



数据湖分析 数据湖管理

文档版本: 20220630



### 法律声明

阿里云提醒您在阅读或使用本文档之前仔细阅读、充分理解本法律声明各条款的内容。 如果您阅读或使用本文档,您的阅读或使用行为将被视为对本声明全部内容的认可。

- 您应当通过阿里云网站或阿里云提供的其他授权通道下载、获取本文档,且仅能用 于自身的合法合规的业务活动。本文档的内容视为阿里云的保密信息,您应当严格 遵守保密义务;未经阿里云事先书面同意,您不得向任何第三方披露本手册内容或 提供给任何第三方使用。
- 未经阿里云事先书面许可,任何单位、公司或个人不得擅自摘抄、翻译、复制本文 档内容的部分或全部,不得以任何方式或途径进行传播和宣传。
- 由于产品版本升级、调整或其他原因,本文档内容有可能变更。阿里云保留在没有 任何通知或者提示下对本文档的内容进行修改的权利,并在阿里云授权通道中不时 发布更新后的用户文档。您应当实时关注用户文档的版本变更并通过阿里云授权渠 道下载、获取最新版的用户文档。
- 4. 本文档仅作为用户使用阿里云产品及服务的参考性指引,阿里云以产品及服务的"现状"、"有缺陷"和"当前功能"的状态提供本文档。阿里云在现有技术的基础上尽最大努力提供相应的介绍及操作指引,但阿里云在此明确声明对本文档内容的准确性、完整性、适用性、可靠性等不作任何明示或暗示的保证。任何单位、公司或个人因为下载、使用或信赖本文档而发生任何差错或经济损失的,阿里云不承担任何法律责任。在任何情况下,阿里云均不对任何间接性、后果性、惩戒性、偶然性、特殊性或刑罚性的损害,包括用户使用或信赖本文档而遭受的利润损失,承担责任(即使阿里云已被告知该等损失的可能性)。
- 5. 阿里云网站上所有内容,包括但不限于著作、产品、图片、档案、资讯、资料、网站架构、网站画面的安排、网页设计,均由阿里云和/或其关联公司依法拥有其知识产权,包括但不限于商标权、专利权、著作权、商业秘密等。非经阿里云和/或其关联公司书面同意,任何人不得擅自使用、修改、复制、公开传播、改变、散布、发行或公开发表阿里云网站、产品程序或内容。此外,未经阿里云事先书面同意,任何人不得为了任何营销、广告、促销或其他目的使用、公布或复制阿里云的名称(包括但不限于单独为或以组合形式包含"阿里云"、"Aliyun"、"万网"等阿里云和/或其关联公司品牌,上述品牌的附属标志及图案或任何类似公司名称、商号、商标、产品或服务名称、域名、图案标示、标志、标识或通过特定描述使第三方能够识别阿里云和/或其关联公司)。
- 6. 如若发现本文档存在任何错误,请与阿里云取得直接联系。

## 通用约定

格式	说明	样例
⚠ 危险	该类警示信息将导致系统重大变更甚至故 障,或者导致人身伤害等结果。	⚠ 危险 重置操作将丢失用户配置数据。
▲ 警告	该类警示信息可能会导致系统重大变更甚 至故障,或者导致人身伤害等结果。	警告 重启操作将导致业务中断,恢复业务 时间约十分钟。
〔〕 注意	用于警示信息、补充说明等,是用户必须 了解的内容。	大) 注意 权重设置为0,该服务器不会再接受新 请求。
⑦ 说明	用于补充说明、最佳实践、窍门等,不是 用户必须了解的内容。	⑦ 说明 您也可以通过按Ctrl+A选中全部文件。
>	多级菜单递进。	单击设置> 网络> 设置网络类型。
粗体	表示按键、菜单、页面名称等UI元素。	在 <b>结果确认</b> 页面,单击 <b>确定</b> 。
Courier字体	命令或代码。	执行    cd /d C:/window    命令,进入 Windows系统文件夹。
斜体	表示参数、变量。	bae log listinstanceid
[] 或者 [alb]	表示可选项,至多选择一个。	ipconfig [-all -t]
{} 或者 {alb}	表示必选项,至多选择一个。	switch {act ive st and}

## 目录

1.元信息管理	05
1.1. 查询Schema详情	05
2.元信息发现	06
2.1. OSS数据源	06
2.2. Tablestore数据源	10
2.3. SLS的OSS投递数据源	13
3.数据入湖	17
3.1. DLA Lakehouse实时入湖	17
3.2. Kafka实时入湖建仓分析	28
3.3. Lindorm实时入湖建仓分析	36
3.4. T+1全量同步一键建仓	43
3.4.1. 概述	43
3.4.2. 如何使用T+1全量同步一键建仓	44
3.4.3. 一键创建OSS数据仓库	45
3.4.4. 授权DLA删除OSS文件	48
3.5. T+1多库合并建仓	51
3.6. ActionTrail日志清洗	54

# 1.元信息管理 1.1. 查询Schema详情

您可以在Schema详情页面管理表,例如查询表数据、删除表等。

#### 操作步骤

- 1. 登录Data Lake Analytics管理控制台。
- 2. 在页面左上角,选择集群所在地域。
- 3. 单击左侧导航栏的数据湖管理 > 元数据管理。
- 4. 单击目标Schema单击右侧的**详细信息**,您可以进行以下操作。
  - **返回**:返回Schema列表。
  - schema搜索: 支持模糊搜索表名。
  - 查询:跳转至SQL执行页面。
  - ○删除:删除表。

# 2.元信息发现

## 2.1. OSS数据源

企业的业务服务所产生的一些标准化表单、日志等数据文件,会被周期性的直接上传到OSS。但是这些存储 在OSS的文件缺少元数据管理,导致难以分析和计算。元数据发现任务可以在单次运行中自动为OSS上面的 数据文件创建和更新数据湖元数据(一张或多张表),具有自动探索文件数据字段及类型、自动映射目录和 分区、自动感知新增列及分区、自动对文件进行分组建表的能力。

#### OSS数据源配置模式

OSS数据源配置支持数仓模式和自由模式,两种模式差异如下:

OSS数据源配置	使用场景	OSS路径格式要求	识别精度	性能
数仓模式	用户直接上传数据 到OSS,并期望构建 可分析与计算的标 准数据仓库。	库/表/文件" 或 者 "库/表/分区/…/ 分区/文件	高	高
自由模式	已存在OSS数据,但 OSS的路径不清晰。 期望通过元信息发 现,构建可分析的 库表分区。	无要求	一般	一般

#### 数仓模式的OSS路径格式要求

OSS是一个开放的文件系统,为了高效的在OSS上面构建数据仓库,数据源路径格式需要有一定的规范性。 OSS数据源数仓模式的元信息发现只支持库/表/文件或者库/表/分区/.../分区/文件的路径格式(根目录对应 schema; 第二级子目录对应Table,且子目录名需要映射到表名;第三级以上如果还有目录,就对应为分 区)。如下图所示:



上图对应的OSS路径如下:

oss://xxxxxx/schema1/Table1/file1.csv
oss://xxxxxx/schema1/Table1/file2.json
oss://xxxxxx/schema1/Table2/file3.csv
oss://xxxxxx/schema1/Table2/file4.csv
oss://xxxxxx/schema1/Table3/year=2020/month=03/day=01/file5.json
oss://xxxxxx/schema1/Table3/year=2020/month=03/day=01/file6.json
oss://xxxxxx/schema1/Table3/year=2020/month=03/day=02/file7.json
oss://xxxxxx/schema1/Table3/year=2020/month=03/day=02/file8.json
oss://xxxxxx/schema1/Table3/year=2020/month=04/day=29/file9.json
<pre>oss://xxxxxx/schemal/Table3/year=2020/month=04/day=29/file10.json</pre>
<pre>oss://xxxxxx/schemal/Table3/year=2020/month=04/day=30/file11.json</pre>
<pre>oss://xxxxxx/schemal/Table3/year=2020/month=04/day=30/file12.json</pre>
oss://xxxxxx/schemal/Table4/age=20/file13.json
oss://xxxxxx/schemal/Table4/age=20/file14.json
oss://xxxxxx/schema1/Table5/2020/03/29/file15.csv
oss://xxxxxx/schema1/Table5/2020/03/29/file16.csv
oss://xxxxxx/schema1/Table5/2020/03/30/file17.csv
oss://xxxxxx/schemal/Table5/2020/03/30/file18.csv

#### 上述数仓模式的OSS数据源进行元信息发现后,在DLA中自动映射的表如下所述:

OSS数据源目录 名称	DLA中自动映射的表名称	映射说明
		Table1目录下的文件类型为csv和json,不一致,故而不 能进行映射创建表。同一目录下的文件类型必须一致才能 进行映射。
Table1	不创建表	⑦ 说明 如果文件类型相同,但是文件里面的字段不是同一种类型也不能进行映射。
Table2	创建表Table2	Table2目录下的文件类型都为csv,可以进行映射。
Table3	创建分区表Table3	分区格式为year=xx/month=xx/day=xx/。自动映射了 以下分区: • year=2020/month=03/day=01 • year=2020/month=03/day=02 • year=2020/month=04/day=29 • year=2020/month=04/day=30
Table4	创建分区表Table4	分区格式为age=xx。自动映射了age=20分区。

OSS数据源目录 名称	DLA中自动映射的表名称	映射说明
T able5	创建分区表Table5	<pre>分区格式为 partition_0=xx/partition_1=xx/partition_2=xx/。自动 映射了以下分区: 2020/03/29 2020/03/30 </pre> ⑦ 说明 由于没有分区键,这里使用 partition_num来补充。

#### 操作步骤

- 1. 登录Dat a Lake Analytics管理控制台。
- 2. 在左侧导航栏,单击数据湖管理 > 元信息发现。
- 3. 在元信息发现页面的OSS数据源区域,单击进入向导。

一元信息发现					
<b>0000</b> <b>SLS的OSS</b> 自动为SLS股激	3) 建数据源 🚯 eloss上面的数据创建和限新数据被元数编。	进入向导	000	OSS数据源 () 自初为OSS存益及場益文件创建和更新数编组元数编	进入向导
<b>Tablestore</b> 目动为Tablesto	数据源 <b>①</b> ▲上型的政治创建风更新政策和元政绩,方使分析和计算。	进入向导			

4. 在OSS数据源页签,根据页面提示进行参数配置,配置说明如下:

参数	说明	
数仓模式和自由模式	<ul> <li>您可以选择数仓模式或自由模式:</li> <li>数仓模式:为"基于OSS而构建的标准数据仓库"的场景构建自动化元信息发现,识别精度高。OSS路径数据布局要求为"库/表/文件"或者"库/表/分区//分区/文件"。</li> <li>自由模式:为"探索OSS上的数据进行分析"的场景构建自动化元信息发现。对OSS数据布局没有要求,可能会产生差异化的表。</li> </ul>	
	文件在OSS中的存储地址,以/结尾。系统会根据您选择的文件夹路径,自动设置OSS路径。	
OSS目录位置	⑦ 说明 系统会自动拉取与DLA同地域的OSS Bucket,您可以根据业务需要从下拉列表中选择Bucket。选择Bucket后,系统会自动列出该Bucket下所有的Object和文件;选中目标Object和文件后,系统会自动将其添加到右侧的OSS路径处。	

参数	说明
格式解析器	格式解析器会读取数据文件内容,从而确定文件的数据格式。默认自动解 析,即按照顺序调用所有内置解析器,也可指定特定文件类型的格式解析 器,比如json、parquet、avro、orc、csv。 o json: 需要读取文件开头以确定文件格式。 o parquet: 需要读取文件结尾处的schema以确定文件格式。 o avro: 需要读取文件开头处的shema以确定文件格式。 o orc: 要读取文件元数据以确定文件格式。 o csv: 检查以下分隔符: 逗号(,)、竖线()、制表符(\t)、分号 (;)、空格()、(\u0001)。
配置选项	高级自定义设置项,如更新,删除规则等。
调度频率	您可以根据需要定期计划运行元信息发现任务。
Schema名称	设置Schema名称,即映射到DLA中的数据库名称(默认每个发现任务会新 创建一个独立的Schema)。

#### 5. 完成上述参数配置后,单击创建,开始创建元信息发现任务。

② 说明 元信息发现任务创建完成后, DLA自动在您设定的时间周期运行发现任务, 如果您想立即同步数据, 也可以在任务列表选择立即执行。

6. 任务开始运行后,会在实例列表显示任务实例的当前运行状态。也可以在任务列表界面管理任务的运行 情况,支持查看任务的运行状态、配置的修改、跳转到DLA的SQL窗口进行快速的数据查询。

云原生数据湖分析	一元信息发现			
概选账号管理	SLS的OSS投递数据源 OSS数据源 Tablestore数据源			
虚拟集群管理数据湖管理 HOT 个				
元信息发现	为"基于OSS而构建的场准数据公库"的场景构建自动化元信息发现,识别精度高。OSS路径数据布局要求为"库/表/文件"或者"库/表/分区//分区/文件"具体参考			
元数据管理 数据入湖	置位录目220 *			
实时数据湖	格式解析器 csv V			
Serverless Presto ^	> 配置选项 (问法)			
SQL执行 SQL监控	调度配置            调度频率         手动执行			
Serverless Spark へ 作业管理	目标元数据配置			
优惠铁餐 系统设置	* Schema名称 您新疆的元信息发现的名称 > 配置选项 (可选)			

#### 注意事项

数仓模式的注意事项如下:

• 如果OSS数据源路径没有被DLA识别出来,您需要查看路径下的文件类型是否相同。如果是CSV文件,您

可以在解析器CSV中配置具体的参数比如分隔符、转义字符、是否有表头等。

- 由于元信息发现通过采样的方式并不能覆盖所有的记录,如果不同行的字段变化很大,会出现生成的表字 段减少的情况。
- 在识别分区及表的时候,只支持子目录下只有文件的场景。如果目录下既有子目录又有文件,则该目录会 被忽略掉,从而导致分区没有生成。

自由模式的注意事项如下:

• 元信息发现会如何生成表名

元信息发现会自动为它创建的表生成名称。存储在元数据管理schema目录中的表的名称遵循以下规则:

- 默认使用最后一级目录名作为表名(针对OSS数据文件)。
- 仅允许使用字母、数字、字符和下划线(\_)。
- 。 表名的最大长度不能超过 128 个字符。发现程序会截断生成的名称以适应限制范围。
- 。如果遇到重复的表名,则元信息发现会在表名后添加MD5字符串后缀。
- 元信息发现如何创建分区

当元信息发现扫描OSS目录文件并检测到多个文件时,它会在目录结构中确定表的根目录,以及哪些目录 是表的分区。

表的名称基于OSS目录前缀或目录名,当某个目录级别下大部分的目录结构和文件格式都相同时,发现程 序会创建一张分区表。例如,对于以下OSS目录结构:

```
oss://bucket01/folder1/table1/partition1/fiile.txt
oss://bucket01/folder1/table1/partition2/fiile.txt
oss://bucket01/folder1/table2/partition3/fiile.txt
oss://bucket01/folder1/table2/partition4/fiile.txt
```

因为table1和table2下的目录和文件内容都是相似的,所以发现程序将创建一个具有两个分区列的表。分区列分别为partition\_0(table这一级目录)、partition\_1(partition这一级目录)。

对于以下OSS目录结构:

```
oss://bucket01/folder1/table1/partition1/fiile.csv
oss://bucket01/folder1/table1/partition2/fiile.csv
oss://bucket01/folder1/table2/partition3/fiile.json
oss://bucket01/folder1/table2/partition4/fiile.json
```

因为table1和table2下的文件格式不同,所以发现程序将创建两张具有一个分区列的表。table1分区列包含partition1和partition2,table2分区列包含partition3和partition4。

对于采用key=value样式的Hive风格分区路径,发现程序会使用键名自动填充列名称。否则,它使用默认 名称,如partition\_0、partition\_1等。

### 2.2. Tablestore数据源

本文介绍如何通过向导创建Tablestore数据源发现任务,自动发现您当前区域下所有的Tablestore实例和表 (包括新增实例和新增表),并自动创建DLA库表映射到Tablestore的实例和表上。

#### 前提条件

当前区域下已经创建了Tablestore实例数据。具体操作请参考创建实例。

#### 操作步骤

- 1. 登录Data Lake Analytics管理控制台。
- 2. 在左侧导航栏,单击数据湖管理 > 元信息发现。
- 3. 在元信息发现页面的Tablestore数据源区域,单击进入向导。

一元信息发现	C.				
	SLS的oss投递数据源 () 自动为sls投递到oss上围的数据创建印度新数据稳元数据。	进入向导	000	OSS数据源 () 自动为OSS存在及理量文件创建和更新数据能元数据	进入向导
***	Tablestore数据源 () 自动为Tablestore上面的数据创建双更新数据输元数据,方便分析和计算。	进入向导			
任务列表	历史列表				

4. 在Tablestore数据源页签,根据页面提示进行参数配置,配置说明如下:

Īπ						
C	OSS数据源 SLS的OSS投递数据源 Tablestore数据源					
	N-COTTON					
	DLA会目动发现您当前region的所有Tablestore实例和表(包括新增实例和新增表),	并且政治論でした元億息				
	调整改革 等小时 新的一轮调度会自动对新增的Tabl	estoref设施识行使取				
	建(G3円) 00	✓ 17				
	目标无数编程器					
<	* schema前缀 tablestore20201230					
	按照Tablestore的实例名,元信息; Tablestore实例的类的关系对象名	2環生成为不同的DLA的Schema,其中Schema的名称规则为"自湿义崩覆_Tablestore实例名"。生成到DLA 中的录名与 目前。				
	》 副置远项 (可选)					
	Tablestore所有实例建表预览					
	支持標糊匹配实例名称 Q					
	实例名称					
	✓ benchmark01					
	主表名称 关联主表的对象名称 对象类型 🚺	DLA的Schema名称 () DLA的Table名称 () schema的规则ocation				
		· 返回 · · · · · · · · · · · · · · · · · ·				
-	参数	说明				
	调度频率	你可以根据季要设置调度Tablestore数据源发现任务的周期				
	制及频十	心可以依据而安议重购及中国的约约代数正方的内别。				
具体的间		设直 I ablestore 数据源反现调度 任务的执行时间。				
-						
	设置Schema的前缀,即映射到DLA中的数据库名称的前缀。Schema的名					
	schema前缀 称规则为"自定义前缀 Tablestore实例名"。生成到DLA中的表名与					
		「あっにってって大学」の文化学大学が学校では同じ。				
配直选项		局级目定义设直项,如米杆米源、米杆条数等。				

5. 完成上述参数配置后,单击创建,开始创建Tablestore数据源发现任务。

6. Tablestore数据源发现任务创建成功后,单击**立即发现**开始运行自动发现Tablestore数据源任务。

立即发现 返回任务管理列表

您也可以在**任务列表**中看到创建成功的任务信息。Tablestore数据源发现任务将根据您设置的**调度频** 率,需要您手动执行或者定期自动调度该任务。

元信息发现								
	SLS的OSS均进数据源 ① 回为为SLS税差到GS上期的数据创建和现象和度相元数据。			进入向导		创建和更新数据湖元数据		进入向导
<b>***</b>	Tablestore数据源 () 目动为Tablestore上面的政旗创建及更新政旗	湖元款据,方便分析和计算。		进入向导				
任务列表	历史列表							
全部类型	$\sim$							周新
< ID	名称	schema名称/前缀	调度计划	调度状态	最近运行状态	最近运行耗时	创建时间	摄作
358	all sls project of this region -> muyuant estonline	muyuantestonline	手动执行	开启调度	• SUCCESS	306₩9	2020-12-30 15:31:37	执行   编辑   历史
349	all ots instances of this region -> tablest ore20201231	tablestore20201231	每小时00分	开启调度	SUCCESS Ø	20秒	2020-12-30 11:10:08	执行   编辑   历史   删除
348	all ots instances of this region -> tablest ore20201228	tablestore20201228	手动执行	开启调度		2010	2020-12-30 10:17:08	执行   编辑   历史   删除

Tablestore数据源发现任务执行成功后,单击schema名称/前缀列下面的数据库名称链接(如单击tablestore20201231),跳转到Serverless Presto > SQL执行页面。您可以看到DLA自动发现创建成功的库、表、列信息。



7. 在Serverless Presto > SQL执行页面编写SQL语句,单击同步执行或者异步执行,执行SQL语句。
例如在tablestore20201231\_benchmark01下执行 select \* from `tablestore20201231\_benchma rk01`.`test000` where `key` = '' limit 20; 。

SQL执行		· 唐法手册 · 函数手册
tablestore20201231 © C	SQL执行集群 public(共享集部) 〜 同参执行(F8) 昇参执行(F9) 推式(以下10) 主題 〜 新建SQL	畫拟集群 👔 登录DMS执行SQL
"双击"切换Schema	1 select * from `tablestore20201231_benchmark01`.`test000` _ where `key` = '' limit 20; '	
tablestore20201231_benchmark01 (current)		
V 🖬 test000		
III key		
i≣ col000		
i col001		
col002		
■ col003		
i≣ col004		
i≣ col005		
≤ col006		
i≣ col007	執行历史 执行结果 SQUE控 1	日出結果集 💙 降離
col008		

## 2.3. SLS的OSS投递数据源

本文介绍如何通过向导创建SLS的OSS投递数据源发现任务,自动发现您当前区域下所有的Logstore投递到 OSS的日志数据(包括新增的投递到OSS的日志数据,以及增量的分区及数据),并自动创建DLA库表映射 到投递的OSS日志数据上。

#### 前提条件

当前区域下的SLS已经将日志服务数据投递到OSS。具体操作请参考将日志服务数据投递到OSS。

#### 业务场景

企业会把服务日志、行为日志等数据存储在日志服务SLS中,当数据量较大时,可以通过投递的方式把全量数据归档到OSS。以前这些数据不可进行分析计算,使用DLA的元信息发现功能,能够一键构建数据湖的元数据,接入DLA的统一数据分析平台。使用DLA的Serverless Spark及Presto引擎能够计算和分析DLA元数据管理的全域数据,可以支持ETL后数据交付、低频全量日志数据分析、日志数据关联DB数据分析等业务场景。



#### 操作步骤

- 1. 登录Data Lake Analytics管理控制台。
- 2. 在左侧导航栏, 单击数据湖管理 > 元信息发现。
- 3. 在元信息发现页面的SLS的OSS投递数据源区域,单击进入向导。

元信息发现					
	SLS的OSS均涵教派演 自动为545税通到OSS上型的改進物違い更新改進為元政選。	进入向导	000	CSS数据版 目の为らSS等量が増量文件创建和更新数据執元数据	进入向导
***	Tablestore数据版 () 目动为Tablestore上版的数据创建风度新数据唱元数据、方便分析和计算。	进入向导			

4. 在SLS的OSS投递数据源页签,根据页面提示进行参数配置,配置说明如下:

١л	信息发现	Q.					
s	SLS的OSS投過較個標準 OSS較過原準 Tablestore較過原源						
	数据源配	<ul> <li></li></ul>	当有新的OSS投递设置时,在下一次执行时可	以自动发现	R		
	调度配置						
		调度频率	毎日	~			
	新的一轮调度会自动对新增的OSS投递进行发现			现			
		具体时间	00 ~ Bil: 00	~	分		
	目标元数据配置						
		* schema前缀	sls20210129				
			按照Logstore投递的Bucket名,聚合到不同DL 会替换为"_"。具体可以参考"SLS投递建表效果	A的Schen 预览"	na,其中Schema名称为"前缀	_Bucket名"; DLA的表名为"ProjectName_StoreName"	特殊字符
		> 配置选项 (可选)					
	SLS投递建	皇表效果预览					
	支持模糊	相匹配Project名称 Q					
	Project名和	R					
	⊻ k	:678e					
	Logstor	e名称	OSS投递路径			DLA的Schema名称 🚺	DLA的Table名称 🚺

参数	说明
数据源配置	<ul> <li>您可以通过以下两种方式选择数据源:</li> <li>自动发现:您无需做任何配置,自动发现所有Project的OSS投递,当 有新的OSS投递设置时,在下一次执行时可以自动发现。</li> <li>手动选择:当选择该方式时,您需要手动选择logstore。</li> </ul>
调度频率	您可以根据需要设置调度SLS的OSS投递数据源发现任务的周期。
具体时间	设置SLS的OSS投递数据源发现调度任务的执行时间。
schema前缀	设置Schema的前缀,即映射到DLA中的数据库名称的前缀。Schema的名称规则为"前缀Logstore投递的Bucket名称"。
配置选项	高级自定义设置项,如文件字段变更规则。

- 5. 完成上述参数配置后,单击创建,开始创建SLS的OSS投递数据源发现任务。
- 6. SLS的OSS投递数据源发现任务创建成功后,单击**立即发现**开始运行自动发现SLS的OSS投递数据源任务。

<b></b> 创建成功					
立即发现	返回任务管理列表				

您也可以在**任务列表**中看到创建成功的任务信息。SLS的OSS投递数据源发现任务将根据您设置的**调度** 频率,需要您手动执行或者定期自动调度该任务。

任务列表	历史列表								
全部类型	$\sim$							R	刷新
ID	名称	schema名称/前缀	调度计划	调度状态	最近运行状态	最近运行耗时	创建时间	操作	
358	all sls project of this region -> muyuant estonline	muyuantestonline	手动执行	开启调度 💽	• SUCCESS 🕜	306€)	2020-12-30 15:31:37	执行  编辑  历史  删除	88
349	all ots instances of this region -> tablest ore20201231	tablestore20201231	每小时00分	开启调度	SUCCESS	20#9	2020-12-30 11:10:08	执行 编辑 历史  删除	

SLS的OSS投递数据源发现任务执行成功后,单击schema名称/前缀列下面的数据库名称链接(如单击muyuantestonline),跳转到Serverless Presto > SQL执行页面。您可以看到DLA自动发现创建成功的库、表、列信息。

SQL执行		请法手册 函数手册
muyuantestonline O C "双击"切换Schema	SQL执行集群         public(共享集群)         回步执行(F8)         异步执行(F9)         格式化(F10)         主題 ∨           新速SQL集队集群	● 登录DMS执行SQL
<ul> <li>muyuantestonline_dla_crawle</li> <li>m k8s_log_c159a750471d2451</li> <li>m k8s_log_c25491812f63b48</li> <li>m k8s_log_cfb01e83f003e4b2</li> <li>m sls_crawler_test_csv_full_ty</li> <li>m sls_crawler_test_full_types</li> <li>m sls_crawler_test_partition</li> <li>m sls_crawler_test_partition</li> </ul>	1 <u>执行历史</u> 执行结果 SQL监控 3	SHARE V BE
>  sls_crawler_test_shipper_	Os /*+ cluster=public */ SELECT count(*) FROM `changqing_tpc_h_3_test`.`orders_base_spark_2` LIMIT 500;	2021年1月4日 11:22:33
>  sls_crawler_test_shipper_p	Os /*+ cluster=public */ show databases;	2021年1月4日 11:22:25
> 🖩 sis_crawier_test_sis_partiti	Os /*+ cluster=public */ show databases;	2021年1月4日 11:21:49
>  sis_crawler_test_testnotpa	0 Os /*+ cluster=public */ show create database changqing_tpc_h_3_test;	2021年1月4日 11:21:40
	Os /*+ cluster=public */ show create database changging_tpc_h_3_test;	2021年1月4日 11:21:10

7. 在Serverless Presto > SQL执行页面编写SQL语句,单击同步执行或者异步执行,执行SQL语句。
例如在muyuantestonline\_\_dla\_crawler\_hangzhou下执行 select \* from `muyuantestonline\_ dla\_crawler\_hangzhou`.`sls\_crawler\_test\_csv\_full\_types` limit 20;

SQL执行					语法手册	函数手册
muyuantestonline 🛛 😋	SQL执行集群	public(共享集群) ~	同步执行(F8) 异步执行(F9)	格式化(F10) 主题 ~		
"双击"切换Schema				新娃SQL	虚拟兼群 し 登求	DMSHATSQL
∨ 🛢 muyuantestonline_dla_crawl¢	1 2 select	* from `muvuantestonline	dla crawler hangzhou`.	`sls crawler test _csv f	ull types` limit :	20.
> 🖩 k8s_log_c159a75947fd2451	3	110m mayaanooboon11no				,
> m k8s_log_c25491812f63b48						
> 🖩 k8s_log_cfb01e83f003e4b:						
> msls_crawler_test_csv_full_ty						
>  sis_crawler_test_full_types	执行历史	执行结果 SQL监控	0		导出结果集	~ 隐藏 □◇
>  sls_crawler_test_json_full_t						
> m sls_crawler_test_partitionm	序号	keyint	keybigint	keyfloat	keydouble	detail
> 🖩 sls_crawler_test_shipper_	1	1	10000000	0.1	0.0000000000001	详情
>  sls_crawler_test_shipper_js	2	1	10000000	0.1	0.0000000000001	详情
> III sis_crawler_test_shipper_p	3	1	10000000	0.1	0.0000000000001	详情
>  sls_crawler_test_testnotpa	4	1	10000000	0.1	0.0000000000001	详情
	5	1	10000000	0.1	0.0000000000001	详情
	6	1	10000000	0.1	0.0000000000001	详情

# 3.数据入湖

### 3.1. DLA Lakehouse实时入湖

DLA Lakehouse实时入湖方案利用数据湖技术,重构数仓语义;分析数据湖数据,实现数仓的应用。本文以 RDS MySQL数据源为例介绍了RDS MySQL从入湖到分析的操作步骤。

#### 背景信息

数据湖分析(Data Lake Analytics)是目前炙手可热的方向,主要是以对象存储系统为核心,构建海量、低成本的结构化、半结构化、非结构化对象文件的入湖、存储和分析业务。目前各大云厂商都在积极跟进,布局相关的业务能力,阿里云数据湖分析团队在这个方向也很早就投入相关产品的研发。随着数据湖的应用越来越多,大家发现依赖数据湖最原始的能力,仅仅做简单的存储和分析,往往会遇到很多的问题。比较典型的痛点如下:

- 多源头数据需要统一存储管理,并需要便捷的融合分析。
- 源头数据元信息不确定或变化大,需要自动识别和管理;简单的元信息发现功能时效性不够。
- 全量建仓或直连数据库进行分析对源库造成的压力较大,需要卸载线上压力规避故障。
- 建仓延迟较长(T+1天),需要T+10m的低延迟入湖。
- 更新频繁致小文件多,分析性能差,需要Upsert自动合并。
- 海量数据在事务库或传统数仓中存储成本高,需要低成本归档。
- 源库行存储格式或非分析型格式,分析能力弱,需要支持列式存储格式。
- 自建大数据平台运维成本高,需要产品化、云原生、一体化的方案。
- 常见数仓的存储不开放,需要自建能力、开源可控。

Lakehouse是一种更先进的范式(Paradigm)和方案,用来解决上述简单入湖分析遇到的各种痛点问题。在 Lakehouse技术中,非常关键的技术就是多版本的文件管理协议,它提供入湖和分析过程中的增量数据实时 写入、ACID事务和多版本、小文件自动合并优化、元信息校验和自动进化、高效的列式分析格式、高效的索 引优化、超大分区表存储等能力。目前开源社区有Hudi、Delta、Iceberg等数据湖方案,阿里云数据湖分析 团队选择了比较成熟的Hudi作为DLA Lakehouse的湖仓一体化格式。关于Lakehouse的更多介绍,请参 见Lakehouse介绍。

#### DLA Lakehouse核心概念和相关约束说明

- Lakehouse (湖仓) 有两重含义:
  - 范式:即解决简单入湖分析所遇到的痛点问题的一种解决方案。
  - 存储空间:用来提供一个从其他地方入湖写入数据的空间,后续所有相关操作都围绕着这个湖仓来进行。
    - 不同的Lakehouse有完全不同的路径,路径之间不可以相互有前缀关系(防止数据覆盖)。
    - Lakehouse不能轻易进行修改。
- Workload (工作负载)是围绕湖仓一体化而展开的核心工作的编排调度(由DLA Lakehouse统一调度), 包括如下功能特点:
  - 。 入湖建仓
    - 为了将其他源头的数据,汇总到整个湖仓内构建一个统一的数据平台,例如有DB类型的入湖建仓, 也有Kafka的入湖建仓,还有OSS的数据转换建仓。
    - 不同的入湖建仓,涉及到全量、增量等多个阶段,会统一编排并统一协调调度,简化用户管理成本。

○ 查询优化

为了提升分析能力,构建各种查询优化方面的工作负载,比如自动构建索引、自动清理历史数据、自动构建物化视图等。

- 。 管理
  - 成本优化:自动生命周期、冷热分层存储等。
  - 数据互通:跨域建仓等。
  - 数据安全: 备份恢复等能力。
  - 数据质量: DQC自动校验等。
- Job作业对于Workload的实际作业拆分和执行,以及调度到不同的计算平台上执行,对用户不可见;目前 DLA只支持调度作业到DLA Serverless Spark上执行。核心单元概念如下:
  - 全量作业(从某个Workload中拆分出来)
  - 增量作业(从某个Workload中拆分出来)
  - 。 Clustering: 小文件聚合
  - Indexing: 自动索引构建
  - Compaction: 自动日志合并
  - Tier: 自动分层存储
  - Lifecycle: 自动生命周期管理
  - MaterializedView:物化视图
- DB ( 库 ): DLA 的库
- Table (表): DLA的表
- Partition (分区): DLA的分区
- Column (列): DLA的列

#### DLA Lakehouse方案介绍

DLA Lakehouse实时入湖是分钟级近实时的数据入湖方案,它能够构建统一、低成本、海量数据、自动元信 息同步的湖仓平台,并支持高性能的DLA Spark计算和DLA Presto分析。DLA Lakehouse实时入湖方案的存储 与计算完全分离,写、存、读完全弹性,它的方案架构如下图所示:



#### 准备工作

#### 您需要在DLA中进行以下操作:

1. 开通云原生数据湖分析服务

⑦ 说明 如果未注册阿里云账号,请先注册账号。具体请参见。

#### 2. 创建虚拟集群

⑦ 说明 DLA基于Spark引擎来运行DLA Lakehouse,因此创建虚拟集群的时候需要选择Spark引擎。

#### 您需要在RDS中进行以下操作:

1. 创建RDS MySQL实例

⑦ 说明 由于DLA Lakehouse只支持专有网络,故创建RDS MySQL实例时,网络类型请选择专有网络。

- 2. 创建数据库和账号
- 3. 通过DMS登录RDS数据库
- 4. 在SQLConsole窗口中执行SQL语句创建库表并插入数据。

#### 您需要在DTS中进行以下操作:

⑦ 说明 目前DLA中RDS数据源的入湖分析工作负载,会先利用RDS做数据的全量同步,然后依赖DTS 数据订阅功能做增量同步,最终实现完整的RDS数据入湖。

#### 1. 创建RDS MySQL数据订阅通道

? 说明

- 由于DLA Lakehouse只支持专有网络,故订阅任务的网络类型请选择专有网络。
- 由于DLA Lakehouse无法自动更新元数据信息,故需要订阅的数据类型请选择数据更新和结构更新。

#### 2. 新增消费组

- 3. 查看订阅Topic和消费者ID。后续的创建RDS入湖负载的增量同步配置中需要使用这2个参数。
  - 在订阅任务的订阅配置中可以查看订阅Topic。

.

认用印度	基本信息	
订阅数据	订阅ID: dtsf5906xn717zzbcs	订阅名称: fengze-benchmark-test
性能监控	订阅的数据决型: 数据更新	
监控报警	遊览例关型: RDS 实例	RDS实例ID: rm- 50
数据消费	数据库美型: MySQL	
任务日志	账号: muyuan	修改实例密码 登录数据库
	订阅T0pic: cn_hangzhou_rm_bp15w0sa2q26d115o_muyuan	

○ 在订阅任务的数据消费中可以查看消费者ⅠD。

订阅配置	数据消费	
订阅数据		
性能监控	消费组ID/名称	消费位点
监控报警	dtszsq10o4715k4bbo group_test_3	
数据消费	dtsdxpo9ar317i6bdv	2020年09月23日 18:25:58
任务日志	group_test_2	
	dtsdcwq9147163ub5h group_test_1	2020年09月23日 18:25:58

#### 您需要在OSS中进行以下操作:

- 1. 开通OSS服务
- 2. 创建存储空间

⑦ 说明 由于目前DLA Lakehouse都是将数据写入OSS的某个空目录内,为了方便进行独立的数据管理,请您尽量选择独立的Bucket。

确保您的数据流在RDS中部署的区域与DTS、DLA、OSS的区域相同。

#### 操作步骤

- 1. 创建湖仓。
  - i. 登录Data Lake Analytics控制台。
  - ii. 在左侧导航栏单击数据湖管理 > 湖仓一体化。
  - iii. 在Lakehouse湖仓一体化页面的湖仓列表页签,单击创建新湖仓。

#### iv. 在**新建湖仓**页面进行参数配置。参数说明如下表所示:

新建湖仓 构建统一的数据湖存储,汇聚各个异构数据源的数据形成数仓体系		
*湖仓名称	填写新建Lakehouse的名字	
描述信息	填写一些关于湖仓备注描述,如此湖仓应用;	场景,应用业务限制等
* 存储类型	OSS 选择Lakehouse数据存储介质,当前仅支持OSS	★型
* 存储路径	: 请仔细规划存储路径,创建后不允许修改;建议 数据被要盖	议选择一个空目录,且不能与之前创建的湖仓目录有互相包含关系,防止历史
编码	UTF8 当前仅支持目标存储数据编码为UTF8	~
参数名称		参数说明
湖仓名称		DLA Lakehouse的名称。
描述信息		湖仓备注描述,例如湖仓应用场景、应用业务限制 等。
存储类型		DLA Lakehouse数据的存储介质,当前仅支持OSS 类型。
		DLA Lakehouse数据在OSS中的存储路径。
存储路径		⑦ 说明 请谨慎规划存储路径,创建后不允 许修改。建议选择一个空目录,且不能与之前

许修改。建议选择一个空目录,且不能与之前 创建的湖仓目录有互相包含关系,防止历史数 据被覆盖。

存储数据的编码类型,当前仅**支**持目标存储数据编码为UT F8。

#### v. 参数配置完成后, 单击**创建**。

湖仓创建成功后,湖仓列表页签中将展示创建成功的湖仓任务。

Lakehouse湖仓一体化								
湖仓列表	工作负载列表							
								创建新湖仓 刷新
ID	名称	描述信息	存储类型	存结路径	工作负载概况	创建时间	操作	
4	Lakehouse123		OSS	oss://123-abcb/Lakehouse123/	0个负载	2021-04-30 15:42:04	创建入湖负载	

2. 创建入湖负载。

编码

i. 在Lakehouse湖仓一体化页面的湖仓列表中, 单击操作列的创建入湖负载。

#### ii. 在新建工作负载页面,进行数据源的基础配置、全量同步配置、增量同步配置、生成目标数据 规则配置。

⑦ 说明 当前仅支持RDS数据源和PolarDB数据源。

■ 基础配置的参数说明如下:

基础配置	
* 名称	请填写工作负载的名称
* 选择湖仓	Lakehouse123
	工作负载将数据输出到所选的湖仓存储空间内,如果还未创建空间,请点击创建新湖仓
* 数据存储格式	HUDI
	HUD!是目前开源很流行的湖仓一体化存储格式,阿里云DLA团队针对HUDI做了大量优化,并持续贡献开源。点击Apache Hudi <b>了解更多</b>
* 源端实例引擎类型	MySQL 🗸
	选择数据来源的引擎类型,例如RDS MySQL或者PolarDB MySQL都是MySQL引擎
* 任务执行Spark虚拟集群	daily-test 🗸 🗸
	目前入湖工作负载是在DLA Spark的VC环境内运行,如果还未创建任何DLA Spark的VC环境,请点击创建DLA Spark虚拟 集群

参数名称	参数说明
名称	工作负载的名称。
选择湖仓	工作负载将数据输出到所选的湖仓存储空间内。可下拉选择已经创建 的湖仓。
数据存储格式	数据的存储格式固定为HUDI。
源端实例引擎类型	数据源的引擎类型。当前仅支持MySQL引擎。
	执行Spark作业的虚拟集群。目前入湖工作负载在DLA Spark的虚拟 集群中运行。如果您还未创建虚拟集群,请进行创建,具体请参见 <mark>创</mark> <mark>建虚拟集群</mark> 。
任务执行Spark虚拟集群	⑦ 说明 请确保您选择的Spark虚拟集群处于正常运行状态,如果您选择的Spark虚拟集群处于非正常运行状态,启动工作负载时将失败。

#### 全量同步配置的参数说明如下:

全量同步配置		
* 实例名称	rm-bp17)	
	选择需要入湖的数据源的实例ID名称。详见rm-bp17u79345w7iq564	
* 用户名	请填写数据库链接的jdbc访问用户名	
* 密码	请填写数据库链接的jdbc访问密码	
* Spark运行所需私有网络ID	vpc- m1	$\sim$
	选择VPC ID, DLA Spark访问数据源时利用ENI技术网络打通该VPC。请查看Spark网络配置	
* Spark运行所需交换机ID	vsw-l >4n	$\checkmark$
	选择交换机ID号,DLA Spark访问数据源时利用ENI技术进行网络打通所需	
* Spark运行所需安全组	sg-l D3x	$\sim$
	选择安全组ID, DLA Spark安全访问数据源时进行网络安全授权所需就转到关联安全组	
* Spark运行所需CU数	24	$\checkmark$
	指定任务执行的DLA Spark所使用的CU数,建议多保留一些CU数,让入湖性能更好、任务更稳定。详细计量计费方题 看	全重

参数名称	参数说明
实例名称	选择需要入湖的数据源的实例ID名称。
用户名	需要入湖的数据源实例的访问用户名。
密码	需要入湖的数据源实例的访问密码。
Spark运行所需私有网络ID	DLA Spark利用ENI技术配置该VPC网络来访问数据 源。关于DLA Spark如何配置数据源VPC网络,请 参见 <mark>配置数据源网络</mark> 。
Spark运行所需交换机ID	DLA Spark运行所需VPC网络下的交换机ID。
Spark运行所需安全组	DLA Spark访问数据源时进行网络安全授权的安全 组ID。您可以到RDS数据源实例的 <b>数据安全性</b> 页面 中获取安全组ID,如未设置安全组请进行添加,具 体操作请参见 <mark>设置安全组</mark> 。
Spark运行所需CU数	指定执行DLA Spark作业所使用的CU数,建议多保 留一些CU数,让入湖性能更好、作业任务更稳 定。

增量同步配置的参数说明如下:

增量同步配置	
* 同步方式	DTS
	选择增量同步通道类型,当前仅支持DTS方式
* 订阅配置	cn_hangzhoutest_version2
	dtsynr 6bak 🗸
	选择增量同步所使用的DTS订阅通道配置(确保可正常使用),分别选择订阅Topic和消费组ID <mark>。请先创建DTS订阅消费此实</mark> 例
* DTS用户名	请填写DTS订阅消费组的用户名
* DTS密码	请填写DTS订阅消费组的密码
* Spark运行所需私有网络ID	vpcm1 ~~~
	选择VPC ID,DLA Spark访问数据源时利用ENI技术网络打通该VPC。请查看Spark网络配置
* Spark运行所需交换机ID	vsw-lo4n 🗸
	选择交换机ID号,DLA Spark访问数据源时利用ENI技术进行网络打通所需
* Spark运行所需安全组	✓ x€0t -p₂
	选择安全组ID,DLA Spark访问数据源时进行网络安全授权所需
* Spark运行所需CU数	24 🗸
	指定任务执行的DLA Spark所使用的CU数,建议多保留一些CU数,让入湖性能更好、任务更稳定。详细计量计费方案查 看
∨ 高级规则配置	置(可选)
消费位	点 earliest V
每批次消费记录等	数 1999998

参数名称	参数说明
同步方式	增量同步的通道类型。当前仅支持DTS方式。
订阅配置	增量同步所使用的DTS订阅通道配置,分别选择订阅Topic和消费组 ID。
DTS用户名	增量同步DTS数据订阅消费组的账号信息。
DTS密码	增量同步DTS数据订阅消费组账号对应的密码信息。
Spark运行所需私有网络ID	DLA Spark利用ENI技术配置该VPC网络来访问数据源。关于DLA Spark如何配置数据源VPC网络,请参见 <mark>配置数据源网络</mark> 。
Spark运行所需交换机ID	DLA Spark运行所需VPC网络下的交换机ID。
Spark运行所需安全组	DLA Spark访问数据源时进行网络安全授权的安全组ID。您可以到 RDS数据源实例的 <b>数据安全性</b> 页面中获取安全组ID,如未设置安全组 请进行添加,具体操作请参见 <mark>设置安全组</mark> 。
Spark运行所需CU数	指定执行DLA Spark作业所使用的CU数,建议多保留一些CU数,让入 湖性能更好、作业任务更稳定。

参数名称	参数说明
高级规则配置	<ul> <li>消费位点:数据消费的时间点。当前取值固定为earliest,表示自动从最开始的时间点获取数据。</li> <li>每批次消费记录条数:表示每次通过DTS拉取的数据量。</li> </ul>

#### ■ **生成目标数据规则配置**的参数说明如下:

生成	目标数据规则配置			
	* 库名前缀 🐧	lakehouse202105071116_		
		生成目标库的数据路径和元信息名称时,会自动加前缀。为了避免海量数据和元数据在DLA中冲突, <mark>请确保输入的前缀在</mark> 当 <mark>前用户下所有工作负数内都不重复</mark>		
	* 前缀应用范围 🐧	数据目录加前缀, 元数据中库表名加前缀 🗸		
		表示上述前缀的应用范围,可以在库的数据目录名中加前缀,也可以在元数据的库名中加前缀		
	库表过滤 🚺	包含 库表选择 库表表达式,逗号分隔,例如db2.*		
		过逾逗号表达式,可以使用模糊符号"*.",例如db1*,db2.table2 表示包含db1的所有表和db2的table2表		
		<b>排除</b> 库表选择   库表表达式,逗号分隔,例如db2.*		
		过海逞号表达式,可以使用模糊符号**.**,例如db3.*, 表示不包含db3的所有表		
	目标端主键字段 🐧	库表选择 可以使用模糊符号**,例如db3.* 主键字段 按顺序多值填写 💙 🕂		
		为指定库表设置主键字段,例如 do Luser_*并设置11.2 表示do 1的所有"user_"前缀的表都使用11.2两个字段为组合主键,若不设置,则系统放大尝试选择主键/唯一 键条作为目标主键字段,若表都没有这些键,这此表现为忽略同步		
	✔ 高级配置(可选)			
	目标端分区字段 🚯	库表选择 可以使用模糊符号***,例如db3.* 分区字段 按顺序多值填写 💙 🕇		
		为指定库表设置分区字段,例如 db1.user_* 并设置gmt_create 表示db1的所有"user_*前缀的表都使用gmt_create字段作为分区字段,若没有设置,则生成的表默认没有 分区		
	源端时区 🐧	系统根据数据演变例所在时区自动解析		
	源端编码 🚺	▶ 系统根据数据源实例设置的编码自动解析		
	目标端版本字段 🚺	3 系统根据通道数据的时间戳自动解析生成目标表的版本字段		
	字段类型转换 🐧	▶ 系统自动根据源端引擎类型和字段类型自动解析缺射到DLA的对应字段类型		
	库/表/列名映射 🐧	目标满库表列的名称生成规则为: 源满库表列名称特殊字符替换为下划线 _ 后的名称		
参	数名称	参数说明		

参数名称	参数说明
库名前缀	生成目标库的数据路径和元信息名称时,会自动添加该前缀。为了避免海量数据和元数据在DLA中冲突,请确保输入的前缀在当前阿里云账号下的所有工作负载内都不重复。
前缀应用范围	设置库名前缀的应用范围。包括: <ul> <li>数据目录加前缀,元数据中库表名加前缀</li> <li>数据目录不加前缀,元数据中库表名加前缀</li> </ul>
库表过滤	设置需要同步的库和表范围。排除的优先级高于包 含。

参数名称	参数说明
	为指定库表设置主键字段。例如: 库表选择输入 db1.user_* , 主键字段输入 f1,f2 , 表示 db1 的所有 user_ 前缀的表都使用 f1,f2 两个字段作为组合主键。
目标端主键字段	⑦ 说明 如果不设置该参数,则系统依次尝试选择表中的主键或唯一键来作为目标端主键字段;如果表中不存在主键或唯一键,则视为忽略同步。
高级配置	目标端分区字段:为指定库表设置分区字段。例 如:库表选择输入 dbl.user_* ,分区字段输 入 gmt_create ,表示 dbl 的所有 user_ 前缀的表都使用 gmt_create 字段作为分区 字段。如果不设置该参数,则生成的表默认没有分 区。

#### iii. 上述参数配置完成后,单击**创建**。

入湖负载创建成功后,在工作负载列表页签中将展示创建成功的工作负载。

Lakehouse	湖仓一体化						
湖仓列表	工作负载列表						
						创建入湖	负载    刷新
ID	名称	关联湖仓	英型	状态	创建时间	操作 🊺	
2	efg_123	abc123	D8实时入湖(全量+增量)	NO STATUS	2021-05-07 15:49:37	洋情 启动	
1	abc_123	abc123	DB实时入湖(全量+增量)	NO STATUS	2021-05-07 15:40:29	洋情 二启动	

#### 3. 启动工作负载。

#### 在工作负载列表页签中,定位到创建成功的入湖负载,在操作列单击启动。

Lakehouse	湖仓一体化					
湖仓列表	工作负载列表					
						创建入湖负载
ID	名称	关联期合	类型	状态	创建时间	操作 🌗
4	baobao12	baobao12	DB实时入湖(全量+增量)	RUNNING	2021-05-08 17:04:55	洋情   停止
3	abc1234	abc	DB实时入湖(全量+增量)	STOPPED	2021-05-08 16:54:44	洋橋   启动   校正
2	efg_123	abc123	DB实时入湖(全量+增量)	NO STATUS	2021-05-07 15:49:37	洋情 启动

工作负载任务启动成功后,状态将由NO STATUS(未启动)变为RUNNING(运行中)。

湖仓列表	工作负载列表					
						创建入湖负载 刷新
ID	名称	关联划合	後型	状态	创建时间	提作 🌗
4	baobao12	baobao12	DB实时入湖(全量+增量)	RUNNING	2021-05-08 17:04:55	详情   停止
3	abc1234	abc	DB宾时入湖(全量+増量)	STOPPED	2021-05-08 16:54:44	详情   启动   校正

您还可以在操作列停止和校正工作负载任务、查看Spark日志。具体说明如下:

操作按钮	含义
详情	单击该按钮,可以查看Spark日志或者UI,并定位工作负载任务启动失败原因。

操作按钮	含义
停止	单击该按钮,可以停止工作负载任务。
	单击该按钮,可以对启动失败的工作负载任务进行数据校正。
校正	⑦ 说明 校正一般使用在库表变更、字段格式不兼容等场景下。校正过程会 重新进行部分存量数据的全量同步,请慎重填写库表筛选表达式,建议使用精 确匹配表达式筛选,避免校正一些不必要的数据。如果未填写需要校正的库 表,则校正失败。

工作负载任务启动成功后,在湖仓列表页签单击存储路径下的OSS路径链接,可以跳转到OSS控制台查 看已经从RDS数据源同步过来的库表路径以及表文件。

○ 数据库路径

abc234 / 华东1 (杭州)	~		
概览		上传文件	新建目录 碎片管理(1) 授权 批量操作 > 刷新
用量查询	>		文件名
文件管理	>	5	/ baobao12/
权限管理	>	🗆 📁	.dla-lakehouse/
基础设置	>		lakehouse20210508173_abc123/
冗余与容错	>		lakehouse20210508173_mysql/

○ 数据表路径

对象存储 / abc234 / 3	文件管理							
abc234 / 华东1 (杭州	b ~							
概览		上传文件	新建目录	碎片管理 (1)	授权	批量操作 🗸	刷新	
用量查询	>		文件名					
文件管理	>	<u>ب</u>	<u>/ baobao</u> i	<u>12/</u> lakehouse2021	0508173_a	bc123/		
权限管理	>		students/					
基础设置	>							
冗余与容错	>							

○ 数据表文件

abc234 / 华东1 (杭州	) ~				版本控制	川 未开通 读写权限 私有 类型 标衡
概范		上传文件	KP 新建目录 荷片管理(1) 接収 加泉協介 > 刷新			请输入文件名前缀匹配
用量查询	>		文件名	文件大小	存储类型	更新时间
文件管理	>		6 / baobao12/ lakehouse20210508173_abc123/ students/			
权限管理	>		🧀 .hoodie/			
基础设置	>		hoodie_partition_metadata	0.091KB	标准存储	2021年5月08日17:06:43
冗余与容错	>		40e80a20-aa9c-4d5f-ae9f-1bb4cd470a9e-0_5-3-8_20210508170631.parquet	424.771KB	标准存储	2021年5月08日17:06:44
传输管理	>		498caec2-4600-49e2-8fd6-9dfa1f9c2224-0_2-3-5_20210508170631.parquet	424.771KB	标准存储	2021年5月08日17:06:44
日志管理	>		70f31792-8555-447e-bc8d-c726c27bc3ef-0_3-3-6_20210508170631.parquet	424.771KB	标准存储	2021年5月08日17:06:44
数据外理	>		7fb6ffce-3b19-47dc-93fe-fc95e31eb7a5-0_1-3-4_20210508170631.parquet	424.771KB	标准存储	2021年5月08日17:06:44
数据安全			83cc1054-16c7-4885-9fc5-26ebd089ecaa-0_0-3-3_20210508170631.parquet	424.772KB	标准存储	2021年5月08日17:06:44
			fdcee6f9-1ca3-4e49-81d2-bd1b2b2c40c6-0_4-3-7_20210508170631.parquet	424.771KB	标准存储	2021年5月08日17:06:44

4. 进行数据分析。

工作负载任务启动成功后,在数据湖管理 > 元数据管理页面中,查看从RDS数据源同步过来的元数据 信息。

云原生数据潮分析	SQL执	行								清法手册	函数手册
概范											
账号管理	testlake	house2324 🛛 😋	执行集群 pu	blic(共享集群) × 国	步执行(F8) 异步执行(F9) 《	和主 格式(化(F10) 主題 ~	]		新建Presto虚拟集群	(i) 255	DMS执行SQL
虚拟集群管理		"双击"切换Schema	1 select	* from `testlakehouse2324`.`d	latest limit 20;		-				
数据动管理 New ^	× 8 1	estlakehouse2324 (current)	2								
元信思发现	~	a diatest									
元数据管理		kanka_timestamptype									
数据入湖		_hoodie_commit_time									
湖合一体化 New		_hoodie_record_key									
Sequerierr Presto		_kafka_timestamp									
500818#		_kafka_offset									
SQL00MM		_hoodie_commit_seqno	执行历史	ANTIAR SOLUTE						11111111111	× be
SQUART		· India have									
SQL监控		_kafka_partition	序号	_kafka_timestamptype	value	_hoodie_commit_time	_hoodie_record_key	_kafka_timestamp	_kafka_offset		detail
Serverless Spark		_kafka_topic	1	0	{age=19, id=1, name=lisi}	20210811155259	99328a3e-62fc-4918-a8df-30	1628667704347000	0		洋街
作业管理	<	_hoodie_partition_path	4								÷
QuickBI报表描述		III day									
注重客餐		i hour									

单击**操作**列的**查询数据**,在Serverless Presto > SQL执行页面,查看从RDS数据源同步过来的全量表数据。

云原生数据潮分析	1 50	山地行																382+32.00		-
概法	130	E1744 1																國法子の	明秋于日	-
张易管理	lakı	house2	2105081	0	c	执行集群	诸这样		> ■歩执行	(F8) 🛱	步统(〒(F9)		格式(k/F10)	=腰 ∨			新聞Prestol思報	1111 <b>()</b> ()	录DMS执行SO	
41/00/00		-75	击"切换Sc	hema	Ŭ			1. het	0508172		1.1.1.1									
	~	🛢 lake	nouse2021	0508173	l_abc123	2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2	eet - 11	on issenouse2021	0306113_006123	. student:	s 11011 2	0.								
AUMORENE New		× = :	tudents																	
元信息发现			💼 _hoodi	ie_comm	it_time															
元数据管理			_hoodi	ie_comm	it_seqno															
数据入湖			hoodi	ie_record	_key															
潮会一体化 New			hoodi	ie_partiti	on_path															
Serverless Presto			id id	e_file_na	me															
SQL访问点			name																	
SOLINE	1		hight																	
coultra			🖬 age																	
200765	1		weight			ļ.														
Serverless Spark						-														
作业管理																				
优惠查督						执行	历史	执行结果 SQL型	± ()									导出结果集	~ 隐蔽	
系统设置						序号		_hoodie_file_na	ime	id			name		hight	 age	 weight		d.	-
						1		498caec2-4600	-49e2-8fd6-9d	з			王虎		154	15	47		<b>;</b> #tr	1¢
						2		83cc1054-16c7	-4885-9fc5-26	1			₩E		151	15	43		洋街	
						3		70f31792-8555	-447e-bc8d-c7	4			柯南		152	17	50		洋橋	
						4		7fb6ffce-3b19-	47dc-93fe-fc9	2			张四		153	16	44		洋情	
						5		40e80a20-aa9c	-4d5f-ae9f-1b	6			柯西		157	19	52		洋悟	
						6		fdcee6f9-1ca3-	4e49-81d2-bd	5			柯北		156	18	51		洋橋	

如果您在数据源RDS中变更了原始数据,在Serverless Presto > SQL执行页面进行查询时,数据会同 步进行更新。

### 3.2. Kafka实时入湖建仓分析

DLA Lakehouse实时入湖方案利用数据湖技术,重构数仓语义,分析数据湖数据,实现数仓的应用。本文介 绍Kafka实时入湖建仓分析的操作步骤。

#### 前提条件

- 已在DLA中开通云原生数据湖分析服务。更多信息,请参见开通云原生数据湖分析服务。
- 已创建Spark引擎的虚拟集群。更多信息,请参见创建虚拟集群。
- 如果您使用RAM子账号登录,还需要进行如下操作:
  - 已授予RAM子账号AliyunDLAFullAccess权限。更多信息,请参见为RAM账号授权。
  - 已将DLA子账号绑定到RAM子账号。更多信息,请参见DLA子账号绑定RAM账号。
- 已在Kafka中创建实例。更多信息,请参见创建实例。
- 已在Kafka中创建Topic。更多信息,请参见创建Topic。
- 已在Kafka中发送消息。更多信息,请参见发送消息。

#### 方案介绍

DLA Lakehouse的Kaf ka实时入湖建仓分析助力企业实现"业务数据化"+"建仓"+"数据业务化"的数据闭 环建设,主要包括三方面。

- Kafka实时入湖建仓引擎:支持T+10min近实时入湖,同时支持Schema推断及变更、嵌套打平、分区管理、小文件合并及Clustering等能力。
- OSS存储:入湖建仓基于OSS的低成本存储介质,有效地降低存储层成本,同时格式层使用HUDI格式来支持Append写入。
- 完全弹性的分析: DLA支持Serverless Presto和Serverless Spark的分析与计算能力,完全按需计费。

方案架构如下图所示。



#### 注意事项

Kafka中创建的Topic数据超过一定的时间会被自动清理,如果Topic数据过期,同时入湖任务失败,再重新 启动时读取不到被清理掉的数据,会有丢失数据的风险。因此请适当调大Topic数据的生命周期并在入湖任 务失败时及时提交工单或者钉钉咨询DLA答疑获得支持。

#### 授予RAM子账号创建库表的权限

如果您使用RAM子账号登录,在开始Kafka实时入湖的操作前,需要先授予RAM子账号创建库表的权限。

- 1. 登录Data Lake Analytics控制台。
- 2. 在左侧导航栏单击Serverless Presto > SQL执行。
- 3. 在右侧运行框输入如下语句,单击同步执行(F8)。

grant create,alter on \*.\* to user1\_s1041577795224301;

#### 操作步骤

- 1. 创建湖仓。
  - i. 登录Data Lake Analytics控制台。
  - ii. 在左侧导航栏单击数据湖管理 > 湖仓一体化。
  - iii. 在Lakehouse湖仓一体化页面的湖仓列表页签,单击创建新湖仓。

#### iv. 在**新建湖仓**页面进行参数配置。参数说明如下表所示:

新建湖仓 构建统一的	的数据湖存储,汇聚各个异构数据源的数据形	形成数仓体系
* 湖仓名称	填写新建Lakehouse的名字	
描述信息	填写一些关于湖仓备注描述,如此湖仓应用:	场景,应用业务限制等
* 存储类型	OSS	~
	选择Lakehouse数据存储介质,当前仅支持OSS	类型
* 存储路径		
	请仔细规划存储路径,创建后不允许修改;建议 数据被要盖	X选择一个空目录,且不能与之前创建的湖仓目录有互相包含关系,防止历史
编码	UTF8	~
	当前仅支持目标存储数据编码为UTF8	
会物夕む		会粉记旧
学致口你		
湖仓名称		DLA Lakehouse的名称。
描述信息		湖仓备注描述,例如湖仓应用场景、应用业务限制 等。
存储类型		DLA Lakehouse数据的存储介质,当前仅支持OSS 类型。
		DLA Lakehouse数据在OSS中的存储路径。
存储路径		⑦ 说明 请谨慎规划存储路径,创建后不允 许修改。建议选择一个空目录,且不能与之前

创建的湖仓目录有互相包含关系,防止历史数 据被覆盖。

存储数据的编码类型,当前仅支持目标存储数据编码为UT F8。

#### v. 参数配置完成后, 单击**创建**。

湖仓创建成功后,湖仓列表页签中将展示创建成功的湖仓任务。

Lakehouse湖仓一体化													
湖仓列表	工作负载列表												
								创建新湖仓 刷新					
ID	名称	描述信息	存储类型	存结路径	工作负载概况	创建时间	操作						
4	Lakehouse123		OSS	oss://123-abcb/Lakehouse123/	0个负载	2021-04-30 15:42:04	创建入湖负载						

2. 创建入湖负载。

编码

i. 在Lakehouse湖仓一体化页面的湖仓列表中, 单击操作列的创建入湖负载。

- ii. 在新建工作负载页面,选择数据源为Kafka数据源。
- iii. 进行数据源的基础配置、增量同步配置、数据解析配置、生成目标数据规则配置。
  - 基础配置的参数说明如下:

新建工作负载	
RDS数据源 PolarDB数	据源 Kafka数据源
基础配置	
* 名称	testlakehouse
* 选择湖仓	Lakehouse123
	工作负载将数据输出到所选的湖仓存储空间内,如果还未创建空间,请点击创建新湖仓
* 数据存储格式	нирі 🗸
	HUDI是目前开源很流行的湖仓一体化存储格式,阿里云DLA团队针对HUDI做了大量优化,并持续贡献开源。点击Apache Hudi,了解更多
* 任务执行Spark虚拟集群	st
	目前入湖工作负载是在DLA Spark的VC环境内运行,如果还未创建任何DLA Spark的VC环境,请点击 创建DLA Spark虚拟集群
* 子账户执行所需RAM角色	377729437257872923(AliyunDLASparkProcessingDataRole)
c	子账户提交Spark 作业时所使用的RAM角色,详见Spark RAM子账户授权

参数名称	参数说明			
名称	工作负载的名称。			
选择湖仓	下拉选择已经创建的湖仓,工作负载将数据输出到所选的湖仓内。			
数据存储格式	数据的存储格式固定为HUDI。			
任务执行Spark虚拟集群	执行Spark作业的虚拟集群。目前入湖工作负载在DLA Spark的虚拟 集群中运行。如果您还未创建虚拟集群,请进行创建,具体请参见创 建虚拟集群。 ⑦ 说明 请确保您选择的Spark虚拟集群处于正常运行状 态,如果您选择的Spark虚拟集群处于非正常运行状态,启动工			
子账号执行所需RAM角色	子账号提交Spark作业时所使用的RAM角色,固定选 择 <b>AliyunDLASparkProcessingDataRole</b> 。更多信息,请参见 <mark>细</mark> 粒度配置RAM子账号权限。			

#### 增量同步配置的参数说明如下:

增量同步配置							
* 实例名称	7(dtstest)	7(dtstest)					
	选择需要入湖的数据源的实例ID名称。详见ali	kafka_post-cn-mp91gcgu2007					
* topic选择	dlatest	$\checkmark$					
消费位点	earliest	~					
	入潮消费数据的位点,latest表示最近位点; ex <mark>详见</mark>	arliest表示最老的位点,如果命中kafka数据清理逻辑可能会读取不到位点					
* Spark运行所需CU数	24	$\sim$					
	指定任务执行的DLA Spark所使用的CU数,建 详细计量计费方案查看	议多保留一些CU数,让入湖性能更好、任务更稳定。					
参数名称		参数说明					
实例名称		选择需要入湖的数据源的实例ID和实例名称。					
topic选择		在Kafka中创建的Topic名称。					
消费位点		入湖消费数据的位点,latest表示最近位点; earliest表示最老的位点。更多信息,请参见 <mark>消息</mark> <mark>队列Kafka版何时删除旧消息?</mark> 。					
Spark运行所需CU数	文	指定Kafka实时入湖任务所需DLA Spark的CU数。					

#### ■ 数据解析配置的参数说明如下:

数据解析配置		
* 消息数据格式	JSON	•
	选择kafka value存储数据的格式类型,对于JSON类型支持json展开	
* Schema设置模式	● 手动设置 ● 自动识别	
	手动模式会从kafka的topic获取一条样例数据进行schema的推断设置	_
* 消息样例数据	"key", "value":{ "id": 1, "name": "lisi", "age": 19 }	
	可以通过调整样例数据进行schema调整	_
* 嵌套打平层数	0	-
	支持json的嵌套打平,0表示不打平,1表示打平一层,默认添加上一层名字作为前缀,_分割	
* Schema预览	root value: struct (nullable = true)   age: long (nullable = true)   id: long (nullable = true)   name: string (nullable = true)   kafka topic: string (nullable = true)	•
	kafka的value字段进行schema处理后在dla表中的schema,其中"_kafka_"开头的为kafka消息的内置字段。具体的消息如果	R
	wttschema中的子般或目小为json,对应的子校立攻重为hull。	

参数名称	参数说明
消息数据格式	Kafka的Value字段的数据存储格式,固定为 JSON。
Schema设置模式	默认为 <b>手动设置</b> ,会从Kafka的Topic获取一条样 例数据进行Schema的推断设置。
消息样例数据	通过调整样例数据对Schema进行调整。
嵌套打平层数	设置JSON的嵌套打平层数,取值如下。 • 0:不打平。 • 1:打平一层。 • 2:打平两层。 • 3:打平三层。 • 4:打平四层。 • 5:打平五层。
Schema预览	Kafka Value进行JSON解析的Schema模板,同时 作为DLA表的Schema。

#### ■ **生成目标数据规则配置**的参数说明如下:

生成目标数据规则配置												
* 库名	t	estlakehouse2324		表名	dlate	est						
	同步 符开	到dla的库名称,不能和已有库同名。) 头,可包含字符、数字、下划线	车名以字		直接使用 母及数书	用kafka的topic名作为表名,其中会把除去号 字的字符替换为_	7					
目标满分区字段	×	源跳字段名kafka_timestamp		格式处	理方法	时间楷式化	/	源端字段格式	微秒级精度	$\sim$		
		目标分区字段名 day		目标分	区配置	yyyy-MM-dd						
	×	源端字段名 _kafka_timestamp		格式处	理方法	时间植式化	/	源端字段格式	微秒级精度	$\sim$		
		目标分区字段名 hour		目标分	区配置	нн						
	对于 格式 十 1	"要同步的kafka topic支持设置多级分区 1,并设置分区目标字段的"名称"及"时 <mark>算加一条</mark>	Z。注意对于"格式处理方法 间格式化方式"比如"yyyy",	去"为"时间 ,详见 jav	明格式化" /a的form	需要选择一个时间字段(也可以使用kafka消 nat	息的时	j间字段'_kafka_ti	imestamp"),之后指定该字段的时间			
主键生成规则	策	IN UUID				$\sim$						
	主領	自动生成UUID,相当于数据APPEND,	而不做更新及删除									
												_
											3500 <del>(</del> 152	重豐

参数名称	参数说明
库名	目标DLA库的名称,库名以字符开头,可包含字 符、数字、下划线,不能和已有库同名否则会报 错。
表名	系统默认使用Kafka的Topic名作为表名,且不允 许用户更改。
目标端分区字段	支持对要同步的Kafka Topic设置多级分区。您可 以按需设置分区字段。
主键生成规则	仅支持UUID。UUID作为主键,数据只能新增插 入,不能做更新和删除。

#### iv. 上述参数配置完成后, 单击**创建**。

入湖负载创建成功后,在工作负载列表页签中将展示创建成功的工作负载。

Lakeho	use湖仓一体化							
湖仓列录	工作负载列表							
								创建入湖负载 刷新
ID	名称	关联湖合	库前缀	类型	状态	创建时间	修改时间	操作 🚺
68	testlakehouse	Lakehouse123	testlakehouse2324	KAFKAI繼星实时入湖	NO STATUS	2021-08-11 15:46:39	2021-08-11 15:46:39	洋情   启动   删除

#### 3. 启动工作负载。

- i. 在工作负载列表页签中,定位到创建成功的入湖负载,在操作列单击启动。
- ii. 在弹出的再次确认窗口单击确定。

工作负载任务启动成功后,状态变为RUNNING。

Lake	house湖仓一体化							
湖仓	列表 工作负载列表							
								创建入湖负载
ID	名称	关联湖合	库前缀	遊	状态	创建时间	修改时间	攝作 🚺
68	testlakehouse	Lakehouse123	testlakehouse2324	KAFKA增量实时入湖	RUNNING	2021-08-11 15:46:39	2021-08-11 15:46:39	洋情   停止   任务UI

- 4. 进行数据分析。
  - i. 在工作负载列表页签单击入湖负载的库前缀。

ii. 在SQL执行页面,系统已经选中了目标库,您可以直接输入SQL语句进行数据分析。

云原生数据端分析	SQL执行								语法手册	函数手册
概范						_		-		
账号管理	testlakehouse2324 🛛 😋 😷	执行集群 publi	ic(共享集剧) > F	同步执行(F8) 异步执行(F9) 例	Ⅰ止 格式化(F10) 主题 >			新建Presto世科集群	() 923	DMS执行SQL
虚拟集群管理	"双击"切换Schema	1 select *	from `testlakehouse2324`.`	dlatest limit 20:						
数据动管理 New へ	testlakehouse2324 (current)	2								
元信意发现	✓ ■ diatest ■ kafka timestamotyne									
元数据管理	a value									
数据入期	_hoodie_commit_time									
御令一体化 New	_hoodie_record_key									
formed as Dente	_kafka_timestamp									
Serveriess Presto	_kafka_offset									
SQL语问点	_hoodie_file_name								_	
SQL执行	_hoodie_commit_sequo	执行历史	SQUEER SQUEER	3				4	9出结果集	∨ 陰嶽
SQL <u>18</u> 39	I _kafka_partition	序号	_kafka_timestamptype	value	_hoodie_commit_time	_hoodie_record_key	_kafka_timestamp	_kafka_offset		detail
Serverless Spark	🗮 _kafka_topic	1	0	{age=19, id=1, name=lisi}	20210811155259	99328a3e-62fc-4918-a8df-30	1628667704347000	0		洋價
作业管理	_hoodie_partition_path	4								×
QuickBI招專搭量	🔳 day									
	t hour									

- 5. (可选) OSS数据存储管理。
  - i. 在湖仓列表页签单击存储路径下的OSS路径链接。
  - ii. 在OSS控制台查看已经从Kafka数据源同步过来的库表路径以及表文件。
    - 数据库路径: /lakehouse123/testlakehouse2324/。
    - 表路径: /lakehousetest123/testlakehouse2324/dlatest。
    - 表文件路径: /lakehousetest123/testlakehouse2324/dlatest/day=2021-8-11。

対象存储 / Buck	et 列表 / alicd	n-log-delivery-14	06926474064770-cn-hangzhou / 义件管理				任务列表
alicdn-log-de	livery-1406920	5474064770 - cn - l	hangzhou / 华乐1(杭州) 🗸 🗸		版本控	制 未开通 读写权限 私有 类型 标准	豊存儲 (本地冗余)
概范		上传文件	新建日表 科州智语 授权 能温度作 > 局衡			请输入文件名前缀匹配	Q ®
用量查询	>		文件名	文件大小	存储类型	更新时间	操作
文件管理	>		/ testlakehouse2324/ dlatest/ day=2021-08-11/ hour=15/				
权限管理	>		.hoodie_partition_metadata	0.091KB	标准存储	2021年8月11日 15:53:10	洋橋 更多 >
基础设置	>		05d/d04d-5dbd-4d10-b04b-e7aa6e908b/e-0_0-7-7_20210811155259.parquet	426.409KB	标准存储	2021年8月11日 15:53:12	洋橋 更多 Y

↓ 注意 请勿删除从Kafka数据源同步过来的库表文件,否则会有丢失数据的风险。

## 3.3. Lindorm实时入湖建仓分析

DLA Lakehouse实时入湖方案利用数据湖技术,重构数仓语义,分析数据湖数据,实现数仓的应用。本文介 绍Lindorm实时入湖建仓分析的操作步骤。

#### 方案介绍

DLA Lakehouse的Lindorm实时入湖建仓分析助力企业构建大数据离在线一体化,主要包括三方面。

- Lindorm实时入湖建仓引擎:支持T+10min近实时入湖,同时支持Schema推断、动态列增加、分区管理、 小文件合并及Clustering等能力。
- Lindorm DFS存储: Lindorm入湖建仓数据回流到Lindorm DFS中,用户无需额外购买其他存储,有效地降低运维管理成本。
- 完全弹性的分析: DLA支持Serverless Presto和Serverless Spark的分析与计算能力,完全按需计费。

#### 方案架构如下图所示。



#### 使用限制

- Lindorm宽表引擎版本必须大于等于2.1.28。
- Lindorm通道服务版本必须大于等于3.5.0。
- Lindorm文件引擎版本必须大于等于3.10.3。

#### 前提条件

- 已在DLA中开通云原生数据湖分析服务。更多信息,请参见开通云原生数据湖分析服务。
- 已创建Spark引擎的虚拟集群。更多信息,请参见创建虚拟集群。
- 已创建Presto CU版虚拟集群。更多信息,请参见DLA Presto CU版本快速入门。
- 如果您使用RAM子账号登录,还需要进行如下操作:
  - 已授予RAM子账号AliyunDLAFullAccess权限。更多信息,请参见为RAM账号授权。
  - 已将DLA子账号绑定到RAM子账号。更多信息,请参见DLA子账号绑定RAM账号。
- 已在Lindorm中开通数据湖分析。

#### 注意事项

- Lindorm通道中的消息数据默认只保留7天,如果数据过期,同时入湖任务失败,再重新启动时读取不到过期的数据,会有丢失数据的风险。因此请在入湖任务失败时及时提交工单或者钉钉咨询DLA答疑获得支持。
- 若您无法在云原生多模数据库Lindorm控制台开通数据湖分析服务,请及时提交工单或者钉钉咨询DLA答 疑获得支持。

#### 操作步骤

- 1. 若使用RAM子账号登录,需授予RAM子账号创建库表的权限。
  - i. 登录Data Lake Analytics控制台。
  - ii. 在左侧导航栏单击Serverless Presto > SQL执行。

iii. 在右侧运行框输入如下语句, 单击**同步执行(F8)**。

grant create,alter on \*.\* to user1\_s1041577795224301;

⑦ 说明 您需要替换 user1 s1041577795224301 为RAM子账号绑定的DLA子账号名称。

- 2. 创建湖仓。
  - i. 登录Data Lake Analytics管理控制台。
  - ii. 在左侧导航栏单击数据湖管理 > 湖仓一体化。
  - iii. 在Lakehouse湖仓一体化页面的湖仓列表页签,单击创建新湖仓。

#### iv. 在新建湖仓页面进行参数配置。参数说明如下:

┃ 新建湖仓 构建统一的	数据湖存储,汇聚	各个异构数据源的数据形成数仓体系					
* 湖仓名称	lindorm_dla_test	t					
描述信息	填写—些关于湖·	填写—些关于湖仓备注描述,如此湖仓应用场景,应用业务限制等					
* 存储类型	LindormDFS(兼署 选择Lakehouse目标	LindormDFS(兼容HDFS) V					
* Lindorm DFS集群	选择Lindorm集群ID	m028n) w 各组件最低版本要求 Lindorm宽表服务(LD):2.1.28 Lindorm传输服务(LTS):3.5.0、Lindorm存储服务					
* 存储路径	(LDFS):3.10.3,名元法 /dla/lindorm_dla 请仔细规划与选择符	320年2020年前,请到经制召开10年33开33 a_test 季储路径,创建后不允许修改;建议选择一个空目录且不能与之前创建的湖仓目录有互相包含关系防止历					
编码	史政攝被覆盖;者在 UTF8 当前仅支持目标存储	?储奕型刀OSS, 请确保选择OSS目 変为非归档类型 ◆ 裁拟据编码为UTF8					
参数名称		参数说明					
湖仓名称		DLA Lakehouse的名称。					
描述信息		湖仓备注描述,例如湖仓应用场景、应用业务限制等。					
存储类型		DLA Lakehouse数据的存储介质,选择 <b>LindormDFS(兼容</b> HDFS)类型。					
Lindorm DFS集群		Lindorm集群的ID。					
存储路径		DLA Lakehouse数据在Lindorm中的存储路径,限定以"/dla"开头。 ⑦ 说明 请谨慎规划存储路径,创建后不允许修改。建议选择 一个空目录,且不能与之前创建的湖仓目录有互相包含关系,防 止历史数据被覆盖。					
编码		存储数据的编码类型,当前仅支持目标存储数据编码为UTF8。					

#### v. 参数配置完成后, 单击**创建**。

湖仓创建成功后,湖仓列表页签中将展示创建成功的湖仓任务。

湖仓列制	<b>湖仓列表</b> 工作负载列表								
								创建新湖合刷新	
ID	名称	描述信息	存储类型	存储路径	工作负载概况	创建时间	修改时间	授作	
74	lindorm_dla_test		LindormDFS	hdfs://ld-bp17j28j2y7pm028n-proxy-haproxy.lindo rm.rds.aliyuncs.com:9000/dla/lindorm_dla_test	1个负载	2021-08-24 15:11:56	2021-08-24 15:11:56	创建入湖负载	
64	Lakehouse123		OSS	oss://alicdn-log-delivery-1406926474064770-cn-h angzhou/	1个负载	2021-08-11 15:35:44	2021-08-11 15:35:44	创建入湖负载	

#### 3. 创建入湖负载。

- i. 在Lakehouse湖仓一体化页面的湖仓列表中, 单击操作列的创建入湖负载。
- ii. 在新建工作负载页面,选择数据源为Lindorm数据源。
- iii. 进行数据源的基础配置、全量和增量同步配置、生成目标数据规则配置。
  - 基础配置的参数说明如下:

基础配置	
* 名称	lindorm-dla
* 选择湖仓	lindorm_dla_test V
	工作负载将数据输出到所选的湖仓存储空间内,如果还未创建空间,请点击创建新湖仓
* 数据存储格式	HUDI
	HUDI是目前开源很流行的湖仓一体化存储格式,阿里云DLA团队针对HUDI做了大量优化,并持续贡献开源。点击A <mark>pache</mark> Hudi , 了解更多
* 任务执行Spark虚拟集群	jobtest 🗸
	目前入湖工作负载是在DLA Spark的VC环境内运行,如果还未创建任何DLA Spark的VC环境,请点击 创建DLA Spark虚拟集群
* 子账户执行所需RAM角色	377729437257872923(AliyunDLASparkProcessingDataRole)

参数名称	参数说明
名称	工作负载的名称。
选择湖仓	下拉选择已经创建的湖仓,工作负载将数据输出到所选的湖仓内。
数据存储格式	数据的存储格式固定为HUDI。
任务执行Spark虚拟集群	执行Spark作业的虚拟集群。目前入湖工作负载在DLA Spark的虚拟 集群中运行。如果您还未创建虚拟集群,请进行创建,具体请参见创 建虚拟集群。 ⑦ 说明 请确保您选择的Spark虚拟集群处于正常运行状态,如果您选择的Spark虚拟集群处于非正常运行状态,启动工 作负载时将失败。
子账号执行所需RAM角色	子账号提交Spark作业时所使用的RAM角色,固定选 择 <b>AliyunDLASparkProcessingDataRole</b> 。若使用主账号登录, 请忽略此参数。更多信息,请参见 <mark>细粒度配置RAM子账号权限</mark> 。

■ **全量和增量同步配置**的参数说明如下:

全量和增量同步配置		
* Lindorm集群	m028n	
	选择数据源关联的Lindorm集群,当前只能与湖仓关联的集群相同。	
* 实时订阅通道ID	ede8	$\checkmark$
	选择Lindorm数据的DLA类型订阅通道ID,若不存在,请点击开通此实例的数据湖分析订阅通道	
* Spark运行所需私有网络ID	52ld	$\sim$
	选择VPC ID,DLA Spark访问数据源时利用ENI技术网络打通该VPC。请查看Spark网络配置	
* Spark运行所需交换机ID	38qh7	$\sim$
	选择交换机ID号,DLA Spark访问数据源时利用ENI技术进行网络打通所需	
* Spark运行所需安全组	3m3	$\checkmark$
	选择安全组ID,DLA Spark安全访问RDS数据源时进行网络安全授权所需,需确保安全组的出方向能访问当前的数据	源。
* Spark运行所需CU数	16	$\sim$
	指定任务执行的DLA Spark所使用的CU数,建议多保留一些CU数,让入湖性能更好、任务更稳定。 <mark>详细计量计费方案查看</mark>	
∨ 高级规则配置	置 (可选)	
消费位	立点 最老点位(earliest)	$\checkmark$
每批次消费记录祭	2数 100	

参数名称	参数说明
Lindorm集群	默认为与数据源关联的Lindorm集群,且不允许用户更改。
实时订阅通道ID	Lindorm数据的DLA类型订阅通道ID。
Spark运行所需私有网络ID	固定值,不支持修改。DLA Spark利用ENI技术配置该VPC网络来 访问数据源,此私有网络ID和Lindorm集群所使用的私有网络ID 相同。更多信息,请参见 <mark>配置数据源网络</mark> 。
Spark运行所需交换机ID	固定值,不支持修改。该参数表示DLA Spark运行所需VPC网络 下的交换机ID。
Spark运行所需安全组	DLA Spark访问数据源时进行网络安全授权的安全组ID。如果未 设置安全组请登录VPC <mark>控制台</mark> 进行添加。
Spark运行所需CU数	Lindorm实时入湖任务所需DLA Spark的CU数。默认情况下无需 修改,若想执行更快,可以增加每批次消费记录条数,并同时增 加CU数。
高级规则配置 (可选)	<ul> <li>消费位点:数据消费的时间点。当前取值固定为earliest,表示自动从最开始的时间点获取数据。</li> <li>每批次消费记录条数:每次通过Lindorm拉取的数据量。</li> </ul>

#### ■ **生成目标数据规则配置**的参数说明如下:

主成目标数据规则配置									
* 库名	testlindorm	表名	test01						
	同步到dia的库名称,不能知已有库同名。库名以字符开头,可包含字符、数字、下划线	直接	直接使用Lindorm入湖通道的表,其中会把非字符数 字的字符替换为"_"						
目标端分区字段	B X 源端字段名 regionid	格式处理方	法 指定分区字段	$\sim$	目标分区字段名	regionid			
∨ 高级配置(可选)									
忽略数据删除									
	在这里填写的表,表示即使源库中数据被删除了,在目标湖仓中,	也不会删除这	些数据,会强制忽略任何对目标数据的删除	余行为(只	已插入和更新)				
目标表的主键生成规则	」 系统根据Lindorm表主键字段自动生成主键								
源端时回	系统根据数据源实例所在时区自动解析								
源端编码	3 系统根据数据源实例设置的编码自动解析								
目标端版本字段	Q 系统根据通道数据的时间戳自动解析生成目标表的版本字段	系统根据通道数据的时间就目动解析生成目标表的版本字段							
字段类型转换	系统自动根据源调引擎英型和字段英型自动解析映射到DLA的对应字段英型								
库/表/列名映	目标编库表列的名称生成规则为: 源编库表列名称特殊字符替换为下划线"」后的名称								
其他配置									

参数名称	参数说明
库名	目标DLA库的名称,库名以字符开头,可包含字符、数字、下划线,不能 和已有库同名否则会报错。
表名	系统默认使用Lindorm入湖通道的表名,会自动把除数字和字母外的字符 替换为"_",且不允许用户更改。
目标端分区字段	对指定库表设置分区字段。不设置则不会进行数据分区。
高级配置(可选)	默认忽略数据删除:即使源库中数据被删除了,在目标湖仓中,也不会删 除这些数据。

#### iv. 上述参数配置完成后, 单击**创建**。

入湖负载创建成功后,在工作负载列表页签中将展示创建成功的工作负载。

- 4. 启动工作负载。
  - i. 在工作负载列表页签中,定位到创建成功的入湖负载,在操作列单击启动。
  - ii. 在弹出的再次确认窗口单击确定。

工作负载任务启动成功后,状态变为RUNNING。

湖仓列	湖台列表 <u>工作负载列表</u>								
	创建入和负载								
ID	名称	关联制合	湖仓存储类型	库前缀	英型	状态	创建时间	修改时间	援作 🌗
93	lindorm-test	lindorm_dla_test	LindormDFS	testlindorm	Lindorm实时入湖(全量+增量)	RUNNING	2021-08-24 17:39:09	2021-08-24 17:39:09	详情   停止   任务UI
68	testlakehouse	Lakehouse123	OSS	testlakehouse2324	KAFKA增量实时入湖	STOPPED	2021-08-11 15:46:39	2021-08-11 15:46:39	编辑 详情 启动 删除

- 5. 待工作负载任务运行一段时间后,进行数据分析。
  - i. 在工作负载列表页签单击入湖负载的库前缀。
  - ii. 在弹出的Lindorm湖仓查询提示窗口单击确定。
  - iii. 在SQL执行页面,系统已经选中了目标库,您可以直接输入SQL语句进行数据分析。

#### Lindorm与DLA的字段类型转换表

目前支持的字段类型转换如下。

Lindorm字段类型	DLA字段类型
long、usigned_long、short、unsigned_short、int、unsigned_integer、 byte	bigint
string、date、unsigned date、time、unsigned time、timestamp、 unsigned timestamp、decimal	string
float、double、unsigned double	double
boolean	boolean
varbinary、binary、encoded binary	binary

⑦ 说明 不支持的Lindorm字段类型不会同步到DLA中。

#### 常见问题

Q:为什么Lindorm中开通数据湖分析后,下拉选择入湖表没有显示已创建的表?

A:目前数据湖分析仅支持有Schema的表,即Lindorm通过CQL创建的表。更多信息,请参见Lindorm CQL操 作文档。

### 3.4. T+1全量同步一键建仓

### 3.4.1. 概述

本文简要介绍了T+1全量同步一键建仓方案。

#### 业务挑战

大部分阿里云用户会将阿里云关系型数据库RDS、PolarDB for MySQL、MongoDB或者云服务器ECS自建数据 库作为业务系统数据库,随着业务数据的增长,业务数据库存储的数据也越来越多。但RDS、PolarDB for MySQL、MongoDB或者ECS自建数据库的计算能力有限,如果直接使用RDS、PolarDB for MySQL、 MongoDB或者ECS自建数据库搭建数据仓库,会占用线上业务的计算资源从而会影响业务的运行。如果使用 自建开源大数据生态体系,例如Hive、Spark等,需要专门的大数据工程师来操作和运维,且操作流程也不 像使用MySQL一样简单,成本极高。

#### 解决方案

T+1全量同步一键建仓是指通过DLA控制台配置数据源(RDS、PolarDB for MySQL、MongoDB数据源、ECS 自建数据库数据)和目标OSS数据仓库,系统按照您设定的数据同步时间自动、无缝的帮您把数据源中的数 据同步到目标数据仓库OSS中,同时在数据仓库和DLA中创建与数据源表相同的表结构,基于目标数据仓库 进行数据分析,不影响数据源端的线上业务运行。



#### 方案优势

T+1全量同步一键建仓方案有以下优势:

- 一键同步数据源(RDS、PolarDB for MySQL、MongoDB数据源、ECS自建数据库数据)中的上千张表数 据,无需其他额外的配置工作。
- 提供Serverless的服务,无需维护任何实例,零运维成本。
- 数据同步过程中,数据源中的数据存储在目标数据仓库OSS中,搭建数据仓库过程中不会对数据源端的业务系统产生任何压力。
- 通过DLA控制台配置建仓任务,支持自定义设置定时数据投递。
- 极致的计算能力,充分发挥DLA的计算能力,通过DLA的大内存、并发计算完成复杂的多表JOIN以及数据仓 库搭建过程中的其他必要操作。

### 3.4.2. 如何使用T+1全量同步一键建仓

本文档主要介绍了T+1全量同步一键建仓功能的使用流程。

#### 前提条件

使用T+1全量同步一键建仓功能之前,您需要完成以下准备工作:

⑦ 说明 创建OSS数据仓库时,要求DLA、RDS、OSS所属Region相同,否则无法使用T+1全量同步一键建仓功能。

- 根据不同的数据源进行如下操作:
  - 如果是云RDS数据源,请完成接入RDS数据源的准备工作,具体请参见快速入门。
  - 如果是PolarDB for MySQL数据源,请完成接入PolarDB for MySQL数据源的准备工作,具体请参见快速入门。
  - 如果是MongoDB数据源,请完成接入MongoDB数据源的准备工作,具体请参见入门概述。

- 在OSS中进行如下操作:
  - i. 开通OSS服务。
  - ii. 创建存储空间。
  - iii. 创建目录。

② 说明 根据业务需求,判断是否需要创建目录来存储RDS、PolarDB for MySQL、MongoDB 数据源或ECS自建数据库数据。

- 在DLA中进行如下操作:
  - i. 开通DLA服务, 具体请参见开通云原生数据湖分析服务。
  - ii. 初始化DLA数据库主账号密码,具体请参见重置数据库账号密码。

#### 操作步骤

• 一键创建OSS数据仓库

创建一键建仓任务后,根据业务需要,您也可以在**数据入湖**页面的**任务列表**页签中,在操作列单击执 行随时手动发起数据同步操作,在目标数据仓库OSS中创建与数据源(RDS数据源、ECS自建数据库数据) 相同的表结构、在DLA中创建对应的数据仓库表结构。

数据入湖										
	- 健建合 - 健康合誉									
> 更多										
任务列表	历史列表									
全部类型	$\sim$								用的	鏩
ID	名称	Schema名称	类型	调度计划	调度状态	最近运行状态	最近运行耗时	创建时间	操作	
425	dla_test -> dla_test_polar	dia lar	一键建合	每天00时30分	停止调度	• FAILED 🍞 🛄	3580	2021-03-20 19:41:59	执行 編輯 历史 删除	
424	dis_it_db -> dis_it_db	d	一键建仓	每天00时30分	开启调度	SUCCESS ⑦	17後	2021-03-19 14:13:26	执行   编辑   历史   删除	

 一键建仓任务进行数据同步时,您可以在数据入湖页面的历史列表页签中,在操作列单击详情实时查看 同步任务的执行状态等信息。

数据入湖								
	一键建合 一键建合是一种目动把在线数据同步到OSS并建立3 建立数合的目标	故焉仓库的服务, 只需少量配置即可完成数据库1	司步到OSS并 进入向导		多库合并建仓 多库合并建仓额助您把分库分表在多个R 行统一的数据分析	DS实例、数据库里面的数据聚合到数合里面的-	一个库里面来,方便忽进	进入向导
〉 更多								
任务列表	万史列表							
								刷新
ID	任务英型	任务名:	执行状态	5	期运行时间	最后更新时间	操作	
172383	一键建合	dis_it_db -> dis_it_db	SUCCESS (?)	21	021-03-29 16:59:45	2021-03-29 17:00:08	洋情 重跑 一删	*
172345	一键建仓	regdb01 -> testexec01	FAILED ⑦	2	021-03-29 15:39:57	2021-03-29 15:40:27	洋情  重跑   删	*

### 3.4.3. 一键创建OSS数据仓库

本文介绍如何通过向导将RDS、PolarDB for MySQL、MongoDB数据库或者ECS自建数据库数据同步到OSS 中,即创建OSS数据仓库(创建Schema)。

#### 操作步骤

- 1. 登录Data Lake Analytics管理控制台。
- 2. 在页面左上角,选择DLA所在地域。

- 3. 在左侧导航栏单击数据湖管理 > 数据入湖。
- 4. 在数据入湖页面,单击一键建仓中的进入向导。
- 5. 根据系统提示进行授权,授权完成后单击下一步。

DLA访问OSS、RDS、PolarDB for MySQL、MongoDB时,需要您将OSS、RDS、PolarDB for MySQL、 MongoDB的只读权限授予DLA。授权操作只需一次,后续使用一键建仓时无需再次授权。

创建Schema		
1 授权	2 配置	3 确认
角色名 DLA访问OSS授权角色 AliyunOpenAnalyticsAccessingOSSRole	已授权	
角色名 DLA访问RDS授权角色 AliyunOpenAnalyticsAccessingRDSRole	未授权 点击这里进行授权	
	下一步	

#### 6. 根据页面提示,进行参数配置。

② 说明 您可以根据实际业务数据的存储方式,选择将RDS、PolarDB for MySQL、MongoDB数 据或者ECS自建数据库数据同步到OSS。

51 mm						数据来源			
小衣				支持模糊搜索	Q	*服务器	m	com	
	类型	实例自定义名称	实例ID			• <b></b>	3306		
	MySQL	DL -1	rm-		-				
	SQLServer		rm-			* 用户名			
	MySQL	DI	rm-						
	MySQL	hi	rm-			- 8054			
	MySQL	51	rm-			* Schema名称	请输入database名称		1073524 <b>8</b>
	MySQL		rm-						
	MySQL		rm-						
	PPAS	rd	rm-			建仓配置			
	MySQL	fe	rm-			* Schema名称	请输入schema名称		
1	SQLServer	te: er	rm-						
	MySQL	m	rm- )			* 数据位置			
	MySQL		rm-						
	MariaDB		rm-uS			* 调度时间	00:30	©	
	MySQL	rr0	rr-b						
	MySQL	102	rm-		*	> 高级选项			

类别	参数	说明
	类型	数据源的类型为RDS、PolarDB for MySQL、MongoDB。 通过单击实例前的圆圈,将实例添 加到 <b>数据来源</b> 中。
云RDS、PolarDB for MySQL、	实例自定义名称	实例的备注名称。
MongoDB	实例ID	实例的ID,系统将自动拉取与DLA 同地域的实例。 支持模糊搜索实例。

类别	参数	说明	
自建数据库	ECS ID	ECS自建数据库中ECS实例的ID。 ⑦ 说明 对于ECS自建数 据库,您需要手动将反向访问 地址段100.104.0.0/16添加到 ECS白名单中。	
	VPC ID	ECS实例中的VPC ID。	
	Engine	ECS自建数据库的类型。	
	服务器	设置一键建仓中RDS、PolarDB for MySQL、MongoDB实例或者ECS自 建数据库数据源。	
	端口	RDS、PolarDB for MySQL、 MongoDB实例或者ECS自建数据库 的连接端口。	
数据来源	用户名	RDS、PolarDB for MySQL、 MongoDB实例或者ECS自建数据库 的数据库账号名。	
	密码	用户名对应的密码。	
	Schema名称	RDS、PolarDB for MySQL、 MongoDB实例或者ECS自建数据库 中的数据库名称。 配置好数据源后,单击 <b>测试连接</b> , 进行连通性测试。	
	Schema名称	设置Schema的名称,即RDS、 PolarDB for MySQL、MongoDB数 据库或者ECS自建数据库在DLA中 的映射数据库名称。	

类别	参数	说明
建仓配置		建仓时,RDS、PolarDB for MySQL、MongoDB数据库或者ECS 自建数据库数据在OSS中的详细存 储地址。
	新记行等	系统将自动拉取与DLA同地域的 OSS Bucket。根据业务需要,选择 Bucket和Object。
	<u> </u>	使用一键建仓功能时,DLA需要有 删除OSS数据的权限,以便进行从 OSS数据到RDS、PolarDB for MySQL、MongoDB数据库或者ECS 自建数据库的ETL(Extract Transform Load)操作,具体请 参见授权DLA删除OSS文件。
	调度时间	设置将RDS、PolarDB for MySQL、MongoDB数据库或者ECS 自建数据库数据同步到OSS的时间。 系统默认的数据同步时间是 00:30,您可以根据业务规律,将 数据同步时间设置在业务低峰期, 以免同步过程中可能对业务造成的 影响。
	高级选项	自定义设置项,例如过滤字段等。

#### 7. 完成上述参数配置后,单击创建,创建OSS数据仓库。

⑦ 说明 数仓创建成功后,DLA自动在您设定的同步时间将RDS、PolarDB for MySQL、MongoDB 数据库或者ECS自建数据库数据同步到OSS中,同时在OSS中创建与RDS、PolarDB for MySQL、 MongoDB或者ECS自建数据库相同的表结构,在DLA中创建对应的OSS表。

### 3.4.4. 授权DLA删除OSS文件

使用一键建仓功能时,如果您需要在DLA中进行从OSS数据到RDS、PolarDB for MySQL、MongoDB、ECS自 建数据库数据的ETL(Extract Transform Load)操作,需要为DLA授予删除OSS数据的权限。

#### 步骤一: 创建自定义授权策略

- 1. 使用DLA服务所属的阿里云账号登录RAM控制台。
- 2. 在左侧导航栏单击权限策略管理。
- 3. 在权限策略管理页面,单击创建权限策略。
- 4. 在新建自定义权限策略页面进行参数配置。
  - 配置模式:选择脚本配置。
  - 策略内容:复制以下内容,替换您的OSS Bucket名称。如果您需要更细粒度的DLA权限控制,可填写

待删除OSS文件的完整路径。

```
{
    "Version": "1",
    "Statement": [
        {
            "Action": [
               "oss:DeleteObject"
             ],
            "Resource": "acs:oss:*:*:<your-bucket-name>/*",
             "Effect": "Allow"
        }
    ]
}
```

RAM访问控制		← 新建自定义权限策略
概览		
人员管理 用户组	^	策略名称 DLADeleteOssObject
用户		备注
设置		删除OSS文件
权限管理	^	
授权		<ul> <li>可视化配置</li> <li>脚本配置</li> </ul>
权限策略管理	«	策略内容
RAM角色管理		导入已有系统策略
OAuth应用管理		1       {         2       "Version": "1",         3       "Statement": [         4       {         5                 6                 7                 8                 9                 10                 11       1

### 步骤二:将授权策略授权给DLA的角色

1. 在左侧导航栏单击RAM角色管理。

2. 在RAM角色管理页面,找到目标RAM角色名称AliyunOpenAnalyticsAccessingOSSRole。

RAM访问控制		RAM访问控制 / RAM角色管理					
概览		RAM角色管理					
人员管理 用户组	^	<b>什么是RAM角色?</b> RAM植色机制造向惊信在的实体(eg, RAM用户、某个应用或阿里云服务)进 ≪完成中下的一个RAM用户(可能是代集一个核动App的后端服务)	行授权的一种安全方法。根据不同应用场景,受信任的实体可能有如下一	些例子:			
设置权限管理	~	- 其他云账户中的RAM用户(需要进行跨账户的资源访问) - ECS条则上运行的应用程序代码(需要对达规源执行操作) - 革些阿里云服务(需要对地账户中的资源进行通信才能提供服务) - 企业的多份提供物心中,可以用于角色集合登录					
授权	«	RAM角色颁发短时有效的访问今碑(STS夺碑),使其成为一种更全全的授予访问权限的方法。 特别说知: RAM角色不同于传统的教科书式角色(其含义是指一组权限集)。如果您需要使用教科书式角色的功能,请参考RAM授权择题(Policy)。					
RAM角色管理							
OAuth应用管理		RAM角色名称	备注	创建时间	操作		
		AliyunOpenAnalyticsAccessingOSSRole	OpenAnalytics默认使用此角色来访问OSS	2018年12月4日 15:07:10	添加权限删除		
		添加收码					

3. 单击RAM角色名称链接,进入角色详情页面,然后单击添加权限。

RAM 访问控制 / RAM角色管理 / A	RAM 说明的 / RAM 编色管理 / AlyesCpenturbyticsAccessingDDSStole									
← AliyunOpenA	← AliyunOpenAnalyticsAccessingDBSRole									
基本信息										
RAM 角色名称 Aliyu	nOpenAnalyticsAccessingDBSRole		创建时间	2021年3月19日16:12:44						
脅注 Oper	Analytics默认使用此角色来访问DBS 编辑		ARN	acs:ram::104157	role 邙 复制					
最大会话时间 3600	わ 調識									
<b>权限管理</b> 信任策略管理										
漆加权限 精确授权						0				
权限应用范围	权限策略名称	权限策略类型	备注		授权时间	摸作				
整个云账号	AliyunOpenAnalyticsAccessingDBSRolePolicy	系统策略	用于开放分析服务的授权策略,包括D	BS的部分访问权限	2021年3月19日16:12:45	移除权限				

4. 在**添加权限**页面,选择**自定义策略**,单击<mark>步骤一:创建自定义授权策略</mark>中的权限策略,将其授权给DLA 的角色AliyunOpenAnalyticsAccessingOSSRole,单击**确定**完成授权。

时和风				
系统策略 自定义策略 -	+ 新建权限策略		已选择 (1)	ň
请输入权限策略名称进行模糊搜索		B	k8sWorkerRolePolicy-f76	e9206c-9264
权限策略名称	备注			
k8sWorkerRolePolicy-f7e920				
dla-delete-luhao				
dla-delete-yunfu				
DlaFsTestPolicy				
dlfs-ut	dlfs-ut			
k8s_c71f7e7a06f7e4f95a2c10	测试			
OssDelete-minghui				
DLADeleteOssObject				
dla_fengshen_del_create				
k8sMasterRolePolicy-134538				

### 3.5. T+1多库合并建仓

本文档主要介绍了多库合并建仓的操作步骤。

#### 背景信息

在数据库应用中,当单个关系型数据库RDS(Relational Database Service)的数据量越来越多时,相应的数据查询时间也会延长,影响用户体验。为保证业务可以继续使用RDS数据库,业务端通常会采用分库分表技术,将一个RDS数据库中的单张表数据拆分到多个数据库的多张表中。上述方案可解决因数据量大而导致的用户体验问题,但在对分库分表数据进行大数据分析时,逻辑上的一个表被拆成了多张表,由于没有类似TDDL中间件来屏蔽物理表的拆分,进行数据分析时变得十分复杂。

#### 解决方案

T+1多库合并建仓是指通过DLA控制台上的多库合并建仓向导将RDS中的分库分表数据聚合到统一的表中,并 以分区表形式存储数据。您可以全局分析所有数据,也可以选择某个分区对分区数据进行分析,进行数据分 析时会非常的方便,并且也不影响RDS端的业务运行。

#### 前提条件

- 多库合并建仓时, DLA将OSS作为存储RDS数据的数据仓库, 您需要在OSS中完成以下准备工作:
  - i. 开通OSS服务,请参见开通OSS服务。
  - ii. 创建存储空间,请参见创建存储空间。
  - iii. 新建目录, 请参见新建目录。

```
? 说明
```

根据业务需求,判断是否需要新建目录存储RDS数据。

• 根据您的业务需要,准备好RDS数据源。具体请参见快速入门。

#### 操作步骤

- 1. 登录Data Lake Analytics管理控制台。
- 2. 在页面左上角,选择DLA所在地域。
- 3. 单击左侧导航栏的数据湖管理>数据入湖,在数据入湖页面单击多库合并建仓中的进入向导。
- 4. DLA首次访问RDS时,需要您将RDS的只读权限授予DLA,授权完成后单击下一步。

如果您之前已经将RDS的只读权限授予DLA,可以忽略该步骤。

<sup>?</sup> 说明

┃ 创建Schema		
1 授权	2 配置	3 确认
角色名 DLA访问OSS授权角色 AliyunOpenAnalyticsAccessingOSSRole	已授权	
角色名 DLA访问RDS授权角色 AliyunOpenAnalyticsAccessingRDSRole	未授权 点击这里进行授权	
	世—不	

#### 5. 根据页面提示,进行参数配置。

类别	参数	说明
手动选择	类型	数据源的类型为RDS。 您可以勾选实例,将RDS实例添 加到数据源中。
<ul> <li>说明</li> <li>通过手动方式指定</li> <li>RDS实例,该方式适</li> <li>用于RDS实例个数不</li> <li>多且实例个数处于</li> <li>静态或者不会频繁</li> <li>动态增加的场景。</li> </ul>	数据库筛选规则	输入您要同步的数据库名字。多个数据库名字之间用 英文逗号(,)分隔。数据库名支持使用通配符%,例 如user_%。
通过查询 ⑦ 说明 指定通过SQL查询方 式指定RDS数据源, 该方式适用于RDS实 例个数较多且实例 个数动态增加的场 景。	-	例如, SELECT 'mysql' AS engine, 'db001' AS db_name, 'rm-111aliyuncs.com' AS host, 3306 AS port, 'rm-123445' AS instance_id, 'vpc-3424555' AS vpc_id FROM tb11
	用户名	为使用方便,DLA要求您选择的所有数据库均使用统一 的用户名和密码。
认证信息	密码	上述用户名对应的密码。输入用户名和密码后,您可 以单击测试连接,进行连通性测试。

类别	参数	说明
	Schema名称	设置Schema的名称,即RDS数据库在DLA中的映射数 据库名称。
	数据位置	建仓时,RDS数据在OSS中的详细存储地址。系统将自动拉取与DLA同地域的OSS Bucket,单击选择位置, 您可以根据业务需要,灵活选取Bucket和Object。使 用多库合并建仓功能时,DLA需要有删除OSS数据的权限,以便进行从OSS数据到RDS数据的ETL(Extract Transform Load)操作,请参见授权DLA删除OSS文件。
	同步时间	设置将RDS数据同步至OSS的时间。系统默认的数据同 步时间是00:30,您可以根据业务规律,将数据同步时 间设置在业务低峰期,以免同步过程中可能对业务造 成的影响。
建仓配置	表名生成规则	<ul> <li>设置DLA建仓时,RDS表在数仓中的映射表名。映射表名将通过以下两种规则自动生成:</li> <li>IdentityResolver:数仓中的表名与RDS表名相同,适用于RDS中有分库但没有分表的场景。</li> <li>RemoveTrailingUnderscoreAndNumberResolver:将RDS表名中最后一次出现的下划线和数字去掉,作为数仓中的表名。例如,RDS表名为tbl_001,则数仓表名为tbl。</li> </ul>
	分区配置	设置数仓的分区字段以及分区字段值的生成方式。分 区字段值为一个包含变量的表达式,例如 \${rdsInstanceId} 。DLA暂时支持以下变量: • rdsEngine: RDS支持的引擎类型,包含MySQL、 SQLServer、PostgreSQL、Oracle。 • rdsDbName: RDS数据库的名字。 • rdsTableName: RDS表的名字。 • rdsInstanceId: RDS实例ID。 • rdsVpcId: RDS实例所属VPC ID。
	高级配置	自定义设置项,例如过滤字段等。

		✓ ────────────────────────────────────			2 配置			3 确认	
选择数据库实例						建仓配置			
手动选择	通过查	间指定				Schema名称	test_adb_schema		
RDS列表 华东1	1(杭州)			支持模糊搜索 RDS	Q	数据位置	oss:// 'test_adb_file/		选择位置
☑ 类型		RDS名称	实例ID			同步时间	00:30	3	
		rm- bp	rm-bp	-		表名生成规则	IdentityResolver	Ň	· ()
数据库筛选规则	test_adb			0		分区配置	添加 您RDS表里面的数据会对应到数仓 个分区的信息,包括分区的名字以 绍。	表里面的一个分区 , 这里请告诉我 及取值方式。点击这里展开更多乡	划们关于如果构建这 关于分区配置的介
用户名	1000						分区名	分区表达式	操作
密码				测试连接				没有数据	
						> 高级选项			

6. 完成上述参数配置后,单击创建,创建数据仓库。

数据仓库创建成功后,DLA自动在您设定的同步时间将RDS数据同步到OSS中,同时在OSS中创建与RDS相同的表结构、在DLA中创建对应的OSS表。

## 3.6. ActionTrail日志清洗

DLA提供ActionTrail日志自动清洗解决方案,可以将ActionTrail投递到OSS的日志文件转换为DLA中可以直接 查询的数据表,同时自动对数据进行分区和压缩,方便您分析和审计对云产品的操作日志。

#### 日志分析痛点

ActionTrail是阿里云提供的云账号资源操作记录的查询和投递服务,可用于安全分析、资源变更追踪以及合规性审计等场景。您可以通过ActionTrail控制台查看各个云产品的操作日志。对于30天以内的日志,ActionTrail支持投递到日志服务进行分析;对于30天以外的数据可以投递到OSS上,但直接分析OSS中的数据有以下痛点。

• 日志数据格式复杂,不利于直接分析

ActionTrail中保存的是JSON格式的数据,一行内有多条数据,数据以一个Array的形式保存,例如 [{"eventId":"event0"...},{"eventId":"event1"...}] 。

理论上可以分析上述格式的JSON数据,但非常不便,需要先把每行数据拆分成多条记录,然后再对拆分后 的记录进行分析。

• 小文件多,分析数据耗时且占用大量系统资源

当您通过账号(阿里云账号和RAM子账号)频繁操作云产品时,每天产生的操作日志文件数非常多。以操作DLA的账号为例,该账号下每天会产生几千个数据文件,一个月的文件数将达到几十万个,大量的数据 文件对大数据分析非常不便,分析数据耗时,且需要足够大的集群资源才能进行大数据分析。

#### 前提条件

使用ActionTrail日志清洗之前,您需要按照以下步骤做好准备工作。

② 说明 使用ActionTrail日志清洗功能时,要求ActionTrail、OSS、DLA所属Region相同,否则无法使用该功能。

● 在ActionTrail中完成以下操作:

在ActionTrail中创建跟踪,请参见创建单账号跟踪或创建多账号跟踪。

- 在OSS中完成以下操作:
  - 开通OSS服务,请参见开通OSS服务。
  - 创建Bucket,请参见创建存储空间。
  - 新建文件夹,请参见创建目录。

⑦ 说明 根据业务需求,判断是否需要新建文件夹,将ActionTrail投递过来的数据存储在新建文件夹中。

- 在DLA中完成以下操作:
  - 开通DLA服务,请参见开通云原生数据湖分析服务。
  - 初始化DLA数据库主账号密码,请参见重置数据库账号密码。

#### 步骤一: 创建Schema

- 1. 登录Data Lake Analytics管理控制台。
- 2. 在页面左上角,选择DLA所在地域。
- 3. 在左侧导航栏单击数据湖管理 > 数据入湖。
- 4. 在数据入湖页面单击ActionTrail日志清洗中的进入向导。
- 5. 在ActionTrail日志清洗页面,根据页面提示进行参数配置。

* ActionTrail文件根目录	oss://e//Aliyur	nLogs/Actiontrail/	Ĵ	选择位置	自动发现		
* Schema名称:	ActionTrail_test_schem	la	0				
* 清洗后数据保存位置	oss://dlaossfile1/ActionTrail_test_schema			自定义 🗸			
* 数据清洗时间	00:30		© ()				
> 高级选项							
		创建					
参数名称		参数描述					
ActionTrail文件根目录		ActionTrail投递到OSS中日志数据的存储目录。目录以 AliyunLog s/Actiontrail/ 结尾。					
		<ul> <li>选择位置:自定义ActionTrail投递到OSS中的日志数据的存储目录。</li> </ul>					
		<ul> <li>● 自动发现: DLA自动设置ActionTrail投递到OSS中的日志数据的存储目录。</li> </ul>					
Schema名称		设置Schema的名称,即OSS在DLA中的映射数据库名称。					

参数名称	参数描述
清洗后数据保存位置	DLA清洗OSS数据后,将结果数据回写入OSS即数据清洗后的存储位 置。DLA会默认指定存储位置。您也可以自定义存储位置。
数据清洗时间	设置每天DLA清洗OSS数据的时间。系统默认的数据清洗时间是 00:30,您可以根据业务规律,将数据清洗时间设置在业务低峰期, 以免清洗过程中可能对业务造成的影响。

#### 6. 完成上述参数配置后单击创建,创建Schema。

Schema创建成功后, DLA自动在您设定的同步时间将ActionTrail投递到OSS中的日志数据同步到DLA中,并在DLA中创建OSS日志文件对应的表。

您也可以在数据入湖页面的任务列表页签中,在操作列单击执行随时手动发起数据同步操作,将 ActionTrail投递到OSS中的日志数据同步到DLA中,并在DLA中创建OSS日志文件对应的表。

任务列表	历史列表									
全部快型										
ID	名称	Schema名称	後型	调度计划	调度状态	最近运行状态	最近运行耗时	创建时间	操作	
424	dis_it_db b	dit	一键建仓	每天00时30分	开启调度	SUCCESS ⑦	22秒	2021-03-19 14:13:26	19407 MHMAA 历史 創除	
423	admin -> tana bab	testc d	一键建合	每天02时00分	停止调度			2021-03-17 19:06:16	执行   编辑   历史   删除	
422	test9928. er	werr	一键建仓	每天00时30分	停止调度	• FAILED 🕐 🗐	333秒	2021-03-17 18:26:48	执行   编辑   历史   删除	
81	Action 1	a	ActionTrail日志清洗	每天00时30分	开启调度	• SUCCESS ⑦	0#9	2021-03-01 14:23:07	执行   编辑   历史   删除	
369	regdb01 - 01	t 1	一键建仓	每天11时30分	停止调度	• FAILED (?) 🗐	99Æ)	2021-01-18 19:59:47	执行   编辑   历史   删除	

数据同步到DLA以后,您就可以在DLA中使用标准SQL语法对ActionTrail日志数据进行分析。