

ALIBABA CLOUD

阿里云

数据湖分析
数据湖管理

文档版本：20201020

 阿里云

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。

1. 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。
2. 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。
3. 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。	 注意 权重设置为0，该服务器不会再接受新请求。
 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。	 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。	在结果确认页面，单击确定。
<code>Courier</code> 字体	命令或代码。	执行 <code>cd /d C:/window</code> 命令，进入Windows系统文件夹。
<i>斜体</i>	表示参数、变量。	<code>bae log list --instanceid</code> <i>Instance_ID</i>
[] 或者 [a b]	表示可选项，至多选择一个。	<code>ipconfig [-all -t]</code>
{ } 或者 {a b}	表示必选项，至多选择一个。	<code>switch {active stand}</code>

目录

1.元数据爬取	05
2.数据入湖	07
2.1. 多库合并建仓	07
2.2. ActionTrail日志清洗	09
3.实时数据湖	12

1.元数据爬取


本文介绍如何通过向导创建元数据爬取任务，爬取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。

操作步骤

1. 进入[元数据爬取](#)页面。
2. 在元数据爬取页面，单击进入向导按钮。
3. 在创建元数据爬取页面左侧，指定要爬取的OSS的数据目录。
4. 在创建元数据爬取页面右侧，根据页面提示进行参数配置，配置说明如下：

参数	说明
格式解析器	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。
爬取频率	您可以根据需要定期计划运行元数据爬取任务。
Schema名称	设置Schema名称，即映射到DLA中的数据库名称（默认每个爬取任务会新创建一个独立的Schema）。
配置选项	高级自定义设置项，如更新，删除规则等。

5. 完成上述参数配置后，单击创建，开始创建元数据爬取任务。

 **说明** 元数据爬取任务创建完成后，DLA自动在您设定的时间周期运行爬取任务，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。



注意事项

● 元数据爬取会如何生成表名

元数据爬取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。爬取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，则元数据爬取会在表名后添加MD5字符串后缀。

● 元数据爬取会如何创建分区

当元数据爬取扫描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，爬取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，所以爬取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，所以爬取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，爬取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。

- 元数据爬取目前支持的格式解析器

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。

DLA 为以下类型数据文件提供了一组内置格式解析器，如果用户没有指定特定的格式解析器，元数据爬取会按照以下顺序调用内置解析器。

解析器类型	备注
JSON	需要读取文件开头以确定文件格式。
Parquet	需要读取文件结尾处的 schema 以确定文件格式。
CSV	检查以下分隔符：逗号(,)、竖线()、制表符(\t)、分号(;)、空格()、(\u0001)。
ORC	需要读取文件元数据以确定文件格式。
AVRO	需要读取文件开头处的 schema 以确定文件格式。

2. 数据入湖

2.1. 多库合并建仓

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，然后修改业务层代码，并配合使用类似TDDL的分库分表中间件。

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

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，不影响RDS端的业务运行。

前提条件

使用多库合并建仓前，您需要完成以下准备工作：

- OSS

背景信息

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Bucket，请参见[创建Bucket](#)。
- iii. 新建文件夹，请参见[新建文件夹](#)。

 说明

根据业务需求，判断是否需要新建文件夹存储RDS数据。

- RDS

背景信息

根据您的业务需要，参照[RDS for MySQL快速入门](#)、[RDS for SQL Server快速入门](#)、[RDS for PostgreSQL快速入门](#)或者[MySQL for PPAS快速入门](#)准备好RDS数据源。

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构建 > 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。

说明
如果您之前已经将RDS的只读权限授予DLA，可以忽略该步骤。



5. 根据页面提示，进行参数配置。

类别	参数	说明
手动选择，通过手动方式指定RDS实例，该方式适用于RDS实例个数不多且实例个数处于静态或者不会频繁动态增加的场景。	类型	数据源的类型为RDS。通过单击实例前的方框，将RDS实例添加到数据源中。
	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，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-111..aliyuncs.com' AS host, 3306 AS port, 'rm-123445' AS instance_id, 'vpc-3424555' AS vpc_id FROM tbl1
认证信息	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。
	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。
	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	<p>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：</p> <p>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。</p> <p>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。</p>
	分区配置	<p>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 <code>\${rdsInstanceId}</code>。DLA暂时支持以下变量：<code>rdsEngine</code>，RDS支持的引擎类型，包含MySQL、SQLServer、PostgreSQL、Oracle。<code>rdsDbName</code>，RDS数据库的名字。<code>rdsTableName</code>，RDS表的名称。<code>rdsInstanceId</code>，RDS实例ID。<code>rdsVpcId</code>，RDS实例所属VPC ID。建议您同时填写RDS实例ID以及数据库名，例如：分区名 <code>rds_instance_id</code>，分区值 <code>\${rdsInstanceId}</code>。分区名 <code>rds_db_name</code>，分区值 <code>\${rdsDbName}</code>。</p>
	高级配置	<p>自定义设置项，例如过滤字段等，请参见高级选线功能。</p>

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。

2.2. 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。

日志分析痛点

操作审计 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SLS](#) 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的帐号为例，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，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。

🔍 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。

- ActionTrail

在ActionTrail中创建跟踪，请参见[创建跟踪](#)

- OSS

- 开通OSS服务，请参见[开通OSS服务](#)
- 创建Bucket，请参见[创建Bucket](#)
- 新建文件夹，请参见[新建文件夹](#)

🔍 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。

- DLA

- 开通DLA服务，请参见 [开通DLA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. 单击左侧导航栏的数据湖构建 > 数据入湖，在数据入湖页面单击ActionTrail日志清洗中的进入向导。
4. 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	<p>ActionTrail投递到OSS中日志数据的存储目录。</p> <p>目录以AliyunLogs/Actiontrail/结尾。</p> <ul style="list-style-type: none"> ○ 选择位置：自定义ActionTrail投递到OSS中的日志数据的存储目录。 ○ 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。
------------------	---

Schema名称	设置Schema的名称，即OSS在DLA中的映射数据库名称。
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。 <ul style="list-style-type: none"> DLA默认指定存储位置。 支持您自定义存储位置。
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。

5. 完成上述参数配置后单击**创建**，创建Schema。

Schema创建成功后，ActionTrail投递到OSS中的日志数据并未同步到DLA中，即DLA中没有创建OSS日志文件对应的表，您需要通过单击**立即同步**来创建表同步表数据。

步骤二：同步数据

创建Schema后，单击**立即同步**同步数据，也可以在任何需要的时候通过以下步骤同步数据。

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**数据湖构建 > 元数据管理**。
4. 在**元数据管理**页面，单击目标数据源右侧的**详细信息**。
5. 在**元数据管理**页面，单击**配置**页签。
 -
6. 单击**立即同步**启动数据同步任务。
在**配置**页签下，单击**更新更新Schema配置**。
7. 单击**表**页签，查看数据同步情况。
 -

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。

3. 实时数据湖