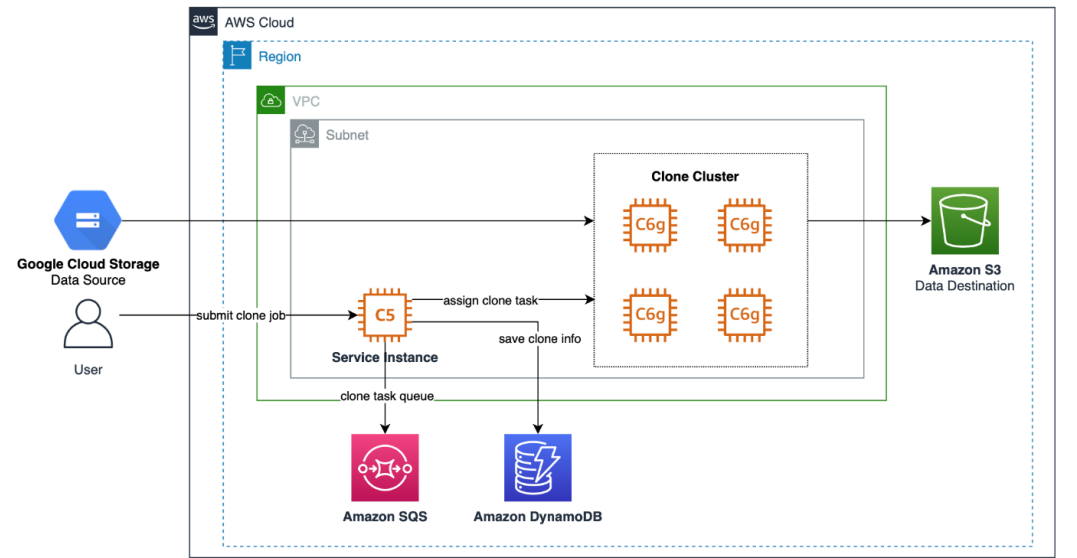


数据传输服务使用手册

本服务用于将Google GCS及其他云厂商（Rclone支持的云厂商）存储服务的数据传输到AWS S3上。可以让您批量启动和管理传输EC2实例。

架构设计

基于Google GCS为数据源的架构图如下：



组件作用：

Amazon SQS：传输任务下发后存储与排序。

Amazon DynamoDB：系统在初始化时会自动创建三张表，用于存储传输设备、任务状态表与配置。

部署目录

1. 创建VPC
2. 修改子网配置
3. 创建部署所需的IAM role
4. 创建SQS Queue
5. 验证邮箱
6. 传输EC2实例的AMI制作
7. 部署传输服务
8. 传输任务配置
9. 启动传输集群
10. 添加传输任务
11. 查看传输实例
12. 查看传输任务
13. 关闭传输集群
14. 清除资源

1. 创建VPC

1. 进入AWS 控制台，搜索栏中输入“VPC”，点击“创建VPC”按钮， 选择vpc与子网一同创建。

VPC > 您的 VPC > 创建 VPC

创建 VPC 信息

VPC 是由 AWS 对象(如 Amazon EC2 实例)填充的 AWS 云的隔离部分。

VPC 设置

要创建的资源 信息
仅创建 VPC 资源或创建 VPC、子网等。

☒ 仅 VPC ☐ VPC、子网等。

名称标签 – 可选
使用“Name”键和您指定的值创建一个标签。

IPv4 CIDR 块 信息
☒ IPv4 CIDR 手动输入
☐ IPAM 分配的 IPv4 CIDR 块 –

2. 使用默认配置即可，点击确认。

VPC / 高级 VPC / 创建 VPC / 创建 VPC 资源

Create VPC workflow

Creating VPC Resources
Thank you for using the new create VPC experience. Let us know what you think.

Associate route table 87%

详细信息

- 创建 VPC: vpc-071185e68c88d1151
- Disable DNS hostnames
- Enable DNS resolution
- Verifying VPC creation: vpc-071185e68c88d1151
- Create S3 endpoint: vpce-0249cd6277f744ba1
- Create subnet: subnet-0796fb0454687ab28
- Create subnet: subnet-06f23680369246dc9
- Create subnet: subnet-04584d1cfc5be5fac
- Create subnet: subnet-04d1e70549cae06c5
- Create internet gateway: igw-0c73a5f32621ec757
- Attach internet gateway to the VPC
- Create route table: rtb-05b5001fb86160105

2. 修改子网配置

1. 选中要修改的子网，点击操作“编辑子网”

search: vpc-071185e68c88d1151 清除筛选条件

	Name	子网 ID	状态	VPC
<input checked="" type="checkbox"/>	项目-subnet-private2-us-east-1b	subnet-04d1e70549cae06c5	Available	vpc-071185e68c88d1151
<input type="checkbox"/>	项目-subnet-public1-us-east-1a	subnet-0796fb0454687ab28	Available	vpc-071185e68c88d1151
<input type="checkbox"/>	项目-subnet-private1-us-east-1a	subnet-04584d1cfc5be5fac	Available	vpc-071185e68c88d1151
<input type="checkbox"/>	项目-subnet-public2-us-east-1b	subnet-06f23680369246dc9	Available	vpc-071185e68c88d1151

46dc9 / 项目-subnet-public2-us-east-1b

操作

创建流日志

编辑子网设置

编辑 VPC CIDR

子网 ARN
arn:aws:ec2:us-east-1:071737308255:subnet/subnet-06f23680369246dc9

IPv6 CIDR
-

VPC
vpc-071185e68c88d1151 | 项目-vpc

状态
Available

可用区
us-east-1b

路由表
rtb-05b5001fb86160105 | 项目-rtb-public2-us-east-1b

IPv4 CIDR
10.0.16.0/20

可用区 ID
use1-az2

网络 ACL
acl-035e570d4e

自动分配客户拥有的 IPv4 地址

编辑网络 ACL 关联

编辑路由表关联

编辑 CIDR 预留

共享子网

管理标签

删除

2. 勾选，自动分配IPv4



3. 创建部署所需的IAM ROLE

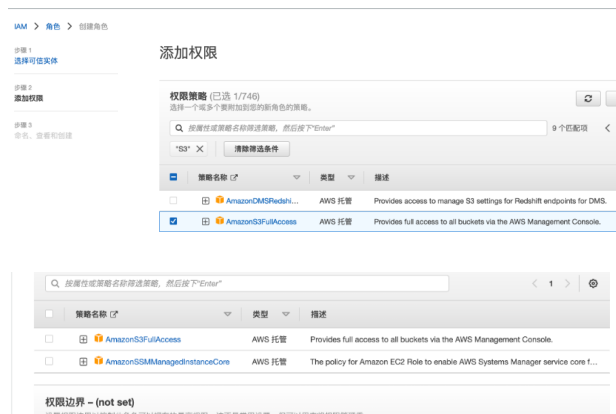
1. 打开AWS console 进入 IAM service



2. 点击左侧“角色”将页面切换到角色，点击右上角的“创建角色”按钮



3. 添加权限，此处需要添加S3的权限和AmazonSSMManagedInstanceCore的权限



4. 为角色命名:Rclone-EC2-Instance-Role，并点击右下角创建按钮(注意此角色名称必须为：Rclone-EC2-Instance-Role)

配置
设置最大消息大小、对其他使用者的可见性和消息保留周期。 [信息](#)

可见性超时 [信息](#)
 小时
应介于 0 秒至 12 小时之间。

消息保留周期 [信息](#)
 天
应介于 1 分钟至 14 天之间。

交付延迟 [信息](#)
 秒
应介于 0 秒至 15 分钟之间。

最大消息大小 [信息](#)
 KB
应介于 1 KB 和 256 KB 之间。

接收消息等待时间 [信息](#)
 秒
应介于 0 至 20 秒之间。

☒ **基于内容的重复数据删除**
已启用基于内容的重复数据删除时，消息重复数据删除 ID 为可选项。

- 配置完成后，点击页面右下角橙色按钮“创建队列”
- 队列创建完成之后，会生成URL，点击复制按钮即可，我们会在后面的应用页面中使用到该URL

Amazon SQS > 队列 > queue_name.fifo

queue_name.fifo [编辑](#) [删除](#) [清除](#) [发送和接收消息](#) [开始 DLQ 重新路由](#)

详细信息 [信息](#)

名称	类型	ARN
queue_name.fifo	FIFO	arn:aws:sqs:us-east-1:071737308255:queue_name.fifo
加密	URL	死信队列
已禁用	https://sq.us-east-1.amazonaws.com/071737308255/queue_name.fifo	-
属性		
已创建	最大消息大小	上次更新时间

5. 验证邮箱

邮箱用来接收数据copy过程中的异常告警，需要您提供两个邮箱，一个From，一个To（From和To可以为同一个邮箱地址），两个邮箱均需要在对应的区域的[Amazon Simple Email Service](#)服务中进行验证，以下为验证步骤：

- 进入AWS管理控制台，将区域切换到us-east-1
- 在上方的搜索栏中输入SES，点击进入[Amazon Simple Email Service](#)
- 点击右侧的橙色按钮“Create Identity”

Customer engagement

Amazon SES

Highly-scalable inbound and outbound email service

Amazon Simple Email Service (SES) is a cloud-based email service that provides cost-effective, flexible and scalable way for businesses of all sizes to keep in contact with their

Send your first email

Get started by creating and verifying a *sender identity* - a domain or email address you use to send email through Amazon SES.

[Create identity](#)

4. 选择Email Address

Amazon SES > Configuration: Verified identities > Create identity

Create identity

A *verified identity* is a domain, subdomain, or email address you use to send email through Amazon SES. Identity verification at the domain level extends to all email addresses under one verified domain identity.

Identity details [Info](#)

Identity type

☐ **Domain**
To verify ownership of a domain, you must have access to its DNS settings to add the necessary records.

☒ **Email address**
To verify ownership of an email address, you must have access to its inbox to open the verification email.

Email address

Email address can contain up to 320 characters, including plus signs (+), equals signs (=) and underscores (_).

- 点击右下方的“Create identity”按钮
- 在您填写的邮箱中，会收到一封邮件，格式如下图：

Amazon Web Services – Email Address Verification Request in region US East (N. Virginia)

Amazon Web Services <no-reply-aws@amazon.com> Today at 8:11
To: Jia, Ting

Dear Amazon Web Services Customer,

We have received a request to authorize this email address for use with Amazon SES and Amazon Pinpoint in region US East (N. Virginia). If you requested this verification, please go to the following URL to confirm that you are authorized to use this email address:

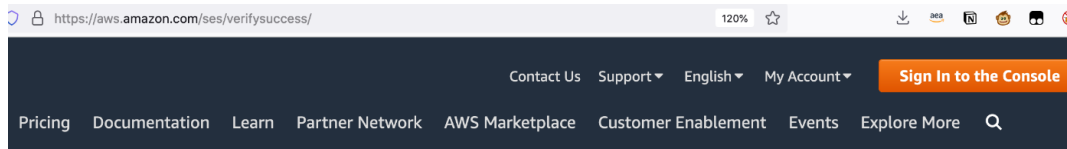
<https://email-verification.us-east-1.amazonaws.com/?Context=071737308255&X-Amz-Date=20200321T121925Z&Identity-Id=arn:aws:iam::40amazon.com:iam:User/AmazonSES-User&Identity-Id=arn:aws:iam::40amazon.com:iam:User/AmazonSES-User&Signature=af6c5280b3006144598e51a009a6222ebc2f7d67164e5270e1940a3e483c75d5>

Your request will not be processed unless you confirm the address using this URL. This link expires 24 hours after your original verification request.

If you did NOT request to verify this email address, do not click on the link. Please note that many times, the situation isn't a phishing attempt, but either a misunderstanding of how to use our service, or someone setting up email-sending capabilities on your behalf as part of a legitimate service, but without having communicated the procedure first. If you are still concerned, please forward this notification to aws-email-domain-verification@amazon.com and let us know if forward that you did not request the verification.

To learn more about sending email from Amazon Web Services, please refer to the Amazon SES Developer Guide at <http://docs.aws.amazon.com/ses/latest/DeveloperGuide/Welcome.html> and Amazon Pinpoint Developer Guide at <http://docs.aws.amazon.com/pinpoint/latest/userguide/welcome.html>.

7. 点击邮件中的 <https://email-verification.us-east-1.amazonaws.com/?Context=071737308255&X-Amz-Date=20200321T121925Z&Identity-Id=arn:aws:iam::40amazon.com:iam:User/AmazonSES-User&Identity-Id=arn:aws:iam::40amazon.com:iam:User/AmazonSES-User&Signature=af6c5280b3006144598e51a009a6222ebc2f7d67164e5270e1940a3e483c75d5> 开头的长链接，进行验证即可，点开链接后无需其他操作，页面如下图：



Congratulations!

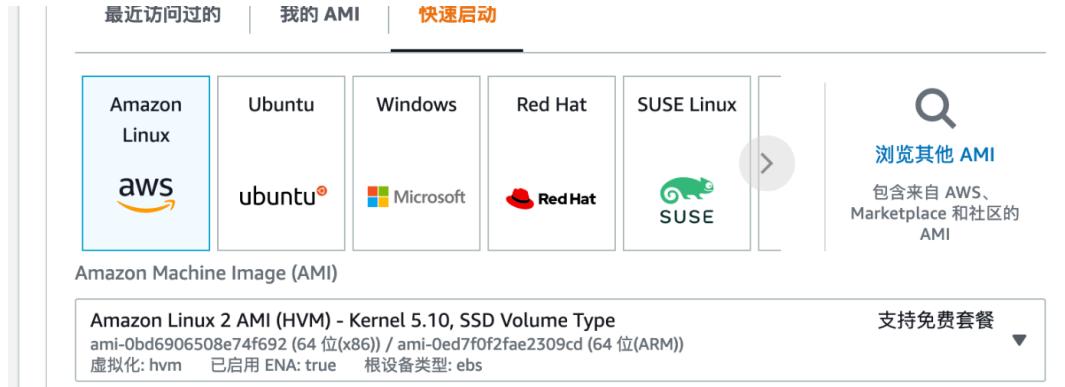
You have successfully verified an email address. You can now start sending email from this address.

For new Amazon SES users—If you have not yet applied for a sending limit increase, then you are still in the [sandbox environment](#), and you can only send email to addresses that have been verified. To verify a new email address or domain, see the **Identity Management** section of the [Amazon SES console](#).

6. 传输EC2实例的AMI制作

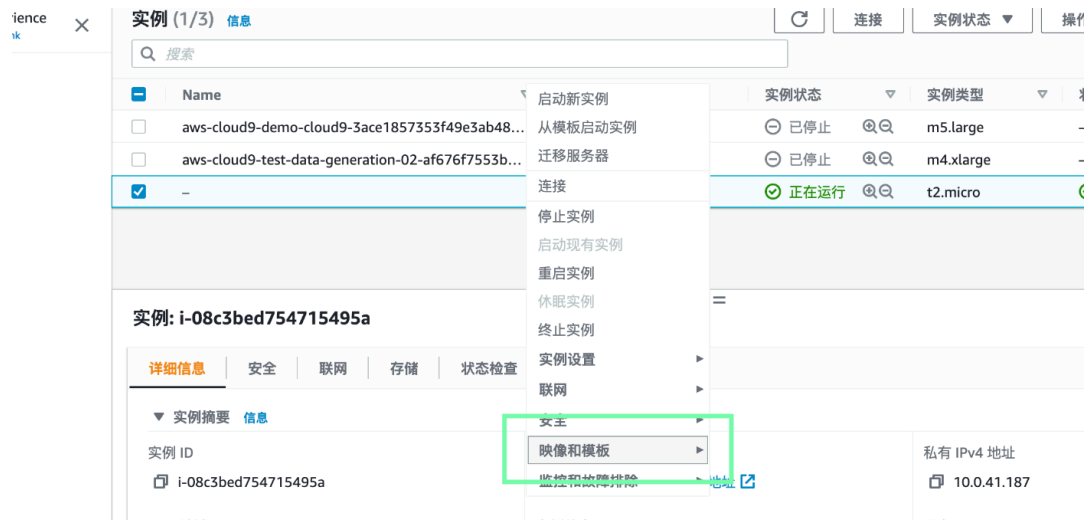
制作的AMI用来启动传输数据时用的EC2实例。

1. 使用系统默认的AMI启动机器



2. 将GCP的密钥放到启动的机器中，可将密钥存放于用户在机器上创建的文件夹（示例/root/auth-key/）下。

3. 生成AMI（可参考[文档](#)）

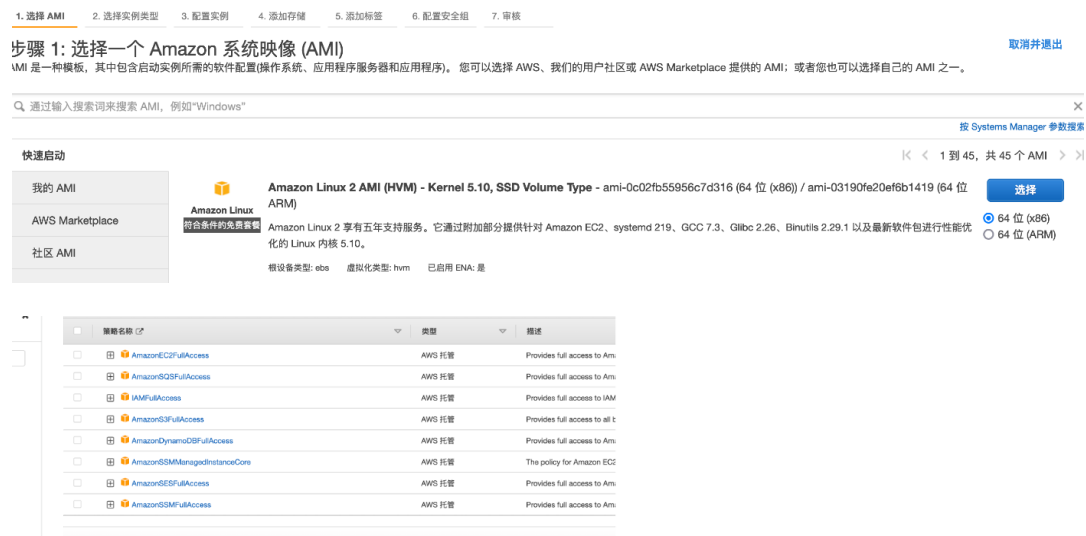


在AMI中查询刚刚生成的AMI的ID，之后会在传输机器的config中配置，请看下图位置。

7. 部署传输服务

如果您没有用于通过SSH登录EC2实例的密钥对，请[参考文档](#)进行创建。

1. 进入AWS管理控制台，搜索EC2并进入，在“实例”页面中点击“启动新实例”按钮，点击右侧“Amazon Linux2 AMI”对应的“选择”按钮



2. 选择机型



3. 在“配置实例详细信息”的步骤中，选择我们在第1步和第2步创建的VPC与修改过后的公有子网

3: 配置实例详细信息

以满足您的需求。您可以从同一AMI上启动多个实例，请求Spot实例以利用其低价优势，向实例分配访问管理角色等等。

实例的数量 [启动至 Auto Scaling 组](#)

购买选项 ☐ 请求 Spot 实例

网络 [新建 VPC](#)

子网 [新建子网](#)
4091 个 IP 地址可用

自动分配公有 IP

主机名称类型

DNS Hostname ☐ Enable IP name IPv4 (A record) DNS requests
☒ 启用基于资源的 IPv4 (A 记录) DNS 请求
☐ 启用基于资源的 IPv6 (AAAA 记录) DNS 请求

置放群组 ☐ 将实例添加到置放群组

4. 在“添加存储”中，调整硬盘大小，调到40GB即可

4: 添加存储

将您使用以下存储设备设置启动，您可以将其他 EBS 卷和实例存储卷附加到您的实例，以满足您的需求。您还可以在启动实例后附加其他 EBS 卷而无需实例。详细了解有关 Amazon EC2 存储选项的信息。

设备	快照	大小 (GiB)	卷类型	IOPS	吞吐量 (MiB/s)	停止时删除
/dev/xvda	snap-0c1ac78a014204c	<input type="text" value="40"/>	通用型 SSD (gp2)	125/3000	不适用	<input checked="" type="checkbox"/>

价格使用免费套餐的用户最多可获得 30GB 的 EBS 通用型(SSD)或磁存储卷。有关免费使用套餐资格和使用限制的信息，请参阅“了解更多”。

5. 在“配置安全组”步骤中，创建一个新的安全组，设置安全组规则，将可访问IP限制为您自己的IP即可，如下图

步骤 6: 配置安全组

安全组是一组防火墙规则，用于控制您的实例的流量。在此页面上，您可以指定规则并允许特定流量到达您的实例。例如，如果您希望设置一个 Web 服务器，并允许 Internet 流量到达您的实例，请添加相应的规则来允许来自任何地址的 HTTP 和 HTTPS 端口。您可以创建一个新安全组或从下面选择一个现有安全组。有关 Amazon EC2 安全组的信息，请参阅“了解更多”。

分配安全组 ☒ 添加一个新安全组

安全组名称

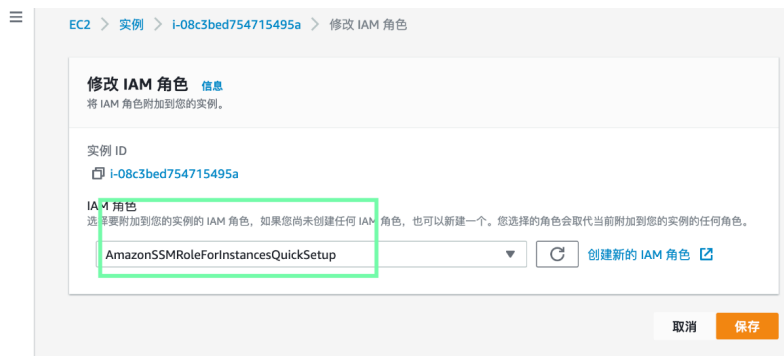
描述

类型	协议	端口范围	来源	策略
所有流量	全部	0 - 65535	<input type="text" value="64.292.61.34/32"/>	拒绝 (Close for Admin Desktop)

[添加规则](#)

6. 选择您的密钥对（Key Pair），启动实例。

7. 设置机器的权限，我们将第3步创建好的Rclone-ec2-controller-role角色赋予这个机器



8. 连接到机器，选择您创建好的实例，点击右上角第一个“连接”按钮



9. 选择“会话管理器”，并点击“连接”进入



9. 查询是否有默认的安装，如过无，则执行安装。

```
java -version
sudo yum install java-1.8.0-openjdk.x86_64
```

10. 从我们您给我们put代码的S3桶拉取代码包

```
aws s3 cp [s3存储路径] [本地路径]
```

11. 启动程序

```
java -jar ***.jar > transfer.log &
```

8. 传输任务配置

1. 打开「EC2的外网IP」:8080 web页面，如下图，按照“说明”一列进行配置。

← → ↺

your-ec2-public-ip:8080

↓

传输任务配置

启动传输集群

传输任务列表

请按照说明填写以下配置

配置名称	内容	说明
GCS Path File:	<div>Browse...</div> No file selected.	文件格式csv，内容举例：warehouse/business/project/dt=2022-02-28/
GCS Bucket:	s3://[redacted]ket-us-east-1/	源GCS的bucket名称,请参考默认例子输入
AWS S3 Bucket:	s3://data-clone-testing[redacted]-delete-after-c	目的AWS S3的bucket名称，请参考默认例子输入
AWS Region:	us-east-1	目的AWS S3所在的区域
AWS EC2 AMI ID:	ami-0c02fb55956c7d316	目的region对应的EC2 AMI ID，region若不变，保持默认值即可
Cloen Log Bucket:	s3://transfer-log-rcloen/	存储数据传输过程中产生的log的S3 bucket名称
Rclone 命令：	rclone copy -v {SourceLocation} {DestinationLocation}	Rclone的传输命令，可根据默认格式进行修改
Alarm Email From:	[redacted]@amazon.com	发送告警的邮箱
Alarm Email To:	[redacted]@amazon.com	接收告警的邮箱
SQS URL:	https://sqs.us-east-1.amazonaws.com/[redacted]	用来传输copy数据的消息队列的URL

保存配置

2. 其中“GCS Path File ”上传的文件需要是csv文件，文件中只有一列，即您需要copy的文件路径，最好按照大小进行拆分为天级别的文件路径， 注意文件中不需要包含bucket的名字，如果您的某个存储路径为gs://your-bucket-name/app1/project1/date=20220501/file.gz+'工作表1'D3, 那么放入csv的内容为：app1/project1/date=20220501/，每个路径对应一个具体的传输任务，保存配置后会存储到DynamoDB对应的表中。文档格式亦可参考：demo.csv ，也可以参考下方csv文件截图：

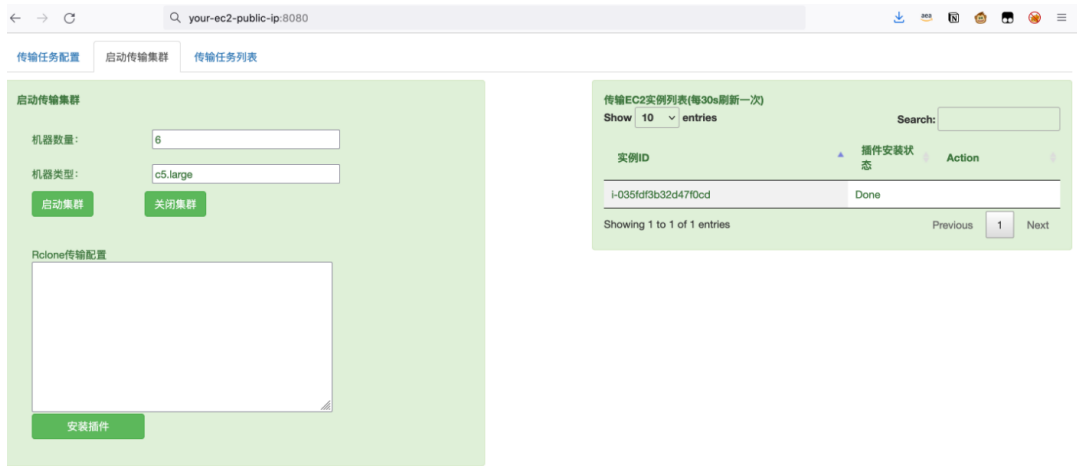
	A
1	app2/project2/business2/2022-01-01/
2	app2/project2/business2/2022-01-02/
3	app2/project2/business2/2022-01-03/
4	app2/project2/business2/2022-01-04/
5	app2/project2/business2/2022-01-05/
6	app2/project2/business2/2022-01-06/
7	app2/project2/business2/2022-01-07/
8	app2/project2/business2/2022-01-08/
9	app2/project2/business2/2022-01-09/
10	app2/project2/business2/2022-01-10/
11	app2/project2/business2/2022-01-11/
12	app2/project2/business2/2022-01-12/
13	app2/project2/business2/2022-01-13/
14	app2/project2/business2/2022-01-14/
15	
16	
17	
18	
19	
20	
21	
22	
23	
24	
25	
26	
27	
28	
29	
30	
31	
32	
33	
34	
35	
36	
37	
38	
39	
40	
41	
42	
43	
44	
45	
46	
47	
48	
49	
50	
51	
52	
53	
54	
55	
56	
57	
58	
59	
60	
61	
62	
63	
64	
65	
66	
67	
68	
69	
70	
71	
72	
73	
74	
75	
76	
77	
78	
79	
80	
81	
82	
83	
84	
85	
86	
87	
88	
89	
90	
91	
92	
93	
94	
95	
96	
97	
98	
99	
100	

- 其中“AWS EC2 AMI ID”为第6步创建的AMI的ID。
- “Alarm Email From”和“Alarm Email To” 为第5步验证过的邮箱地址。
- “SQS URL”为第4步创建的SQS的队列的URL，在AWS管理控制台SQS队列的页面上即可以复制。
- “Rclone命令”默认使用copy命令，如果需要添加传输参数，请参考附录中的Rclone文档参考，并在此修改，注意 {SourceLocation}、{DestinationLocation}与rclone -v需保持默认的格式，示例如下。

```
rclone copy -v --transfers 8 --checkers 4 --s3-chunk-size 10M {SourceLocation} {DestinationLocation}
```
- 所有信息填写好后， 点击“保存配置”即可。

9. 启动传输集群

- 在web页面中点击“启动传输集群”



- 填写机器数量和机器类型，**注意**机器类型需要严格和AWS EC2的机型保持一致， 点击“启动集群”。

以下是验证过的两种机型上的传输效果对比，以供参考：

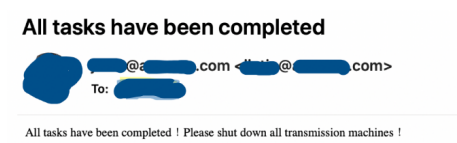
	A	B	C	D
1	机型	价格	网络性能	传输性能
2	C5.large	0.097 USD	最高可达10GB	最高300Mib/s，通常在传输400GB数据后会降低至85Mib/s
3	C6gn.large	0.0986 USD	最高可达25GB	最高300Mib/s，通常在传输600GB数据后会在200Mib/s上下波动。最低可至30Mib/s，这种情况比较少。

- 2. 等待3-5分钟，上一步启动的EC2实例需要进行初始化。
- 3. “Rclone传输配置”可以在已有rclone的服务器上配置好后直接复制过来（可通过rclone config show查看）。
- 4. “Rclone传输配置”格式参考如下：

```
[gcp]
type = google cloud storage
client_id = your-clienid
client_secret = your-client-secret
project_number = your-number
service_account_file = /file-path/gcs.json # GCS证书
object_acl = authenticatedRead
bucket_acl = authenticatedRead

[s3]
type = s3
provider = AWS
env_auth = true
region = us-east-1
location_constraint = us-east-1
```

- 5. 其他Rclone支持的云厂商对应的配置可以参考[Rclone文档](#)。
- 6. 点击“安装插件”，大约1分钟左右，插件安装完后，右侧EC2实例列表中的“插件安装状态”一列会变为“Done”。
- 7. 此时在第7步上传的文件路径生成的传输任务已经在传输的SQS队列中，因此数据自动开始传输。
- 8. 所有任务传输完成后您填写的To邮箱会收到以下邮件：



10. 添加传输任务

如果您的一批传输任务已经完成，需要添加新的传输任务，则将文件路径放在一个新的文件中，参考第8步，直接上传新的文件，点击“保存配置”即可，其他信息如果没有需要可不进行修改（第一次写入的其他配置信息已保存）。

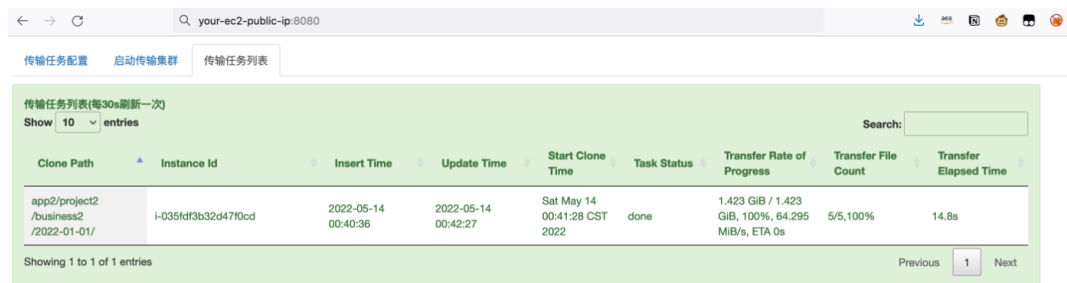
11. 查看传输实例

页面右侧可以看到您当前启动的所有用来传输数据的实例。



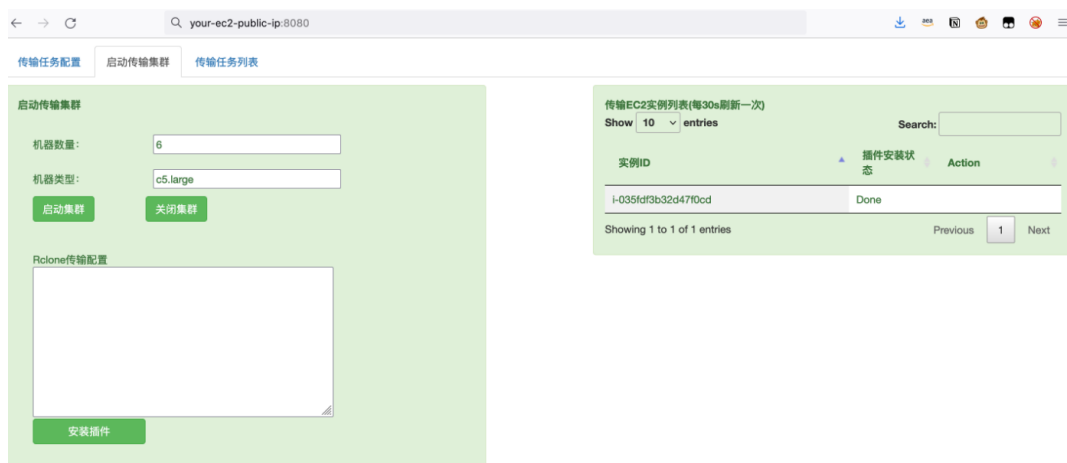
12. 查看传输任务

在web页面中，点击“传输任务列表”，进行查看；最后的三列如果数据还未传输完成，显示为“N/A”。



13. 关闭传输集群

点击以下页面中的“关闭集群”的按钮即可，注意：关闭集群意味着所有用来传输数据的实例都将被Terminate（终止），请务必确认所有实例上的传输任务已完成（所有任务完成后会有邮件通知）。



14. 清除资源

需要手动删除包括VPC、IAM role、SQS、SES与所创建的AMI。

DynamoDB中的消息task列表的内容是传输过程中的完成状态，如不需要直接删除表即可。

附录

1. [AWS system manager介绍](#)
2. [Rclone文档](#)
3. [DynamoDB介绍](#)