



E-MapReduce 数据开发

文档版本: 20210916



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 {act ive st and} |

目录

| 1.概述 | 05 |
|---------------------------|----|
| 2.项目管理 | 06 |
| 3.作业编辑 | 09 |
| 4.工作流编辑 | 14 |
| 5.临时查询 | 18 |
| 6.运维中心 | 20 |
| 7.创建集群模板 | 22 |
| 8.云监控事件编码 | 23 |
| 9.作业配置 | 24 |
| 9.1. 作业日期设置 | 24 |
| 9.2. Shell作业配置 | 25 |
| 9.3. Hive作业配置 | 26 |
| 9.4. Hive SQL作业配置 | 26 |
| 9.5. Spark作业配置 | 27 |
| 9.6. Spark SQL作业配置 | 28 |
| 9.7. Spark Shell作业配置 | 29 |
| 9.8. Spark Streaming作业配置 | 30 |
| 9.9. Hadoop MapReduce作业配置 | 30 |
| 9.10. Sqoop作业配置 | 31 |
| 9.11. Pig作业配置 | 32 |
| 9.12. Flink(VVR)作业配置 | 33 |
| 9.13. Streaming SQL作业配置 | 34 |
| 9.14. Presto SQL作业配置 | 35 |
| 9.15. Impala SQL作业配置 | 36 |
| 10.数据开发常见问题 | 38 |

1.概述

创建E-MapReduce集群后,您可以在E-MapReduce数据开发中创建项目。数据开发是可编程、调度和监控的工作流平台,基于有向无环图(DAG),可以定义一组有依赖的 作业,按照依赖依次执行。您可以通过控制台管控作业调度任务,并监控各类作业的运行状态,以便管理和运维工作流。

↓ 注意 如果您的EMR高安全集群对接的是外部的MIT Kerberos,则无法使用数据开发功能。

E-MapReduce数据开发的具体功能包括:

- 数据开发项目管理:为项目关联集群资源和添加项目成员,详情请参见项目管理。
- 大数据作业开发和编辑:支持Hive、HiveSQL、MapReduce、Spark和Shell等作业类型的开发,详情请参见作业编辑。
- 工作流开发和调度:通过拖拽构建工作流,并设置时间调度策略和工作流之间依赖,详情请参见工作流编辑。
- 临时查询:支持HiveSQL、SparkSQL、Spark和Shell四种类型的临时查询作业,详情请参见临时查询。
- 运行记录:查看任务、工作流的运行记录和日志,可以重新运行失败的作业和工作流,并可以查看项目成员在项目中的操作历史,详情请参见运维中心。

2.项目管理

创建E-MapReduce集群后,您可以在数据开发中创建项目,并在项目中进行作业的编辑和工作流的调度。新建项目之后,您可以对项目进行管理,为项目关联集群资源、添加 项目成员以及设置全局变量。

前提条件

已创建集群,详情请参见创建集群。

使用限制

只有阿里云账号才能创建项目、添加项目成员和添加集群资源,即控制台上**新建项目、用户管理和集群设置**功能只对阿里云账号管理员可见,RAM用户不可见。

新建项目

- 1. 进入数据开发的项目列表页面。
 - i. 通过阿里云账号登录阿里云E-MapReduce控制台。
 - ii. 在顶部菜单栏处,根据实际情况选择地域和资源组。
 - iii. 单击上方的**数据开发**页签。

阿里云账号可以查看该账号下的所有项目列表,RAM用户仅可以查看具有开发权限的项目列表。如需为RAM用户添加项目开发权限,则需要阿里云账号来配置,详情 请参见<mark>用户管理</mark>。

- 2. 在**项目列表**页面,单击右上角的新建项目。
- 3. 在新建项目对话框中,输入项目名称和项目描述,从资源组选择列表中,选择已有的资源组。

⑦ 说明 如果不指定资源组,项目会加入默认的资源组,资源组使用详情请参见使用资源组。

4. 单击**创建**。

在**项目列表**页面,可查看或者操作新增的项目。

查看项目基本信息

- 1. 进入数据开发的项目列表页面。
 - i. 通过阿里云账号登录阿里云E-MapReduce控制台。
 - ii. 在顶部菜单栏处,根据实际情况选择地域和资源组。
 - iii. 单击上方的**数据开发**页签。
- 2. 进入目标项目的项目管理页面。
 - i. 在**项目列表**页面,单击目标项目的项目ID。
 - ii. 单击上方的**项目管理**页签。
- 3. 查看项目基本信息。

在基本信息页面,可以查看该项目的项目名称、创建时间、Created User和项目描述信息。

通用配置

通用配置中的安全模式适用于需要对E-MapReduce数据开发运行的作业进行权限管理的场景。

开启安全模式后,需要在运行集群的用户管理中添加提交作业的EMR用户,详情请参见<mark>管理用户</mark>。在开启安全模式的项目中提交作业时,阿里云账号以默认hadoop用户执 行,RAM用户则默认以当前RAM用户同名的EMR用户执行作业。

- 1. 进入数据开发的项目列表页面。
 - i. 通过阿里云账号登录阿里云E-MapReduce控制台。
 - ii. 在顶部菜单栏处,根据实际情况选择地域和资源组。
 - iii. 单击上方的数据开发页签。
- 2. 进入目标项目的项目管理页面。
 - i. 在**项目列表**页面,单击目标项目的项目ID。
 - ii. 单击上方的**项目管理**页签。
- 3. 配置安全模式。
 - i. 单击左侧菜单**通用配置**。
 - ii. 根据作业情况,选择开启或者关闭安全模式。

↓ 注意 开启安全模式后, Shell和Hive类型的作业无法运行。

用户管理

您可以通过以下步骤为RAM用户添加或删除某个项目的开发权限。

- 1. 进入数据开发的项目列表页面。
 - i. 通过阿里云账号登录阿里云E-MapReduce控制台。
 - ii. 在顶部菜单栏处,根据实际情况选择地域和资源组。
 - ⅲ. 单击上方的**数据开发**页签。
- 2. 进入目标项目的项目管理页面。
 - i. 在**项目列表**页面,单击目标项目的项目ID。
- ii. 单击上方的**项目管理**页签。
- 3. 在左侧导航栏中,单击**用户管理**。
- 4. 在**用户管理**页面,根据业务诉求添加用户或者删除用户。

| ○ 添加 a. b. | 用户。 单击右上角的 添加用户 。 在 添加用户 对话框中,选择需要添加的RAM用户,然后 用户添加成功后,就可以在 用户管理 页面查看新增的用 | 5单击 添加 。 月户信息。 |
|--------------------|--|--|
| | ⑦ 说明 被添加的RAM用户将成为该项目的成员, | 拥有查看、开发该项目下的作业和工作流的权限。 |
| ○ 删除 在 用 | 用户。 户管理 页面,单击待删除用户所在行的 删除 。 | |
| 集群设置 | ł | |
| 通过以下步 | 骤,可以为项目设置集群资源,使该项目中的作业运行 | 在关联的集群上。 |
| 1. 进入数 | 据开发的项目列表页面。 | |
| i. 通 | 过阿里云账号登录 <mark>阿里云E-MapReduce控制台</mark> 。 | |
| ii. 在 | 顶部菜单栏处,根据实际情况选择地域和资源组。 | |
| iii. 单 | 击上方的 数据开发 页签。 | |
| 2. 进入目 | 标项目的项目管理页面。 | |
| i. 在 | 项目列表 页面,单击目标项目的项目ID。 | |
| ii. 单 | 击上方的 项目管理 页签。 | |
| 3. 在左侧 | 导航栏中,单击 集群设置 。 | |
| 4. 在 集群 | 设置 页面,可以进行如下操作: | |
| ∘ 添加 | 集群。 | |
| a. | 单击右上角的 添加集群 | |
| b. | 在 添加集群 对话框中,选择资源组和集群信息。 | |
| | 从选择集群下拉列表中选择已购买的包年包月或按量作 | J费的集群(不包括通过集群模板创建的集群)。 |
| с. | 单击 确定 。 | |
| | 在 集群设置 页面,可以查看新增的集群信息。 | |
| ○ 修改 | 集群配置。 | |
| a. | 单击目标集群所在行的 修改配置 。 | |
| b. | 在 修改配置 对话框中,设置提交作业到该集群的队列、 | 用户信息和客户端白名单。 |
| | 配置项 | 描述 |
| | 提交作业默认用户 | 设置项目使用所选集群提交作业时的默认用户,默认值是hadoop,默认用户只能有一个。 |
| | 提交作业默认队列 | 设置项目使用所选集群提交作业时的默认队列,默认作业提交到default队列。 |
| | 提交作业用户白名单 | 设置可以提交作业的用户,如果有多个用户,可以通过英文半角逗号(,)分隔。 |
| | 提交作业队列白名单 | 设置项目中的作业可以运行在所选集群的队列,如果有多个队列,可以通过英文半角逗号(,)分隔。 |
| | 客户端白名单 | 配置可以提交作业的客户端,您可以使用Master节点或Gateway节点。通过ECS自建的Gateway暂不支持在此处配置。 |

- c. 单击**确定**。
- 取消关联的集群资源。

在**集群设置**页面,单击目标集群所在行的**删除**,取消关联的集群资源。

变量定义

您可以通过以下步骤设置项目级别的自定义变量,这些变量可以被本项目中的作业项目作为全局变量调用。

- 1. 进入数据开发的项目列表页面。
 - i. 通过阿里云账号登录阿里云E-MapReduce控制台。
 - ii. 在顶部菜单栏处,根据实际情况选择地域和资源组。
 - ⅲ. 单击上方的**数据开发**页签。
- 2. 进入目标项目的项目管理页面。
 - i. 在**项目列表**页面,单击目标项目的项目ID。
 - ii. 单击上方的**项目管理**页签。
- 3. 在左侧导航栏中,单击**变量定义**。
- 4. 在**变量定义**页面,可以根据业务诉求添加自定义变量或者删除自定义变量。
- 。 添加自定义变量。
 - a. 单击右上角的**添加**。

```
b. 在添加自定义变量对话框中,添加变量名称和变量值,根据需要选择是否为变量名的变量值加密。
```

作业中以 \${VariableName} 的形式调用变量。例如, 您添加变量名为ENV_ABC, 变量值为12345, 不开启是否为密码。Shell类型作业内容示例如下。

```
echo ${ENV_ABC}
```

返回结果如下。

12345

这里环境变量的设置相当于执行了如下的脚本。

export ENV_ABC=12345

c. 单击**确认**。

在**变量定义**页面,可查看新增的变量信息。

```
。 删除自定义变量。
```

单击目标变量所在行的**删除**,就可以删除对应的变量。

问题反馈

如果您在使用阿里云E-MapReduce过程中有任何疑问,欢迎您扫描下面的二维码加入钉钉群进行反馈。



3.作业编辑

在项目中,您可以通过创建作业来进行任务开发。本文为您介绍作业编辑相关的创建、设置和运行等操作。

背景信息

本文为您提供作业编辑的相关操作,具体如下:

- 新建作业
- 设置作业
- 在作业中添加注解
- 运行作业
- 作业可执行操作
- 作业提交模式说明

前提条件

已创建项目或已被加入到项目中,详情请参见项目管理。

新建作业

- 1. 进入数据开发的项目列表页面。
 - i. 通过阿里云账号登录阿里云E-MapReduce控制台。
 - ii. 在顶部菜单栏处,根据实际情况选择地域和资源组。
 - iii. 单击上方的**数据开发**页签。
- 2. 单击待编辑项目所在行的**作业编辑**。
- 3. 新建作业。
 - i. 在页面左侧,在需要操作的文件夹上单击右键,选择**新建作业**。

⑦ 说明 您还可以通过在文件夹上单击右键,进行新建子文件夹、重命名文件夹和删除文件夹操作。

ii. 在**新建作业**对话框中,输入**作业名称**和**作业描述**,从**作业类型**列表中,选择新建的作业类型。

目前E-MapReduce数据开发支持的作业类型有: Shell、Hive、Hive SQL、Spark、SparkSQL、Spark Shell、Spark Streaming、MR、Sqoop、Pig 、Flink、Streaming SQL、Presto SQL和Impala SQL。

⑦ 说明 创建作业时作业类型一经确定,不能修改。

ⅲ. 单击**确定**。

作业创建成功后,就可以做相应的作业设置、作业编辑等操作了。

设置作业

各类作业类型的开发与设置,请参见<mark>作业</mark>部分。以下内容介绍的是作业的**基础设置、高级设置、共享库**和告警设置。

- 1. 在**作业编辑**页面,单击右上角的**作业设置**。
- 2. 在**作业设置**面板,设置基础信息。

| 配置项 | | 说明 |
|------|--------|---|
| | 作业名称 | 您创建作业的名称。 |
| | 作业类型 | 您创建作业的类型。 |
| | 失败重试次数 | 作业运行失败后的重试次数,可以选择的重试次数范围为: 0~5次。 |
| 作业概要 | 失败策略 | 作业运行失败后支持的策略如下: • 暂停当前工作流:作业运行失败后,不再继续执行当前工作流。 • 继续执行下一个作业:作业运行失败后,继续执行下一个作业。 根据业务情况,可以打开或者关闭使用最新作业内容和参数开关。 • 关闭:作业失败后重新执行时,使用初始作业内容和参数生成作业实例。 • 打开:作业失败后重新执行时,使用最新的作业内容和参数生成作业实例。 |
| | 作业描述 | 单击右侧的 编辑 ,可以修改作业的描述。 |
| 运行资源 | | 单击右侧的 <mark>十</mark> 图标,添加作业执行所依赖的JAR包或UDF等资源。 您需要将资源先上传至OSS,然后在运行资源中直接添加即可。 |
| 配置参数 | | 指定作业代码中所引用的变量的值。您可以在代码中引用变量,格式为 <i>\$/变量名]。</i> 单击右侧的 <mark>十</mark> 图标,添加Key和Value,根据需要选择是否为Value进行加密。其中,Key为变量名,Value为变 量的值。另外,您还可以根据调度启动时间在此配置时间变量,详情请参见 <mark>作业日期设置</mark> 。 |

3. 在**作业设置**面板,单击**高级设置**页签。

| 配置项 | 说明 | |
|------|---|--|
| 模式 | 提交节点包括以下两种模式,详情请参见作业提交模式说明。 在Worker节点提交: 作业通过Launcher在YARN上分配资源进行提交。 在Header/Gateway节点提交: 作业在分配的机器上直接运行。 預期最大运行时长: 0~10800秒。 | |
| 环境变量 | 添加作业执行的环境变量,也可以在作业脚本中直接export环境变量。 示例一:一个Shell类型的任务,内容是 echo \${ENV_ABC} 。如果此处设置了一个环境变量 ENV_ABC=12345 ,则 echo 命令的输出结果为 12345 。 示例二:一个Shell类型的作业,内容是 java-jarabc.jar ,其中<i>abc.jar</i>的内容如下: public static void main(String[] args) {System.out.println(System.getEnv("ENV_ABC"));} 返回的结果是 12345 。此处环境变量的设置相当于执行了以下脚本。 export ENV_ABC=12345 java - jarabc.jar | |
| 调度参数 | 设置作业运行YARN队列、内存、虚拟核数、优先级和执行用户等信息。当未设置这些参数时,作业会直接采用Hadoop集群的默认值。 ⑦ 说明 内存设置用于设置启动器Launcher的内存配额。 | |

4. 在**作业设置**面板,单击**共享库**页签。

在**依赖库**区域,填写**库列表**。

执行作业需要依赖一些数据源相关的库文件。E-MapReduce将这些库以依赖库的形式发布在调度服务的仓库中,在创建作业时需要指定使用哪个版本的依赖库。您只需设 置相应的依赖库版本,例如 sharedlibs:streamingsql:datasources-bundle:2.0.0 。

5. 在**作业设置**面板,单击**告警设置**页签。

| 配置项 | 说明 |
|--------|-----------------------------|
| 执行失败 | 设置作业执行失败时,是否通知到用户告警组或钉钉告警组。 |
| 启动超时 | 设置作业启动超时时,是否通知到用户告警组或钉钉告警组。 |
| 作业执行超时 | 设置作业执行超时时,是否通知到用户告警组或钉钉告警组。 |

在作业中添加注解

进行数据开发时,您可以通过在作业内容里添加特定的注解来添加作业参数。注解的格式如下。

!!! @<注解名称>: <注解内容>

⑦ 说明 … 必须顶格,并且每行一个注解。

当前支持的注解如下。

| 注解名称 | 说明 | 示例 |
|-----------------|--|--|
| rem | 表示一行注释。 | !!! @rem: 这是一行注释 |
| env | 添加一个环境变量。 | !!!@env:ENV_1=ABC |
| var | 添加一个自定义变量。 | !!! @var: var1="value1 and \"one string end with 3 spaces\" " !!! @var: var2=\${yyyy-MM-dd} |
| resource | 添加一个资源文件。 | <pre>!!! @resource: oss://bucket1/dir1/file.jar</pre> |
| sharedlibs | 添加依赖库,仅对Streaming SQL作业有效。包含多个依赖库时,依赖库间 用英文半角逗号(,)隔开。 | !!! @sharedlibs: sharedlibs:streamingsql:datasources- bundle:1.7.0, |
| scheduler.queue | 设置提交队列。 | !!! @scheduler.queue: default |

E-MapReduce

| 注解名称 | 说明 | 示例 |
|--------------------|----------------------|----------------------------|
| scheduler.vmem | 设置申请内存,单位MB。 | !!! @scheduler.vmem: 1024 |
| scheduler.vcores | 设置申请的核数。 | !!! @scheduler.vcores: 1 |
| scheduler.priority | 设置申请的优先级,取值范围为1~100。 | !!! @scheduler.priority: 1 |
| scheduler.user | 设置提交用户名。 | !!! @scheduler.user: root |

♫ 注意

使用注解时,需要注意以下事项:

- 无效注解将被自动跳过。例如,设置未知注解、注解内容不符合预期等。
- 注解中的作业参数优先级高于作业配置中的参数,如果作业注解和作业配置中有相同的参数,则以作业注解为准。

作业注解示例如下:

| ₽ 333 | | × | | Ξ |
|-------|-------|--|--|-----------------------------------|
| > SHE | LL F. | 作业内容: 🕐 適待止 | 保存 作业设置 | 帮助 🗗 |
| 1 | !!! | @rem: 只是一行注释,不会有任何作用 | B BURGHAR | |
| 2 | | | HERE'S AND A CONTRACT OF A CON | and Character menod for table (1) |
| 3 | 111 | @rem: 添加自定义变量,可以在作业内容中引用 | an feast should be the | |
| 4 | !!! | @var: var1 = value1 32 | | |
| 5 | 111 | @var: var2 = \${yyyy-hh-mm} | | |
| 6 | | | | |
| / | 111 | @rem: 添加目定义环境变量 | | |
| 8 | | Geun: ENNTT =YYYY-AAAA. | | |
| 10 | | | | |
| 11 | | | | |
| 12 | | Scheduler vmen 1074 | | |
| 13 | | Scheduler.vcores: 1 | | - |
| 14 | 111 | Gscheduler.priority: 1 | | |
| 15 | | | | |
| 16 | 111 | @rem: 添加资源文件 | | |
| 17 | 111 | @resource: oss://emr-demo/resource/jars/1.jar | | |
| 18 | | | | |
| 19 | 111 | @rem: 添加依赖库(streaming-sql作业) | | |
| 20 | 111 | @ shared libs: shared libs: streaming sql: data sources - bundle: 1.7.0, shared libs: streaming sql: flink-bundle: 1.7.0 | 个 | |
| 21 | | | 1 | |
| 22 | | @rem:在作业中引用相关变量 | + | |
| 23 | ech | o "varl=ş{varl}, var2=ş{var2} ENV_1=ş{ENV_1}, ş{var3}, ş{var4}" | | |
| | | | _ | |
| | | | | |
| | | | 57 | |
| | 实 | 远行(仅供参考) | ^ | |
| 日志 | 运行 | C录 所属工作流 + 摄入OSS路径 ↔ | 去OSS控制台上传 🗗 | |
| | | | | |

运行作业

1. 执行作业。

- i. 在新建的作业页面,单击右上方的**运行**来执行作业。
- ii. 在**运行作业**对话框中,选择资源组和执行集群。
- iii. 单击**确定**。
- 2. 查看作业运行日志。

i. 作业运行后,您可以在**日志**页签中查看作业运行的日志。

| 实际运行(仅供参考) | × ", |
|---|--|
| <pre>bash -c "echo \${TEST}"</pre> | £ |
| 日志 运行记录 所属工作流 审计日志 版本控制 + 插入OSS路径 | ∂ 去OSS控制台上传 🗗 🔥 🗸 🗸 |
| 2021-09-02 16:37:54.653 [main] INFO c.a.e.f.a.j.l.impl.CommonshellJobLauncherImpl - [COMMAND][FJI-]FINFTIJLCES 9] envs(override): (BRK_FLOW_AGENT 2021-09-02 16:37:54.654 [main] INFO c.a.e.f.a.j.l.impl.CommonshellJobLauncherImpl - [COMMAND][FJI-]FINFTIJLCES 9] envs(override): (BRK_FLOW_AGENT PRTH=/mmt/disX)yarn/usercahe/haddop/Agench/Application_163031177326_0022.01 [main] 1000001 _ERK_FLOW_GLOBERT_DIC- 2021-09-02 16:37:54.654 [main] INFO c.a.e.f.a.j.l.impl.CommonshellJobLauncherImpl - [COMMAND][FJI-36'N#*'J;'m.MI *_0] Shell Executor type: com.aliyun ashellExecutor. | 详细日志内容 JOB_ID=FJI-I F+IN :TJINICOC, ¹ &I ⁴ G, FLOW_SKIP_SCI_ANALYZE-fa -c, echo 234] .emr.flow.agent.common.shell.Jav |
| JOB OUTPUT BEGIN | |
| 234 | |
| JOB CUTPUT END************************************ | |
| 2021-09-02 16:37:55.159 [main] HHTO c.a.e.f.a.j.l.impl.CommonShellJobLauncherImpl - [COMMAND][FJJ-3E7880 ⁺⁺⁺ , # +# +# +# +# +# +# +# +# +# +# +# +# + | e=0. wun hook. her is closed already, skip. |
| #####END_OF_LOG##### | |

ii. 单击运行记录页签,可以查看作业实例的运行情况。

iii. 单击目标运行记录右侧的**详情**,跳转到运维中心,可以查看作业实例的详细信息。

作业可执行操作

在**作业编辑**区域,您可以在作业名称上单击右键,执行如下操作。

| 操作 | 说明 |
|-------|---------------------------------------|
| 克隆作业 | 在相同文件夹下,克隆当前作业的配置,生成一个新的作业。 |
| 重命名作业 | 重新命名作业名称。 |
| 删除作业 | 只有在作业没有关联工作流,或关联的工作流没有在运行或调度时,才可以被删除。 |

作业提交模式说明

Spark-Submit进程(在数据开发模块中为启动器Launcher)是Spark的作业提交命令,用于提交Spark作业,一般占用600 MB以上内存。作业设置面板中的内存设置,用于设 置Launcher的内存配额。

新版作业提交模式包括以下两种。

| 作业提交模式 | 描述 |
|---------------------|--|
| 在Header/Gateway节点提交 | Spark-Submit进程运行在Header节点上,不受YARN监控。Spark-Submit内存消耗大,作业过多会造成Header节点资源紧张, 导致整个集群不稳定。 |
| 在Worker节点提交 | Spark-Submit进程运行在Worker节点上,占用YARN的一个Container,受YARN监控。此模式可以缓解Header节点的资源使 用。 |

在E-MapReduce集群中,作业实例消耗内存计算方式如下。

作业实例消耗内存 = Launcher消耗内存 + 用户作业(Job)消耗内存

在Spark作业中,用户作业(Job)消耗内存又可以进一步细分,计算方式如下。

Job消耗内存 = Spark-Submit(指逻辑模块,非进程)消耗内存 + Driver端消耗内存 + Executor端消耗内存

作业配置不同,Driver端消耗的物理内存的位置也不同,详细内容如下表。

| Spark使用模式 | | | Spark-Submit和Driver端 |
|-------------------|-----------------|-------------------------------|---|
| | 作业提交进程使用LOCAL模式 | | 作业提交进程是Header节点上的一个进程,不受YARN监控。 |
| | 作业提交进程使用YARN模式 | | 作业提交进程是Worker节点上的一个进程,占用YARN的一个Container, 受YARN监控。 |
| Yarn-Client模 式 | | Spark-Submit和Driver端是在同一个进程中。 | |
| | | | |

| Spark使用模式 | Spark-Submit和Driver端 |
|----------------|---------------------------------------|
| | |
| | |
| Yarn-Cluster模式 | Driver端是独立的一个进程,与Spark-Submit不在一个进程中。 |
| | |
| | |

问题反馈

如果您在使用阿里云E-MapReduce过程中有任何疑问,欢迎您扫描下面的二维码加入钉钉群进行反馈。



4.工作流编辑

在E-MapReduce数据开发项目中,您可以在作业编辑中定义一组有依赖的作业,然后创建工作流,使作业按照依赖依次执行。E-MapReduce工作流支持基于有向无环图 (DAG)并行执行大数据作业,您可以通过E-MapReduce控制台管控工作流调度以及查看工作流的执行状态。

背景信息

本文为您提供工作流编辑的相关操作介绍,具体如下:

- 新建工作流
- 编辑工作流
- 配置工作流调度
- 执行工作流
- 查看工作流运行记录
- 工作流可执行操作

前提条件

- 已创建项目,详情请参见项目管理。
- 已编辑好作业,详情请参见作业编辑。

新建工作流

- 通过以下步骤,可以新建工作流。
- 1. 进入数据开发的项目列表页面。
 - i. 通过阿里云账号登录阿里云E-MapReduce控制台。
 - ii. 在顶部菜单栏处,根据实际情况选择地域和资源组。
 - iii. 单击上方的**数据开发**页签。
- 2. 在**项目列表**页面,单击目标项目右侧的**工作流设计**。
- 3. 新建工作流。
 - i. 在**工作流设计**区域, 在需要操作的文件夹上单击右键, 选择**新建工作流**。
 - ii. 在新建工作流对话框中,填写工作流名称和工作流描述,选择资源组和执行集群。 在选择执行集群时,可以做如下选择:
 - ■选择当前已存在集群:表示工作流执行时,相关任务会下发到该集群中。
 - 选择集群模板:表示调度系统在工作流启动时先按模版创建一个集群,然后将作业下发到该集群上执行。在工作流结束后,调度系统会自动释放该集群,详情请参见创建集群模板。

⑦ 说明 选择当前已存在集群的下拉列表中只会出现当前项目已关联的集群,如需选择其他集群,您需要取消项目的关联集群,详情请参见项目管理。

ⅲ. 单击**确定**。

工作流创建成功后,就可以进行工作流编辑和配置等操作。

编辑工作流

1. 在工作流编辑画布上,根据业务情况拖拽作业节点到画布上。

每个作业节点被拖入画布后,在**节点编辑**面板可以做如下配置。

| 配置项 | 描述 |
|---------|---|
| 关联作业 | 需要关联对应作业类型的作业。 |
| 自定义作业配置 | 根据业务情况打开或者关闭 自定义作业配置。 • 打开:您可以选择该作业节点的执行集群。 • 关闭(默认):该作业节点的任务会下发到在工作流的执行集群上。 |

2. 设置作业间的依赖关系。

按照作业间的依赖关系,从每个作业节点底部中心拖拽出连线对作业节点进行关联,其中箭头所指方向为工作流的运行方向。

3. 设置控制节点,完成整个工作流的设计。

从**控制节点**区域拖拽END节点到画布,将START节点与工作流开始的作业节点关联,END节点与工作流结束的作业节点关联,完成整个工作流的设计。您可以单击右上 角**自动布局**,调整工作流节点的展示。

| r a l | Ê û Q Q 🛟 🛛 🤧 🌮 🧮 🗒 🗉 |
|----------|-----------------------|
| | |
| 节点组件 | |
| 控制节点 | START |
| START | |
| END | |
| 作业节点 | |
| | V jobwcr 💭 MRjob |
| | |
| | |
| Spark | |
| SparkSQL | \sim |
| MR | END |
| | |
| | |
| | |

在编辑工作流时,您可以单击右上角的**上锁**来为工作流加上编辑锁,此时只有您可以编辑,其他项目成员无法编辑和运行该工作流。只有解锁之后,项目中其他成员才可 以编辑该作业。

⑦ 说明 上锁之后,只有上锁的成员和阿里云账号能够解锁。

配置工作流调度

您可以打开工作流调度配置面板,配置工作流调度参数,调度系统会按照参数定时运行相关工作流,并将作业下发到指定集群上执行。以下介绍如何配置工作流的基本属性、 调度属性和告警设置。

- 1. 进入数据开发的项目列表页面。
 - i. 通过阿里云账号登录阿里云E-MapReduce控制台。
 - ii. 在顶部菜单栏处,根据实际情况选择地域和资源组。
 - iii. 单击上方的**数据开发**页签。
- 2. 在**项目列表**页面,单击目标项目右侧的**工作流设计**。
- 3. 在右侧**工作流**区域,单击**配置**。
- 4. 在工作流调度配置面板的基本属性页签下,您可以修改工作流描述,选择资源组和执行集群。
- 5. 在**工作流调度配置**面板,单击调度属性页签,设置工作流的调度属性。

| 配置项 | | 说明 |
|--------|---------|--|
| 调度状态 | | 支持的调度状态如下: • 启动:可以启动工作流调度。调度开启后,工作流编辑画布上方会出现 调度中 的状态提示。 • 停止:调度状态为停止。 |
| | 开始时间 | 工作流调度的开始时间。 |
| 时间层类调度 | 结束时间 | 可选,工作流调度的结束时间。 |
| 时间周生间皮 | 调度周期 | 工作流调度的周期。 |
| | CRON表达式 | 工作流调度周期的CRON表达式。 |
| 依赖调度属性 | 所属项目 | 可选,当前工作流的前序工作流所属的项目。 |
| | 依赖工作流 | 可选,当前工作流的前序工作流,即前序工作流执行完成后,当前工作流才会被调度执行。 |

6. 在**工作流调度配置**面板,单击**告警设置**页签,设置工作流的告警配置项。

| 配置项 | 说明 |
|------|--------------------------------|
| 执行失败 | 设置工作流执行失败时,是否通知到用户告警组或钉钉告警组。 |
| 节点失败 | 设置工作流节点执行失败时,是否通知到用户告警组或钉钉告警组。 |
| 执行成功 | 设置工作流执行成功时,是否通知到用户告警组或钉钉告警组。 |

| 配置项 | 说明 |
|--------|---|
| 启动超时 | 设置如果工作流中有节点在下发到集群后30分钟内还没有启动时,是否通知到用户告警组或钉钉告警组。 |
| 节点执行超时 | 设置如果节点执行时长超过作业配置里的预期最大运行时长时,是否通知到用户告警组或钉钉告警组。 |

执行工作流

您也可以指定工作流的业务时间,此时工作流作业中的时间相关变量将使用指定的业务时间进行计算,一般用于重跑某个时间段的工作流实例,可以设置单次重跑或批量重 跑。如果您的作业中没有任何时间相关变量,可以选择在当前时间立即执行,即可运行工作流。

- 1. 进入数据开发的项目列表页面。
 - i. 通过阿里云账号登录阿里云E-MapReduce控制台。
 - ii. 在顶部菜单栏处,根据实际情况选择地域和资源组。
 - ⅲ. 单击上方的**数据开发**页签。
- 2. 在**项目列表**页面,单击目标项目右侧的**工作流设计**。
- 3. 运行工作流。
 - i. 在**工作流设计**页面,单击运行。
 - ii. 在运行工作流对话框中,配置相关运行信息。

您可以根据业务情况,选择运行方式。支持的运行方式有**立即执行**和设置时间调度运行,两种运行方式的具体情况如下表。

| 说明 |
|---|
| 立即运行一个工作流,可以将 指定运行时间 作为本工作流的业务时间,时间相关的变量将使用该时间进行计算。 |
| 立即运行一批工作流,将指定调度规则的触发时间作为本工作流的业务时间,时间相关的变量将使用该时间进行计算。一次最多支持100个触发点。您需要设置的信息如下: |
| ■ 开始时间:工作流调度的开始时间。 |
| ■ 结束时间:可选,工作流调度的结束时间。 |
| ■ 调度周期:工作流调度的周期。 |
| ■ CRON表达式:设置调度周期后的CRON表达式。 |
| 跳过成功节点:您可以根据业务情况选择是否打开跳过成功节点。打开该开关后,如果某个业务时间对应的工作流实 例是成功的,将会跳过该业务时间的实例,继续运行其他业务时间的工作流实例。 |
| |

ⅲ. 单击**确定**。

查看工作流运行记录

运行工作流后,可通过以下步骤查看工作流运行记录。

- 1. 在**工作流设计**页面,单击下方运行记录页签。
- 您可以查看工作流实例的运行状态。 2. 单击工作流实例所在行的**详情**,跳转至运维中心。

您可以查看工作流实例的详细情况,也可以暂停、恢复、停止和重跑工作流实例,详情请参见运维中心。

| 功能 | 说明 |
|---------|---|
| 详情 | 查看工作流实例的详细信息,包括工作流节点实例的详细信息和运行状态。 |
| 停止工作流 | 终止正在运行的工作流实例,所有正在运行的作业节点立即停止。 |
| 暂停工作流 | 暂停正在运行的工作流实例,正在运行的作业节点会继续执行,但后续的作业节点不再执行。 |
| 恢复工作流 | 恢复已被暂停的工作流实例。 |
| 重跑工作流实例 | 重新运行已经结束的工作流实例。单击 重跑工作流实例 后,可以选择只重试失败节点,也可以从头重跑所有节点。 |

工作流可执行操作

在工作流设计区域,您可以在工作流名称上单击右键,执行如下操作。

| 操作 | 说明 | |
|--------|-----------------------|--|
| | 在同一文件夹下克隆出相同图形的工作流。 | |
| 克隆工作流 | ⑦ 说明 工作流的调度参数目前无法克隆。 | |
| 重命名工作流 | 重新命名工作流的名称。 | |
| 删除工作流 | 删除工作流。当工作流在运行状态时无法删除。 | |

问题反馈

如果您在使用阿里云E-MapReduce过程中有任何疑问,欢迎您扫描下面的二维码加入钉钉群进行反馈。



5.临时查询

临时查询主要针对即席查询(Ad Hoc)的场景,面向数据科学家和数据分析师,其主要使用工具为SQL。运行临时查询作业时,将会在页面下方显示日志和查询结果。本文为 您介绍如何在临时查询页面新建作业、设置作业、运行作业和编辑锁操作。

背景信息

本文为您介绍临时查询的相关操作,具体如下:

- 新建作业
- 设置作业
- 运行作业
- 编辑锁

前提条件

已创建项目或已被加入到项目中,详情请参见项目管理。

新建作业

1. 进入数据开发的项目列表页面。

- i. 通过阿里云账号登录阿里云E-MapReduce控制台。
- ii. 在顶部菜单栏处,根据实际情况选择地域和资源组。
- ⅲ. 单击上方的**数据开发**页签。
- 2. 在**项目列表**页面,单击目标项目所在行的**作业编辑**。

3. 新建临时查询作业。

- i. 単击最左侧的 Q 图标。
- ii. 在临时查询区域,在需要操作的文件夹上单击右键,选择新建作业。

⑦ 说明 您还可以通过在文件夹上单击右键,进行新建子文件夹、重命名文件夹和删除文件夹操作。

iii. 在新建交互式作业对话框中,填写作业名称和作业描述,从作业类型列表中,选择新建的作业类型。
 E-MapReduce数据开发支持Shell、SparkSQL、Spark Shell和HiveSQL四种类型的临时查询作业。

↓ 注意 创建作业时作业类型一经确定,不能修改。

```
iv. 单击确定。
```

设置作业

各个具体作业类型的开发与设置,请参见<mark>作业</mark>部分。以下内容介绍的是作业的**基础设置、高级设置、共享库**和告警设置。

- 在临时查询页面,单击右上角的作业设置。
- 2. 在**作业设置**面板,设置基础信息。

| 参数 | 描述 |
|------|---|
| 作业概要 | 作业名称:您创建作业的名称。 作业类型:您创建作业的类型。 作业描述:单击右侧的编辑,可以修改作业的描述。 |
| 运行资源 | 单击右侧的 <mark>,</mark> 图标,添加作业执行所依赖的JAR包或UDF等资源。 您需要将资源先上传至OSS,然后在 运行资源 中直接添加即可。 |
| 配置参数 | 指定作业代码中所引用的变量的值。您可以在代码中引用变量,格式为 <i>¥变量名)。</i> 单击右侧的 <mark>-+</mark> 图标,添加Key和Value,根据需要选择是否为Value进行加密。其中,Key为变量名,Value为变量的值。另外,您还可以根据 调度启动时间在此配置时间变量,详情请参见 <mark>作业日期设置。</mark> |

3. 在**作业设置**面板,单击**高级设置**页签。

| 配置项 | 说明 |
|-----|---|
| 模式 | 提交节点包括以下两种模式,详情请参见作业提交模式说明。 在Worker节点提交:作业通过Launcher在YARN上分配资源进行提交。 在Header/Gateway节点提交:作业在分配的机器上直接运行。 预期最大运行时长: 0~10800秒。 |

E-MapReduce

| 配置项 | 说明 |
|------|---|
| 环境变量 | 添加作业执行的环境变量,也可以在作业脚本中直接export环境变量。 示例一:一个Shell类型的任务,内容是 echo \${ENV_ABC} 。如果此处设置了一个环境变量 ENV_ABC=12345 ,则 echo 命令的输出 结果为 12345 。 示例二:一个Shell类型的作业,内容是 java-jarabc.jar ,其中.abc.jar的内容如下: public static void main(String[] args) {System.out.println(System.getEnv("ENV_ABC"));} 返回的结果是 12345 。此处环境变量的设置相当于执行了以下脚本。 export ENV_ABC=12345 java - jara bc.jar |
| 调度参数 | 设置作业运行YARN队列、内存、虚拟核数、优先级和执行用户等信息。当未设置这些参数时,作业会直接采用Hadoop集群的默认值。 ⑦ 说明 内存设置用于设置启动器Launcher的内存配额。 |

4. 在**作业设置**面板,单击**共享库**页签。

在**依赖库**区域,填写**库列表**。

执行作业需要依赖一些数据源相关的库文件。E-MapReduce将这些库以依赖库的形式发布在调度服务的仓库中,在创建作业时需要指定使用哪个版本的依赖库。您只需设 置相应的依赖库版本,例如 sharedlibs:streamingsql:datasources-bundle:2.0.0 。

5. 在**作业设置**面板,单击**告警设置**页签。

| 配置项 | 说明 |
|--------|-----------------------------|
| 执行失败 | 设置作业执行失败时,是否通知到用户告警组或钉钉告警组。 |
| 启动超时 | 设置作业启动超时时,是否通知到用户告警组或钉钉告警组。 |
| 作业执行超时 | 设置作业执行超时时,是否通知到用户告警组或钉钉告警组。 |

运行作业

1. 执行作业。

- i. 在**临时查询**页面,单击右上方的运行执行作业。
- ii. 在**运行作业**对话框中,选择资源组和执行集群。
- ⅲ. 单击**确定**。
- 2. 查看作业运行日志。
 - i. 作业运行后,您可以在**日志**页签下查看作业运行的日志。
 - ii. 单击运行记录页签,可以查看作业实例的运行情况。
 - iii. 单击目标运行记录右侧的**详情**,跳转到运维中心,可以查看作业实例的详细信息。

编辑锁

在编辑作业时,您可以单击右上方的**上锁**,为该作业加上编辑锁,保证此时只有您可以编辑作业,项目其他成员无法编辑该作业。只有解锁之后,项目中其他成员才可以编辑 该作业。

⑦ 说明 上锁之后,只有上锁的成员和阿里云账号管理员能够解锁。

问题反馈

如果您在使用阿里云E-MapReduce过程中有任何疑问,欢迎您扫描下面的二维码加入钉钉群进行反馈。



6.运维中心

本文介绍通过运维中心管控工作流调度任务、监控任务运行状态、查看工作流记录和审计日志,便于您对工作流的管理和运维。

背景信息

本文为您介绍运维中心的相关操作,具体如下:

- 查看项目任务概览
- 管理工作流记录
- 查看审计日志

前提条件

已创建项目,具体请参见项目管理。

查看项目任务概览

- 1. 进入数据开发的项目列表页面。
 - i. 通过阿里云账号登录<mark>阿里云E-MapReduce控制台</mark>。
 - ii. 在顶部菜单栏处,根据实际情况选择地域和资源组。
 - iii. 单击上方的**数据开发**页签。
- 2. 在**项目列表**页面,单击目标任务所在行的运行记录。
- 3. 在左侧导航栏中,单击**概览**。
 - 您可以查看项目任务的概览信息。

管理工作流记录

- 1. 进入数据开发的项目列表页面。
 - i. 通过阿里云账号登录阿里云E-MapReduce控制台。
 - ii. 在顶部菜单栏处,根据实际情况选择地域和资源组。
 - ⅲ. 单击上方的**数据开发**页签。
- 2. 在**项目列表**页面,单击目标任务所在行的**运行记录**。
- 3. 工作流记录管理。
 - 管理工作流记录信息的方式如下:
 - 。 工作流记录

在工作流记录页面,您可以查看工作流实例的相关信息,同时可以对工作流实例进行停止、暂停和恢复操作。

| 功能 | 描述 |
|----|--|
| 详情 | 单击目标工作流实例所在行的 详情 ,可以查看工作流实例的详细信息,包括工作流实例的详细信息和运行状态。 |
| 停止 | 单击目标工作流实例所在行的 停止 ,可以终止正在运行的工作流实例。 |
| 暂停 | 单击目标工作流实例所在行的 暂停 ,可以暂停正在运行的工作流实例。 |
| 恢复 | 单击目标工作流实例所在行的 恢复 ,可以恢复已被暂停的工作流实例。 |

○ 手动运行作业记录

a. 在左侧导航栏中,选择**工作流记录 > 手动运行作业记录**。

b. 在**手动运行作业记录**页面,可以查看作业实例运行详情,同时可以对作业实例进行停止操作。

| 功能 | 描述 |
|----|--|
| 详情 | 单击目标作业实例所在行的 详情 ,可以查看作业实例的详细信息,包括执行参数、作业内容以及作业日志。 |
| 停止 | 单击目标作业实例所在行的 停止 ,可以终止正在运行的作业实例。 |

- 。 流式作业
 - a. 在左侧导航栏中,选择**工作流记录 > 流式作业**。
 - b. 在流式作业页面,可以查看流式作业实例运行详情,同时可以对流式作业实例进行编辑、启动、停止和查看历史记录操作。

| 功能 | 描述 |
|------|--|
| 详情 | 单击目标作业所在行的 详情 ,可以查看流式作业实例的详细信息,包括执行参数、作业内容以及作业日志。 |
| 编辑 | 单击目标作业所在行的编辑,可以进入 作业编辑 页面,修改作业内容。 |
| 启动 | 单击目标作业所在行的 启动 ,可以执行流式作业。 |
| 停止 | 单击目标作业所在行的 停止 ,可以终止正在运行的流式作业。 |
| 历史记录 | 单击目标作业所在行的 历史记录 ,可以查看流式作业的运行记录。 |

查看审计日志

您可以通过以下步骤查看项目成员在项目中的操作历史。

1. 进入数据开发的项目列表页面。

- i. 通过阿里云账号登录阿里云E-MapReduce控制台。
- ii. 在顶部菜单栏处,根据实际情况选择地域和资源组。
- ⅲ. 单击上方的**数据开发**页签。
- 2. 在**项目列表**页面,单击目标任务所在行的运行记录。
- 3. 在左侧导航栏中,单击**审计日志**。
 - 在**审计日志**页面,可以查看项目成员的操作历史。

问题反馈

如果您在使用阿里云E-MapReduce过程中有任何疑问,欢迎您扫描下面的二维码加入钉钉群进行反馈。



7.创建集群模板

集群模版是为快速创建集群而保存的配置。本文为您介绍如何创建集群模板。

背景信息

集群模版主要用于数据开发工作流自动创建临时集群。在使用数据开发工作流完成作业任务时,如果您只关注作业任务是否完成,可以使用集群模版功能来快速建立集群,调 度系统会在工作流启动时按照模版创建一个集群,然后将作业下发到该集群上执行。当工作流结束后,调度系统会自动释放该集群。

使用限制

集群模版目前仅支持Hadoop和Dataflow两种集群类型,如需其他集群类型,您可以<mark>提交工单</mark>处理。

操作步骤

- 1. 进入数据开发的项目列表页面。
 - i. 通过阿里云账号登录阿里云E-MapReduce控制台。
 - ii. 在顶部菜单栏处,根据实际情况选择地域和资源组。
 - iii. 单击上方的**数据开发**页签。
- 2. 在**项目列表**页面,单击右上角的**集群模板**。

在集群模板列表页面,您可以对已创建的模板进行编辑和删除操作。

| 功能 | 说明 |
|----|---|
| 编辑 | 单击目标模板所在行的 编辑 ,可以修改集群模板。修改完成后,会立即生效到使用该模板的工作流。 |
| | 单击目标模板所在行的 删除 ,可以删除该集群模板。 |
| 删除 | ⑦ 说明 系统不会检查此模版是否被引用,删除之后,通过集群模版自动创建集群的工作流会失败。 |

3. 创建集群模板。

- i. 单击右上角的**创建集群模板**。
- ii. 在创建集群模板页面,配置相关信息。

```
创建集群模版的过程与创建集群基本一致,详情请参见创建集群。
```

在创建集群模板的硬件配置中,您可以为实例设置多机型实例,可以避免单一机型库存不足造成集群创建失败最终影响作业执行。

| ☑ 实例 ● 洗型配置 ┏ | 单机型实例 | 多机型实例 | | | | | |
|---------------|--------|-------|-------|------------------|-----|-----------------|--|
| | Master | | VCPU: | 8 | ~ | | |
| | Core | | | Core | | | |
| | Task | | 内存: | 16 | ~ | GB | |
| | | | 实例: | 备选实例 | 已选到 | 实例 (最多选择三条) | |
| | | | | ecs.hfc6.2xlarge | | ecs.c6e.2xlarge | |
| | | | | ecs.n4.2xlarge | | ecs.c6.2xlarge | |
| | | | | | | | |
| | | | | | | | |
| | | | | 2 10 | |) 顶 | |
| | | | | 乙 坝 | | 4 坝 | |

iii. 完成上述参数配置后,选中E-MapReduce服务条款后,单击保存模板。 模板创建成功后,就可以在模板列表中查看到。

问题反馈

如果您在使用阿里云E-MapReduce过程中有任何疑问,欢迎您扫描下面的二维码加入钉钉群进行反馈。



8.云监控事件编码

在云监控的事件监控模块中,您可以订阅E-MapReduce数据开发相关的系统事件,监控集群的核心组件服务状态。

云监控系统事件编码及其含义如下。

| 事件编码 | 事件描述 | 事件类型 |
|---------------|-------------|------|
| EMR-110401002 | 工作流已成功。 | FLOW |
| EMR-110401003 | 工作流已提交。 | FLOW |
| EMR-110401004 | 作业已提交。 | FLOW |
| EMR-110401005 | 工作流节点已启动。 | FLOW |
| EMR-110401006 | 工作流节点状态已检查。 | FLOW |
| EMR-110401007 | 工作流节点已完成。 | FLOW |
| EMR-110401008 | 工作流节点已结束。 | FLOW |
| EMR-110401009 | 工作流节点已取消。 | FLOW |
| EMR-110401010 | 工作流已取消。 | FLOW |
| EMR-110401011 | 工作流已重跑。 | FLOW |
| EMR-110401012 | 工作流已恢复。 | FLOW |
| EMR-110401013 | 工作流已暂停。 | FLOW |
| EMR-110401014 | 工作流已结束。 | FLOW |
| EMR-110401015 | 工作流节点已失败。 | FLOW |
| EMR-110401016 | 作业已失败。 | FLOW |
| EMR-210401001 | 工作流已失败。 | FLOW |
| EMR-210401003 | 工作流节点启动超时。 | FLOW |
| EMR-210401004 | 作业启动超时。 | FLOW |

9.作业配置 9.1. 作业日期设置

在作业编辑的过程中,支持在作业参数中设置时间变量通配符。

变量通配符格式

E-MapReduce所支持的变量通配符的格式为≴(dateexpr-1d)或者≸(dateexpr-1h)。其中dateexpr表示标准的时间格式表达式,对应的规则如下。

↓ 注意 请注意时间格式的大小写。

| 格式 | 描述 |
|------|--------------------|
| уууу | 表示4位的年份。 |
| ММ | 表示月份。 |
| dd | 表示天。 |
| нн | 表示24小时制,12小时制使用hh。 |
| mm | 表示分钟。 |
| SS | 表示秒。 |

时间变量可以是包含*yyyy*年份的任意时间组合,同时支持用加号(+)和减号(-)来分别表示延后和提前。例如,变量*≴(yyyy-MM-dd)*表示当前日期,则:

- 后1年的表示方式: *\${yyyy+1y}*或者*\${yyyy-MM-dd hh:mm:ss+1y*}。
- 后3月的表示方式: \${yyyyMM+3m}或者\${yyyy-MM-dd hh:mm:ss+3m}。
- 前5天的表示方式: \${yyyyMMdd-5d}或者\${yyyy-MM-dd hh:mm:ss-5d}。
- 例如,假设当前时间为20160427 12:08:01:
- 如果在作业参数中写成\$(yyyyMMdd HH:mm:ss-1d),那么这个参数通配符在真正执行的时候会被替换成20160426 12:08:01,即在当前日期上减了一天并精确到了秒。
- 如果写成 \${yyyyMMdd-1d},则执行时会替换成20160426,表示当前日期的前一天。
- 如果写成*\${yyyyMMdd},*则会被替换成20160427,直接表示当前的日期。

? 说明

- 目前E-MapReduce仅支持小时和天维度的加减,即只支持在dateexpr后面+Nd、-Nd、+Nh、-Nh的形式(dateexpr为时间格式表达式,N为整数)。
- 时间变量参数必须以yyyy开始,如*\${yyyy-MM*}。如果希望单独获取月份等特定时间区域的值,可以在作业内容中使用如下两个函数提取:
 - parseDate(<参数名称>, <时间格式>):将给定参数转换为Date对象。其中,参数名称为上述配置参数中设置的一个变量名,时间格式为设置该变量时所使用的时间格式。如设置一个变量 current_time = \${yyyy/MMddHHmmss-1d},则此处时间格式应设置为yyyy/MdddHHmmss。
 - formatDate(<Date对象>, <时间格式>):将给定Date对象转换为给定格式的时间字符串。

函数使用示例:

- 获取current_time变量的小时字面值: \${formatDate(parseDate(current_time, 'yyyyMMddHHmmss'), 'HH')}
- 。 获取current_time变量的年字面值: \${formatDate(parseDate(current_time, 'yyyy/MMddHHmmss'), 'yyyy')}

操作示例

- 1. 进入数据开发的项目列表页面。
 - i. 通过阿里云账号登录阿里云E-MapReduce控制台。
 - ii. 在顶部菜单栏处,根据实际情况选择地域和资源组。
 - ⅲ. 单击上方的**数据开发**页签。
- 2. 在**项目列表**页面,单击目标项目所在行的**作业编辑**。
- 3. 进行作业设置。
 - i. 在作业编辑页面,选择需要操作的作业名称,单击右上角的作业设置。

ii. 在配置参数区域,单击-图标,新增参数,并按照上文介绍的变量通配符格式填写参数。

| | 11-11天主. | T II V L |
|---|---|---|
| -f rankings_uservisits_join_hdfs.hivehiveconf date{ <mark>\${dy_date}</mark> } | 失败重试次数: 失败策略: | 0次 ~ 继续执行下一个作业 ~ |
| | 作业描述: | 测试用例: 集成测试, Hive SQL QueryHdfs |
| | | 0 编辑 |
| | 运行资源 | |
| | oss://em | /rankings_userv |
| | 配置参数 | ? |
| 实际运行(汉供参考) hive -f rankings_uservisits_join_hdfs.hivehiveconf date=\${dy_date} | ▲ 查码 参数1: | dy_date - |
| | -f rankings_uservisits_join_hdfs.hivehiveconf date 実际运行(仅供参考) hive -f rankings_uservisits_join_hdfs.hivehiveconf date=\${dy_date} | 中国大学部 中国大学部 中国大学部 中国大学部 中国大学部 中国大学部 中国大学部 中国大学部 中国大学 中国大学 中国大学 中国大学 中国大学 中国大学 中国大学 中国大学 |

配置完成后就可以在作业中引用配置参数的Key了。

问题反馈

如果您在使用阿里云E-MapReduce过程中有任何疑问,欢迎您扫描下面的二维码加入钉钉群进行反馈。



9.2. Shell作业配置

本文介绍如何配置Shell类型的作业。

前提条件

已创建好项目,详情请参见<mark>项目管理</mark>。

操作步骤

- 1. 进入数据开发的项目列表页面。
 - i. 通过阿里云账号登录阿里云E-MapReduce控制台。
 - ii. 在顶部菜单栏处,根据实际情况选择地域和资源组。
 - iii. 单击上方的**数据开发**页签。
- 2. 单击待编辑项目所在行的**作业编辑**。
- 3. 新建Shell类型作业。
 - i. 在页面左侧,在需要操作的文件夹上单击右键,选择**新建作业**。
 - ii. 在新建作业对话框中,输入作业名称和作业描述,从作业类型下拉列表中选择Shell作业类型。
 - 表示创建的作业是一个Bash Shell作业。
- ⅲ. 单击确定。
- 4. 编辑作业内容。
 - i. 在**作业内容**中,填写提交该作业需要提供的命令行参数。

示例如下。

DD=`date`;

echo "hello world, \$DD"

```
ii. 单击保存,作业内容编辑完成。
```

问题反馈

如果您在使用阿里云E-MapReduce过程中有任何疑问,欢迎您扫描下面的二维码加入钉钉群进行反馈。



9.3. Hive作业配置

E-MapReduce默认提供了Hive环境,您可以直接使用Hive来创建和操作创建的表和数据。

前提条件

- 已创建好项目,详情请参见项目管理。
- 已准备好Hive SQL的脚本,并上传到OSS的某个目录中(例如*oss://path/to/uservisits_aggre_hdfs.hive*)。
- uservisits_aggre_hdfs.hive内容如下。

USE DEFAULT;

```
DROP TABLE uservisits;
CREATE EXTERNAL TABLE IF NOT EXISTS uservisits (sourceIP STRING,destURL STRING,visitDate STRING,adRevenue DOUBLE,userAgent STRING,countryCode STRING,language
Code STRING,searchWord STRING,duration INT) ROW FORMAT DELIMITED FIELDS TERMINATED BY ',' STORED AS SEQUENCEFILE LOCATION '/HiBench/Aggregation/Input/uservi
sits':
```

DROP TABLE uservisits_aggre;

CREATE EXTERNAL TABLE IF NOT EXISTS uservisits_aggre (sourceIP STRING, sumAdRevenue DOUBLE) STORED AS SEQUENCEFILE LOCATION '/HiBench/Aggregation/Output/us ervisits_aggre';

INSERT OVERWRITE TABLE uservisits_aggre SELECT sourceIP, SUM(adRevenue) FROM uservisits GROUP BY sourceIP;

操作步骤

- 1. 进入数据开发的项目列表页面。
 - i. 通过阿里云账号登录阿里云E-MapReduce控制台。
 - ii. 在顶部菜单栏处,根据实际情况选择地域和资源组。
 - iii. 单击上方的**数据开发**页签。
- 2. 单击待编辑项目所在行的**作业编辑**。
- 3. 新建Hive类型作业。
 - i. 在页面左侧,在需要操作的文件夹上单击右键,选择**新建作业**。
 - ii. 在新建作业对话框中,输入作业名称和作业描述,从作业类型下拉列表中选择Hive作业类型。 表示创建的作业是一个Hive作业。这种类型的作业,实际是通过以下方式提交Hive作业运行。

hive [user provided parameters]

- ⅲ. 单击**确定**。
- 4. 编辑作业内容。
 - i. 在**作业内容**中,填写提交该作业需要提供的命令行参数。

例如,如果需要使用刚刚上传到OSS的Hive脚本,则填写的内容如下。

-f ossref://path/to/uservisits_aggre_hdfs.hive

⑦ 说明 path 为 uservisits_aggre_hdfs.hive 在OSS上的路径。

您也可以单击下方的**+插入OSS路径**,从OSS中进行浏览和选择,系统会自动补齐OSS上Hive脚本的路径。请务必将Hive脚本的前缀修改为OSSREF,以保证E-MapReduce可以正确下载该文件。

ii. 单击**保存**,作业内容编辑完成。

问题反馈

如果您在使用阿里云E-MapReduce过程中有任何疑问,欢迎您扫描下面的二维码加入钉钉群进行反馈。



9.4. Hive SQL作业配置

本文介绍如何配置Hive SQL类型的作业。

前提条件

已创建好项目,详情请参见项目管理。

操作步骤

- 1. 进入数据开发的项目列表页面。
 - i. 通过阿里云账号登录阿里云E-MapReduce控制台。
 - ii. 在顶部菜单栏处,根据实际情况选择地域和资源组。
 - iii. 单击上方的**数据开发**页签。

```
2. 单击待编辑项目所在行的作业编辑。
```

```
3. 新建Hive SQL类型作业。
```

```
i. 在页面左侧, 在需要操作的文件夹上单击右键, 选择新建作业。
```

| 其中 SOL CONTENT 为作 | | |
|---|---------------------|-------------|
| | 业编辑器中填写的SQL语句。 | |
| 新建作业 | × | |
| * 所属项目: | | |
| * 所属文件夹: | | |
| * 作业名称: | Hive SQL test | |
| * 作业描述: | this is a test | |
| | | |
| * 作业类型 | tiveSQL \vee | |
| | 教徒を取得 | |
| show databases; show tables; 系统会自动为SELECT语作 select * from test1; | 1加上'limit 2000'的限制。 | |
| | A Ø (の)返行 | ◎停止 保存 作业设置 |
| 1 SQL语句示例 2 SQL语句最大不能超过64KB 3 show databases; | 上'limit 2000'的限制 | |
| 4 show tables; 5 系统会自动为SELECT语句加 6 select * from test1; | | |
| 4 shou tables; 5 新始全自动为SELECT语句加 6 select * from test1; | | |
| 4 show tables; 5 系统全自动为SELECT语句加 6 select * from test1; | | |
| 4 show tables; 5 系统合目动为seter语句讷 6 select * from test1; | | |

ii. 单击**保存**,作业内容编辑完成。

问题反馈

如果您在使用阿里云E-MapReduce过程中有任何疑问,欢迎您扫描下面的二维码加入钉钉群进行反馈。



9.5. Spark作业配置

本文介绍如何配置Spark类型的作业。

前提条件 已创建好项目, 详情请参见项目管理。 操作步骤 1. 进入数据开发的项目列表页面。 i. 通过阿里云账号登录阿里云E-MapReduce控制台。 ii. 在顶部菜单栏处,根据实际情况选择地域和资源组。 iii. 单击上方的**数据开发**页签。 2. 单击待编辑项目所在行的**作业编辑**。 3. 新建Spark类型作业。 i. 在页面左侧, 在需要操作的文件夹上单击右键, 选择**新建作业**。 ii. 在新建作业对话框中,输入作业名称和作业描述,从作业类型下拉列表中选择Spark作业类型。 表示创建的作业是一个Spark作业。这种类型的作业,实际是通过以下方式提交的Spark作业运行。 spark-submit [options] --class [MainClass] xxx.jar args ⅲ. 单击**确定**。 4. 编辑作业内容。 i. 在**作业内容**中,填写提交该作业需要提供的命令行参数。 只需要填写spark-submit之后的参数即可。 以下分别展示如何填写创建Spark作业和Pyspark作业的参数: ■ 创建Spark作业。 新建一个Spark作业,作业名称为Wordcount,应用参数填写示例如下: ■ 在命令行下提交完整的命令。 spark-submit --master yarn-client --driver-memory 7G --executor-memory 5G --executor-cores 1 --num-executors 32 --class com.aliyun.emr.checklist.benchmark.S parkWordCount emr-checklist_2.10-0.1.0.jar oss://emr/checklist/data/wc oss://emr/checklist/data/wc-counts 32 ■ 在E-MapReduce作业的作业内容输入框中填写如下命令。 --master yarn-client --driver-memory 7G --executor-memory 5G --executor-cores 1 --num-executors 32 --class com.aliyun.emr.checklist.benchmark.SparkWordCou nt ossref://emr/checklist/jars/emr-checklist_2.10-0.1.0.jar oss://emr/checklist/data/wc oss://emr/checklist/data/wc-counts 32 ↓ 注意 JAR包保存在OSS中,引用这个JAR包的方式是ossref://emr/checklist/jars/emr-checklist_2.10-0.1.0.jar。您可以单击下方的+插入OSS路径,文 件前缀选择OSSREF,从文件路径中进行浏览和选择,系统会自动补齐OSS上Spark脚本的路径。 ■ 创建Pyspark作业。 E-MapReduce除了支持Scala或者Java类型作业外,还支持Python类型Spark作业。新建一个Python脚本的Spark作业,作业名称为Python-Kmeans,应用参数填写 示例如下: --master yarn-client --driver-memory 7g --num-executors 10 --executor-memory 5g --executor-cores 1 ossref://emr/checklist/python/kmeans.py oss://emr/checklist/d ata/kddb 5 32 ↓ 注意 ■ 支持Python脚本资源的引用,同样使用ossref协议。 ■ Pyspark不支持通过作业方式安装Python工具包。 ii. 单击**保存**, 作业内容编辑完成。

问题反馈

如果您在使用阿里云E-MapReduce过程中有任何疑问,欢迎您扫描下面的二维码加入钉钉群进行反馈。



9.6. Spark SQL作业配置

本文介绍如何配置Spark SQL类型的作业。

前提条件

已创建好项目,详情请参见<mark>项目管理</mark>。

操作步骤

1. 进入数据开发的项目列表页面。

```
i. 通过阿里云账号登录<mark>阿里云E-MapReduce控制台</mark>。
```

```
ii. 在顶部菜单栏处,根据实际情况选择地域和资源组。
```

ⅲ. 单击上方的**数据开发**页签。

2. 单击待编辑项目所在行的**作业编辑**。

3. 新建Spark SQL类型作业。

- i. 在页面左侧,在需要操作的文件夹上单击右键,选择**新建作业**。
- ii. 在新建作业对话框中,输入作业名称和作业描述,从作业类型下拉列表中选择Spark SQL作业类型。

⑦ 说明 Spark SQL提交作业的模式默认是Yarn-client模式。

此类型的作业,实际是通过以下方式提交的Spark SQL作业运行。

spark-sql [options] [cli options] {SQL_CONTENT}

参数描述如下表。

| 参数 | 说明 |
|--------------|---|
| options | 在作业设置面板的高级设置页签,单击环境变量所在行的 <mark>,</mark> 图标,添加环境变量SPARK_CLI_PARAMS,例如 SPARK_CLI_PARAMS="- -executor-memory 1gexecutor-cores" 。 |
| cli options | 示例如下: • -e <quoted-query-string> : 表示运行引号内的SQL查询语句。 • -f <filename> : 表示运行文件中的SQL语句。</filename></quoted-query-string> |
| SQL_CONT ENT | 填写的SQL语句。 |

ⅲ. 单击确定。

4. 编辑作业内容。

```
i. 在作业内容中,输入Spark SQL语句。
示例如下。
-- SQL语句示例。
-- SQL语句最大不能超过64 KB。
show databases;
```

```
show tables;
-- 系统会自动为SELECT语句加上'limit 2000'的限制。
select * from test1;
```

ii. 单击**保存**,作业内容编辑完成。

问题反馈

如果您在使用阿里云E-MapReduce过程中有任何疑问,欢迎您扫描下面的二维码加入钉钉群进行反馈。



9.7. Spark Shell作业配置

本文介绍如何配置Spark Shell类型的作业。

前提条件

已创建好项目,详情请参见项目管理。

操作步骤

- 1. 进入数据开发的项目列表页面。
 - i. 通过阿里云账号登录<mark>阿里云E-MapReduce控制台</mark>。
 - ii. 在顶部菜单栏处,根据实际情况选择地域和资源组。
 - ⅲ. 单击上方的**数据开发**页签。
- 2. 单击待编辑项目所在行的**作业编辑**。
- 3. 新建Spark Shell类型作业。
 - i. 在页面左侧,在需要操作的文件夹上单击右键,选择**新建作业**。
 - ii. 在新建作业对话框中,输入作业名称和作业描述,从作业类型下拉列表中选择Spark Shell作业类型。
 - ⅲ. 单击**确定**。
- 4. 编辑作业内容。

i. 在**作业内容**中,输入Spark Shell命令后续的参数。

示例如下。

```
val count = sc.parallelize(1 to 100).filter { _=>
val x = math.random
val y = math.random
x*x + y*y < 1
}.count();
println("Pi is roughly ${4.0 * count / 100}")</pre>
```

ii. 单击**保存**,作业内容编辑完成。

问题反馈

如果您在使用阿里云E-MapReduce过程中有任何疑问,欢迎您扫描下面的二维码加入钉钉群进行反馈。



9.8. Spark Streaming作业配置

本文介绍如何配置Spark Streaming类型的作业。

前提条件

- 已创建好项目,详情请参见项目管理。
- 已准备好作业所需的资源,以及作业要处理的数据。

操作步骤

- 1. 进入数据开发的项目列表页面。
 - i. 通过阿里云账号登录阿里云E-MapReduce控制台。
 - ii. 在顶部菜单栏处,根据实际情况选择地域和资源组。
 - iii. 单击上方的**数据开发**页签。
- 2. 单击待编辑项目所在行的**作业编辑**。
- 3. 新建Spark Streaming类型作业。
 - i. 在页面左侧,在需要操作的文件夹上单击右键,选择**新建作业**。
 - ii. 在新建作业对话框中,输入作业名称和作业描述,从作业类型下拉列表中选择Spark Streaming作业类型。

ⅲ. 单击**确定**。

- 4. 编辑作业内容。
 - i. 在**作业内容**中,填写提交该作业需要提供的命令行参数。

Spark Streaming作业提交命令的格式如下。

spark-submit [options] --class [MainClass] xxx.jar args

```
作业名称以SlsStreaming为例,作业内容示例如下。
```

--master yarn-client --driver-memory 7G --executor-memory 5G --executor-cores 1 --num-executors 32 --class com.aliyun.emr.checklist.benchmark.SlsStreaming emr-checklist_2.10-0.1.0.jar <project> <logstore> <accessKey> <secretKey>

↓ 注意

- 如果作业JAR包保存在OSS中,则引用这个JAR包的方式是ossref://xxx/.../xxx.jar。
- 您可以单击下方的+插入OSS路径,选择文件前缀为OSSREF,从文件路径中进行浏览和选择,系统会自动补齐OSS上Spark Streaming脚本的路径。

ii. 单击**保存**,作业内容编辑完成。

问题反馈

如果您在使用阿里云E-MapReduce过程中有任何疑问,欢迎您扫描下面的二维码加入钉钉群进行反馈。



9.9. Hadoop MapReduce作业配置

本文介绍如何配置Hadoop MapReduce类型的作业。

前提条件

已创建好项目,详情请参见项目管理。

操作步骤

- 1. 进入数据开发的项目列表页面。
 - i. 通过阿里云账号登录阿里云E-MapReduce控制台。
 - ii. 在顶部菜单栏处,根据实际情况选择地域和资源组。
 - iii. 单击上方的数据开发页签。
- 2. 单击待编辑项目所在行的**作业编辑**。
- 3. 新建Hadoop MapReduce类型作业。
 - i. 在页面左侧,在需要操作的文件夹上单击右键,选择**新建作业**。
 - ii. 在新建作业对话框中,输入作业名称和作业描述,从作业类型下拉列表中选择MR作业类型。
 - 表示创建的作业是一个Hadoop MapReduce作业。这种类型的作业,实际是通过以下方式提交运行。

hadoop jar xxx.jar [MainClass] -D xxx

```
ⅲ. 单击确定。
```

4. 编辑作业内容。

```
i. 在作业内容中,填写提交该作业需要提供的命令行参数。
```

填写的命令行参数需要从 hadoop jar 命令后的第一个参数开始填写,即在输入框中首先填写运行该作业所需JAR包的所在路径,再填写 [MainClass] 和其它您想要设 置的命令行参数。

例如,您想要提交一个Hadoop的sleep作业,该作业不读写任何数据,只提交一些mapper和reducertask到集群中,且每个task执行时需要 sleep一段时间。在 Hadoop(以hadoop-2.6.0版本为例)中,该作业处于Hadoop发行版的 hadoop-mapreduce-client-jobclient-2.6.0-tests.jar包文件中。如果您通过命令行的方式提 交该作业,需要执行以下命令。

hadoop jar /path/to/hadoop-mapreduce-client-jobclient-2.6.0-tests.jar sleep -m 3 -r 3 -mt 100 -rt 100

而在E-MapReduce中配置这个作业,则应在作业内容输入框中填写以下内容。

/path/to/hadoop-mapreduce-client-jobclient-2.6.0-tests.jar sleep -m 3 -r 3 -mt 100 -rt 100

⑦ 说明 您也可以单击下方的+插入OSS路径,选择文件前缀为OSSREF,从文件路径中进行浏览和选择,系统会自动补齐OSS上Hadoop MapReduce脚本的路径。

ii. 单击**保存**,作业内容编辑完成。

上面示例中,sleep作业并没有数据的输入输出,如果作业要读取数据,并输出处理结果(例如Wordcount),则需要指定数据的Input和Output路径。 您可以读写E-MapReduce集群HDFS或OSS上的数据。如果需要读写OSS上的数据,只需要在填写Input和Output路径时,将数据路径写成OSS上的路径地址即可。

jar ossref://emr/checklist/jars/chengtao/hadoop/hadoop-mapreduce-examples-2.6.0.jar randomtextwriter -D mapreduce.randomtextwriter.totalbytes=320000 oss://emr/checklist/data/chengtao/hadoop/Wordcount/Input

问题反馈

如果您在使用阿里云E-MapReduce过程中有任何疑问,欢迎您扫描下面的二维码加入钉钉群进行反馈。



9.10. Sqoop作业配置

本文介绍如何配置Sqoop类型的作业。

前提条件

已创建好项目,详情请参见项目管理。

使用限制

EMR-1.3.0及后续版本支持Sqoop作业类型。在低版本集群上运行Sqoop作业会失败,errlog会报不支持的错误。参数详细请请参见Sqoop。

操作步骤

- 1. 进入数据开发的项目列表页面。
 - i. 通过阿里云账号登录阿里云E-MapReduce控制台。
 - ii. 在顶部菜单栏处,根据实际情况选择地域和资源组。
- ⅲ. 单击上方的**数据开发**页签。
- 2. 单击待编辑项目所在行的**作业编辑**。
- 3. 新建Sqoop类型作业。
 - i. 在页面左侧,在需要操作的文件夹上单击右键,选择**新建作业**。
 - ii. 在**新建作业**对话框中,输入**作业名称**和**作业描述**,从**作业类型**下拉列表中选择Sqoop作业类型。
 - ⅲ. 单击确定。
- 4. 编辑作业内容。

i. 在作业内容中,输入Sqoop命令后续的参数。

示例如下所示。

sqoop [args]

ii. 单击**保存**,作业内容编辑完成。

问题反馈

如果您在使用阿里云E-MapReduce过程中有任何疑问,欢迎您扫描下面的二维码加入钉钉群进行反馈。



9.11. Pig作业配置

本文介绍如何配置Pig类型的作业。

前提条件

- 已创建好项目,详情请参见项目管理。
- 已准备好Pig的脚本,示例如下。

/ * Licensed to the Apache Software Foundation (ASF) under one

* or more contributor license agreements. See the NOTICE file

* distributed with this work for additional information

* regarding copyright ownership. The ASF licenses this file

- * to you under the Apache License, Version 2.0 (the * "License"); you may not use this file except in compliance
- * with the License. You may obtain a copy of the License at

* http://www.apache.org/licenses/LICENSE-2.0

* Unless required by applicable law or agreed to in writing, software * distributed under the License is distributed on an "AS IS" BASIS.

* WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.

* See the License for the specific language governing permissions and

* limitations under the License.

-- Query Phrase Popularity (Hadoop cluster)

-- This script processes a search query log file from the Excite search engine and finds search phrases that occur with particular high frequency during certain times of the day.

-- Register the tutorial JAR file so that the included UDFs can be called in the script.

REGISTER oss://emr/checklist/jars/chengtao/pig/tutorial.jar;

-- Use the PigStorage function to load the excite log file into the "raw" bag as an array of records.

-- Input: (user,time,query)

raw = LOAD 'oss://emr/checklist/data/chengtao/pig/excite.log.bz2' USING PigStorage('\t') AS (user, time, query);

-- Call the NonURLDetector UDF to remove records if the query field is empty or a URL

clean1 = FILTER raw BY org.apache.pig.tutorial.NonURLDetector(query);

-- Call the ToLower UDF to change the query field to lowercase.

clean2 = FOREACH clean1 GENERATE user, time, org.apache.pig.tutorial.ToLower(query) as query; -- Because the log file only contains queries for a single day, we are only interested in the hour.

-- The excite query log timestamp format is YYMMDDHHMMSS.

-- Call the ExtractHour UDF to extract the hour (HH) from the time field.

houred = FOREACH clean2 GENERATE user, org.apache.pig.tutorial.ExtractHour(time) as hour, query;

-- Call the NGramGenerator UDF to compose the n-grams of the query. ngramed1 = FOREACH houred GENERATE user, hour, flatten(org.apache.pig.tutorial.NGramGenerator(query)) as ngram;

-- Use the DISTINCT command to get the unique n-grams for all records.

ngramed2 = DISTINCT ngramed1;

-- Use the GROUP command to group records by n-gram and hour

hour_frequency1 = GROUP ngramed2 BY (ngram, hour);

-- Use the COUNT function to get the count (occurrences) of each n-gram.

hour_frequency2 = FOREACH hour_frequency1 GENERATE flatten(\$0), COUNT(\$1) as count;

-- Use the GROUP command to group records by n-gram only. -- Each group now corresponds to a distinct n-gram and has the count for each hour.

uniq_frequency1 = GROUP hour_frequency2 BY group::ngram;

-- For each group, identify the hour in which this n-gram is used with a particularly high frequency.

-- Call the ScoreGenerator UDF to calculate a "popularity" score for the n-gram.

uniq_frequency2 = FOREACH uniq_frequency1 GENERATE flatten(\$0), flatten(org.apache.pig.tutorial.ScoreGenerator(\$1));

-- Use the FOREACH-GENERATE command to assign names to the fields.

uniq_frequency3 = FOREACH uniq_frequency2 GENERATE \$1 as hour, \$0 as ngram, \$2 as score, \$3 as count, \$4 as mean;

-- Use the FILTER command to move all records with a score less than or equal to 2.0.

filtered_uniq_frequency = FILTER uniq_frequency3 BY score > 2.0;

-- Use the ORDER command to sort the remaining records by hour and score. ordered_uniq_frequency = ORDER filtered_uniq_frequency BY hour, score;

-- Use the PigStorage function to store the results.

-- Output: (hour, n-gram, score, count, average_counts_among_all_hours) STORE ordered_uniq_frequency INTO 'oss://emr/checklist/data/chengtao/pig/script1-hadoop-results' USING PigStorage();

• 已保存该脚本文件script1-hadoop-oss.pig, 并上传到OSS的某个目录中,例如oss://path/to/script1-hadoop-oss.pig。

操作步骤

| 1. 进, | 入数据开发的项目列表页面。 |
|-------|---|
| i | .通过阿里云账号登录 <mark>阿里云E-MapReduce控制台。</mark> |
| ii | . 在顶部菜单栏处,根据实际情况选择地域和资源组。 |
| iii | . 单击上方的 数据开发 页签。 |
| 2. 单i | 击待编辑项目所在行的 作业编辑 。 |
| 3. 新發 | 建Pig类型作业。 |
| i | .在页面左侧,在需要操作的文件夹上单击右键,选择 新建作业 。 |
| ii | .在 新建作业 对话框中,输入 作业名称 和 作业描述 ,从 作业类型 下拉列表中选择Pig作业类型。 |
| | 表示创建的作业是一个Pig作业。这种类型的作业,实际是通过以下方式提交的Pig作业运行。 |
| | pig [user provided parameters] |
| iii | . 单击 确定 。 |
| 4. 编辑 | 辑作业内容。 |
| i | .在 作业内容 中,填写提交该作业需要提供的命令行参数。 |
| | 例如,如果需要使用刚刚上传到OSS的Pig脚本,则填写的内容如下。 |
| | -x mapreduce ossref://emr/checklist/jars/chengtao/pig/script1-hadoop-oss.pig |
| | |

⑦ 说明 您也可以单击下方的+插入OSS路径,选择文件前缀为OSSREF,从文件路径中进行浏览和选择,系统会自动补齐OSS上Pig脚本的路径。

ii. 单击**保存**,作业内容编辑完成。

问题反馈

如果您在使用阿里云E-MapReduce过程中有任何疑问,欢迎您扫描下面的二维码加入钉钉群进行反馈。



9.12. Flink (VVR) 作业配置

EMR-3.27.x及之前版本使用Flink社区开源版本,EMR-3.27.x之后版本使用完全兼容开源Flink的企业版(VVR)。本文介绍如何配置Flink(VVR)类型的作业。

背景信息

Flink企业版由Apache Flink创始团队官方出品,拥有全球统一商业化品牌。

VVR提供企业版StateBackend,性能是开源版本的3~5倍。在EMR Hadoop集群中,您可使用VVR引擎和EMR数据开发功能提交作业。VVR支持开源Flink 1.10版本,默认使用商 业GeminiStateBackend,具备以下特性:

- 采用创新的数据结构,提高随机查询、降低读磁盘I/O的性能。
- 优化Cache策略,内存充足情况下热数据不落盘,并且Compaction后Cache不会失效。
- 完全使用Java实现, 消除RocksDB的JNI开销。
- 使用堆外内存,并基于GeminiDB的特点实现高效的内存分配器,消除JVM GC带来的影响。
- 支持异步增量Checkpoint,同步阶段只进行内存索引的拷贝,相较于RocksDB可以避免I/O带来的抖动。
- 支持Local Recovery和Timer落盘。

② 说明 如果您想使用GeminiStateBackend,请不要在代码中指定StateBackend类型。使用GeminiStateBackend启动时,TM的内存不少于1728 MB。

Flink中Checkpoint和StateBackend的基础配置同样适用于GeminiStateBackend,具体请参见Configuration。

您可以根据具体需求配置参数,部分特殊参数设置如下。

| 参数 | 说明 |
|--|---|
| state.backend.gemini.memory.managed | 默认值为true,表示将自动根据Managed Memory以及Task Slot数计算每个Backend的内存。取值如下: • true • false |
| state.backend.gemini.offheap.size | 默认值为2 GB,当state.backend.gemini.memory.managed为false时,设置每个Backend的 内存。 |
| state.backend.gemini.local.dir | 表示GeminiDB本地数据文件的存放目录。 |
| state.backend.gemini.timer-service.factory | 默认值为HEAP,表示timer-service state的存储位置。取值如下: • HEAP • GEMINI |

⑦ 说明 参数配置方法请参见管理组件参数。

前提条件

- 已创建Hadoop集群,详情请参见创建集群。
- 已创建项目,详情请参见项目管理。
- 已获取作业所需的资源,以及作业需要处理的数据文件,例如,JAR包、数据文件名称及其保存路径。

操作步骤

- 1. 进入数据开发的项目列表页面。
 - i. 通过阿里云账号登录阿里云E-MapReduce控制台。
 - ii. 在顶部菜单栏处,根据实际情况选择地域和资源组。
 - iii. 单击上方的**数据开发**页签。
- 2. 单击待编辑项目所在行的**作业编辑**。
- 3. 新建Flink类型作业。
 - i. 在页面左侧, 在需要操作的文件夹上单击右键, 选择新建作业。
 - ii. 在新建作业对话框中,输入作业名称和作业描述,从作业类型下拉列表中选择Flink作业类型。
 - iii. 单击**确定**。
- 4. 编辑作业内容。
 - i. 在**作业内容**中,填写提交该作业需要提供的命令行参数。
 - Flink类型作业支持JAR包形式的Flink Datastream、Table和SQL作业,示例如下。

run -m yarn-cluster -yjm 1024 -ytm 2048 ossref://path/to/oss/of/WordCount.jar --input oss://path/to/oss/to/data --output oss://path/to/oss/to/result

■ EMR-3.x版本自EMR-3.28.2版本开始, Flink类型作业同时支持PyFlink作业, 示例如下。

run -m yarn-cluster -yjm 1024 -ytm 2048 -py ossref://path/to/oss/of/word_count.py

PyFlink作业其它可用参数,请参见Apache Flink官方文档。

```
ii. 单击保存。
```

- ⑦ 说明 您可以根据集群的版本来访问Flink的Web UI:
 - EMR-3.29.0之前版本
 - 仅支持通过SSH隧道方式访问Web Ul时,请参见<mark>通过SSH隧道方式访问开源组件Web</mark> Ul。
 - EMR-3.29.0及后续版本
 - (推荐)您可以通过EMR控制台的方式访问Web Ul时,请参见访问链接与端口。
 - 您可以通过SSH隧道方式访问Web Ul时,请参见通过SSH隧道方式访问开源组件Web Ul。

问题反馈

如果您在使用阿里云E-MapReduce过程中有任何疑问,欢迎您扫描下面的二维码加入钉钉群进行反馈。



9.13. Streaming SQL作业配置

本文介绍Streaming SQL作业配置的操作步骤。

背景信息

Streaming SQL的详细信息请参见Spark Streaming SQL。

在Streaming SQL作业配置过程中,您需要设置依赖库。以下列出了Spark Streaming SQL提供的数据源依赖包的版本信息和使用说明,建议使用最新版本。

| 库名称 | 版本 | 发布日期 | 引用字符串 | 详细信息 |
|--------------------|-----------|------------|--|--|
| | 2.0.0(推荐) | 2020/02/26 | sharedlibs:streamingsql:datasources- bundle:2.0.0 | 支持数据源:Kafka、Loghub、Druid、 TableStore、HBase、JDBC、DataHub、Redis、 Kudu和DTS。 |
| datasources-bundle | 1.9.0 | 2019/11/20 | sharedlibs:streamingsql:datasources- bundle:1.9.0 | 支持数据源: Kafka、Loghub、Druid、 TableStore、HBase、JDBC、DataHub、Redis和 Kudu。 |
| | 1.8.0 | 2019/10/17 | sharedlibs:streamingsql:datasources- bundle:1.8.0 | 支持数据源: Kafka、Loghub、Druid、 TableStore、HBase、JDBC、DataHub和Redis。 |
| | 1.7.0 | 2019/07/29 | sharedlibs:streamingsql:datasources- bundle:1.7.0 | 支持数据源:Kafka、Loghub、Druid、 TableStore、HBase和JDBC。 |

如果需要了解更详细的使用方法,请参见数据源。

前提条件

- 已创建项目,详情请参见项目管理。
- 已获取作业所需的资源和数据文件。例如, JAR包、数据文件名称以及两者的保存路径。

```
操作步骤
 1. 进入数据开发的项目列表页面。
    i. 通过阿里云账号登录阿里云E-MapReduce控制台。
    ii. 在顶部菜单栏处,根据实际情况选择地域和资源组。
    iii. 单击上方的数据开发页签。
 2. 单击待编辑项目所在行的作业编辑。
 3. 新建Streaming SQL类型作业。
    i. 在页面左侧, 在需要操作的文件夹上单击右键, 选择新建作业。
    ii. 在新建作业对话框中,输入作业名称和作业描述,从作业类型下拉列表中选择Streaming SQL作业类型。
    ⅲ. 单击确定。
 4. 编辑作业内容。
   在作业内容中,填写提交该作业需要提供的命令行参数。示例如下。
     ---- 创建SLS数据表。
    CREATE TABLE IF NOT EXISTS ${slsTableName}
     USING loghub
     OPTIONS (
      sls.project = '${logProjectName}',
      sls.store = '${logStoreName}',
      access.key.id = '${accessKeyId}',
       access.key.secret = '${accessKeySecret}',
      endpoint = '${endpoint}'
     );
     ,,,
---- 导入数据至HDFS。
    INSERT INTO
     ${hdfsTableName}
    SELECT
     col1. col2
    FROM ${slsTableName}
    WHERE ${condition}
```

⑦ 说明 此类型的作业是通过 streaming-sql-f {sql_script 提交的。 sql_script 中保存着作业编辑器中填写的SQL语句。

```
5. 配置依赖库和失败策略。
```

```
i. 单击右上方的作业设置。
```

```
ii. 分别在共享库和流任务设置页签下,配置依赖库和失败处理策略。
```

| 区域 | 配置项 | 说明 |
|--------|-----------|---|
| 依赖库 | 库列表 | 执行作业需要依赖一些数据源相关的库文件。E-MapReduce将这些库以依赖库的形式发布在调度服务的仓库 中,在创建作业时需要指定使用哪个版本的依赖库。 您只需设置相应的依赖库版本,例如 <i>sharedlibs:streamingsql:datasources-bundle:2.0.0</i> 。 |
| 失败处理策略 | 当前语句执行失败时 | 当前语句执行失败时,支持如下策略: 继续执行下一条语句:如果查询语句执行失败,继续执行下一条语句。 终止当前作业:如果查询语句执行失败,终止当前作业。 |

iii. 单击**保存**,完成作业内容及相关配置。

问题反馈

如果您在使用阿里云E-MapReduce过程中有任何疑问,欢迎您扫描下面的二维码加入钉钉群进行反馈。



9.14. Presto SQL作业配置

在数据开发过程中如果您需要使用Presto SQL,可以在E-MapReduce中配置Presto SQL作业。本文介绍如何配置Presto SQL作业。

操作步骤

```
1. 进入数据开发的项目列表页面。
```

- i. 通过阿里云账号登录<mark>阿里云E-MapReduce控制台</mark>。
- ii. 在顶部菜单栏处,根据实际情况选择地域和资源组。

```
iii. 单击上方的数据开发页签。
```

```
2. 单击待编辑项目所在行的作业编辑。
```

```
3. 新建Presto SQL类型作业。
```

```
i. 在页面左侧, 在需要操作的文件夹上单击右键, 选择新建作业。
```

```
ii. 在新建作业对话框中,输入作业名称和作业描述,从作业类型下拉列表中选择Presto SQL作业类型。
```

表示创建的作业是一个Presto SQL作业。这种类型的作业,其运行实际是通过以下方式提交的Presto SQL作业。

presto <options> -f {SQL_SCRIPT}

```
⑦ 说明 其中 SQL_SCRIPT 中保存着作业编辑器中填写的SQL语句。
```

```
iii. 单击确定。
```

4. 编辑作业内容。

- i. 在**作业内容**中,填写提交该作业需要提供的命令行参数。
 - 示例如下。
 - SELECT * from table1;

ii. 单击**保存**, 作业内容编辑完成。

设置Presto Cli参数

默认情况下,Presto查询catalog=hive,schema=default下的数据表,可以通过设置Presto Cli参数来指定不同的Catalog和Schema。Presto SQL作业支持如下两种方式设置 Presto Cli参数:

- 通过环境变量设置
 - 设置密码:如果Presto服务开启了密码认证,可以通过添加名为 PRESTO_PASSWORD 的环境变量来传入密码。
- 设置其他参数:可以将参数设置到名为 PRESTO_CLI_PARAMS 的环境变量中,如 PRESTO_CLI_PARAMS="--catalog mysql --schema db1"。
- 通过自定义变量
 - 设置密码:在作业自定义变量中添加名为 presto.password 的变量,即可设置Presto认证密码。
 - 设置其他参数:在作业自定义变量中添加如 _presto.xxx 的变量,都会被添加到Presto Cli参数列表中,对应的选项为 --xxx 。
 - 支持如下自定义变量。
 - ## 基本参数
 - _presto.schema <schema>
 - * _presto.catalog <catalog>
 - ## 控制/调试参数
 - * _presto.trace-token <trace token>
 - *_presto.session <session>.
 - * _presto.source <source>
 - * _presto.resource-estimate <resource-estimate>...
 - * _presto.log-levels-file <log levels file>
 - ## 连接参数
 - * _presto.server <server>
 - * _presto.http-proxy <http-proxy> * ignore-errors
 - * _presto.socks-proxy <socks-proxy> ## 认证相关参数
 - ## 以血怕大学奴
 - * _presto.user <user> * _presto.password <password>
 - * _presto.client-info <client-info>
 - *_presto.client-request-timeout <client request timeout>
 - * _presto.client-tags <client tags>
 - * _presto.access-token <access token>
 - * _presto.truststore-password <truststore password>
 - * _presto.truststore-path <truststore path>
 - * _presto.keystore-password <keystore password>
 - * _presto.keystore-path <keystore path>
 - * _presto.extra-credential <extra-credential>..
 - ## 高安全相关参数
 - * presto.krb5-config-path <krb5 config path>
 - * _presto.krb5-credential-cache-path <krb5 credential cache path>

 - * _presto.krb5-keytab-path <krb5 keytab path>
 - * _presto.krb5-principal <krb5 principal>
 - * _presto.krb5-remote-service-name <krb5 remote service name>
 - *_presto.krb5-service-principal-pattern <krb5 remote service principal pattern>

问题反馈

如果您在使用阿里云E-MapReduce过程中有任何疑问,欢迎您扫描下面的二维码加入钉钉群进行反馈。



9.15. Impala SQL作业配置

在数据开发过程中如果您需要使用Impala SQL,可以在E-MapReduce中配置Impala SQL作业。本文介绍如何配置Impala SQL作业。

前提条件

已创建好项目,详情请参见项目管理。

操作步骤

- 1. 进入数据开发的项目列表页面。
 - i. 通过阿里云账号登录阿里云E-MapReduce控制台。
 - ii. 在顶部菜单栏处,根据实际情况选择地域和资源组。
 - iii. 单击上方的**数据开发**页签。
- 2. 单击待编辑项目所在行的**作业编辑**。
- 3. 新建Impala SQL类型作业。
 - i. 在页面左侧,在需要操作的文件夹上单击右键,选择**新建作业**。
 - ii. 在**新建作业**对话框中,输入**作业名称**和**作业描述**,从**作业类型**下拉列表中选择Impala SQL作业类型。

此类型作业,实际是通过以下方式提交的Impala SQL作业。

impala-shell -f {SQL_CONTENT} [options];

参数描述如下表。

| 参数 | 说明 |
|--------------|--|
| SQL_CONT ENT | 填写的SQL语句。 |
| options | 在作业设置面板的高级设置页签,单击环境变量所在行的 <mark>,</mark> 图标,添加环境变量IMPALA_CLI_PARAMS,例如 IMAPAL_CLI_PARAMS= "-u hive" 。 |

ⅲ. 单击确定。

4. 编辑作业内容。

i. 在作业内容中,输入Impala SQL语句。
 示例如下。

show databases; show tables; select * from test1;

ii. 单击**保存**,作业内容编辑完成。

问题反馈

如果您在使用阿里云E-MapReduce过程中有任何疑问,欢迎您扫描下面的二维码加入钉钉群进行反馈。



10.数据开发常见问题

本文汇总了数据开发时的常见问题。

作业问题:

- 作业和执行计划的区别是什么?
- 多个ConsumerID消费同一个Topic时为什么TPS不一致?
- 为什么Hive创建的外部表没有数据?
- 为什么Spark Streaming作业运行一段时间后无故结束?
- 为什么Spark Streaming作业已经结束,但是E-MapReduce控制台显示作业还处于"运行中"状态?
- 如何在MR作业中使用本地共享库?
- 如何在MR或Spark作业中指定OSS数据源文件路径?
- Beeline如何访问Kerberos安全集群?
- Spark接收Flume数据时为什么出现内存超用?
- 为什么作业运行较慢?
- 为什么AppMaster调度启动Task的时间过长?
- 导入RDS数据至EMR时,时间字段显示延迟8小时如何处理?
- 使用数据开发提交的作业一直处于Submit状态,该怎么办?

日志问题:

- 如何查看作业日志?
- 如何在OSS上查看日志?
- E-MapReduce中是否可以查看作业的Worker上日志?
- 如何查看E-MapReduce服务的日志?
- 如何清理已经完成作业的日志数据?

异常诊断:

- "Error: Could not find or load main class"
- Spark SQL连RDS出现 "Invalid authorization specification, message from server: ip not in whitelist"
- 读写MaxCompute时, 抛出java.lang.RuntimeException.Parse response failed: '<!DOCTYPE html>...'
- Spark SQL报错 "Exception in thread "main" java.sql.SQLException: No suitable driver found for jdbc:mysql:xxx"
- Hive或Impala作业读取SparKSQL导入的Parquet表报错(表包含Decimal格式的列) "Failed with exception java.io.IOException:org.apache.parquet.io.ParquetDecodingException: Can not read value at 0 in block -1 in file hdfs://.../.../part-00000-xxx.snappy.parquet"
- ThriftServer进程正常,但链接出现异常,报错 "Connection refused telnet emr-header-1 10001"无法连接
- Spark作业报错 "Container killed by YARN for exceeding memory limits." 或者MR作业报错 "Container is running beyond physical memory limits."
- "Error: Java heap space"
- "No space left on device"
- 访问OSS或LogService时报错ConnectTimeoutException或ConnectionException
- 读取Snappy文件时报错Out Of MemoryError
- "Exception in thread main java.lang.RuntimeException: java.lang.ClassNotFoundException: Class com.aliyun.fs.oss.nat.NativeOssFileSystem not found"
- Spark中使用OSS SDK出现
- "java.lang.NoSuchMethodError:org.apache.http.conn.ssl.SSLConnetionSocketFactory.init(Ljavax/net/ssl/SSLContext;Ljavax/net/ssl/HostnameVerifier)"
- "java.lang.lllegalArgumentException: Wrong FS: oss://xxxxx, expected: hdfs://ip:9000
- Spark作业报错 "java.lang.lllegalArgumentException: Size exceeds Integer.MAX_VALUE"

功能使用:

- E-MapReduce是否提供实时计算的功能?
- 导入RDS数据至EMR时,时间字段显示延迟8小时如何处理?
- 如何修改Spark服务的spark-env配置?
- 作业参数传递至脚本文件该如何处理?
- 如何设置HiveServer2的认证方式为LDAP?
- 如何使用阿里云E-MapReduce HDFS的Balancer功能以及参数调优?
- 如何使用standlone模式提交Spark任务?
- 如果E-MapReduce控制台上没有自定义配置选项,该如何处理?

作业和执行计划的区别是什么?

• 创建作业

在E-MapReduce中创建作业,实际只是创建了作业如何运行的配置,该配置中包括该作业要运行的JAR包、数据的输入输出地址以及一些运行参数。该配置创建好后,给它 命名即定义了一个作业。

- 执行计划
- 执行计划是将作业与集群关联起来的一个纽带:
- 可以把多个作业组合成一个作业序列。
- 。 可以为作业准备一个运行集群(或者自动创建出一个临时集群或者关联一个已存在的集群)。
- 可以为这个作业序列设置周期执行计划,并在完成任务后自动释放集群。
- 可以在执行记录列表上查看每一次执行的执行成功情况与日志。

如何查看作业日志?

您可以直接在EMR控制台上查看作业日志。如果您是登录到Master节点提交作业和运行脚本,则您可以根据脚本自行规划。

如何在OSS上查看日志?

- 1. 在E-MapReduce数据开发的页面,找到对应的工作流实例,单击运行记录。
- 2. 在运行记录区域,单击待查看工作流实例所在行的详情,在作业实例信息页面查看执行集群ID。
- 3. 在日志保存目录 OSS://mybucket/emr/spark下,查找执行集群ID目录。
- 4. 在OS5://mybucket/emr/spark/clusterID/jobs目录下会按照作业的执行ID存放多个目录,每个目录下存放了这个作业的运行日志文件。

读写MaxCompute时,抛出java.lang.RuntimeException.Parse response failed: '<!DOCTYPE html>...'

问题分析:可能是MaxCompute Tunnel Endpoint填写错误。

解决方法:输入正确的MaxCompute Tunnel Endpoint。

多个ConsumerID消费同一个Topic时为什么TPS不一致?

有可能这个Topic在公测或其他环境创建过,导致某些Consumer组消费数据不一致。请将对应的Topic和ConsumerID提交工单处理。

E-MapReduce中是否可以查看作业的Worker上日志?

可以。

为什么Hive创建的外部表没有数据?

问题描述:创建完外部表后查询没有数据返回。

外部表创建语句举例如下。

CREATE EXTERNAL TABLE storage_log(content STRING) PARTITIONED BY (ds STRING) ROW FORMAT DELIMITED FIELDS TERMINATED BY 'lt' STORED AS TEXTFILE LOCATION 'oss://log-12453****/your-logs/airtake/pro/storage';

查询没有数据返回。

hive> select * from storage_log;

问题分析: Hive不会自动关联指定Partitions目录。

解决方法:需要您手动指定Partitions目录。

alter table storage_log add partition(ds=123); Time taken: 0.137 seconds hive> select * from storage_log;

返回如下数据。

OK abcd 123 efgh 123

为什么Spark Streaming作业运行一段时间后无故结束?

- 首先检查Spark版本是否是1.6之前版本,如果是的话更新Spark版本。
- Spark 1.6之前版本存在内存泄漏的BUG,此BUG会导致Container被中止掉。
- 检查自己的代码在内存使用上有没有做好优化。

为什么Spark Streaming作业已经结束,但是E-MapReduce控制台显示作业还处于"运行中"状态?

ок

问题分析: Spark Streaming作业的运行模式是Yarn-Client。

解决方法:因为E-MapReduce对Yarn-Client模式的Spark Streaming作业的状态监控存在问题,所以请修改为Yarn-Cluster模式。

"Error: Could not find or load main class"

检查作业配置中作业JAR包的路径协议头是否是 ossref , 如果不是请改为 ossref 。

如何在MR作业中使用本地共享库?

您可以在阿里云E-MapReduce控制台,YARN服务的配置页面,修改mapred-site.xml页签如下参数。

<property>
<name>mapred.child.java.opts</name>
<value>-Xmx1024m -Djava.library.path=/usr/local/share/</value>
</property>
<property>
<name>mapreduce.admin.user.env</name>
<value>LD_LIBRARY_PATH=\$HADOOP_COMMON_HOME/lib/native:/usr/local/lib</value>
</property>

如何在MR或Spark作业中指定OSS数据源文件路径?

您可以在作业中指定输入输出数据源时使用 OSS URL: oss://[accessKeyId:accessKeySecret@]bucket[.endpoint]/object/path 形式,类似 hdfs://。

您在操作OSS数据时:

• (建议)E-MapReduce提供了MetaService服务,支持免AccessKey访问OSS数据,直接写*oss://bucket/object/path*。

 (不建议)可以将AccessKey ID, AccessKey Secret以及Endpoint配置到Configuration (Spark作业是SparkConf, MR作业是Configuration)中,也可以在URI中直接指定 AccessKey ID、AccessKey Secret以及Endpoint。详情请参见开发准备。

Spark SQL报错 "Exception in thread "main" java.sql.SQLException: No suitable driver found for jdbc:mysql:xxx"

问题分析: mysql-connector-java版本过低。

解决方法:更新*mysql-connector-java*至最新版本。

Spark SQL连RDS出现 "Invalid authorization specification, message from server: ip not in whitelist" 检查RDS的白名单设置,将集群机器的内网地址添加到RDS的白名单中。

Hive或Impala作业读取SparkSQL导入的Parquet表报错(表包含Decimal格式的列) "Failed with exception java.io.IOException:org.apache.parquet.io.ParquetDecodingException: Can not read value at 0 in block -1 in file hdfs://.../part-00000-xxx.snappy.parquet"

由于Hive和SparkSQL在Decimal类型上使用了不同的转换方式写入Parquet,导致Hive无法正确读取SparkSQL导入的数据。对于已有的使用SparkSQL导入的数据,如果有被Hive 或Impala使用的需求,建议加上spark.sql.parquet.writeLegacyFormat=true重新导入数据。

Beeline如何访问Kerberos安全集群?

● HA集群 (Discovery模式)

!connect jdbc:hive2://emr-header-1:2181,emr-header-2:2181,emr-header-3:2181/;serviceDiscoveryMode=zooKeeper;zooKeeperNamespace=hiveserver2;principal=hive/_HOST @EMR.\${clusterld).COM

- HA集群
 - 连接emr-header-1

!connect jdbc:hive2://emr-header-1:10000/;principal=hive/emr-header-1@EMR.\${clusterId}.COM

◦ 连接emr-header-2

!connect jdbc:hive2://emr-header-2:10000/;principal=hive/emr-header-2@EMR.\${clusterId}.COM

非HA集群

!connect jdbc:hive2://emr-header-1:10000/;principal=hive/emr-header-1@EMR.\${clusterId}.COM

ThriftServer进程正常,但链接出现异常,报错"Connection refused telnet emr-header-1 10001"无法连接

您可以查看/mnt/disk1/log/spark日志。该问题是由于thrift server oom需要扩大内存,因此调大spark.driver.memory的值即可。

如何查看E-MapReduce服务的日志?

登录Master节点在/mnt/disk1/log中查看对应服务的日志。

Spark作业报错 "Container killed by YARN for exceeding memory limits." 或者MR作业报错 "Container is running beyond physical memory limits."

问题分析:提交App时申请的内存量较低,但JVM启动占用了更多的内存,超过了自身的申请量,导致被NodeManager异常终止。特别是Spark类型作业,可能会占用多的堆 外内存,很容易被异常终止。

解决方法:

- Spark作业,在阿里云E-MapReduce控制台,YARN服务的配置页面,调大spark.yarn.driver.memoryOverhead或spark.yarn.executor.memoryOverhead的值。
- MapReduce作业,在阿里云E-MapReduce控制台,YARN服务的配置页面,调大mapreduce.map.memory.mb或mapreduce.reduce.memory.mb的值。

"Error: Java heap space"

问题分析:作业Task处理的数据量较大,但Task JVM申请的内存量不足从而抛出Out Of MemoryError。 解决方法:

- Spark 作业,在阿里云E-MapReduce控制台,YARN服务的配置页面,调大spark.execut or.memory或 spark.driver.memory的值。
- MapReduce作业, 在阿里云E-MapReduce控制台, YARN服务的配置页面, 调大mapreduce.map java.opts或mapreduce.reduce.java.opts的值。

"No space left on device"

问题分析:

- Master或Worker节点空间不足,导致作业失败。
- 磁盘空间满导致本地Hive元数据库(MySQL Server)异常,Hive Met ast ore连接报错。

解决方法:清理Master节点磁盘空间、系统盘的空间以及HDFS空间。

访问OSS或LogService时报错ConnectTimeoutException或ConnectionException

问题分析:OSS Endpoint需要配置为公网地址,但EMR Worker节点并无公网IP,所以无法访问。

解决方法:

- 修改OSS Endpoint 地址修为内网地址。
- 使用EMR met aservice功能,不指定Endpoint。

例如 select * from tbl limit 10 可以正常运行,但是执行 Hive SQL: select count(1) from tbl 时报错。修改OSS Endpoint地址为内网地址。

alter table tbl set location "oss://bucket.oss-cn-hangzhou-internal.aliyuncs.com/xxx" alter table tbl partition (pt = 'xxx-xx-xx') set location "oss://bucket.oss-cn-hangzhou-internal.aliyuncs.com/xxx"

读取Snappy文件时报错OutOfMemoryError

问题分析:LogService等服务写入的标准Snappy文件和Hadoop的Snappy文件格式不同,EMR默认处理的是Hadoop修改过的Snappy格式,处理标准格式时会抛出 Out Of MemoryError。

解决方法,在阿里云E-MapReduce控制台,相应服务的配置页面,配置如下:

- Hive作业: 配置 set io.compression.codec.snappy.native=true 。
- MR作业: 配置 Dio.compression.codec.snappy.native=true 。
- Spark作业: 配置 spark.hadoop.io.compression.codec.snappy.native=true 。

"Exception in thread main java.lang.RuntimeException: java.lang.ClassNotFoundException: Class com.aliyun.fs.oss.nat.NativeOssFileSystem not found"

在Spark作业中读写OSS数据时,需要安装E-MapReduce SDK,详情请参见准备工作。

Spark接收Flume数据时为什么出现内存超用?

检查接收数据方式是否是Push-based。如果不是,请修改为Push-based方式接收数据,详情请参见Spark Streaming + Flume Integration Guide。

"Caused by: java.io.IOException: Input stream cannot be reset as 5242880 bytes have been written, exceeding the available buffer size of 524288"

(OSS)网络连接重试时缓存不足,请使用1.1.0版本以上的aliyun-java-sdk-emr包。

Spark中使用OSS SDK出现

"java.lang.NoSuchMethodError:org.apache.http.conn.ssl.SSLConnetionSocketFactory.init(Ljavax/net/ssl/SSLCor 因为OSS SDK与Spark和Hadoop的运行环境存在版本依赖冲突,所以不建议在代码中使用OSS SDK。

"java.lang.IllegalArgumentException: Wrong FS: oss://xxxxx, expected: hdfs://ip:9000"

因为在操作OSS数据时,使用HDFS的默认fs,所以在初始化时,需要使用OSS的路径来初始化fs,以便于使用fs来操作OSS源上的数据。

Path outputPath = new Path(EMapReduceOSSUtil.buildOSSCompleteUri("oss://bucket/path", conf)); org.apache.hadoop.fs.FileSystem fs = org.apache.hadoop.fs.FileSystem.get(outputPath.toUri(), conf); if (fs.exists(outputPath)) {

fs.delete(outputPath, true);

}

如何清理已经完成作业的日志数据?

问题描述:集群的HDFS容量被写满,发现/spark-history下有大量的数据。

解决方法:

- 1. 在Spark**配置**页面的**服务配置**区域,查看是否有spark_history_fs_cleaner_enabled参数:
 - 是:修改参数值为true,可以周期性清理已经完成的作业的日志数据。
- 否:在spark-defaults页签下,单击自定义配置,新增spark_history_fs_cleaner_enabled为true。
- 2. 单击右上角的操作 > 重启 All Components
- 3. 在执行集群操作对话框,输入执行原因,单击确定。
- 4. 在弹出的确认对话框中,单击确定。

为什么作业运行较慢?

问题分析:作业的JVM Heap Size设置过小,可能会引起长时间的GC,影响作业性能。 解决方法:

- Tez: 在阿里云E-MapReduce控制台, Tez服务的配置页面, 可以调大Hive参数hive.tez.java.opts的值。
- Spark: 在阿里云E-MapReduce控制台, YARN服务的配置页面, 可以调大spark.execut.or.memory或spark.driver.memory的值。

• Mapreduce: 在阿里云E-MapReduce控制台, YARN服务的配置页面,可以调大mapreduce.map.java.opts或mapreduce.reduce.java.opts的值。

为什么AppMaster调度启动Task的时间过长?

问题分析:作业Task数目过多或Spark Executor数目过多,导致AppMasten调度启动Task的时间过长,单个Task运行时间较短,作业调度的Overhead较大。 解决方法:

- 减少Task数目,使用CombinedInputFormat。
- 提高前序作业产出数据的Block Size(dfs.blocksize)。
- 提高mapreduce.input.fileinputformat.split.maxsize。
- 对于Spark作业,在阿里云E-MapReduce控制台,Spark服务的配置页面,调节spark.executor.instances减少Executor数目,或者调节spark.default.parallelism降低并发数。

Spark作业报错 "java.lang.IllegalArgumentException: Size exceeds Integer.MAX_VALUE"

在Shuffle时,Partition数量过少使得Block Size超过Integer.MAX_VALUE最大值。您可以尝试增大Partition数目,在阿里云E-MapReduce控制台,YARN服务的配置页面,调 大spark.default.parallelism和spark.sql.shuffle.partitions,或者在Shuffle前执行Repartition。

E-MapReduce是否提供实时计算的功能?

E-MapReduce提供Spark Streaming、Storm和Flink三种实时计算服务。

导入RDS数据至EMR时,时间字段显示延迟8小时如何处理?

问题描述:

1. 例如,在云数据库RDS数据源中,数据表Test_Table中包含时间戳(TIMESTAMP)字段。

| 单行 | 详情 | ③ 新 | 建 🥥 删除 | ◎ 提交修改 | → 导出数据 ▼ | 🔕 生成报表 | 【表格数据可以编辑】 |
|----|----|------|------------|-------------|----------|--------|------------|
| | | id 🤻 | applied_at | | - | | |
| 1 | | 1 | 2018-12-2 | 21 12:15:09 | | | |
| 2 | | 2 | 2018-12-2 | 21 12:17:22 | | | |
| 3 | | 3 | 2018-12-2 | 21 12:17:22 | | | |
| 4 | | 4 | 2018-12-2 | 21 12:17:22 | | | |
| 5 | | 5 | 2018-12-2 | 21 12:17:23 | | | |
| | | | | | | | |

- 2. 您可以执行以下命令,导入Test_Table中的数据至HDFS。
 - sqoop import \

--connect jdbc:mysql://rm-2ze****341.mysql.rds.aliyuncs.com:3306/s***o_sqoopp_db \ --username s***o \ --password ****** --table play_evolutions \ --target-dir /user/hadoop/output \

- --delete-target-dir \
- --direct \setminus

--split-by id \setminus

- --fields-terminated-by '|' \
- -m 1

3. 查询导入结果。

查询结果显示,源数据的时间字段显示延迟8小时。

解决方法:在使用TIMESTAMP字段导入数据至HDFS时,请删除--direct参数。

```
sqoop import \
```

--connect jdbc:mysql://rm-2ze****341.mysql.rds.aliyuncs.com:3306/s***o_sqoopp_db \

- --username s***o \ --password ****** \
- --table play_evolutions \
- --target-dir /user/hadoop/output \
- --delete-target-dir \
- --split-by id \ --fields-terminated-by '|' \

```
-m 1
```

查询结果显示正常。

| [root@emr-header-1 ~]# hadoop fs -cat /user/hadoop/output1/part-m-00000 | | | | | |
|---|------------|-----------|--------------|--|--|
| 1 a1 | 2018-12-21 | 12:15:09. | b1 c1 d1 f1 | | |
| 21a2 | 2018-12-21 | 12:17:22. | lb21c21d21f2 | | |
| 31a3 | 2018-12-21 | 12:17:22. | lb3lc3ld3lf3 | | |
| 41a1 | 2018-12-21 | 12:17:22. | 1b41c41d41f4 | | |
| 51a1 | 2018-12-21 | 12:17:23. | lb51c51d51f5 | | |

如何修改Spark服务的spark-env配置?

登录集群的Header节点,修改/etc/ecm/spark-conf/spark-env.sh和/var/lib/ecm-agent/cache/ecm/service/SPARK/<版本号>/package/templates/spark-env.sh中的配 置.

```
⑦ 说明 如果您在Worker节点提交任务,则需要同步修改Worker节点相关配置。
```

作业参数传递至脚本文件该如何处理?

在Hive作业中,您可以通过-hivevar选项,传递作业中配置的参至脚本中。

```
1. 准备脚本文件。
```

```
脚本文件中引用变量的方式为 ${varname} (例如 ${rating} )。本示例中脚本的相关信息如下:
```

- 脚本名称: hivesql.hive
- 脚本的OSS路径: oss://bucket_name/path/to/hivesql.hive
- ∘ 脚本内容

```
use default;
drop table demo;
create table demo (userid int, username string, rating int);
insert into demo values(100,"john",3),(200,"tom",4);
select * from demo where rating=${rating};
```

```
2. 登录阿里云E-MapReduce控制台。
```

- 3. 在顶部菜单栏处,根据实际情况选择地域和资源组。
- 4. 单击上方的**数据开发**页签。
- 5. 在**项目列表**页面,单击对应项目所在行的**作业编辑**。
- 6. 在**作业编辑**区域,右键单击需要操作的文件夹,选择**新建作业**。
- 7. 输入**作业名称和作业描述**,选择*Hive*作业类型。

```
8. 配置作业。
```

E-MapReduce

| i. 在 基础设置 页面,设置参数的 Key和Value ,其中Key为脚本文件中的变量名,必须与脚本一致,例如rating。 | | | | | | |
|--|--|------------------------------|---------------------|--|--|--|
| 配置参数 🕜 🕴 🕇 | | | | | | |
| | | | | | | |
| 2519 夢致は: rating | 3 | | | | | |
| ii. 作业内容中必须添加 -hivevar 选项,以便 | 传递作业中配置的参数值至脚本变量。 | | | | | |
| -hivevar rating=\${rating} -f ossref://bucket | _name/path/to/hivesql.hive | | | | | |
| 9. 执行作业。 | | | | | | |
| 本示例执行结果如下。 | | | | | | |
| 作业实例信息 提交日志 YARN容器列表 | | | | | | |
| | | | | | | |
| Stage-4 is selected by condition resolver. | | | | | | |
| Stage-3 is filtered out by condition resolver | | | | | | |
| Moving data to directory hdfs://emr-header-1. | clu /user/hive/warehouse/demo/.hive-s | taging_hive_2019-09-11_ | | | | |
| Loading data to table default.demo | | | | | | |
| MapReduce Jobs Launched: Stage-Stage-1: Map: 1 Cumulative CPU: 2.45 | sec HDFS Read: 5572 HDFS Write: 89 SUCCESS | | | | | |
| Total MapReduce CPU Time Spent: 2 seconds 450 col1 col2 col3 | msec | | | | | |
| OK Time taken: 16.619 seconds | | | | | | |
| OK demo userid demo username demo rating | | | | | | |
| 19/09/11 10:38:12 INFO Configuration.deprecat | ion: mapred.input.dir is deprecated. Instead, use ma | preduce.input.fileinput | | | | |
| 19/09/11 10:38:12 INFO mapred.FileInputFormat | : Total input files to process : 1 | | | | | |
| 100 john 3 Time taken: 0.192 seconds, Fetched: 1 row(s) | | | | | | |
| JOB OUTPUT END | | | | | | |
| | | | | | | |
| 如何设置HiveServer2的认证方式为L | DAP? | | | | | |
| 1. 登录阿里云E-MapReduce控制台。 | | | | | | |
| 2. 在顶部菜单栏处,根据实际情况选择地域和资源 | 〔组。 | | | | | |
| 3. 単击上万的集群管理负益。 4. 左集群管理页面 单土相应集群所左行的详持 | | | | | | |
| 4. 在集研目違贝面, 半面相应集研剂在11的序角。 5. 在左侧导航栏, 单击集群服务 > Hive。 | , | | | | | |
| 6. 新增LDAP认证配置项并重启HiveServer2。 | | | | | | |
| i. 单击 配置 页签,在 服务配置 区域单击hives | server2-site。 | | | | | |
| < 返回 正常 🔞 Hive - 当前集群: | CHadoop | ● 宣看操作历史 | 唐链接 × 日 操作 × | | | |
| | | | | | | |
| 状态 部署拓扑 配置 配置修改历史 | | | | | | |
| 配置过滤 | 記置过滤 服务配置 ② 部署客户编述 | | | | | |
| 配置搜索: | 全部 hive-site hiveserver2-site hive-env | hivemetastore-site | 自定义配置 | | | |
| 靖輸入 Q | | | | | | |
| 配置范围: | hive.service.metrics.file.frequency | 30000 | | | | |
| 集群默认配置 ⋎ | hive.server2.metrics.enabled | true | | | | |
| 配置类型 | hive.server2.session.check.interval | 1h | | | | |
| 基础配置 高级配置 只读配置 | hive.server2.idle.operation.timeout | 6h | | | | |
| 数据路径目志路径目志相关 | hive.server2.logging.operation.enabled | true | 0 | | | |
| | | <i>C</i> b | | | | |
| | nive.server2.idle.session.timeout | on | | | | |
| | hive.service.metrics.file.location | /tmp/hiveserver2_metric.json | | | | |

hive.server2.enable.impersonation true

每页显示: 20 50 100 全部 (1) 共8条

URL或URI

磁盘相关 网络相关 文件路径

ii. 单击自定义配置。

LDAP认证方式需要新增如下三个配置项。

| 配置项 | 值 | 描述 |
|---|---|--|
| hive.server2.authentication | LDAP | 认证方式。 |
| hive.server2.authentication.ldap.url | 格式为ldap://\$[emr-header-1-hostname]:10389 | \$femr-header-1-hostname]是您实际的主机名称,您可 以在集群的emr-header-1上执行 hostname 命令获 取,连接主机的步骤请参见 登录集群 。 |
| hive.server2.authentication.ldap.baseDN | ou=people,o=emr | 无 |

- iii. 完成上述参数配置后,单击右上方的**保存**。
- iv. 在**确认修改**对话框中,配置各项参数,单击**确定**。
- v. 单击右上方的操作 > 重启 HiveServer2。
- vi. 在执行集群操作对话框中,配置各项参数,单击确定。
- vii. 在**确认**对话框中,单击**确定**。
- 7. 在LDAP中添加账号。

在E-MapReduce集群中,OpenLDAP组件是LDAP的服务,默认用于管理Knox的用户账号,HiveServer2的LDAP认证方式可以复用Knox的账号体系。添加账号的方法请参 见用户管理。

本示例新增账号为emr-guest。

8. 测试新增账号是否可以正常登录HiveServer2。

通过/usr/lib/hive-current/bin/beeline登录HiveServer2,正常登录情况如下。

beeline> !connect jdbc:hive2://emr-header-1:10000/ Enter username for jdbc:hive2://emr-header-1:10000/: emr-guest Enter password for jdbc:hive2://emr-header-1:10000/: emr-guest-pwd Transaction isolation: TRANSACTION_REPEATABLE_READ

如果账号密码不正确,则会显示如下异常。

Error: Could not open client transport with JDBC Uri: jdbc:hive2://emr-header-1:10000/: Peer indicated failure: Error validating the login (state=08S01,code=0)

如何使用阿里云E-MapReduce HDFS的Balancer功能以及参数调优?

1. 登录待配置集群任意节点。

- 2. 执行以下命令, 切换到hdfs用户并执行Balancer参数。
 - su hdfs

/usr/lib/hadoop-current/sbin/start-balancer.sh -threshold 10

- 3. 执行以下命令,查看Balancer运行情况:
 - 。 方式一

less /var/log/hadoop-hdfs/hadoop-hdfs-balancer-emr-header-xx.cluster-xxx.log

。 方式二

 $tailf\/var/log/hadoop-hdfs/hadoop-hdfs-balancer-emr-header-xx.cluster-xxx.log$

⑦ 说明 当提示信息包含 Successfully 字样时,表示执行成功。

Balancer的主要参数。

| 参数 | 说明 |
|---|---|
| Threshold | 默认值为10%,表示上下浮动10%。 当集群总使用率较高时,需要调小Threshold,避免阈值过高。 当集群新增节点较多时,您可以适当增加Threshold,使数据从高使用率节点移向低使用率 节点。 |
| dfs.datanode.balance.max.concurrent.moves | 默认值为5。 指定DataNode节点并发移动的最大个数。通常考虑和磁盘数匹配,推荐在DataNode端设 置为 4*磁盘数 作为上限,可以使用Balancer的值进行调节。 例如: 一个DataNode有28块盘,在Balancer端设置为28, DataNode端设置为 28*4 。具 体使用时根据集群负载适当调整。在负载较低时,增加concurrent数;在负载较高时,减少 concurrent数。 ⑦ 说明 DataNode端需要重启来刷新配置。 |
| dfs.balancer.dispatcherT hreads | Balancer在移动Block之前,每次迭代时查询出一个Block列表,分发给Mover线程使用。 ⑦ 说明 dispatcherThreads是该分发线程的个数,默认为200。 |

E-MapReduce

| 参数 | 说明 |
|---------------------------------------|---|
| dfs.balancer.rpc.per.sec | 默认值为20,即每秒发送的rpc数量为20。 因为分发线程调用大量getBlocks的rpc查询,所以为了避免NameNode由于分发线程压力过 大,需要控制分发线程rpc的发送速度。 例如,您可以在负载高的集群调整参数值,减小10或者5,对整体移动进度不会产生特别大 的影响。 |
| dfs.balancer.getBlocks.size | Balancer会在移动Block前,每次迭代时查询出一个Block列表,给Mover线程使用,默认 Block列表中Block的大小为2GB。因为getBlocks过程会对RPC进行加锁,所以您可以根据 NameNode压力进行调整。 |
| dfs.balancer.moverThreads | 默认值为1000。 Balancer处理移动Block的线程数,每个Block移动时会使用一个线程。 |
| dfs.namenode.balancer.request.standby | 默认值为false。 Balancer是否在Standby NameNode上查询要移动的Block。因为此类查询会对NameNode 加锁,导致写文件时间较长,所以HA集群开启后只会在Standby NameNode上进行查询。 |
| dfs.balancer.getBlocks.min-block-size | Balancer查询需要移动的参数时,对于较小Block(默认10 MB)移动效率较低,可以通过此 参数过滤较小的Block,增加查询效率。 |
| dfs.balancer.max-iteration-time | 默认值为1200000,单位毫秒。 Balancer一次迭代的最长时间,超过后将进入下一次迭代。 |
| dfs.balancer.block-move.timeout | 默认值为0,单位毫秒。 Balancer在移动Block时,会出现由于个别数据块没有完成而导致迭代较长的情况,您可以通 过此参数对移动长尾进行控制。 |

DataNode的主要参数。

| 参数 | 说明 |
|---|--|
| dfs.datanode.balance.bandwidthPerSec | 指定DataNode用于Balancer的带宽,通常推荐设置为100 MB/s,您也可以通过dfsadmin - setBalancerBandwidth 参数进行适当调整,无需重启DataNode。 例如,在负载低时,增加Balancer的带宽。在负载高时,减少Balancer的带宽。 |
| dfs.datanode.balance.max.concurrent.moves | 指定DataNode上同时用于Balancer待移动Block的最大线程个数。 |

如何使用standlone模式提交Spark任务?

E-MapReduce默认使用Spark on Yarn模式,暂不支持standlone模式。

如果E-MapReduce控制台上没有自定义配置选项,该如何处理?

1. 登录集群的Master节点,详情请参见<mark>登录集群</mark>。

2. 进入配置模板的目录。

cd /var/lib/ecm-agent/cache/ecm/service/HUE/4.4.0.3.1/package/templates/

| | [root@emr-header- [root@emr-header- total 44 | <pre>1 templates]# d 1 templates]# 1 root 42826 Tul</pre> | d /var/lib/ecm-agent/ 1 | /cache/ecm/service | /HUE/4.4.0.3.1/pac} | age/templates/ | |
|-----------------------------|--|---|----------------------------|--------------------|---------------------|----------------|--------|
| 4 | 本 1 1 1 1 1 1 1 1 0 0 0 0 0 0 0 0 0 0 0 | : | 22 10.30 <u>Inc.1111</u> | | | | |
| 0 | HUE 表示服务的 | 目录。 | | | | | |
| 0 | 4.4.0.3.1 为Hue的 | 的版本。 | | | | | |
| 0 | hue.ini 为配置文 | 件。 | | | | | |
| 3. 抈 | 1行以下命令,添加 | 您需要的配置。 | | | | | |
| | vim hue.ini | | | | | | |
| 뇔 | 自配置项已存在时 <i>,</i> | 您可以根据时间 | 青况修改参数值。 | | | | |
| 4. 荏 | EE-MapReduce控制 | 1台,重启服务以 | 生效配置。 | | | | |
| 使用 ^{问题现} | 数据开发提交 ^{1象如下图所示。} | 的作业一直统 | 处于Submit状态, | 该如何处理? | | | |
| 执 | 行集群 | 作业类型 | 作业提交节点主机 | | 开始时间 11 | 作业完成时间 11 | 执行时长 |
| C- | 53D4 | SPARK | emr-header-2.clust | | 2021-06-18 03:15:43 | | 11分27秒 |

出现此问题,通常是因为EMRFLOW中部分组件状态错误,您需要在控制台重启状态错误的组件。

1. 进入EMRFLOW页面。

操作

详情

执行状态 🏾

i. 进入任意服务页面,修改访问链接后的服务名为EMRFLOW。

| https://emr.con | isole.aliyun.com/?sp | which the state of the state of the | cn-hangzhou/cluster/C- | /service/EMRFL0 | ov |
|-----------------|--|-------------------------------------|------------------------|--|-----|
| 1.00.1 | 19 - 1 - 1 - 1 - 1 - 1 - 1 - 1 - 1 - 1 - | A DECK A DECK A AM A | IN A DOUBLE A DO | 1.44 | |
| 10.000 | A POINT OF A P | 12.000 | | | |
| uce II 概览 | 聶 朱群答理 ○ 操作历史 ■ 事件列表 山 数据开发 ○ 监控大盘 Beta | | | ◎ 系统管理 | ~ 6 |
| 首页 → | 集群管理 → 集群 (C-2D → 服务 → HDFS | | | | |
| < 返回 | When the second seco | | C | · 查看操作历史 · · · · · · · · · · · · · · · · · · · | · ß |
| | | | | | |

⑦ 说明 本示例是先进入HDFS服务页面。

ii. 单击**部署拓扑**页签。

2. 启动组件。

i. 在**部署拓扑**页签,单击组件处于STOPPED状态操作列的启动。

| <返回 EMRFLOW > | | | | | ◎ 查看操作 | F历史 1 | ♪ 快捷遊接 ~ ♪ 操作 ~ |
|-----------------------|-----------|---------|--------------|----------------|--------------|----------------------|------------------------|
| 状态 部署 在 配置 配置 | 修改历史 | | | | | | |
| 组件名: | 服务名: | ECS ID: | 主机名: | | <u>査询</u> 重置 | | |
| 组件名 ↓↑ | 组件状态↓↑ ₽ | 服务名 | ECS ID 11 | 主机名↓↑ | 主机角色 ↓↑ | IP | 操作 |
| Flow Agent Init | INSTALLED | EmrFlow | i-bp1jdil53l | emr-worker-4 📮 | TASK | 内网:192.1 | 配置 |
| Flow Agent Init | INSTALLED | EmrFlow | i-bp1jdil53l | emr-worker-3 📮 | TASK | 内网:192.1 | 配置 |
| Flow Agent Job Server | STARTED | EmrFlow | i-bp1idgirb | emr-header-1 | MASTER | 内网:192.1 外网:116.6 | 重启 停止 配置 |
| Flow Agent Init | INSTALLED | EmrFlow | i-bp1cjcem | emr-worker-1 📮 | CORE | 内网:192.1 | 配置 |
| Flow Agent Init | INSTALLED | EmrFlow | i-bp1cjcem | emr-worker-2 📮 | CORE | 内网:192.1 | 配置 |
| Flow Agent Init | INSTALLED | EmrFlow | i-bp1idgirb | emr-header-1 📮 | MASTER | 内网:192.1 外网:116.6 | 配置 |
| Emr Meta Command | STARTED | EmrFlow | i-bp1idgirb | emr-header-1 📮 | MASTER | 内网:192.1 外网:116.6 | 重启 停止 配置 |
| Flow Agent Daemon | • STOPPED | EmrFlow | i-bp1idgirb | emr-header-1 | MASTER | 内网:192.1 外网:116.6 | 停止 启动 配置 |

ii. 在**执行集群操作**对话中,输入**执行原因**,单击**确定**。

iii. 在**确认**对话中,单击**确定**。

3. 查看日志信息,检查组件是否启动。

i. 单击上方的查看操作历史。

15012994

ii. 在操作历史对话框中,单击操作类型列的Start EMRFLOW FlowAgentDaemon。

START_FlowAge

| | += /= | | | | | | | | | |
|------|-----------------------|-----------------------------------|---------------------|-------|---------|-------|----|----|-------|----|
| | 操"F历史 | | | | | | | | | × |
| | | | | | | | | | | 刷新 |
| | ID | 操作类型 | 开始时间 | 耗时(s) | 状态 | 进度(%) | 备注 | | 管理 | |
| | 171258 | Start EMRFLOW FlowAge ntDaemon | 2021-06-18 11:16:09 | 4 | ⊘ 成功 | 100 | 1 | | 终止 | |
| iii. | 单击 主机名 列的emr-h | neader-1。 | | | | | | | | |
| | 操作历史 | | | | | | | | | × |
| | 攝作ID: 171258 > 主机列表 | | | | | | | | | 刷新 |
| | Host ID | | 主机名 | | 状 | 5 | | | 进度(%) | |
| | 362877 | | emr-header-1 | | \odot | 成功 | | | 100 | |
| iv. | 单击 任务名 列的STAR | T_FlowAgent Daem | on_ON_emr-header-1。 | | | | | | | |
| | 操作历史 | | | | | | | | | × |
| | 操作ID: 171258 > 主机列表 | > 任务列表 | | | | | | | | 刷新 |
| | Task ID | 任务名 | | | | | | 状态 | | |

⊘ 成功

v. 当任务日志区域,提示如下图所示时表示组件启动成功。

| p |
|---|
| Fri, 18 Jun 2021 11:16:11 config_util.py[line:170] INFO step6 compare_write_and_move move remove tmp_file=/etc/ecm/flow-agent-conf/flow-agent.conf.tmp.1623986171, maybe md5sum are same skip |
| Fri, 18 Jun 2021 11:16:11 comment4gevent.py[line:14] INFO cmd=source /etc/profile.d/ecm_env.sh |
| Fri, 18 Jun 2021 11:16:11 comment4gevent.py[line:22] INFO returncode=0 |
| Fri, 18 Jun 2021 11:16:11 comment4gevent.py[line:23] INFO stdout=,stderr= |
| Fri, 18 Jun 2021 11:16:11 comment4gevent.py[line:14] INFO cmd=su -l root -c "mkdir -p /etc/ecm/flow-agent-conf/security" |
| Fri, 18 Jun 2021 11:16:11 comment4gevent.py[line:22] INFO returncode=0 |
| Fri, 18 Jun 2021 11:16:11 comment4gevent.py[line:23] INFO stdout=,stderr= |
| Fri, 18 Jun 2021 11:16:11 comment4gevent.py[line:14] INFO cmd=source /etc/profile.d/ecm_env.sh |
| Fri, 18 Jun 2021 11:16:11 comment4gevent.py[line:22] IMFO returncode=0 |
| Fri, 18 Jun 2021 11:16:11 comment4gevent.py[line:23] INFO stdout=,stderr= |
| Fri, 18 Jun 2021 11:16:11 comment4gevent.py[line:14] INFO cmd=su -l root -c "chown flowagent:hadoop -R /etc/ecm/flow-agent-conf" |
| Fri, 18 Jun 2021 11:16:11 comment4gevent.py[line:22] INFO returncode=0 |
| Fri, 18 Jun 2021 11:16:11 comment4gevent.py[line:23] INFO stdout=,stderr= |
| Fri, 18 Jun 2021 11:16:11 comment4gevent.py[line:14] INFO cmd=source /etc/profile.d/ecm_env.sh |
| Fri, 18 Jun 2021 11:16:11 comment4gevent.py[line:22] INFO returncode=0 |
| Fri, 18 Jun 2021 11:16:11 comment4gevent.py[line:23] INFO stdout=,stderr= |
| Fri, 18 Jun 2021 11:16:11 comment4gevent.py[line:14] INFO_cmd=su -1 flowagent -c "/usr/lib/flow-agent-current/sbin/flow-agentd start" |
| Fri, 18 Jun 2021 11:16:11 comment4gevent.py[line:22] INFC returncode=0 |
| Fri, 18 Jun 2021 11:16:11 comment4gevent.py[line:23] DHFO stdout=Started flow-agentd, Logging at /mnt/disk1/log/flow-agent/flow-agentd.out, PID=6549. |
| ,stderr= |

⑦ 说明 组件启动后.如果还有报错,请根据日志信息排查并处理。如果报错信息提示权限问题,您可以先通过SSH方式登录集群,执行命令 sudo chown flo wagent:hadoop /mnt/disk1/log/flow-agent/*处理,然后按照上述步骤重新操作以重启状态错误的组件。