



数据湖分析 DMS

文档版本: 20211227



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 <i>,</i> 是用户必须 了解的内容。	大意 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。	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 结果确认 页面,单击 确定 。
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 {a b}	表示必选项,至多选择一个。	switch {active stand}

目录

1.调度DLA	Presto任务	 05
2.调度DLA	Spark任务	 10

1.调度DLA Presto任务

您可以在数据管理DMS中创建和调度DLA任务流实现数仓开发,任务流只需定义一次,即可周期性地自动被 调度执行,大大减少人工操作成本。同时执行DLA任务流时产生的中间结果可以被复用,例如用于其他数据 分析任务,从而最大化利用DLA的资源。

任务调度中一个重要的功能是任务之间的依赖,为演示该功能,本文在DMS中创建两个DLA任务,表、任务 之间的关系如下图所示。



● 任务一

从orders表查询出已经完成的订单: o_orderstatus = 'F' ,并将其写入finished_orders表。

● 任务二

从finished_orders中查询出总价大于10000的订单: o_totalprice > 10000 ,并将其写入 high value finished orders表。

前提条件

- 您需要开通DLA服务,然后创建一个Schema (database)以及在Schema中创建测试表。
 - i. 开通DLA服务,请参见<mark>开通云原生数据湖分析服务</mark>。
 - ii. 初始化数据库主账号密码,请参见重置数据库账号密码。
 - iii. 创建Schema和表。

创建Schema

```
CREATE SCHEMA gmall with DBPROPERTIES(
 LOCATION = 'oss://oss-bucket-name/gmall/',
 catalog='oss'
);
```

创建orders表

```
CREATE EXTERNAL TABLE orders (

O_ORDERKEY INT,

O_CUSTKEY INT,

O_ORDERSTATUS STRING,

O_TOTALPRICE DOUBLE,

O_ORDERDATE DATE,

O_ORDERPRIORITY STRING,

O_CLERK STRING,

O_SHIPPRIORITY INT,

O_COMMENT STRING

)

ROW FORMAT DELIMITED FIELDS TERMINATED BY '|'

STORED AS TEXTFILE

LOCATION 'oss://bucket-name/gmall/orders';
```

finished_orders表

high_value_finished_orders表

● 在DMS产品详情页单击**立即购买**,开通DMS服务。

步骤一:新增DLA数据源

- 1. 登录数据管理DMS控制台。
- 2. 在左侧导航栏选择实例管理,单击实例列表。
- 3. 单击目标实例地址所在行更多中的编辑实例。
- 4. 在编辑实例页面,根据页面提示进行参数配置。

编辑实例						×
∨ 基本信息						
* 数据来源	云数据库 ECS自建数据库	本地/他云数据库	公网数据库	VPC专线IDC	文件存储	
* 数据库类型	ADB3.0-MySQL					\sim
* 实例地区	华东1 (杭州)			\sim	跨阿里云账号3	定例
* 录入方式 (● 实例ID ○ 连接串地址					
* 实例ID	am-		-			\sim
数据库账号	请输入 数据库账号					
数据库密码						
* 管控模式	🧧 自由操作 🏼 🧿 稳	定变更 🔤 安全	协同点此了	解		
* 敏感数据保护 (关闭。点此了解					
> 高级信息 (可查看环境类	型、名称、DBA及更多高级特	性)				
测试连接						提交
参数类别	参数名称		47	参数说明		

参数类别	参数名称	参数说明
	数据来源	选择 云数据库 。
	数据库类型	选择DLA-数据湖分析。
	实例地区	选择DLA所在地域。
	录入方式	您可以通过两种方式录入DLA信 息: • 实例ID ,从实例列表中选择DLA 实例ID。 • 连接串地址 ,手动输入DLA连 接地址。
基本信息	实例ID	 录入方式为实例ID时,从实例列表中选择DLA实例ID。 录入方式为连接串地址时,手动输入DLA连接地址。
	数据库账号	DLA中的数据库账号。

参数类别	参数名称	参数说明
	数据库密码	数据库账号对应的密码。
	管控模式	请参见管控模式。
	安全规则	如果您选择的管控模式为 安全协 同,需要填写 安全规则 ,选择DLA 默认规则。
	环境类型	本例选择 生产 。
	实例名称	DLA的实例ID或者名称。
	开启跨库查询	默认选中 开启跨库查询 ,且不支 持用户修改。
	无锁结构变更	默认选中 关闭 ,且不支持用户修 改。
高级信息	开启SSL	您可以选择 默认(DMS与数据库 服务自动协商是否开启)、开 启或关闭。
	实例DBA	您可以根据实际情况选择一个账号 作为 实例DBA 。
	查询超时时间(s)	本例设置为60。
	导出超时时间(s)	本例设置为600。

5. 完成参数配置后单击测试连接,测试成功后单击提交。

步骤二:编排任务流

- 1. 登录数据管理DMS控制台。
- 2. 在上方导航栏,单击传输与加工(DTS)。
- 3. 单击左侧导航栏中的任务编排,然后单击新增任务流。
- 4. 在新建任务流对话框,设置任务流名称和描述,然后单击确认。

新建任务流			×
*任务流名称	dla_demo	٥	
描述	dla任务流		
			确认 取消

5. 拖拽两个DLA节点到中间的空白区域中, 然后单击任务1的下边框连接点, 引出一条到达任务2的带箭头 直线关联任务1和任务2。

- 6. 分别单击任务1和任务2节点,在内容设置中输入或者选择DLA Schema的名字,然后单击保存。
- 7. 单击任务1节点,在内容设置中输入SQL实现任务一,即从orders表查询出已经完成的订单 o_ordersta tus = 'F',并将其写入finished_orders表,单击保存。

```
insert into finished_orders
select O_ORDERKEY, O_TOTALPRICE
from orders
where O ORDERSTATUS = 'F';
```

8. 单击任务2节点,在内容设置中输入SQL实现任务二,即从finished_orders中查询出总价大于10000的订 单 ototalprice > 10000,并将其写入high_value_finished_orders表,单击保存。

```
insert into high_value_finished_orders
select * from finished_orders
where O TOTALPRICE > 10000;
```

步骤三:运行任务

任务流编排完成后,您可以通过如下方式运行任务。

- 单击**试运行**,系统弹出提示窗口,单击**确认**运行任务。
- 开启**调度配置**为任务设置周期运行策略。

步骤四: 查看任务

任务运行结束后,您可以通过如下方式查看任务。

- 单击运维中心查看任务运行状态。
- 登录Data Lake Analytics管理控制台,在DLA中查询任务运行结果。
 - 任务一: 从orders表查询出已经完成的订单 o_orderstatus = 'F'。
 - 任务二:从finished_orders中查询出总价大于10000的订单 o_totalprice > 10000 。

2.调度DLA Spark任务

DLA Serverless Spark目前支持DataWorks和DMS编排调度任务,同时也提供自定义SDK和Spark-Submit工具 包供用户自定义编排调度。本文将介绍如何使用DMS来编排调度Spark任务。

前提条件

- 您已开通DMS服务。
- 您已开通云原生数据湖分析服务,详情请参见开通云原生数据湖分析服务
- 您已开通OSS服务,详情请参见开通OSS服务。
- 您已创建DLA Spark虚拟集群,详情请参见创建虚拟集群。

⑦ 说明 如果您想用RAM用户提交Spark作业,且之前未使用过RAM用户在DLA控制台提交作业,您可以参见细粒度配置RAM子账号权限进行子账号提交作业配置。

操作步骤

任务调度中一个重要的功能是任务之间的依赖,为演示这个功能,本文会在DMS中创建一个DLA Serverless Spark任务和两个 DLA Spark SQL任务,任务之间的依赖关系如下图所示,任务DLA Spark SQL-1和DLA Spark SQL-2 依赖上游任务完成之后,才能执行。



- 1. 登录DMS控制台。
- 2. 在顶部菜单栏中, 单击数据工厂 > 任务编排。
- 3. 在任务编排页面的自由编排任务区域,单击新增任务流。



- 4. 在新建任务流对话框,将任务流名称设置为spark_demo,将描述设置为spark demo,完成后单击确认。
- 5. 在**任务编排**页面,从左侧任务类型中拖拽一个DLA Serverless Spark任务和两个DLA Spark SQL任务,并按照下图中的依赖关系进行连线。



- 6. 依次单击3个节点,在每个节点的右侧面板中选中内容设置页签,并配置以下信息:
 - 在地域列表中,选择目标Spark集群所在的地域。
 - 在Spark 集群列表中,选择目标Spark集群。
 - 在作业配置文本框中已有运行内置Sparkpi的配置参数或者 show databases ,如果您需要运行自定 义作业,可根据您的需求进行修改。

F线 自动布局 保存位置 添加任务 V 添加便签	内容设置	变量设置	高级设置
发布状态:未发布 调度状态:尚未设置周期调度,请在右侧设置	* 地域		提示 🚯
	华东1(杭州)		~
© DLA Serverless Spark	* Spark 集群 jobtest		◇ 去创建
DLA Spark SQL-2	<pre>*作业配置</pre>		格式化 mp/spark- pache.spa ourceSpec nstances' esourceSp
		保存	

完成以上配置后,单击保存按钮,然后单击页面左上方的试运行或指定时间运行或指定时间范围运行任务。

#519 ×	X) 50L janghu X) Efisii K	
+ -l: spark-test X		帮助手册~
· 國語行 へ 文布	下档 自动有两 保存位置 激加性务 ✓ 激励性务	内容设置 交量设置 高级设置
指定时间运行	2.1919년: 朱武帝 國政党本: 長米近夏期期間度, 请在在街设置	*地域 提示 0
指定时间范围运行		年末1 (統州) ~
> ROBERTSHI		· Crack With
>归杨建仓		release-test 公式创建
> 状态检查		
> 备份恢复		*作业配置 格式化
(164	C DLA Span SQ.1 Z	<pre>* """"""""""""""""""""""""""""""""""""</pre>
	新日本 (〒朝秋7) 10.0 Sgark 502 (「11907) 10.0 Sgark 502 ()	1
	8 m/18/3/ [TH18/m/1] (DA Book SOL-1	
	RATES	88

8. 您也可以在点击任务编排界面的空白界面,进行整个任务的调度设置。

T	× 502_5mm/u × 任务编号 ×				
			_		帮助手册~
]	111 日山市県 保存位置 添加任务 V 添加研究	调度配置	基础属性	操作历史	任务流变量
	发布形态: 未定方 调度记念: 未完成 建筑物构成, 确立在 新公理	开启调度 の1 測度类型			
		周期调查	t.	~	
		• 住 6934日			
		1970-0	1-01		
		9999-0	1-01	63	
		注: 调度将 效期外的任	在有效日期内生 务将不会自动词	放并自动调度。 1度。	反之,在有
		• 调度周期			
		в		\sim	
		• 具体时间			
		00:00		O	
	0	oron表达式			
	S DLA Serverses Spark Z	00 00 00 "	?		
		登 看运行5	چ -		
	N/F18.k UMBN/F1_CLASpank SGL-2				
	RFTRR/D TPRMITI Du A Reux SGL-1				
	94735		¢	φ.	
est	Intellatil nov derivere obtain				

查看任务日志

1. 点击任务编排界面, 左侧的运维中心。



2. 单击对应任务左侧的加号按钮,即可看到任务流的所有任务,点击子任务的查看按钮,即可查看运行日志。

83	फ्रेन्द्र स	1889 × 188948 (1889 ×	任务说 spark-test	~								Fill	C RM
		48t	秋志	能发方式	責任人	业务时间	开始时间			结束时间		操作	
-	-	spark-test(1365055)	 atso 	∓i280.30	open, analytica_test	2021-06-03 15:16:40	2021-06-04 15:16:	59		2021-06-04 15:19:35		DAG图 执行历史 更多	
	_	任务名 2					88.2	任务类型	Hittens	e Antonia	结果时间 4	道行日本	
		DLA Serverless Spark-1					9 #33	DLA Serverless Spark	82725	2021-05-04 15:18:39	2021-06-04 15:18:03	2021-06-04 15:16:39 CST INFO - 开始运行任务[D Serveriess Spark-1],任务地动1(28,158),任务1	
		DLA Spark SQL-1					● 成功	DLA Spark SQL	62203	2021-06-04 15:18:02	2021-06-04 15:19:05	2021-06-04 15:18:02 CST INFO - 开始运行任务(D Bperk SQL-1).任务版本id(128,159).任务版例	u și
		DLA Spark SQL-2					्र इस्क्र	DLA Sperk SQL	92906	2021-06-04 15:18:02	2021-06-04 15:19:30	2021-06-04 15:18:02 CST INFO - 开始运行任务[D Spark SQL-2],任务政策间128,160],任务实例	u 👬
		spark-test(1365026)	0 先敗	于动脉发	open_analytics_test	2021-06-03 15:14:10	2021-06-04 15:14:	10		2021-06-04 15:15:12		DAGH MITTE EF	
	+	spark-test(1365024)	0 朱欣	干动脉发	open_analytics_test	2021-06-03 15:13:15	2021-06-04 15:13:	15		2021-08-04 15:12:15		DAG图 我行历史 更多	
												< 上→茂	1 F-R >

3. 从日志中可以查看该作业的JobId和对应的SparkUI,如果任务出错,无法排查,请记下JobId和SparkUI 联系Spark值班。

自定义任务编排调度

DLA Serverless Spark除了上述调度集成之外,还提供了SDK以及Spark-Submit工具用于提交Spark作业、查询作业状态、获取作业日志等功能,详情请参见:

- SDK安装与使用
- Spark-Submit命令行工具

您可以利用上述工具,使用第三方任务编排调度系统(例如Apache Airflow)来打造自己的工作流。