flag` string COMMENT '早产标志',
`fpssc_no` string COMMENT '计划生育服务证号',
`birctrl_matn_date` string COMMENT '计划生育手术或生育日期',
`cop_flag` string COMMENT '伴有并发症标志',
`trt_dcla_detl_sn` string COMMENT '待遇申报明细流水号',
`vali_flag` string COMMENT '有效标志',
`memo` string COMMENT '备注',
`rid` string COMMENT '数据唯一记录号',
`updt_time` string COMMENT '数据更新时间',
`crter_id` string COMMENT '创建人 ID',
`crter_name` string COMMENT '创建人姓名',
`crte_time` string COMMENT '数据创建时间',
`crte_optins_no` string COMMENT '创建机构编号',
`opter_id` string COMMENT '经办人 ID',
`opter_name` string COMMENT '经办人姓名',
`opt_time` string COMMENT '经办时间',
`optins_no` string COMMENT '经办机构编号',
`poolarea_no` string COMMENT '统筹区编号',
`chfpdr_name` string COMMENT '主诊医师姓名',
`dscg_maindiag_name` string COMMENT '住院主诊断名称',
`adm_caty` string COMMENT '入院科别',
`dscg_caty` string COMMENT '出院科别',
`ttp_pay_flag` string COMMENT '第三方赔付标志',
`ttp_pay_prop` string COMMENT '第三方赔付比例',
`dise_type_code` string COMMENT '病种类型代码',
`same_dise_adm_flag` string COMMENT '同病种入院标志',
`quts_type` string COMMENT '编制类型',
`subsys_codg_src` string COMMENT '子系统编码',
`dty_flag` string COMMENT '是否脏数据(0 否 1 是)')
COMMENT ''
PARTITIONED BY (
  `dt` string);

```

## 2. ods\_prd.ods\_cep\_stcdb\_mdtrt\_d

```
CREATE TABLE `ods_prd.ods_cep_stcdb_mdtrt_d`(  
  `mdtrt_id` string COMMENT '就诊 ID',  
  `medins_setl_id` string COMMENT '医药机构结算 ID',  
  `psn_no` string COMMENT '人员编号',  
  `psn_insu_rlts_id` string COMMENT '人员参保关系 ID',  
  `psn_cert_type` string COMMENT '人员证件类型',  
  `certno` string COMMENT '证件号码',  
  `psn_name` string COMMENT '人员姓名',  
  `gend` string COMMENT '性别',  
  `naty` string COMMENT '民族',  
  `brdy` date COMMENT '出生日期',  
  `age` decimal(4,1) COMMENT '年龄',  
  `coner_name` string COMMENT '联系人姓名',  
  `tel` string COMMENT '联系电话',  
  `addr` string COMMENT '联系地址',  
  `insutype` string COMMENT '险种类型',  
  `psn_type` string COMMENT '人员类别',  
  `cvlserv_flag` string COMMENT '公务员标志',  
  `cvlserv_lv` string COMMENT '公务员等级',  
  `sp_psn_type` string COMMENT '特殊人员类型',  
  `sp_psn_type_lv` string COMMENT '特殊人员类型等级',  
  `clct_grde` string COMMENT '缴费档次',  
  `flxempe_flag` string COMMENT '灵活就业标志',  
  `nwb_flag` string COMMENT '新生儿标志',  
  `insu_admdvs` string COMMENT '参保所属医保区划',  
  `emp_no` string COMMENT '单位编号',  
  `emp_name` string COMMENT '单位名称',  
  `emp_type` string COMMENT '单位类型',  
  `econ_type` string COMMENT '经济类型',  
  `afil_indu` string COMMENT '所属行业',  
  `afil_rlts` string COMMENT '隶属关系',  
  `emp_mgt_type` string COMMENT '单位管理类型',  
  `pay_loc` string COMMENT '支付地点类别',  
  `fixmedins_code` string COMMENT '定点医药机构编号',  
  `fixmedins_name` string COMMENT '定点医药机构名称',  
  `hosp_lv` string COMMENT '医院等级',  
  `fix_blng_admdvs` string COMMENT '定点归属医保区划',
```

```

`lmtpric_hosp_lv` string COMMENT '限价医院等级',
`dedc_hosp_lv` string COMMENT '起付线医院等级',
`begntime` timestamp COMMENT '开始时间',
`endtime` timestamp COMMENT '结束时间',
`mdtrt_cert_type` string COMMENT '就诊凭证类型',
`mdtrt_cert_no` string COMMENT '就诊凭证编号',
`med_type` string COMMENT '医疗类别',
`rloc_type` string COMMENT '异地安置类别',
`ars_year_ipt_flag` string COMMENT '跨年度住院标志',
`pre_pay_flag` string COMMENT '先行支付标志',
`year` string COMMENT '年度',
`refl_old_mdtrt_id` string COMMENT '转诊前就诊 ID',
`ipt_otp_no` string COMMENT '住院/门诊号',
`medrcdno` string COMMENT '病历号',
`chfpdr_code` string COMMENT '主治医师代码',
`adm_diag_dscr` string COMMENT '入院诊断描述',
`adm_dept_codg` string COMMENT '入院科室编码',
`adm_dept_name` string COMMENT '入院科室名称',
`adm_bed` string COMMENT '入院床位',
`wardarea_bed` string COMMENT '病区床位',
`traf_dept_flag` string COMMENT '转科室标志',
`dscg_maindiag_code` string COMMENT '住院主诊断代码',
`dscg_dept_codg` string COMMENT '出院科室编码',
`dscg_dept_name` string COMMENT '出院科室名称',
`dscg_bed` string COMMENT '出院床位',
`dscg_way` string COMMENT '离院方式',
`main_cond_dscr` string COMMENT '主要病情描述',
`dise_no` string COMMENT '病种编号',
`dise_name` string COMMENT '病种名称',
`oprn_oprt_code` string COMMENT '手术操作代码',
`oprn_oprt_name` string COMMENT '手术操作名称',
`otp_diag_info` string COMMENT '门诊诊断信息',
`inhosp_stas` string COMMENT '在院状态',
`die_date` date COMMENT '死亡日期',
`ipt_days` decimal(16,0) COMMENT '住院天数',
`geso_val` decimal(2,0) COMMENT '孕周数',
`birctrl_type` string COMMENT '计划生育手术类别',
`fetts` decimal(3,0) COMMENT '胎次',

```



```

`fetus_cnt` decimal(3,0) COMMENT '胎儿数',
`matn_type` string COMMENT '生育类别',
`prey_time` timestamp COMMENT '妊娠时间',
`latechb_flag` string COMMENT '晚育标志',
`pret_flag` string COMMENT '早产标志',
`fpsc_no` string COMMENT '计划生育服务证号',
`birctrl_matn_date` timestamp COMMENT '计划生育手术或生育日期',
`cop_flag` string COMMENT '伴有并发症标志',
`trt_dcla_detl_sn` string COMMENT '待遇申报明细流水号',
`vali_flag` string COMMENT '有效标志',
`memo` string COMMENT '备注',
`rid` string COMMENT '数据唯一记录号',
`updt_time` timestamp COMMENT '数据更新时间',
`crter_id` string COMMENT '创建人 ID',
`crter_name` string COMMENT '创建人姓名',
`crte_time` timestamp COMMENT '数据创建时间',
`crte_optins_no` string COMMENT '创建机构编号',
`opter_id` string COMMENT '经办人 ID',
`opter_name` string COMMENT '经办人姓名',
`opt_time` timestamp COMMENT '经办时间',
`optins_no` string COMMENT '经办机构编号',
`poolarea_no` string COMMENT '统筹区编号',
`chfpdr_name` string COMMENT '主诊医师姓名',
`dscg_maindiag_name` string COMMENT '住院主诊断名称',
`adm_caty` string COMMENT '入院科别',
`dscg_caty` string COMMENT '出院科别',
`ttp_pay_flag` string COMMENT '第三方赔付标志',
`ttp_pay_prop` decimal(5,4) COMMENT '第三方赔付比例',
`dise_type_code` string COMMENT '病种类型代码',
`same_dise_adm_flag` string COMMENT '同病种入院标志',
`quts_type` string COMMENT '编制类型',
`subsys_codg_src` string COMMENT '子系统编码',
`dty_flag` string COMMENT '是否脏数据(0 否 1 是)',
`local_dty_flag` string COMMENT '本地规则-脏数据标识(0 否 1 是)',
`exch_updt_time` timestamp COMMENT '入仓时间')
COMMENT ''
PARTITIONED BY (
  `dt` string,

```

```
`region` string);
```

### 3. dwd\_prd.dwd\_dgn\_mdtrt\_d

```
CREATE TABLE `dwd_prd.dwd_dgn_mdtrt_d`(  
  `mdtrt_id` string COMMENT '就诊 ID',  
  `medins_setl_id` string COMMENT '医药机构结算 ID',  
  `psn_no` string COMMENT '人员编号',  
  `psn_insu_rlts_id` string COMMENT '人员参保关系 ID',  
  `psn_cert_type` string COMMENT '人员证件类型',  
  `certno` string COMMENT '证件号码',  
  `psn_name` string COMMENT '人员姓名',  
  `gend` string COMMENT '性别',  
  `naty` string COMMENT '民族',  
  `brdy` date COMMENT '出生日期',  
  `age` decimal(4,1) COMMENT '年龄',  
  `coner_name` string COMMENT '联系人姓名',  
  `tel` string COMMENT '联系电话',  
  `addr` string COMMENT '联系地址',  
  `insutype` string COMMENT '险种类型',  
  `psn_type` string COMMENT '人员类别',  
  `cvlserv_flag` string COMMENT '公务员标志',  
  `cvlserv_lv` string COMMENT '公务员等级',  
  `sp_psn_type` string COMMENT '特殊人员类型',  
  `sp_psn_type_lv` string COMMENT '特殊人员类型等级',  
  `clct_grde` string COMMENT '缴费档次',  
  `flxempe_flag` string COMMENT '灵活就业标志',  
  `nwb_flag` string COMMENT '新生儿标志',  
  `insu_admdvs` string COMMENT '参保所属医保区划',  
  `emp_no` string COMMENT '单位编号',  
  `emp_name` string COMMENT '单位名称',  
  `emp_type` string COMMENT '单位类型',  
  `econ_type` string COMMENT '经济类型',  
  `afil_indu` string COMMENT '所属行业',  
  `afil_rlts` string COMMENT '隶属关系',  
  `emp_mgt_type` string COMMENT '单位管理类型',  
  `pay_loc` string COMMENT '支付地点类别',  
  `fixmedins_code` string COMMENT '定点医药机构编号',  
  `fixmedins_name` string COMMENT '定点医药机构名称',
```

```

`hosp_lv` string COMMENT '医院等级',
`fix_blng_admdvs` string COMMENT '定点归属医保区划',
`lmtpric_hosp_lv` string COMMENT '限价医院等级',
`dedc_hosp_lv` string COMMENT '起付线医院等级',
`begntime` timestamp COMMENT '开始时间',
`endtime` timestamp COMMENT '结束时间',
`mdtrt_cert_type` string COMMENT '就诊凭证类型',
`mdtrt_cert_no` string COMMENT '就诊凭证编号',
`med_type` string COMMENT '医疗类别',
`rloc_type` string COMMENT '异地安置类别',
`ars_year_ipt_flag` string COMMENT '跨年度住院标志',
`pre_pay_flag` string COMMENT '先行支付标志',
`year` string COMMENT '年度',
`refl_old_mdtrt_id` string COMMENT '转诊前就诊 ID',
`ipt_otp_no` string COMMENT '住院/门诊号',
`medrcdno` string COMMENT '病历号',
`chfprdr_code` string COMMENT '主治医师代码',
`adm_diag_dscr` string COMMENT '入院诊断描述',
`adm_dept_codg` string COMMENT '入院科室编码',
`adm_dept_name` string COMMENT '入院科室名称',
`adm_bed` string COMMENT '入院床位',
`wardarea_bed` string COMMENT '病区床位',
`traf_dept_flag` string COMMENT '转科室标志',
`dscg_maindiag_code` string COMMENT '住院主诊断代码',
`dscg_dept_codg` string COMMENT '出院科室编码',
`dscg_dept_name` string COMMENT '出院科室名称',
`dscg_bed` string COMMENT '出院床位',
`dscg_way` string COMMENT '离院方式',
`main_cond_dscr` string COMMENT '主要病情描述',
`dise_no` string COMMENT '病种编号',
`dise_name` string COMMENT '病种名称',
`oprn_oprt_code` string COMMENT '手术操作代码',
`oprn_oprt_name` string COMMENT '手术操作名称',
`otp_diag_info` string COMMENT '门诊诊断信息',
`inhosp_stas` string COMMENT '在院状态',
`die_date` date COMMENT '死亡日期',
`ipt_days` decimal(16,0) COMMENT '住院天数',
`geso_val` decimal(2,0) COMMENT '孕周数',

```

```

`birctrl_type` string COMMENT '计划生育手术类别',
`fetts` decimal(3,0) COMMENT '胎次',
`fetus_cnt` decimal(3,0) COMMENT '胎儿数',
`matn_type` string COMMENT '生育类别',
`prey_time` timestamp COMMENT '妊娠时间',
`latechb_flag` string COMMENT '晚育标志',
`pret_flag` string COMMENT '早产标志',
`fpssc_no` string COMMENT '计划生育服务证号',
`birctrl_matn_date` timestamp COMMENT '计划生育手术或生育日期',
`cop_flag` string COMMENT '伴有并发症标志',
`trt_dcla_detl_sn` string COMMENT '待遇申报明细流水号',
`vali_flag` string COMMENT '有效标志',
`memo` string COMMENT '备注',
`rid` string COMMENT '数据唯一记录号',
`updt_time` timestamp COMMENT '数据更新时间',
`crter_id` string COMMENT '创建人 ID',
`crter_name` string COMMENT '创建人姓名',
`crte_time` timestamp COMMENT '数据创建时间',
`crte_optins_no` string COMMENT '创建机构编号',
`opter_id` string COMMENT '经办人 ID',
`opter_name` string COMMENT '经办人姓名',
`opt_time` timestamp COMMENT '经办时间',
`optins_no` string COMMENT '经办机构编号',
`poolarea_no` string COMMENT '统筹区编号',
`chfpdr_name` string COMMENT '主诊医师姓名',
`dscg_maindiag_name` string COMMENT '住院主诊断名称',
`adm_caty` string COMMENT '入院科别',
`dscg_caty` string COMMENT '出院科别',
`ttp_pay_flag` string COMMENT '第三方赔付标志',
`ttp_pay_prop` decimal(5,4) COMMENT '第三方赔付比例',
`dise_type_code` string COMMENT '病种类型代码',
`same_dise_adm_flag` string COMMENT '同病种入院标志',
`quts_type` string COMMENT '编制类型',
`subsys_codg_src` string COMMENT '子系统编码',
`dty_flag` string COMMENT '是否脏数据(0 否 1 是)',
`exch_updt_time` timestamp COMMENT '入仓时间')
COMMENT '就诊信息表'
PARTITIONED BY (

```

```
`dt` string,  
`region` string,  
`subs_code` string);
```