

Moving data from CDP Private Cloud Base to Public Cloud with NiFi site-to-site

Date published: 2019-12-16

Date modified: 2024-04-03

The Cloudera logo is displayed in a bold, orange, sans-serif font. The word "CLOUDERA" is written in all caps, with a stylized 'E' that has a horizontal bar extending to the right.

Legal Notice

© Cloudera Inc. 2024. All rights reserved.

The documentation is and contains Cloudera proprietary information protected by copyright and other intellectual property rights. No license under copyright or any other intellectual property right is granted herein.

Unless otherwise noted, scripts and sample code are licensed under the Apache License, Version 2.0.

Copyright information for Cloudera software may be found within the documentation accompanying each component in a particular release.

Cloudera software includes software from various open source or other third party projects, and may be released under the Apache Software License 2.0 (“ASLv2”), the Affero General Public License version 3 (AGPLv3), or other license terms. Other software included may be released under the terms of alternative open source licenses. Please review the license and notice files accompanying the software for additional licensing information.

Please visit the Cloudera software product page for more information on Cloudera software. For more information on Cloudera support services, please visit either the Support or Sales page. Feel free to contact us directly to discuss your specific needs.

Cloudera reserves the right to change any products at any time, and without notice. Cloudera assumes no responsibility nor liability arising from the use of products, except as expressly agreed to in writing by Cloudera.

Cloudera, Cloudera Altus, HUE, Impala, Cloudera Impala, and other Cloudera marks are registered or unregistered trademarks in the United States and other countries. All other trademarks are the property of their respective owners.

Disclaimer: EXCEPT AS EXPRESSLY PROVIDED IN A WRITTEN AGREEMENT WITH CLOUDERA, CLOUDERA DOES NOT MAKE NOR GIVE ANY REPRESENTATION, WARRANTY, NOR COVENANT OF ANY KIND, WHETHER EXPRESS OR IMPLIED, IN CONNECTION WITH CLOUDERA TECHNOLOGY OR RELATED SUPPORT PROVIDED IN CONNECTION THEREWITH. CLOUDERA DOES NOT WARRANT THAT CLOUDERA PRODUCTS NOR SOFTWARE WILL OPERATE UNINTERRUPTED NOR THAT IT WILL BE FREE FROM DEFECTS NOR ERRORS, THAT IT WILL PROTECT YOUR DATA FROM LOSS, CORRUPTION NOR UNAVAILABILITY, NOR THAT IT WILL MEET ALL OF CUSTOMER’S BUSINESS REQUIREMENTS. WITHOUT LIMITING THE FOREGOING, AND TO THE MAXIMUM EXTENT PERMITTED BY APPLICABLE LAW, CLOUDERA EXPRESSLY DISCLAIMS ANY AND ALL IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO IMPLIED WARRANTIES OF MERCHANTABILITY, QUALITY, NON-INFRINGEMENT, TITLE, AND FITNESS FOR A PARTICULAR PURPOSE AND ANY REPRESENTATION, WARRANTY, OR COVENANT BASED ON COURSE OF DEALING OR USAGE IN TRADE.

Contents

Moving data from Private Cloud to Public Cloud with NiFi site-to-site.....	4
Understanding the use case.....	4
Preparing your clusters.....	4
Setting up your network configuration.....	5
Configuring your truststores.....	5
Defining your CDP Public Cloud data flow.....	7
Configuring Ranger policies for site-to-site communication.....	8
Defining your CDP Private Cloud Base data flow.....	10

Moving data from Private Cloud to Public Cloud with NiFi site-to-site

You can use Apache NiFi's site-to-site functionality to design a flow that moves data from your CDP Private Cloud Base environment to CDP Public Cloud to ensure scalability and load balancing.

Understanding the use case

You can use the Apache NiFi site-to-site functionality to move data between a Public Cloud and a CDP Private Cloud Base environment. To do this, set up a cluster in each environment, prepare your network and truststore configurations, and then define your CDP Private Cloud Base and CDP Public Cloud data flows and Apache Ranger configuration for site-to-site functionality.

Moving data between CDP Public Cloud and CDP Private Cloud Base clusters is a common use case when there is a need for a lot of temporary compute resources that can be quickly provisioned in the cloud.

Imagine you have a large dataset on-premises and you wish to perform heavy computations on the dataset. You can use the following workflow to design a data flow that:

- Moves the dataset from your CDP Private Cloud Base environment to your CDP Public Cloud environment
- Pushes the data to the appropriate destination
- Triggers the workload that processes the data while leveraging the auto-scaling capabilities that CDP Public Cloud provides
- Returns the results in your CDP Private Cloud Base environment

All of this is powered by CDP Private Cloud Base and CDP Public Cloud distributions, while ensuring consistent security policies at a fine-grained level with Apache Ranger, and data management and data lineage with Apache Atlas across the environments.

Preparing your clusters

The first step in preparing to move data from a CDP Private Cloud Base cluster to a CDP Public Cloud cluster is to ensure that you have each cluster set up correctly.

Requirements for your CDP Private Cloud Base cluster:

- CFM running on CDP Private Cloud Base
- Three-node NiFi compute cluster, secured with AutoTLS and configured with Apache Ranger

For details on deploying your CFM cluster on CDP Private Cloud Base, see *CFM Deployment to CDP Private Cloud Base*.

Requirements for your CDP Public Cloud Flow Management cluster:

- Flow Management clusters running on CDP Public Cloud
- Three-node NiFi compute cluster, secured with AutoTLS and configured with Apache Ranger

For details on deploying your CFM cluster on CDP Public Cloud, see *Setting up your Flow Management cluster*.

Related Information

[CFM Deployment to CDP Private Cloud Base](#)

[Setting up your Flow Management cluster](#)

Setting up your network configuration

You can use NiFi's site-to-site capabilities over a RAW TCP or over an HTTP network configuration. For this use case, you must configure site-to-site using HTTP over TLS. This has the advantage of using the NiFi port, which is also used to access the NiFi UI and APIs.

For the purpose of this use case, set up your site-to-site network configurations with the following assumptions:

- You are not using site-to-site through any proxy configuration
- You have a direct connection on port 8443 between NiFi nodes on your CDP Private Cloud Base and CDP Private Cloud clusters

Set up your network configuration according to your architecture and requirements. For more information, see your Cloud provider documentation.

In this use case, NiFi on CDP Private Cloud Base is responsible for initiating the site-to-site connection between the two environments to push and pull data to and from the NiFi cluster in Public Cloud. The NiFi nodes in CDP Public Cloud must be reachable on port 8443 from the NiFi nodes in CDP Private Cloud Base, but not necessarily the other way around.



Tip:

The site-to-site connection is bi-directional and depending on the cluster initiating the site-to-site connection, you will be in a push or pull model.

Configuring your truststores

You must configure your truststores so that each cluster is aware of and trusts the other cluster, to support the two-way TLS that is used to initiate the site-to-site communication between clusters. To do this, you need to download and merge the truststores for NiFi in CDP Private Cloud Base and CDP Public Cloud.

Before you begin

- You have set up a CDP Private Cloud Base and CDP Public Cloud cluster, and have the necessary network configurations established.
- You have the necessary administrative permissions to manipulate the truststore files. You require root access, and this is typically done by an Environment Administrator.
- You have the passwords for the Java Keystore (JKS) files available.

Procedure

1. Download the truststore from the clusters.

- a) Create a temporary directory in which you can edit the truststore files.

```
$ mkdir s2s-temp && cd s2s-temp
```

- b) Download the JKS truststore file for CA used by NiFi in your CDP Private Cloud Base cluster.

```
$ scp -i <key>  
root@<nifi_node>:/var/lib/cloudera-scm-agent/agent-cert/cm-auto-global  
_truststore.jks  
privatecloud_cm-auto-global_truststore.jks
```

- c) Download the JKS truststore file for CA used by NiFi in your CDP Public Cloud cluster.

```
$ scp -i <key>
```

```
cloudbreak@<nifi_node_public_cloud>:/var/lib/cloudera-scm-agent/agent-
cert/cm-auto-global_truststore.jks
publiccloud_cm-auto-global_truststore.jks
```

2. Merge the truststores.

a. Make a copy of the CDP Private Cloud Base JKS:

```
$ cp privatecloud_cm-auto-global_truststore.jks
privatecloud_cm-auto-global_truststore.jks.bak
```

b. Merge the CDP Public Cloud JKS into the CDP Private Cloud Base JKS and rename the entries alias to prevent conflict:

```
$ keytool
-importkeystore
-srckeystore publiccloud_cm-auto-global_truststore.jks
-destkeystore privatecloud_cm-auto-global_truststore.jks
```

The result will be similar to:

```
Importing keystore publiccloud_cm-auto-global_truststore.jks to private
cloud_cm-auto-global_truststore.jks...
Enter destination keystore password:
Enter source keystore password:
Entry for alias imported-ca-b379e6601f5ecfbbec2fefc4eb2efd4a successfull
y imported.
[...]
Entry for alias imported-ca-5945bad341623ae14991e09ffe851725 successfu
lly imported.
Entry for alias cmrootca-1 successfully imported.
Existing entry alias cmrootca-0 exists, overwrite? [no]: no
Enter new alias name (RETURN to cancel import for this entry): cmrootc
a-1-bis
Entry for alias cmrootca-0 successfully imported.
Entry for alias imported-ca-10c56ecc972802e53d1b7287ac2d1c6c successfu
lly imported.
[...]
Entry for alias imported-ca-840644351dd523125493ff4c28e694f7 successfull
y imported.
Import command completed: 140 entries successfully imported, 0 entries
failed or cancelled
```

c. Use the copy you made to merge the CDP Private Cloud Base JKS into the CDP Public Cloud JKS and rename the entries alias to prevent conflict:

```
$ keytool
-importkeystore
-srckeystore privatecloud_cm-auto-global_truststore.jks.bak
-destkeystore publiccloud_cm-auto-global_truststore.jks
```

The result will be similar to:

```
Importing keystore privatecloud_cm-auto-global_truststore.jks.bak to pu
bliccloud_cm-auto-global_truststore.jks...
Enter destination keystore password:
Enter source keystore password:
Existing entry alias cmrootca-0 exists, overwrite? [no]: no
```

```
Enter new alias name (RETURN to cancel import for this entry): cmrootc
a-0-bis
Entry for alias cmrootca-0 successfully imported.
Import command completed: 1 entries successfully imported, 0 entries
failed or cancelled
```

3. Deploy the truststores.

Deploy the modified CDP Private Cloud Base and CDP Public Cloud JKS files on each NiFi node of the respective clusters.

**Note:**

Do not change the permissions and owners of the file (chmod/chown).

4. Restart your CDP Private Cloud Base and CDP Public Cloud clusters.

What to do next

After you have configured your truststores, proceed by defining your data flow in your CDP Public Cloud cluster.

Defining your CDP Public Cloud data flow

To move data between cloud environments using NiFi site-to-site communication, you require a data flow in CDP Public Cloud that can receive data from the CDP Private Cloud Base data flow. To create this data flow, configure a process group, and both an input and output port.

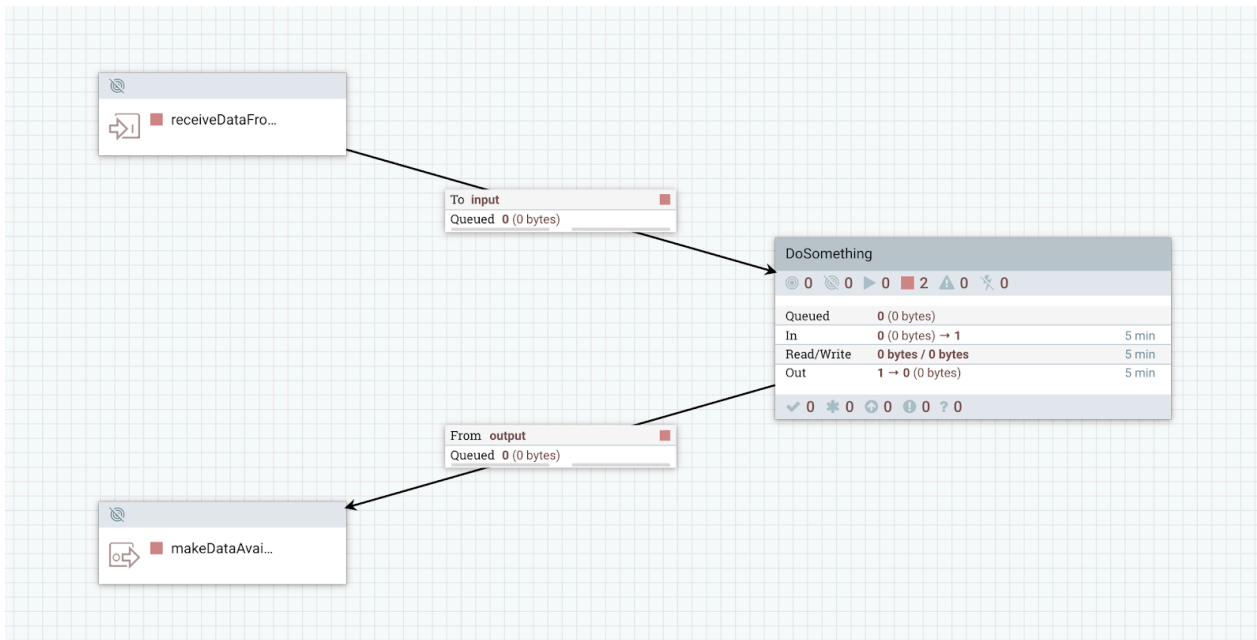
Before you begin

You have prepared your clusters, set up your network configurations, and configured your truststores.

Procedure

1. From your CDP Public Cloud NiFi cluster, create a Process Group to perform the operations you want to complete on the data received from and returned to the CDP Private Cloud Base cluster.
2. Drag an Input Port onto the NiFi canvas.
You must use this port for receiving data from NiFi's CDP Private Cloud Base cluster.
3. Drag an Output Port onto the NiFi canvas.
You must use this port to make data available for download to the CDP Private Cloud Base cluster.
4. Connect your CDP Public Cloud data flow components.
Ensure that you have specified the public endpoints of your NiFi nodes in the CDP Public Cloud cluster.
5. Start your data flow and ensure that both the input and output ports are running.

Example



What to do next

When you have completed your CDP Public Cloud data flow, proceed by configuring Apache Ranger to allow NiFi's site-to-site transmission.

Configuring Ranger policies for site-to-site communication

To allow NiFi's site-to-site communication between CDP Public Cloud and CDP Private Cloud Base clusters, you need to configure Ranger authorization between the two clusters. To do this, create Ranger users in your CDP Public Cloud cluster that correspond to the CDP Private Cloud Base NiFi nodes. Then create a new Ranger policy with site-to-site resources configured, and assign your CDP Private Cloud Base NiFi node users to the policy.

Before you begin

- You have defined your CDP Public Cloud data flow.
- You have a list of your FQDN CDP Private Cloud Base host names. You need the host names to create the Ranger policies in CDP Public Cloud.

Procedure

1. In your CDP Public Cloud environment, launch the Ranger UI, click **Settings Users/Groups/Roles User Create** to add the users corresponding to the nodes of the CDP Private Cloud Base cluster.

- Click User Create to create one user per NiFi node running your CDP Private Cloud Base environment.

You create this user to make Ranger aware of the CDP Private Cloud Base nodes, so that you can create policies by including them. Because this user is not used to authenticate on the Ranger UI, the password can be random.

The screenshot shows the Ranger 'User Create' interface. At the top, there is a navigation bar with 'Ranger', 'Access Manager', 'Audit', 'Security Zone', and 'Settings'. Below this is a breadcrumb trail: 'Users/Groups/Roles' > 'User Create'. The main section is titled 'User Detail' and contains the following fields:

- User Name ***: nifi-d-compute2.field.hortonwo
- New Password ***: [Masked with dots]
- Password Confirm ***: [Masked with dots]
- First Name ***: nifi-d-compute2.field.hortonwo
- Last Name**: [Empty]
- Email Address**: [Empty]
- Select Role ***: User
- Group**: Please select (with a '+' button next to it)

At the bottom of the form, there are two buttons: 'Save' and 'Cancel'.

- Create a new policy in the NiFi Service in Ranger.

You need to enter the following NiFi Resources:

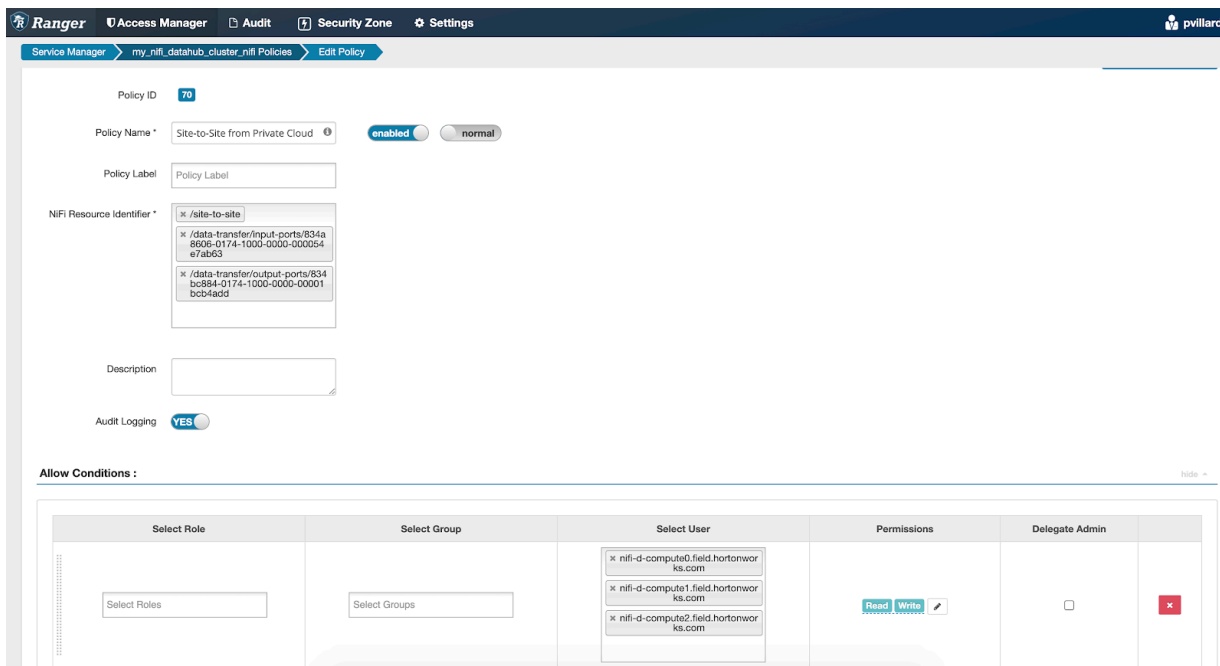
- /site-to-site
- /data-transfer/input-ports/<ID of the Input Port>
- /data-transfer/output-ports/<ID of the Output Port>



Note:

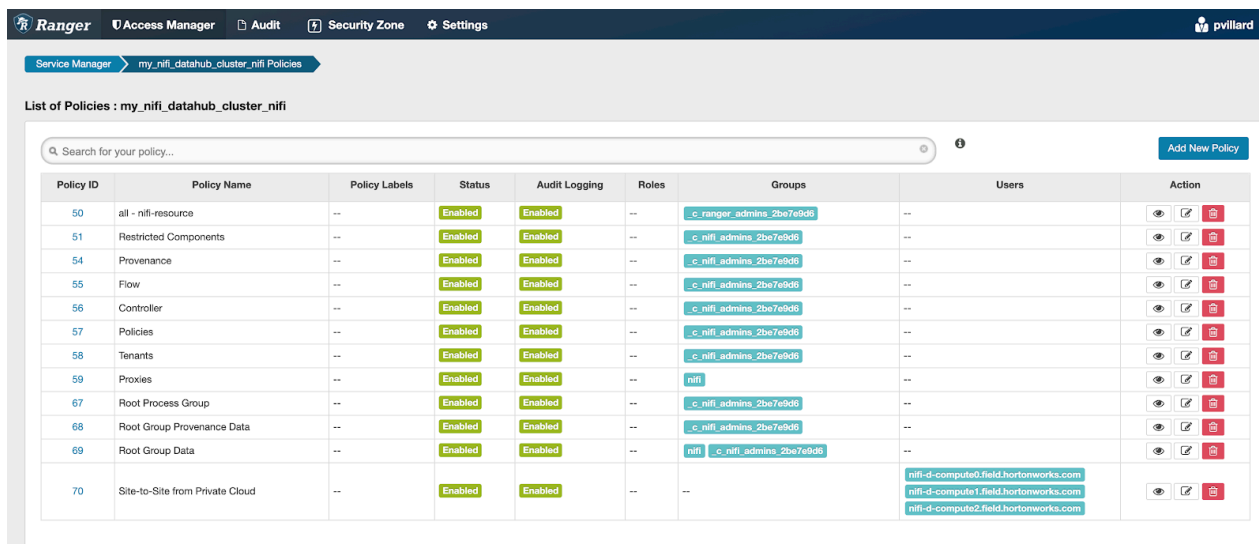
You can retrieve the input and output port IDs by right-clicking the component and reviewing the configuration view.

4. Add the CDP Private Cloud Base users you created in Step 2, and assign Read and Write permissions:



Results

Your policies are now listed.



Defining your CDP Private Cloud Base data flow

To move data between cloud environments using NiFi site-to-site communication, you require a data flow in your CDP Private Cloud Base cluster that can send and receive data from the CDP Public Cloud cluster. To create this data flow, connect a processor to a Remote Process Group configured with HTTP and enable transmission.

Before you begin

- You have defined your CDP Public Cloud data flow and configured Ranger policies for site-to site communication.

- You have the public FQDNs for your CDP Public Cloud cluster nodes.

Procedure

- In your CDP Private Cloud Base cluster, launch the NiFi UI and drag a `GenerateFlowFile` processor onto the canvas.

For this use case, `GenerateFlowFile` creates 1MB files every 10 seconds.

- Drag a Remote Process Group onto the NiFi canvas, configure HTTP protocol, and specify one or more of the NiFi nodes running on your CDP Public Cloud cluster.

After the site-to-site connection is initiated, the source NiFi cluster is aware of the topology of the remote NiFi cluster and of any increase or decrease of the size of the remote cluster. However, it is recommended that you specify at least 2 nodes to ensure higher availability when the site-to-site connection is initiated.

Add Remote Process Group

URLs ?

Transport Protocol ?

Local Network Interface ?

HTTP Proxy Server Hostname ?

HTTP Proxy Server Port ?

HTTP Proxy User ?

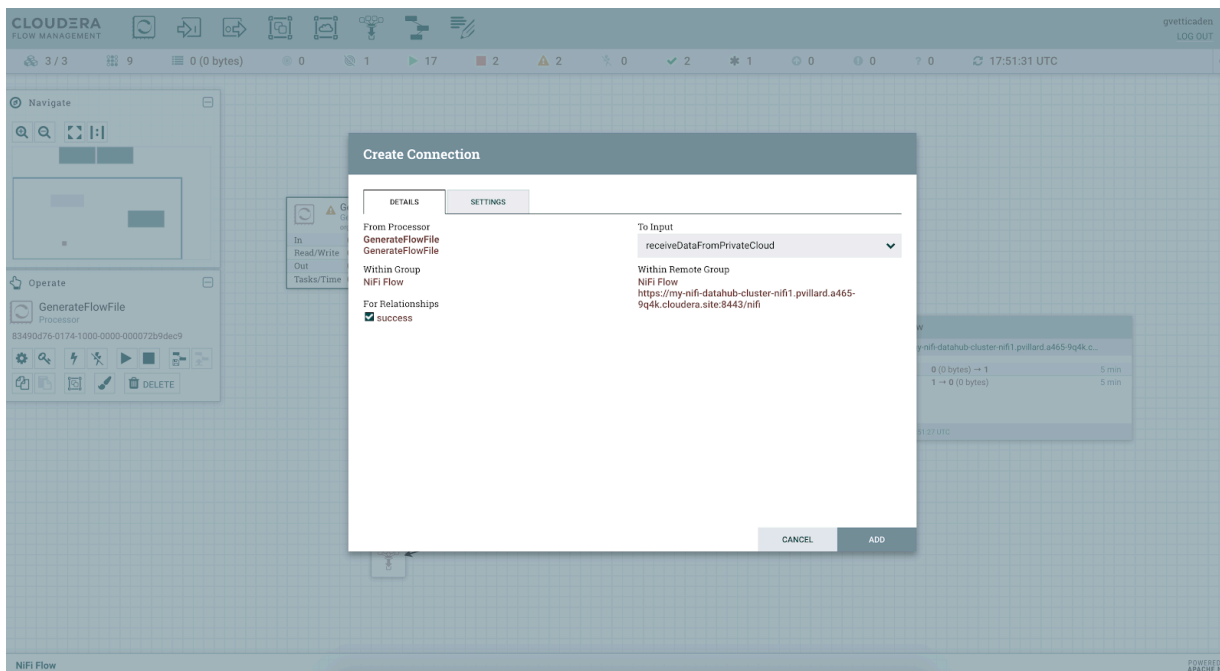
HTTP Proxy Password ?

Communications Timeout ?

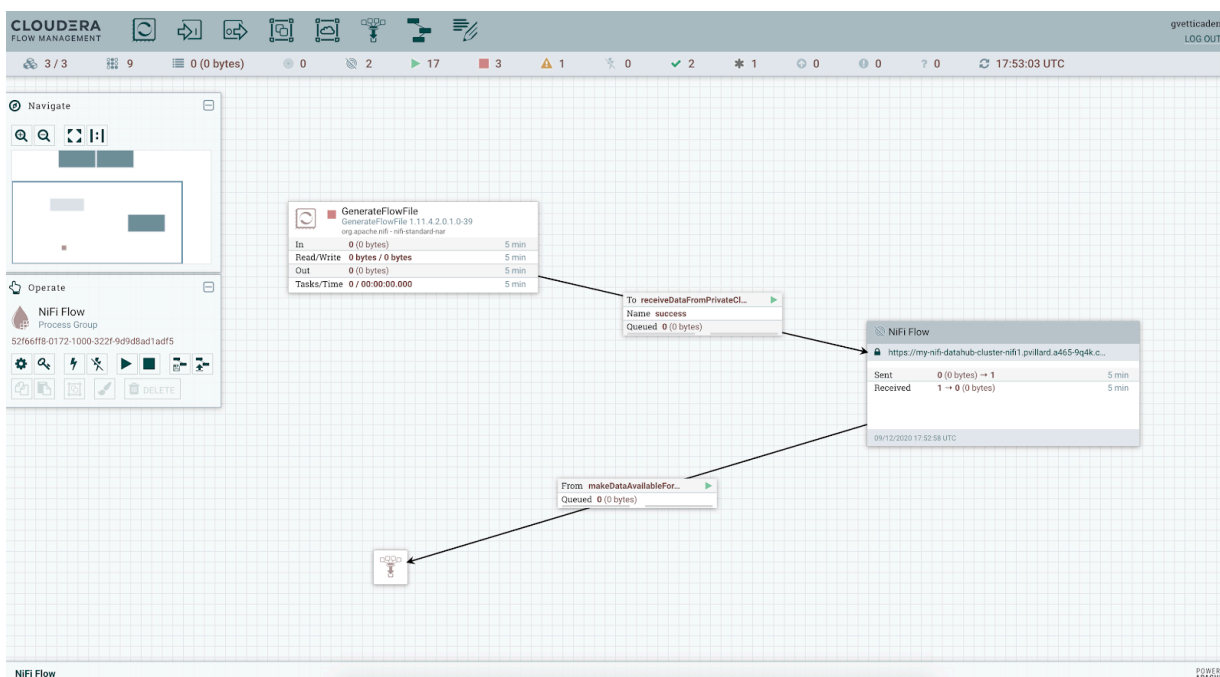
Yield Duration ?

- Right-click the Remote Process group and select Enable transmission.

4. Connect the `GenerateFlowFile` processor to the Remote Process Group and select the Input Port that you created and started on the remote cluster in CDP Public Cloud:



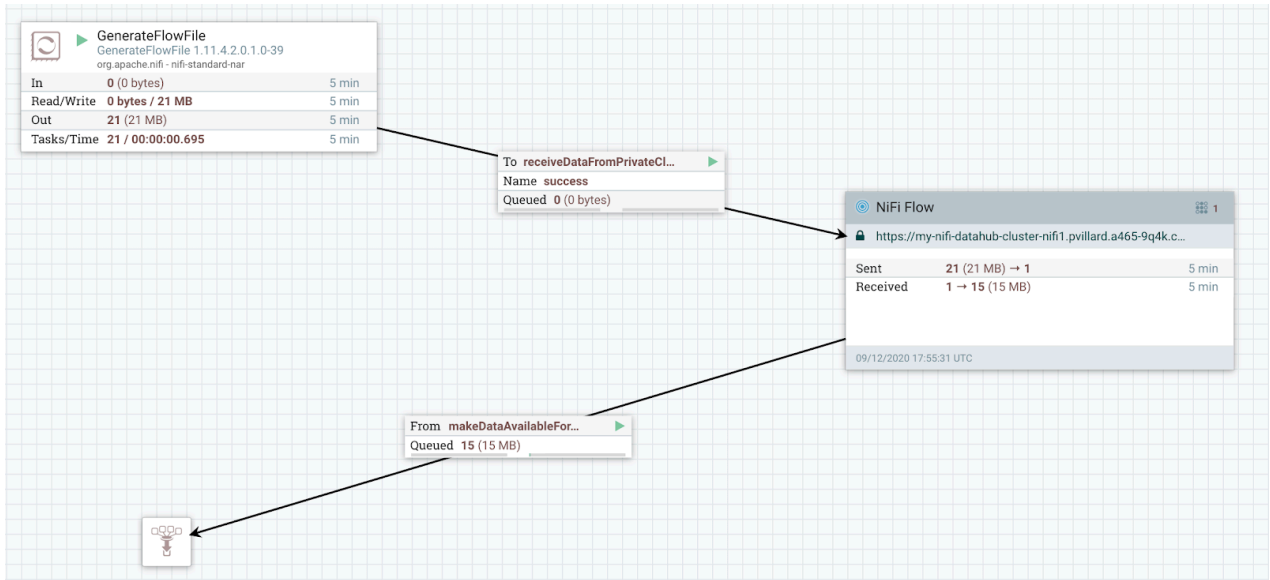
5. You can also define a connection from the Remote Process Group to another component to download data made available by the remote cluster running in the CDP Public Cloud environment. In this example, the Remote Process Group is connected to a funnel.



Results

After you have defined the data flow for your CDP Private Cloud Base cluster, start the CDP Private Cloud Base data flow and confirm that the data is moving back and forth between the environments:

In the CDP Private Cloud Base environment, your data flow looks similar to the following:



In the CDP Public Cloud environment, your data flow will look similar to the following:

