

Kafka实验

1.创建集群

注意选择流式集群以及勾选Kafka组件即可，并绑定公网IP。

集群节点	节点类型 ?	计费模式	实例规格	实例数量
	Master节点 ?	按需计费	通用计算增强型 16 vCPUs 64 GB c6.4xlarge.4 系统盘 高IO 480 GB x 1 数据盘 高IO 600 GB x 1	<div><div></div><div>2</div><div></div></div>
	分析Core节点 ?	按需计费	通用计算增强型 8 vCPUs 32 GB c6.2xlarge.4 系统盘 高IO 480 GB x 1 数据盘 高IO 600 GB x 1	<div><div></div><div>3</div><div></div></div>

Kerberos认证

☐ ?

用户名

root/admin

密码

该密码同时用于远程登录ECS机器或BMS机器和集群管理页面。远程登录ECS机器的用户名为root，集群管理页面的用户名为admin。

确认密码

通信安全授权

☒ 确认授权

授权MRS集群创建和切换子网时开通相应的访问控制规则，从而使用户可以通过MRS管理控制台进行大数据组件部署和后续集群的使用、运维和管理等操作，此时不授权将无法创建集群。[了解更多](#)

需要开通的访问控制规则

2.安装客户端

进入Manager界面



详细信息



客户端生成成功，保存在服务器（192.168.0.3）如下路径：**/tmp/FusionInsight-Client/FusionInsight_Cluster_1_Services_Client.tar**。使用后请及时删除。

关闭

登录对应的服务器（192.168.0.3），安装客户端。

输入指令：

```
tar -xvf FusionInsight_Cluster_1_Services_Client.tar
sha256sum -c FusionInsight_Cluster_1_Services_ClientConfig.tar.sha256
tar -xvf FusionInsight_Cluster_1_Services_ClientConfig.tar
cd /opt/Bigdata/client/FusionInsight_Cluster_1_Services_ClientConfig
./install.sh /opt/hadoopclient #安装到/opt/hadoopclient目录
cd /opt/hadoopclient
source bigdata_env
```

最终效果：

```
[22-10-18 20:07:26]: Install KrbClient begin ...
[22-10-18 20:07:26]: Copy /tmp/FusionInsight-Client/FusionInsight_Cluster_1_Services_ClientConfig/KrbClient/FusionInsight-Client-1.17.tar.gz to /opt/client/KrbClient.
[22-10-18 20:07:26]: Copy /tmp/FusionInsight-Client/FusionInsight_Cluster_1_Services_ClientConfig/KrbClient/FusionInsight-Client-1.17.tar.gz to /opt/client/KrbClient.
[22-10-18 20:07:26]: Copy KRB config files to "/opt/client/KrbClient/kerberos/conf"
[22-10-18 20:07:26]: Copy security script files to "/opt/client/KrbClient/kerberos/bin"
[22-10-18 20:07:26]: Create KRB env file "/opt/client/KrbClient/component_env".
[22-10-18 20:07:26]: KrbClient installation is complete.
[22-10-18 20:07:26]: Install Presto begin ...
[22-10-18 20:07:26]: /opt/client/Presto
[22-10-18 20:07:26]: /tmp/FusionInsight-Client/FusionInsight_Cluster_1_Services_ClientConfig/Presto
[22-10-18 20:07:26]: Copy config files to "presto/conf"
[22-10-18 20:07:26]: cp: cannot stat '/tmp/FusionInsight-Client/FusionInsight_Cluster_1_Services_ClientConfig/Presto/lib': No such file or directory
[22-10-18 20:07:26]: Presto installation is complete.
[22-10-18 20:07:26]: Install Spark2x begin ...
[22-10-18 20:07:26]: Copy /tmp/FusionInsight-Client/FusionInsight_Cluster_1_Services_ClientConfig/Spark2x/FusionInsight-Client-1.17.tar.gz to /opt/client/Spark2x
[22-10-18 20:07:27]: Copy Hadoop config files to "/opt/client/Spark2x/spark/etc/hadoop"
[22-10-18 20:07:27]: Create Spark env file "/opt/client/Spark2x/component_env".
[22-10-18 20:07:27]: add SPARK_LOCAL_HOSTNAME env to spark-env.sh
[22-10-18 20:07:28][Hudi]: Start to install Hudi client.
[22-10-18 20:07:28][Hudi]: check and init success!
[22-10-18 20:07:28][Hudi]: complete copy "/opt/client/Spark2x/hudi" to "/opt/client/Hudi".
[22-10-18 20:07:28][Hudi]: complete copy "/opt/client/Spark2x/spark/conf" to "/opt/client/Hudi/hudi/conf".
[22-10-18 20:07:28][Hudi]: complete modify hudi conf.
[22-10-18 20:07:28][Hudi]: complete generate component_env.
[22-10-18 20:07:28][Hudi]: Successfully installed Hudi client.
[22-10-18 20:07:28]: Spark2x installation is complete.
[22-10-18 20:07:28]: Install Yarn begin ...
[22-10-18 20:07:28]: Copy Yarn config files to "/opt/client/Yarn/config"
[22-10-18 20:07:28]: Yarn installation is complete.
[22-10-18 20:07:28]: Install ZooKeeper begin ...
[22-10-18 20:07:28]: Copy /tmp/FusionInsight-Client/FusionInsight_Cluster_1_Services_ClientConfig/ZooKeeper/FusionInsight-Client-1.17.tar.gz to /opt/client/ZooKeeper.
[22-10-18 20:07:28]: /tmp/FusionInsight-Client/FusionInsight_Cluster_1_Services_ClientConfig/ZooKeeper
[22-10-18 20:07:28]: Create Zookeeper env file "/opt/client/ZooKeeper/component_env".
[22-10-18 20:07:28]: Copy zookeeper config files to /opt/client/ZooKeeper/zookeeper/conf
[22-10-18 20:07:28]: ZooKeeper installation is complete.
[22-10-18 20:07:28]: The component client is installed successfully
[22-10-18 20:07:28]: Persisting client information to the database...
[2022-10-18 20:07:28]: Platform=x86_64
[2022-10-18 20:07:28]: ClientIp=192.168.0.72
[2022-10-18 20:07:28]: InstallUser=root
[2022-10-18 20:07:28]: ClientPath=/opt/client
[2022-10-18 20:07:30]: Do Save client information to database Successfully.
[root@node-master1fxou FusionInsight_Cluster_1_Services_ClientConfig]# cd /opt/client/
[root@node-master1fxou client]# source bigdata_env
[root@node-master1fxou client]#
```

3. 登录Master节点

在“集群列表 > 现有集群”列表中，单击名称“Kafka_test”，在“节点管理”页签中找到类型为“Master1”的节点，并单击其名称，跳转至云服务器控制台上的该弹性云服务器详情页面。

节点管理

组件管理

告警管理

文件管理

作业管理

引导操作

弹性伸缩

标签管理

配置Task节点

节点操作

所有节点组

节点名称

请输入节点名称

Q

C

节点组名称	节点类型	节点数	操作									
master_node_default_group	Master	可用区2: 2	升级规格									
<div><input type="checkbox"/></div> 节点名称	IP	机架	操作状态	健康状态	CPU使用率	内存使用率	磁盘使用率	网络速度	规格名	规格	付费类型	可用区
<div><input checked="" type="checkbox"/></div> ★ node-master1M...	192.168.0.3	/default/rac...	<div>运行中</div>	<div>良好</div>	8% <div></div>	21.84% 49.04 GB/63.00 GB <div></div>	6.64% 990.00 GB/1060.42 C <div></div>	R: 24.4 Byte/s; W: 18.6	c6.4xlarge.4	16 vCPUs 64 GB 600 GB 高L...	按需计费	可用区2
<div><input type="checkbox"/></div> ★ node-master2p...	192.168.0.75	/default/rac...	<div>运行中</div>	<div>良好</div>	8% <div></div>	11.28% 55.67 GB/63.00 GB <div></div>	6.43% 992.21 GB/1060.42 C <div></div>	R: 15.2 Byte/s; W: 13.1	c6.4xlarge.4	16 vCPUs 64 GB 600 GB 高L...	按需计费	可用区2

节点组名称	节点类型	节点数	操作									
core_node_streaming_group	Streaming Core	可用区2: 3	扩容 缩容									
<div><input type="checkbox"/></div> 节点名称	IP	机架	操作状态	健康状态	CPU使用率	内存使用率	磁盘使用率	网络速度	规格名	规格	付费类型	可用区
<div><input type="checkbox"/></div> node-str-coreAHVK...	192.168.0.34	/default/rac...	<div>运行中</div>	<div>良好</div>	4% <div></div>	9.38% 28.31 GB/32.00 GB <div></div>	5.81% 998.76 GB/1060.42 C <div></div>	R: 4.2 Byte/s; W: 6.1 By	c6.2xlarge.4	8 vCPUs 32 GB 600 GB 高IO ...	按需计费	可用区2
<div><input type="checkbox"/></div> node-str-coreEvIG...	192.168.0.1...	/default/rac...	<div>运行中</div>	<div>良好</div>	5% <div></div>	9.61% 28.24 GB/32.00 GB <div></div>	5.81% 998.76 GB/1060.42 C <div></div>	R: 4.3 Byte/s; W: 5.9 By	c6.2xlarge.4	8 vCPUs 32 GB 600 GB 高IO ...	按需计费	可用区2
<div><input type="checkbox"/></div> node-str-coreJluH...	192.168.0.1...	/default/rac...	<div>运行中</div>	<div>良好</div>	5% <div></div>	9.70% 28.21 GB/32.00 GB <div></div>	5.82% 998.66 GB/1060.42 C <div></div>	R: 4.4 Byte/s; W: 6.5 By	c6.2xlarge.4	8 vCPUs 32 GB 600 GB 高IO ...	按需计费	可用区2

4.使用Kafka客户端创建topic

1. 在“集群列表 > 现有集群”列表中，单击名称“Kafka_test”，进入集群“概览”页面。

在“概览”页面单击“IAM用户同步”后的“同步”等待同步完成。

基本信息 [了解更多](#)

集群名称	mrs_Kafka_test ✎
集群状态	运行中
集群管理页面	? 前往 Manager ⇄
付费类型	按需计费
集群版本	MRS 3.1.0
集群类型	分析集群
集群ID	e4d2cf92-b5d0-4807-9b00-fd311e0e9e57
创建时间	2022/10/18 22:25:00 GMT+08:00
可用区	可用区2
虚拟私有云	vpc-default
默认生效子网	subnet-default 切换子网 ?
弹性公网IP	120.46.190.113 ? 添加安全组规则
OBS权限控制	? 单击管理
委托	? -- 管理委托
Kerberos认证	关闭
日志记录	? <input checked="" type="checkbox"/>
安全组	mrs_mrs_Kafka_test_hjac
IAM用户同步	? 已同步 同步
数据连接	? 单击管理
通信安全授权	? <input checked="" type="checkbox"/>

2. 选择“组件管理 > ZooKeeper > 实例”，查看ZooKeeper角色实例的IP地址。记录ZooKeeper角色实例中任意一个的IP地址即可。这里记下 192.168.0.3 。

服务 ZooKeeper / 实例

服务状态

实例

服务配置

更多

每30s刷新一次

C

所有角色

高级搜索

<input type="checkbox"/>	角色	主机名	管理IP	业务IP	机架	运行状态	配置状态
<input type="checkbox"/>	quorumpeer	node-str-coreGUuH.mrs-qhte.com	192.168.0.118	192.168.0.118	/default/rack0	良好	已同步
<input type="checkbox"/>	quorumpeer	node-master1Mb5N.mrs-qhte.com	192.168.0.3	192.168.0.3	/default/rack5916	良好	已同步
<input type="checkbox"/>	quorumpeer	node-master2p1OM.mrs-qhte.com	192.168.0.75	192.168.0.75	/default/rack5916	良好	已同步

3. 执行如下命令，创建kafka topic。
- ```
kafka-topics.sh --create --zookeeper <ZooKeeper角色实例所在节点IP:2181/kafka> --
partitions 3 --replication-factor 3 --topic <Topic名称>
```

相应的，我们输入 `sh kafka-topics.sh --create --zookeeper 192.168.0.3:2181/kafka --
partitions 3 --replication-factor 3 --topic topicTest`

```
[root@node-master1Mb5N bin]# sh kafka-topics.sh --create --topic topicTest --partitions 3 --replication-factor 3 --zookeeper 192
.168.0.3:2181/kafka
Created topic topicTest.
[root@node-master1Mb5N bin]#
```

输入 `sh kafka-topics.sh --list --zookeeper 192.168.0.3:2181/kafka` 查看当前topic

```
Created topic topicTest.
[root@node-master1Mb5N bin]# sh kafka-topics.sh --list --zookeeper 192.168.0.3:2181/kafka
_KafkaMetricReport
_consumer_offsets
_default_metrics
topicTest
[root@node-master1Mb5N bin]#
```

## 5.管理Kafka主题中的消息

1. 选择“组件管理 > Kafka > 实例”，查看Kafka角色实例的IP地址。记录Kafka角色实例中任意一个的IP地址。如 192.168.0.105 。

服务 Kafka / 实例

服务状态

实例

服务配置

更多

每30s刷新一次

C

所有角色

高级搜索

| <input type="checkbox"/> | 角色     | 主机名                            | 管理IP          | 业务IP          | 机架             | 运行状态 | 配置状态 |
|--------------------------|--------|--------------------------------|---------------|---------------|----------------|------|------|
| <input type="checkbox"/> | Broker | node-str-coreEvIG.mrs-qhte.com | 192.168.0.105 | 192.168.0.105 | /default/rack0 | 良好   | 已同步  |
| <input type="checkbox"/> | Broker | node-str-coreGUuH.mrs-qhte.com | 192.168.0.118 | 192.168.0.118 | /default/rack0 | 良好   | 已同步  |
| <input type="checkbox"/> | Broker | node-str-coreAHVK.mrs-qhte.com | 192.168.0.34  | 192.168.0.34  | /default/rack0 | 良好   | 已同步  |

2. 登录Master节点，在topic test中产生消息。首先执行命令 `kafka-console-producer.sh --
broker-list <Kafka角色实例所在节点IP:9092> --topic <Topic名称> --producer.config
/opt/hadoopclient/Kafka/kafka/config/producer.properties` 其中<Topic名称>为Step7创建的Topic名称。然后输入指定的内容作为生产者产生的消息，输入完成后按回车发送消息。如果需要结束产生消息，使用“Ctrl + C”退出任务。

输入 `kafka-console-producer.sh --broker-list 192.168.0.105:9092 --topic topicTest --producer.config /opt/hadoopclient/Kafka/kafka/config/producer.properties`

3. 消费topic test中的消息。输入 `kafka-console-consumer.sh --topic topicTest --bootstrap-server 192.168.0.105:9092 --consumer.config /opt/hadoopclient/Kafka/kafka/config/consumer.properties`

结果:

```
[root@node-master1MbsN bin]# kafka-console-producer.sh --broker-list 192.168.0.105:9092 --topic topicTest --producer.config /opt/hadoopclient/Kafka/kafka/config/producer.properties
[2022-10-19 00:53:31,868] WARN The configuration 'producer.type' was supplied but isn't a known config. (org.apache.kafka.clients.producer.ProducerConfig)
[2022-10-19 00:53:31,872] WARN The configuration 'serializer.class' was supplied but isn't a known config. (org.apache.kafka.clients.producer.ProducerConfig)
>hello
<hello
>^C[root@node-master1MbsN bin]#
[root@node-master1MbsN bin]# kafka-console-consumer.sh --topic topicTest --bootstrap-server 192.168.0.105:9092 --consumer.config /opt/hadoopclient/Kafka/kafka/config/consumer.properties
hello
hello
```

## 6. Python使用Kafka

准备Python环境，导入kafka-python包（这里使用2.0.1版本），注意绑定公网IP下载。

导入producer.py以及consumer.py文件，注意相应的本机信息。其中 `consumer_id` 可查找 `/opt/hadoopclient/Kafka/kafka/config/consumer.properties` 获取。

```
#!/usr/bin/python3

from kafka import KafkaProducer

conf = {
 'bootstrap_servers': ['192.168.0.105:9092'],
 'topic_name': 'topicTest',
}

print('start producer')
producer = KafkaProducer(bootstrap_servers=conf['bootstrap_servers'])

data = bytes("hello kafka!", encoding="utf-8")
producer.send(conf['topic_name'], data)
producer.close()
print('end producer')
```

```
#!/usr/bin/python3

from kafka import KafkaConsumer

conf = {
 'bootstrap_servers': ['192.168.0.105:9092'],
 'topic_name': 'topicTest',
 'consumer_id': 'example-group1'
}

print('start consumer')
consumer = KafkaConsumer(conf['topic_name'],
 bootstrap_servers=conf['bootstrap_servers'],
 group_id=conf['consumer_id'])

for message in consumer:
 print("%s: value=%s" % (message.topic, message.value))

print('end consumer')y运行producer.py
```

依次运行producer.py和consumer.py，结果如图：

```
[root@node-master1MBSN pytest]# python producer.py
start producer
end producer
[root@node-master1MBSN pytest]# python consumer.py
start consumer
topicTest: value=hello kafka
topicTest: value=hello kafka
topicTest: value=hello kafka!
topicTest: value=hello kafka!
```

## 附：Kafka学习笔记

### 1.介绍

Kafka是Apache旗下的一款分布式流媒体平台，Kafka是一种高吞吐量、持久性、分布式的发布订阅的消息队列系统。它最初由LinkedIn(领英)公司发布，使用Scala语言编写，与2010年12月份开源，成为Apache的顶级子项目。它主要用于处理消费者规模网站中的所有动作流数据。动作指(网页浏览、搜索和其它用户行动所产生的数据)。

**活动流数据**是几乎所有站点在对其网站使用情况做报表时都要用到的数据中最常规的部分。活动数据包括页面访问量（Page View）、被查看内容方面的信息以及搜索情况等内容。这种数据通常的处理方式是先把各种活动以日志的形式写入某种文件，然后周期性地对这些文件进行统计分



析。**运营数据**指的是服务器的性能数据（CPU、IO 使用率、请求时间、服务日志等等数据）。运营数据的统计方法种类繁多。

## Kafka三大特点:

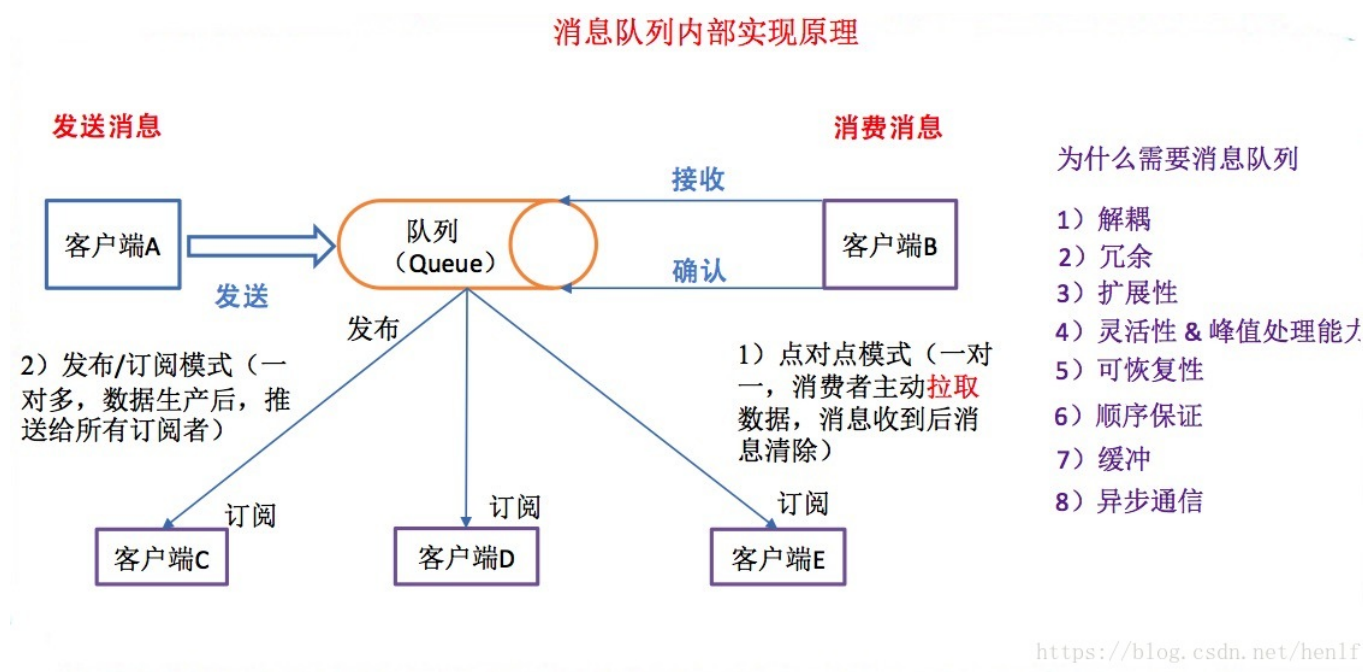
- 1.**高吞吐量**：可以满足每秒百万级别消息的生产和消费。
- 2.**持久性**：有一套完善的消息存储机制，确保数据高效安全且持久化。
- 3.**分布式**：基于分布式的扩展；Kafka的数据都会复制到几台服务器上，当某台故障失效时，生产者和消费者转而使用其它的Kafka。

## 1.关键概念

- topic kafka 把收到的消息按 topic 进行分类，因此可以理解为 topic 是一种类别
- producer 往 kafka 发送消息的用户
- consumer 接收 kafka 消息的用户
- broker kafka 集群可以由多个 kafka 实例组成，每个实例（server）称为 broker

无论是 kafka broker 本身，还是 producer 或者 consumer，都依赖于 zookeeper 集群保存一些 meta 信息，保证系统可用性，以及使用 zookeeper 的选举机制。

## 2.消息队列实现原理



**点对点模式** 一对一，**消费者主动拉取数据，消息收到后消息清除**。点对点模式通常是一个基于拉取或轮询的消息发送模型。此模型中，消费者从队列主动拉取信息，而不是消息系统推送消息给消费者，并且，**消息只能被一个且只有一个消费者接收处理**，即使有多个消息监听者也是如此。

**发布订阅模式** 一对多，数据生产后，**推送给所有订阅者**。发布订阅模型则是一个**基于推送**的消息传送模型。发布订阅模型可以有多种不同的订阅者，临时订阅者只在主动监听主题时才接收消息，而持久订阅者则监听主题的所有消息，即使当前订阅者不可用，处于离线状态。