

NGScloud v0.90

Bioinformatic system for RNA-seq analysis of non-model species using cloud computing

GI Genética, Fisiología e Historia Forestal
Dpto. Sistemas y Recursos Naturales
ETSI Montes, Forestal y del Medio Natural
Universidad Politécnica de Madrid

<http://gfhforestal.com/>
<https://github.com/ggfhf/>

Table of contents

Introduction	1
Installation.....	3
NGScloud installation	3
Additional software installation and dependencies	4
Ubuntu Linux	4
Microsoft Windows	5
Mac OS X	6
First steps	8
Connect to your AWS account	8
Search the Account Id	9
Create an Access Key Id and Secret Access Key	9
Starting NGScloud	10
Configuring your first NGScloud environment.....	11
A step by step example	16
Create volumes	17
Link volumes in cluster templates.....	18
Create cluster with the t2.micro template.....	19
Upload the read files to the cluster	21
Setup bioinformatic applications in the cluster	27
Review the quality of reads using FastQC.....	27
Download the quality analysis results.....	34
Trim reads using Trimmomatic	37
Terminate the cluster with a t2.micro template and create another cluster with a r3.4xlarge template	42
Assembly reads using Trinity.....	45
Evaluate the transcriptome using RSEM-EVAL.....	53
Terminate the cluster with r3.4xlarge template and create another cluster with r3.xlarge template	59
Transcriptome filtering using transcript-filter.....	63
Transcriptome clustering using CD-HIT-EST.....	70
Terminate the cluster with r3.xlarge template and create another cluster with c3.xlarge template	78
Upload the protein database to the cluster.....	82

Add nodes to the cluster with a c3.xlarge template	86
Annotate the filtered and clustered transcriptome using transcriptome-blastx.....	89
Terminate the cluster with c3.xlarge template and create another cluster with t2.micro template	97
Download the transcriptome, evaluation and annotation files.....	101
Terminate the cluster with the t2.micro template	113
How-to.....	115
How to display this manual.....	115
How to recreate the NGScloud config file.....	115
How to create a new environment	115
How to change to another environment	115
How to view characteristics of a cluster template.....	115
How to create a cluster	115
How to terminate a cluster	115
How to list the running clusters	115
How to create a volume	115
How to remove a volume.....	115
How to list the created volumes	115
How to link a volume in cluster templates.....	115
How to add a node in a cluster	115
How to remove a node in a cluster	115
How to open a terminal of a cluster	115
How to set up a bioinformatic software in a cluster.....	116
How to run a RNA-seq bioinformatic software in a cluster	116
How to display datasets of a volume	116
How to display the contents of a dataset	116
How to upload reference files to a cluster.....	116
How to compress/decompress reference files in a cluster.....	116
How to upload database files to a cluster.....	116
How to compress/decompress database files in a cluster.....	116
How to upload read files to a cluster	117
How to compress/decompress read files in a cluster	117
How to download results files from a cluster	117
How to compress/decompress results files in a cluster	117
How to view submission logs in the local computer.....	117
How to view result logs in the cluster	117

Introduction

NGScloud is a bioinformatic system developed to analyze RNA-seq data using the cloud computing services of Amazon - Elastic Compute Cloud (EC2)- that permit the access to ad hoc computing infrastructure scaled according to the complexity of the experiment, so its costs and times can be optimized. The application provides a user-friendly front-end to easily operate Amazon's hardware resources, and to control a workflow of RNA-seq analysis oriented to non-model species, incorporating the cluster concept, which allows parallel runs of common RNA-seq analysis programs in several virtual machines for faster analysis (see Figure 1).

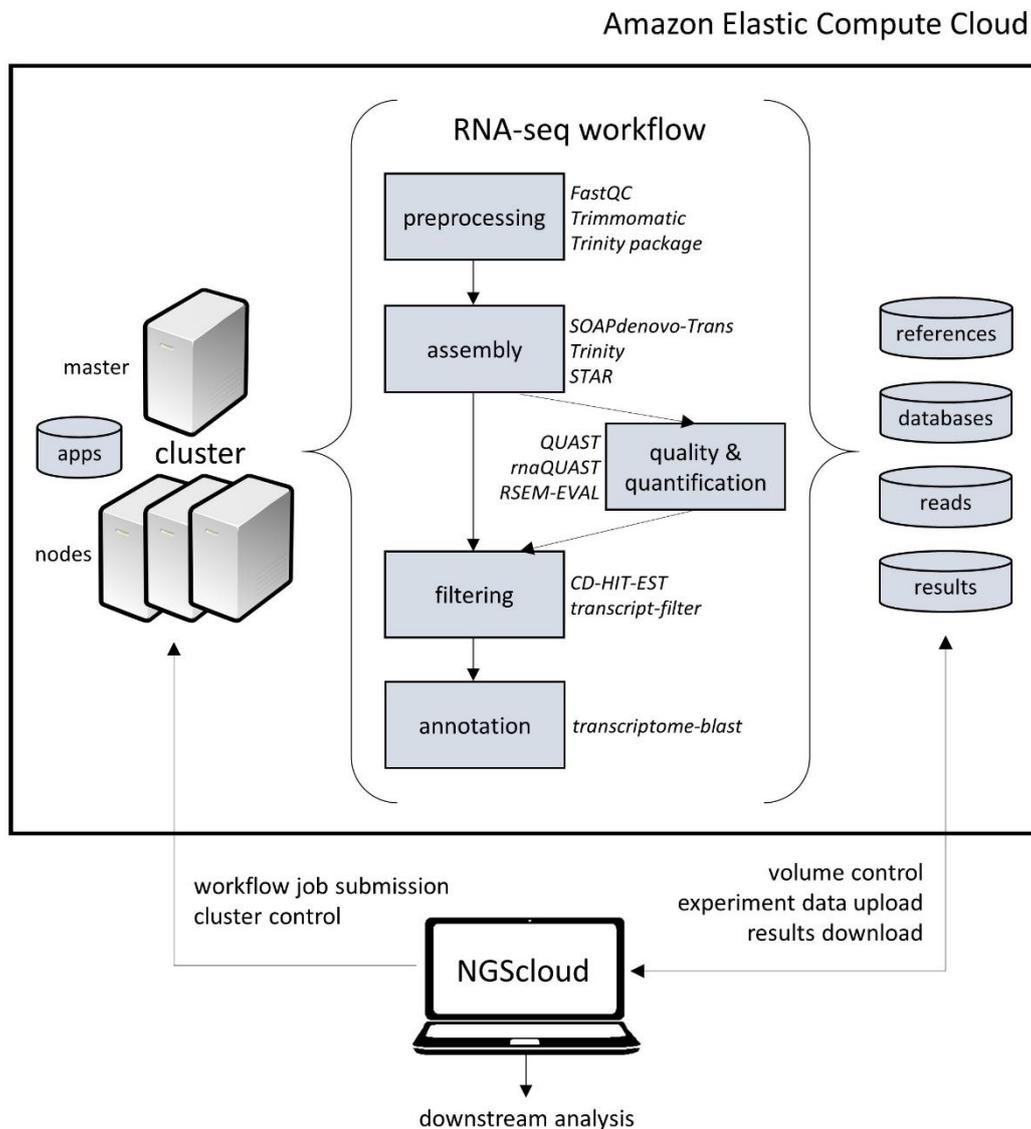


Figure. 1. NGScloud architecture. NGScloud operates EC2 resources, submits workflow and manages datasets from RNA-seq experiments.

The development of NGScloud stems from the needs of specific user-friendly tools for RNA-seq analysis in small laboratories, or by researchers that lack advanced knowledge in the bioinformatic analysis of RNA-seq experiments. NGScloud is specially oriented to RNA-seq analysis in non-model organisms or when large experiments involving many libraries and massive data generation is expected. NGScloud was designed to facilitate RNA-seq analyses since the researcher is guided in the choice of the input files for bioinformatic applications and the parameters to be used, encapsulating the complexity of the command line. In addition, NGScloud takes advantage of the resources provided by the Amazon Web Services, so it can also be considered as an alternative to private clusters, to perform the analysis, and to store the read files, results, and associated databases.

Installation

NGScloud installation

NGScloud was programmed in Python3, and it runs in any computer with an OS that allows for Python 3: Linux, Microsoft Windows, Mac OS X and other platforms.

NGScloud is available from the GitHub software repository of the Forest Genetics and Physiology Research Group (<https://github.com/GGFHF/NGScloud/>), and it is distributed under GNU General Public Licence Version 3.

To download NGScloud, click in *Clone or download* and then in *Download ZIP*:

GGFHF / NGScloud

Watch 2 Star 0 Fork 0

Code Issues 0 Pull requests 0 Projects 0 Insights

NGScloud provides a user-friendly front-end to operate the EC2 resources, and to control the workflow of non-model species oriented RNA-seq experiment in a modular way. The application allows to optimize the cost-efficiency ratio of RNA-seq experiments when appropriate computational facilities are not available.

5 commits 1 branch 0 releases 2 contributors GPL-3.0

Branch: master New pull request Find file Clone or download

fernandomoramarquez v 0.87	Latest commit f3b43d5 4 hours ago
Package	v 0.87 17 hours ago
LICENSE	Initial commit 2 days ago
README.md	v 0.87 4 hours ago

README.md

NGScloud

NGScloud is a bioinformatic system developed to analyze RNA-seq data using the cloud computing services of Amazon that permit the access to ad hoc computing infrastructure scaled according to the complexity of the experiment, so its costs and times can be optimized. The application provides a user-friendly front-end to operate Amazon's hardware resources, and to control a workflow of RNA-seq analysis oriented to non-model species, incorporating the cluster concept, which allows parallel runs of com-mon RNA-seq analysis programs in several virtual machines for faster analysis.

To install NGScloud on Linux and Mac OS X, simply decompress the NGScloud-master.zip into a directory, typing the following command in a terminal window:

```
$ unzip NGScloud-master.zip
```

Then, the execution permissions of the programs must be set by using this command:

```
$ chmod u+x *.py *.sh
```

For Microsoft Windows, simply unzip NGScloud-master.zip in the usual way.

Additional software installation and dependencies

Python 2 and Python 3 are necessary for a correct functioning of NGScloud. Python 2 is necessary because StarCluster, an additional software used to manage clusters of EC2 virtual machines (see “Additional software installation”) has been programmed in Python 2. For Ubuntu Linux both versions are already preinstalled. However, Python is not preinstalled on Microsoft Windows and Mac OS X in any of its versions.

If you use Windows, you can download both Python versions from the official website (<https://www.python.org/>), or use one of the several distributions that include Python along with other software packages for standard bioinformatic analysis. We recommend installing Anaconda (a version corresponding to Python 3.6 or higher). Anaconda is a free cross-platform for Microsoft Windows, Linux and Mac OS X (<https://www.continuum.io/>). The installation instructions for Anaconda are available on its web site.

If you are a Mac OS X user and you are not sure about how to install Python, we recommend installing Anaconda as well.

Next, we present how to install the additional software that is required to run NGScloud (AWS CLI, Boto3, Paramiko and StarCluster) on Ubuntu Linux, Microsoft Windows and Mac OS X.

To work properly, NGScloud needs the following software packages to be installed in the OS:

- AWS CLI (<https://aws.amazon.com/cli/>), the AWS Command Line Interface.
- Boto3 (<https://boto3.readthedocs.io/>), the AWS SDK (Software Development Kit) for Python.
- Paramiko (<http://www.paramiko.org/>), an implementation of the SSHv2 protocol in Python.
- StarCluster (<http://star.mit.edu/cluster/>), an open source cluster-computing toolkit for Amazon EC2 (Elastic Compute Cloud).

Ubuntu Linux

First, you open a terminal window and type the following command to install the Python3 modules Tk, PIL and PIL.ImageTk, if necessary:

```
$ sudo apt-get install python3-tk python3-pil python3-pil.imagetk
```

The additional software may be installed by typing the following commands in the terminal window.

- AWS CLI:

```
$ sudo pip3 install awscli
```

- Boto3:

```
$ sudo pip3 install boto3
```

- Paramiko:

```
$ sudo apt-get install build-essential libssl-dev libffi-dev python3-dev
```

```
$ sudo pip3 install cryptography
```

```
$ sudo pip3 install paramiko
```

- StarCluster:

```
$ sudo pip install starcluster
```

Microsoft Windows

Assuming that Anaconda has been installed in Windows with Python 3.6 or higher as main environment, Python 2.7 must be installed as additional environment identified as *py27* running the following command on a Command Prompt started as Administrator:

```
> conda create --name py27 python=2.7 anaconda
```

Then, the additional software is installed into the same Command Prompt started as Administrator:

- AWS CLI:

```
> pip install awscli
```

- Boto3:

```
> conda install boto3
```

- Paramiko:

```
> pip install paramiko
```

- StarCluster:

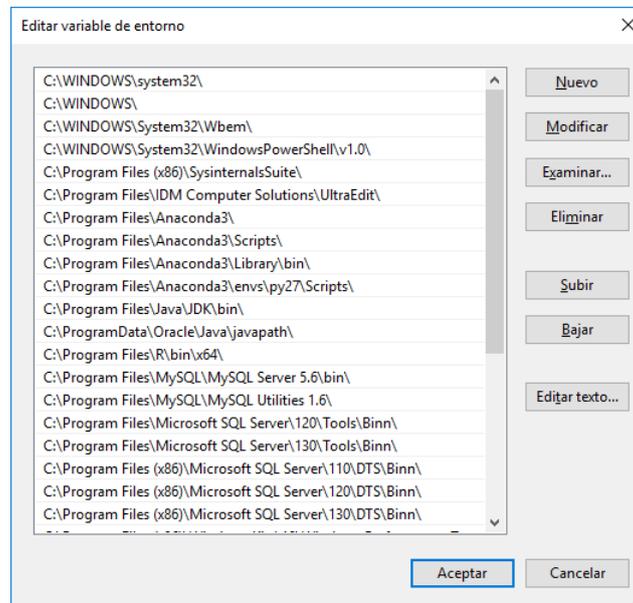
```
> activate py27
```

```
> pip install starcluster
```

```
> deactivate py27
```

If *Anaconda3_path* is the directory where you have installed Anaconda3, you must review the "Environment Variables" in "System Properties" dialog box and verify that the following directories are declared as PATH variables:

- *Anaconda3_path*
- *Anaconda3_path\Scripts*
- *Anaconda3_path\Library\bin*
- *Anaconda3_path\envs\py27\Scripts*



Mac OS X

Assuming that Anaconda distribution has been installed in Windows with Python 3.6 or higher as main environment, the steps are similar to the installation in Microsoft Windows.

First, Python 2.7 must be installed as additional environment identified as *py27* typing the following command on a terminal:

```
$ conda create --name py27 python=2.7 anaconda
```

Then you can install the additional software, by typing:

- AWS CLI:

```
$ pip install awscli
```

- Boto3:

```
$ conda install boto3
```

- Paramiko:

```
$ pip install paramiko
```

- StartCluster:

```
$ activate py27
```

```
$ pip install starcluster
```

```
$ deactivate py27
```

If *Anaconda3_path* is the directory where you have installed Anaconda3, you must review the *.bash_profile* file in your home directory to include in the PATH variable the following directories:

- *Anaconda3_path/bin*
- *Anaconda3_path/envs/py27/bin*

The last line in *.bash_profile* should be something like this:

```
export PATH=Anaconda3_path/bin:Anaconda3_path/envs/py27/bin:$PATH
```

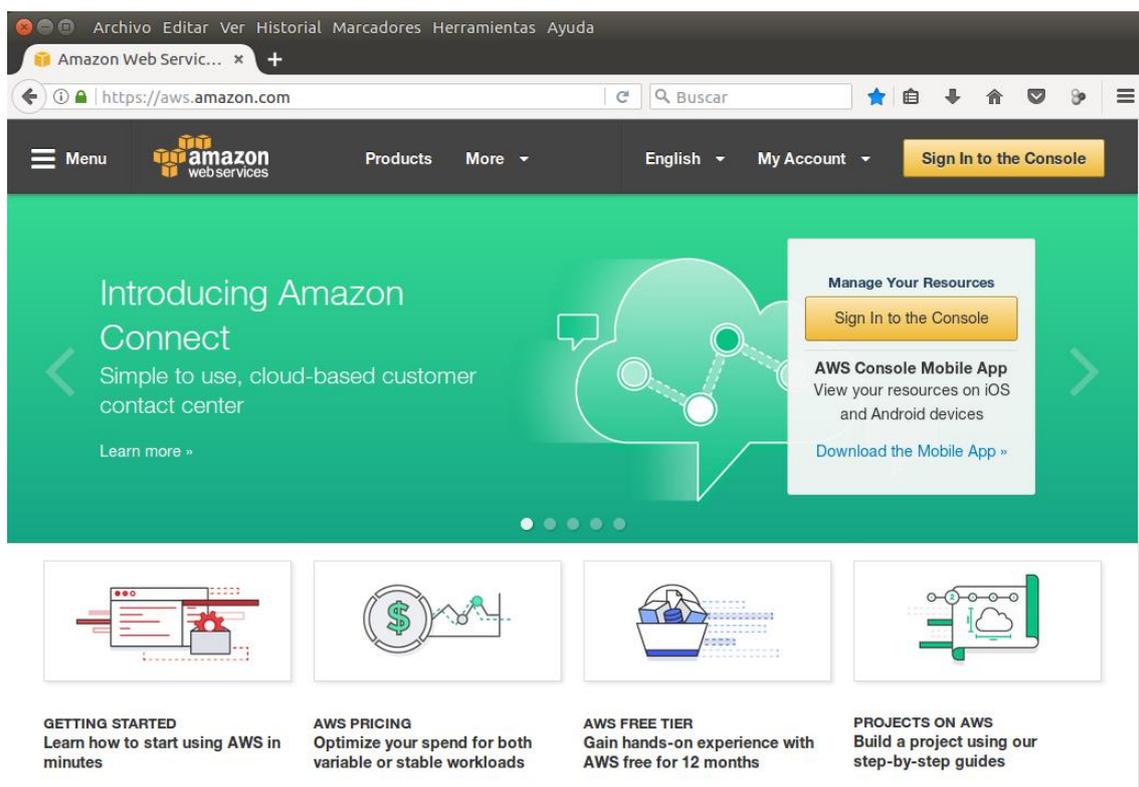
First steps

The following steps are mandatory before you can use NGScloud, after NGScloud and the additional software have been installed:

- Connect to your AWS Account
- Search the Account Id
- Create an Access Key Id and Secret Access Key
- Start NGScloud
- Configuring your first NGScloud environment

Connect to your AWS account

First, you must connect to your AWS Account in the web site <https://aws.amazon.com> clicking in *Sign in to the Console*:



Then, complete your e-mail and password personal information in the corresponding text boxes:



Sign In or Create an AWS Account

What is your email (phone for mobile accounts)?

E-mail or mobile number:

I am a new user.

I am a returning user and my password is:

[Sign in using our secure server](#)

[Forgot your password?](#)

If you don't have an AWS Account, you can create one. Currently, Amazon allows the users the access to restricted services for free for one year. Information about how to use the free tier is properly explained in: <http://docs.aws.amazon.com/awsaccountbilling/latest/aboutv2/billing-free-tier.html>.

Search the Account Id

The Account Id is a 12 digits' number located in *My Account* information:

Dashboard	▼ Account Settings Edit
Bills	Account Id: ←
Cost Explorer	Account Name:
Budgets	Password:
Reports	
Cost Allocation Tags	▼ Contact Information Edit
Payment Methods	Full Name:
Payment History	Address:
Consolidated Billing	City:
Preferences	State:
Credits	Postal Code:
Tax Settings	Country:
DevPay	Phone Number:
	Company Name:
	Website URL:

Create an Access Key Id and Secret Access Key

In *Your Security Credentials*, you must click in *Access Keys (Access Key Id and Secret Access Key)* option. Then, new information will be displayed, and you must click in *Create New Access Key*.

Dashboard

Search IAM

Details

Groups

Users

Roles

Policies

Identity Providers

Account Settings

Credential Report

Encryption Keys

Your Security Credentials

Use this page to manage the credentials for your AWS account. To manage credentials for AWS Identity and Access Management (IAM) users, use the [IAM Console](#).

To learn more about the types of AWS credentials and how they're used, see [AWS Security Credentials](#) in AWS General Reference.

- + Password
- + Multi-Factor Authentication (MFA)
- Access Keys (Access Key ID and Secret Access Key)

You use access keys to sign programmatic requests to AWS services. To learn how to sign requests using your access keys, see the [signing documentation](#). For your protection, store your access keys securely and do not share them. In addition, AWS recommends that you rotate your access keys every 90 days.

Note: You can have a maximum of two access keys (active or inactive) at a time.

Created	Deleted	Access Key ID	Last Used	Last Used Region	Last Used Service	Status	Actions
Create New Access Key							
<p>Important Change - Managing Your AWS Secret Access Keys</p> <p>As described in a previous announcement, you cannot retrieve the existing secret access keys for your AWS root account, though you can still create a new root access key at any time. As a best practice, we recommend creating an IAM user that has access keys rather than relying on root access keys.</p>							
<ul style="list-style-type: none"> + CloudFront Key Pairs + X.509 Certificates + Account Identifiers 							

Next, a dialog box will display to confirm that your Access Key Id and Secret Access Key have been created successfully. Once this has been checked, you must download a file containing your personal Access Key Id and Secret Access key by clicking in *Download Key File* button.

Create Access Key

✔ Your access key (access key ID and secret access key) has been created successfully.

Download your key file now, which contains your new access key ID and secret access key. If you do not download the key file now, you will not be able to retrieve your secret access key again.

To help protect your security, store your secret access key securely and do not share it.

[Show Access Key](#)

[Download Key File](#) [Close](#)

Starting NGScloud

You must set the directory where NGScloud.zip was decompressed as the current directory in a terminal window or command prompt. NGScloud will run in graphical mode using the graphical user interface (GUI), but it can also be run in console mode on server machines without GUI installed. Here, we explain how to run NGScloud in GUI mode. However, the console mode has menus with the same options available in GUI mode.

If you are a Linux or Mac OS X user, you start NGScloud in GUI mode typing the following command in a terminal window in the directory where the package of NGScloud is downloaded:

```
$ ./NGScloud.py
```

Alternatively, you can type also:

```
$ ./NGScloud.py --mode=gui
```

To run NGScloud in console mode:

```
$ ./NGScloud.py --mode=console
```

The file NGScloud.bat allows to execute NGScloud.py to the Microsoft Windows users calling the Python interpreter. Then, type the following command to run NGScloud in a Command Prompt in the directory where the package of NGScloud is downloaded:

```
> NGScloud
```

You can type too:

```
> NGScloud --mode=gui
```

And to run NGScloud in console mode, type the command:

```
> NGScloud --mode=console
```

Configuring your first NGScloud environment

NGScloud philosophy is based on the "cluster" concept. A cluster is a set of virtual machines of an AWS instance type. Each instance type has its hardware features: machine type, CPU number, memory amount, etc. You can consult these features in <https://aws.amazon.com/ec2/instance-types/>.

When a cluster is created, it has only a virtual machine named "master node". After the master node creation, you can add "subsidiary nodes" if they are necessary to run some processes in parallel. In this case, the new job will be run in the node determined according to the workload.

"Data volumes" allow us to save data and keep them even if there is not any cluster created. NGScloud always uses the following volumes:

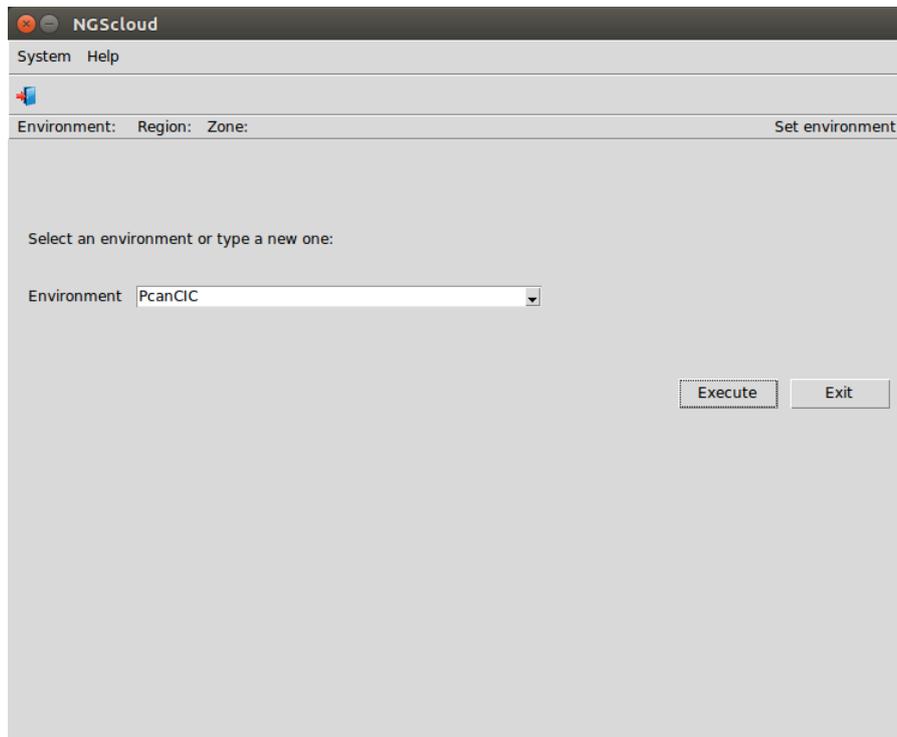
- (1) "application volume": to install the bioinformatic applications; this volume is mandatory.
- (2) "read volume": to upload the read files of the experiments; this volume is mandatory.
- (3) "result volume": to store the results of the experiments; this volume is mandatory.

(4) "reference volume": to hold reference genomes/transcriptomes and information about gene structure that may be used by some applications to refine the results; this volume is optional.

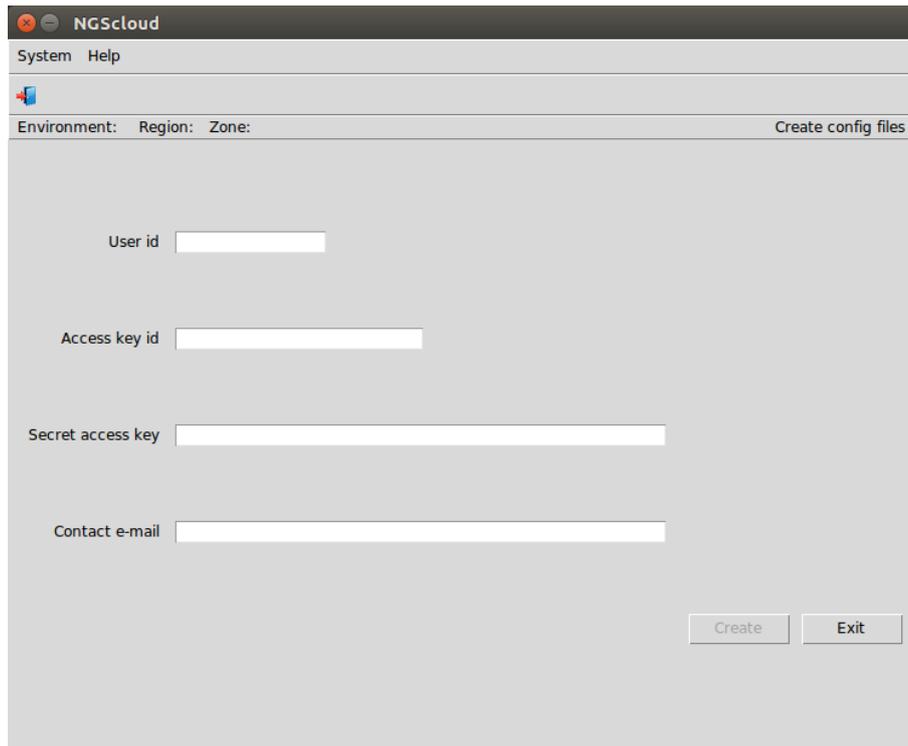
(5) "database volume": to hold data from reference sequence databases (RefSeq) used by some annotation processes; this volume is optional.

Before starting an experiment in NGScloud, it is very important to estimate the sizes each volume will need, particularly for the *reads* and *results volumes*. The *reads volume* must be able to store the uploaded read files and new read dataset obtained after trimming, if needed. The results of running the bioinformatic applications implemented in NGScloud may have big size; therefore, the *results volume* size must be set accordingly. We recommend configuring unique *results* and *reads volumes* for each experiment.

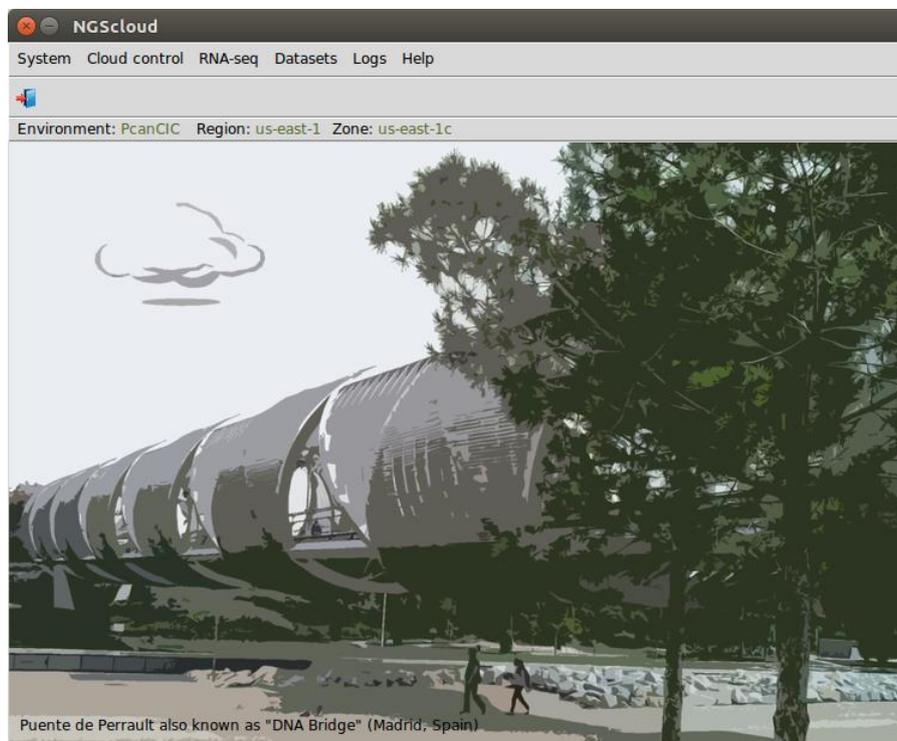
An "environment" identifies a user, a volume set and the AWS zone where processes run and the volumes are stored. When starting NGScloud for the first time, you must type the name of the environment (alphanumeric characters only) in the box *Environment*. E. g. **PcanCIC**:



In the next window, you must type your *AWS user id*, *access key id* and *secret access key*. A *contact e-mail* address is required too. This e-mail address is used to warn you when a submitted job ends.



Once the first steps are completed, the main window of NGScloud is shown, and it is ready to use.



NGScloud is structured in several menus:

System menu

Just to exit the application.

Cloud control menu

This menu contains all the items related to:

- Set an environment
- NGScloud configuration and security
- Creation of clusters, nodes and volumes, and options to operate with them
- Setup of bioinformatic applications in a cluster
- Open a terminal in a cluster node

RNA-seq menu

Here, all options related to RNA-seq experiments are implemented:

- Quality, trimming and digital normalization of reads
- *De-novo* assembly and reference-based assembly
- Assembly quality assessment and transcript quantification
- Transcriptome filtering
- Annotation

Datasets menu

The options included here allow to handle the read, reference, database and result datasets:

- List dataset
- Upload read files from local computers to a cluster.
- Download the results files from a cluster to the local computer.
- Compress and decompress files in a cluster.
- Remove datasets.

Logs menu

This menu allows the access to logs of submissions in the local computer and logs of results in the clusters.

Help menu

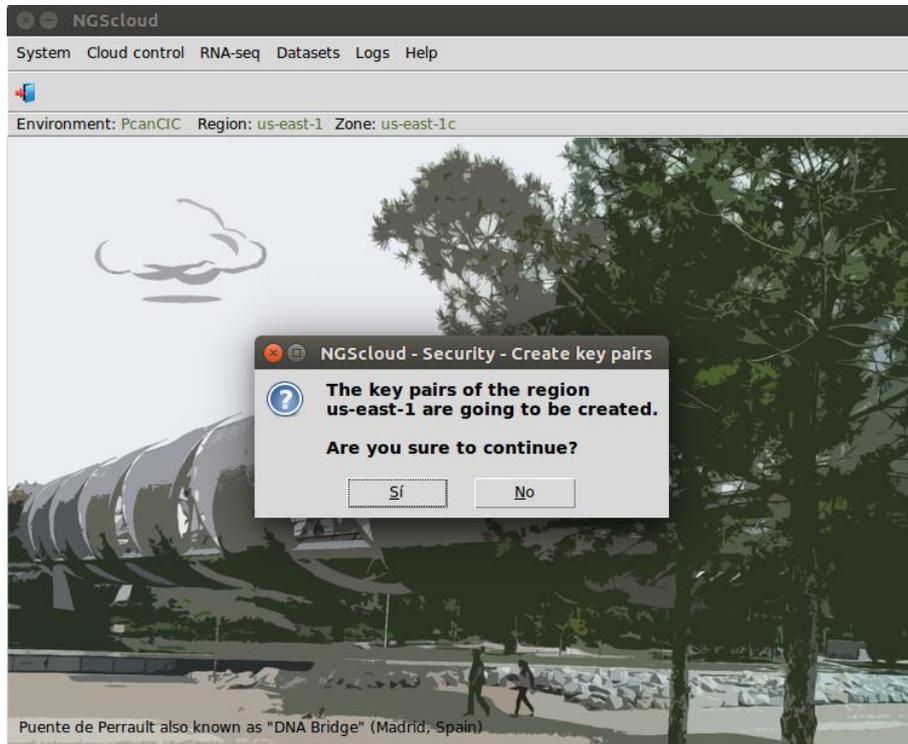
It contains the documentation of the application.

Before using any of the options in the menus, "key pairs" need to be created. Key pairs are used to encrypt and decrypt login information. You can create key pairs, by selecting the menu item with this path:

Main menu > Cloud control > Security > Create key pairs

A dialog box will be raised to confirm the action.

A key pair is valid for all zones within a region. Then if you have created a key pair in a zone, and you change to another zone of the same region, you do not have to create the key pair again.



A step by step example

We have sequencing data corresponding to a RNA-seq Illumina library of an experiment about the process of cicatrization after wounding the xylem of the stem of the Canary Island pine (*Pinus canariensis*). The next table shows the size characteristics of the read files yield by the NGS platform:

Library	File	Read number	Compressed size (in B)	Decompressed size (in B)
Pcan-CIC	Pcan-CIC_1.fastq.gz	21.444.414	1.648.808.931	5.178.672.000
	Pcan-CIC_2.fastq.gz	21.444.414	1.632.988.106	5.178.672.000
		2	42.888.828	10.357.344.000

In this example, we are going to review the quality of reads, to trim read ends with bad scores and to assembly the reads yielding a transcriptome.

The steps are the following:

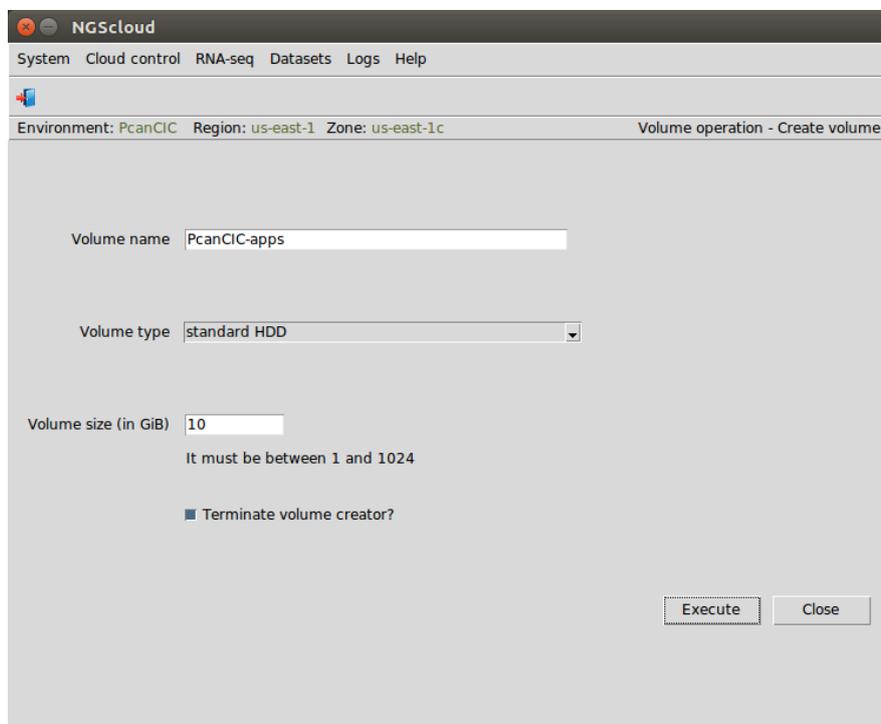
- Create volumes
- Link volumes in cluster templates
- Create a cluster with the t2.micro template
- Upload the read files to the cluster
- Setup the bioinformatic applications in the cluster
- Review the quality of reads using FastQC
- Trim the reads using Trimmomatic
- Terminate the cluster with a t2.micro template and create another cluster with a r3.4xlarge template
- Assembly the reads using Trinity
- Evaluate the transcriptome quality using RSEM-EVAL
- Terminate the cluster with a r3.4xlarge template and create another cluster with r3.xlarge template
- Transcriptome filtering using transcript-filter
- Transcriptome clustering using CD-HIT-EST
- Terminate the cluster with a r3.xlarge template and create another cluster with c3.xlarge template
- Upload the protein database to the cluster
- Add nodes to the cluster with a c3.xlarge template
- Annotate the filtered and clustered transcriptome using transcriptome-blastx
- Terminate the cluster with a c3.xlarge template and create another cluster with t2.micro template
- Download the transcriptome, evaluation and annotation files
- Terminate the cluster with a t2.micro template

Create volumes

First, we need to create the data volumes to have persistent storage of the installed bioinformatic applications, read data, and results. We have to decide the type and the size of each volume. Ten GiB can be enough size for the *app volume*. In this case, we choose a standard HDD type, given the size and the cost per GiB of each volume type. To create the volume, select the menu item with this path:

Main menu > Cloud control > Volume Operation > Create volume

In the raised window, we type **PcanCIC-apps** in *Volume name* textbox, we select **standard HDD** in *Volume type* combobox, and we type **10** in *Volume size (in GiB)* textbox; we untick *Terminate volume creator?* checkbox; and we press the *Execute* button:



A volume creator instance will be started to create and format the volume. When the volume is created and formatted, the volume creator will not be terminated (we have unticked *Terminate volume creator?* checkbox), which will allow us to create the other volumes quickly.

The sizes of read, reference, database and result volumes are 20 GiB, 5GiB, 5GiB and 100 GiB, respectively. The volume type is a standard HDD for both cases. We repeat the steps done to configure the app volume, to configure the reads and results volumes, making sure that the flag of *Terminate volume creator?* checkbox must be ticked when creating our last volume.

We can review the created volumes selecting the menu item as follows:

Main menu > Cloud control > Volume Operation > List volumes

The raised window will show information about the created volumes:

Zone	Volume Name	Volume Id	Size (GiB)	State	Attachments
us-east-1c	PcanCIC-apps	vol-06ca1b476b0410bf1	10	available	0
us-east-1c	PcanCIC-databases	vol-0a91e05f13ba61536	5	available	0
us-east-1c	PcanCIC-reads	vol-00eee59f776a0fd83	20	available	0
us-east-1c	PcanCIC-references	vol-0628624c8d1eeff88	5	available	0
us-east-1c	PcanCIC-results	vol-04392abe34f0e9762	100	available	0

Link volumes in cluster templates

A cluster template identifies the instance type, the machine image and other characteristics of a cluster when it is booted.

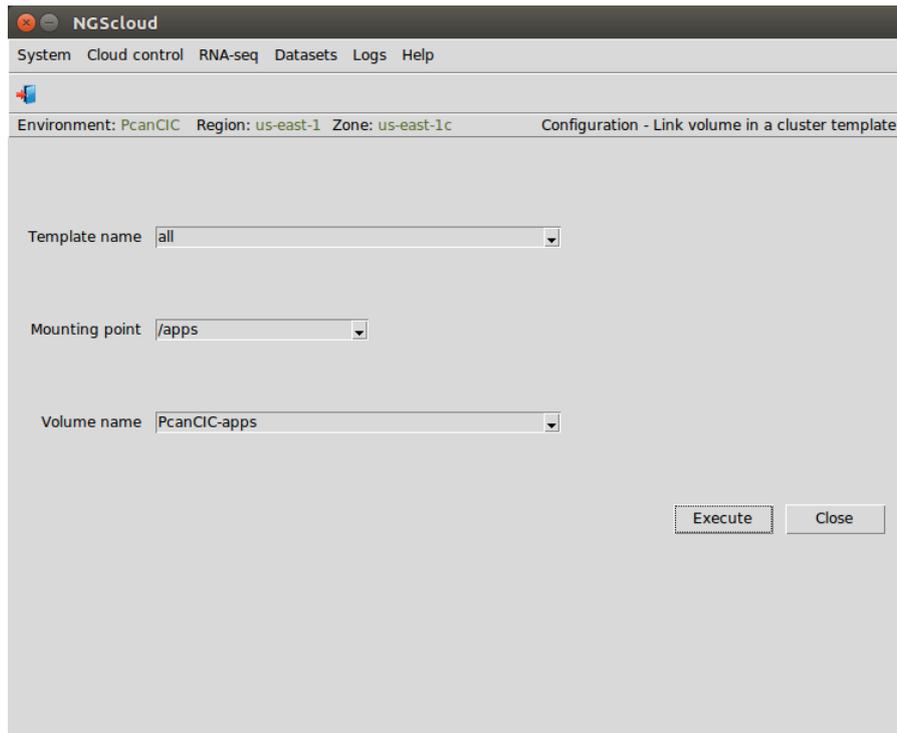
We must link the created volumes to the cluster templates so that the volumes are automatically attached at the start of the cluster. There are five mounting points.

- /apps: to the application volume;
- /references, to the reference volume;
- /databases to database volume;
- /reads to read volume
- /results to result volume

To link a volume to a cluster template, select the menu item with this path:

Main menu > Cloud control > Configuration > Link volume in a cluster template

In the raised window, we must fill in the boxes with the information relative to the template (we can choose a specific template or *all* templates), the mounting point and the name of the volume in *Volume name*. To link the application volume to all the templates, we select **all** in *Template name* combobox, **/apps** in *Mounting point* combobox, and **PcanCIC-apps** in *Volume name* combobox. Then we have to press the *Execute* button.



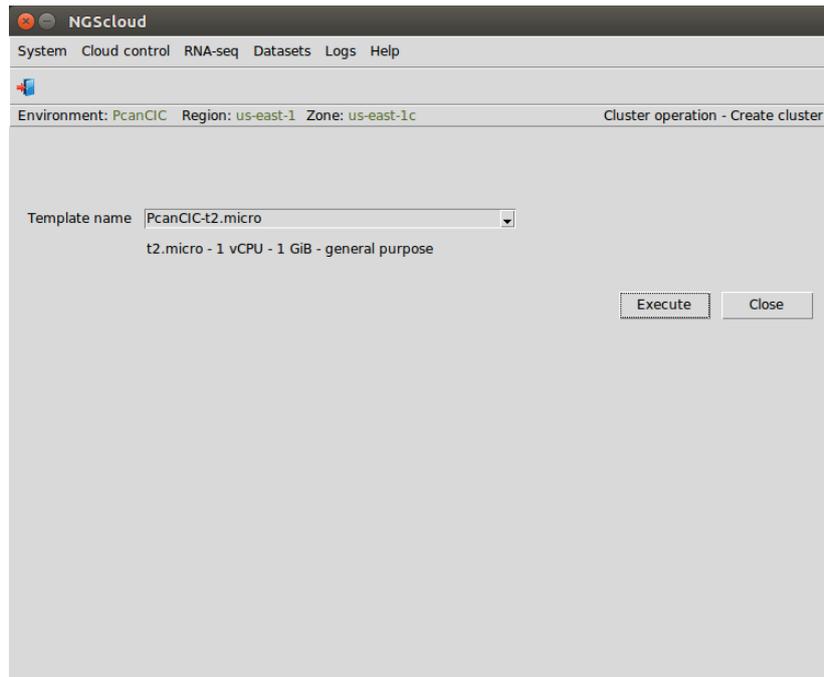
Then we repeat this step for the other two volumes created earlier.

Create cluster with the t2.micro template

Now we create a cluster with a t2.micro template, 1 CPU and 1 GiB of RAM, because the read file upload and the read trimming require few hardware resources. We select the menu item with this path:

Main menu > Cloud control > Cluster operation > Create cluster

In the raised window, we select **PcanCIC-t2.micro**, the template corresponding to a t2.micro instance type, in *Template name* combo-box; and then we press the *Execute* button:



A window is raised displaying the run log:

```

NGScloud - Cluster operation - Create cluster - Log
This process might take several minutes. Do not close this window, please wait!
*****
Verifying process requirements ...
Process requirements are OK.
*****
Creating the cluster PcanCIC-t2.micro using StarCluster ...

StarCluster - (http://star.mit.edu/cluster) (v. 0.95.6)
Software Tools for Academics and Researchers (STAR)
Please submit bug reports to starcluster@mit.edu

>>> Validating cluster template settings...
>>> Cluster template settings are valid
>>> Starting cluster...
>>> Launching a 1-node cluster...
>>> Creating security group @sc-PcanCIC-t2.micro...
Reservation:r-0fb43e9cea93167a5
>>> Waiting for instances to propagate...
0/1 | | 0%
1/1 | | 100%
>>> Waiting for cluster to come up... (updating every 30s)
>>> Waiting for all nodes to be in a 'running' state...
0/1 | | 0%
0/1 | | 0%
1/1 | | 100%
>>> Waiting for SSH to come up on all nodes...
0/1 | | 0%
0/1 | | 0%
0/1 | | 0%
0/1 | | 0%
0/1 | | 0%
0/1 | | 0%
0/1 | | 0%
0/1 | | 0%
0/1 | | 0%
0/1 | | 0%
1/1 | | 100%
>>> Waiting for cluster to come up took 1.106 mins
>>> The master node is ec2-54-152-166-158.compute-1.amazonaws.com
>>> Configuring cluster...
>>> Attaching volume vol-00eee59f776a0fd83 to master node on /dev/sdx ...
>>> Attaching volume vol-06ca1b476b0410bf1 to master node on /dev/sdx

```

When the cluster is started, an infrastructure software will be installed. At the end of the installation, an email is sent, informing of its completion:



Upload the read files to the cluster

Each task related to datasets or to the run of a bioinformatic program has:

- First, a window to help to select datasets or specific files. A config file is created according to the user selection and default values of the parameters of the program.
- Then, a window where the parameters of the program are shown with an explanation of its meaning. Every parameter has a default value that can be changed.
- Finally, a building of a bash script to run the program using the config file and the submission of this script to the cluster.

To select specific files, the first window has a text box where a pattern must be entered. The pattern must be a Python regular expression. A regular expression is used to find a string in other string(s). The pattern is formed by a sequence of characters; some of them have a special meaning, e.g. "." means "any character except newline" and "*" means "0 or more repetitions of the preceding element". You can learn about Python regular expressions at: <https://docs.python.org/3.6/library/re.html>

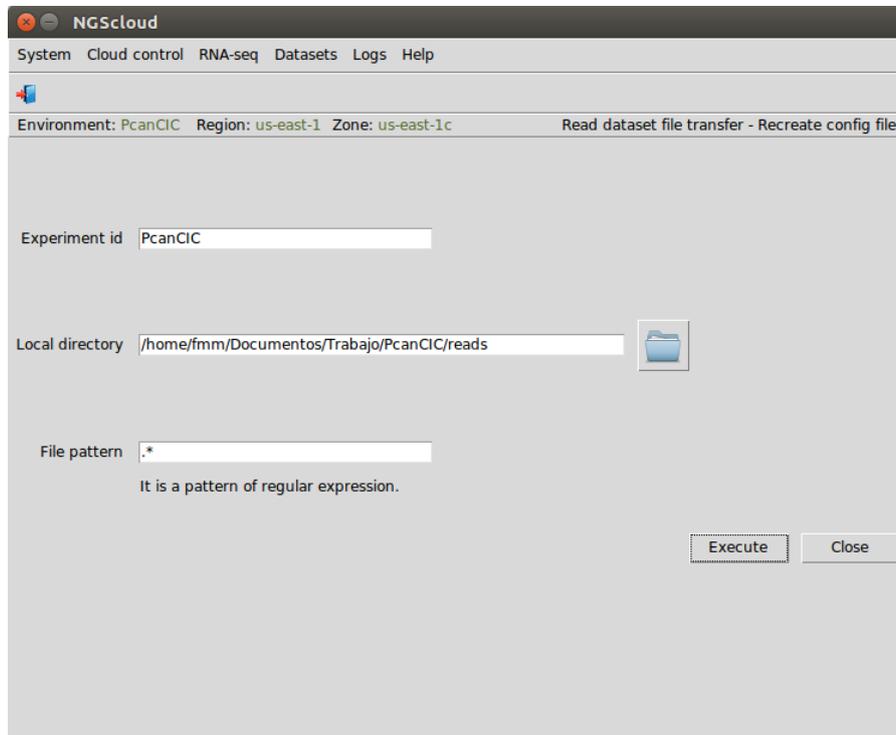
Perhaps, these examples are useful for your selection:

Pattern	Selection
.*	all the files
transcriptome.fasta	the file whose name is "transcriptome.fasta"
.*fastq	the files whose name ended in "fastq"
.*fastq.gz	the files whose name ended in "fastq".gz"
.*Pcan.*	the files whose name contains the characters "Pcan"
.*Pcan.*fastq	the files whose name contains the characters "Pcan" and ends in "fastq"

To create a config file to upload the read files to a cluster, select the menu item with this path:

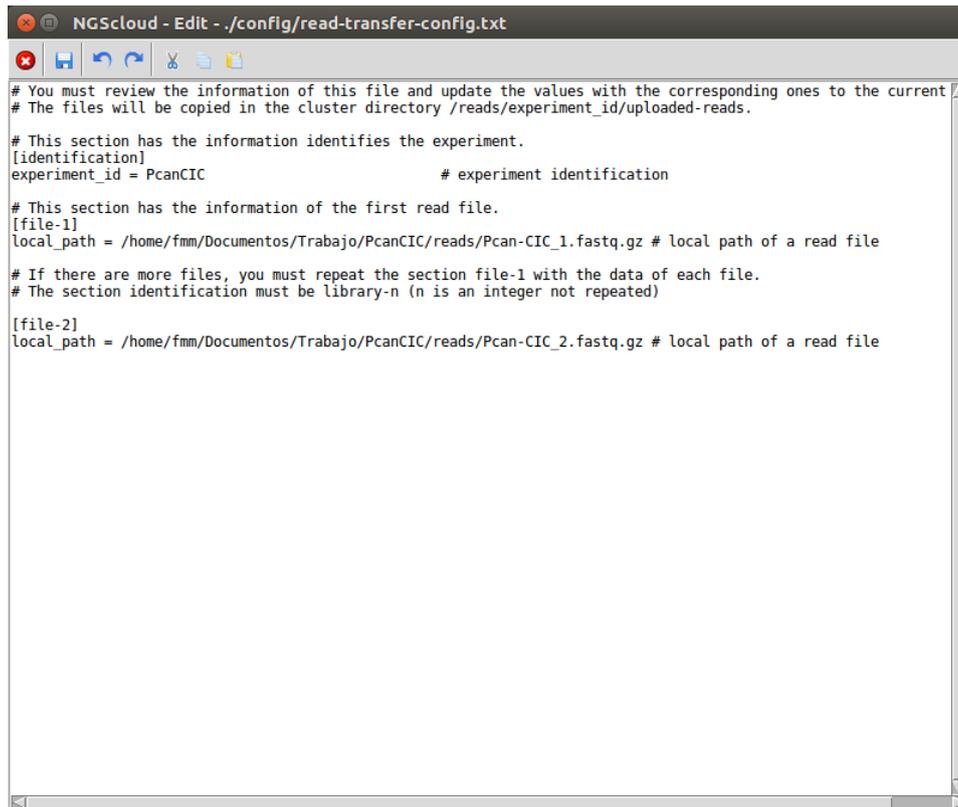
Main menu > Dataset > Read dataset file transfer > Recreate config file

In the raised window, we type **PcanCIC** in the *Experiment id* textbox, the local directory where the files are in the *Local directory* textbox (or we select it using the next button), and we type **.*** as the pattern to select the files in the *File pattern* textbox. Then we press the *Execute* button:



In the next window, we can edit the config file created, and remove files or add new files if the file pattern has not selected the appropriate ones. This window is a text editor, and can be easily modified. When we save the configuration file, the modifications are validated. If there are errors, a list of them is displayed.

In this example, we can notice that the configure file has three sections: *identification*, with the experiment identification; *file-1*, with the local path of the first read file; and *file-2* with the second one:



```

# You must review the information of this file and update the values with the corresponding ones to the current
# The files will be copied in the cluster directory /reads/experiment_id/uploaded-reads.

# This section has the information identifies the experiment.
[identification]
experiment_id = PcanCIC                # experiment identification

# This section has the information of the first read file.
[file-1]
local_path = /home/fmm/Documentos/Trabajo/PcanCIC/reads/Pcan-CIC_1.fastq.gz # local path of a read file

# If there are more files, you must repeat the section file-1 with the data of each file.
# The section identification must be library-n (n is an integer not repeated)

[file-2]
local_path = /home/fmm/Documentos/Trabajo/PcanCIC/reads/Pcan-CIC_2.fastq.gz # local path of a read file

```

This window is a text editor, and can be easily modified. When we save the configuration file, the modifications are validated. If there are errors, a list of them is displayed.

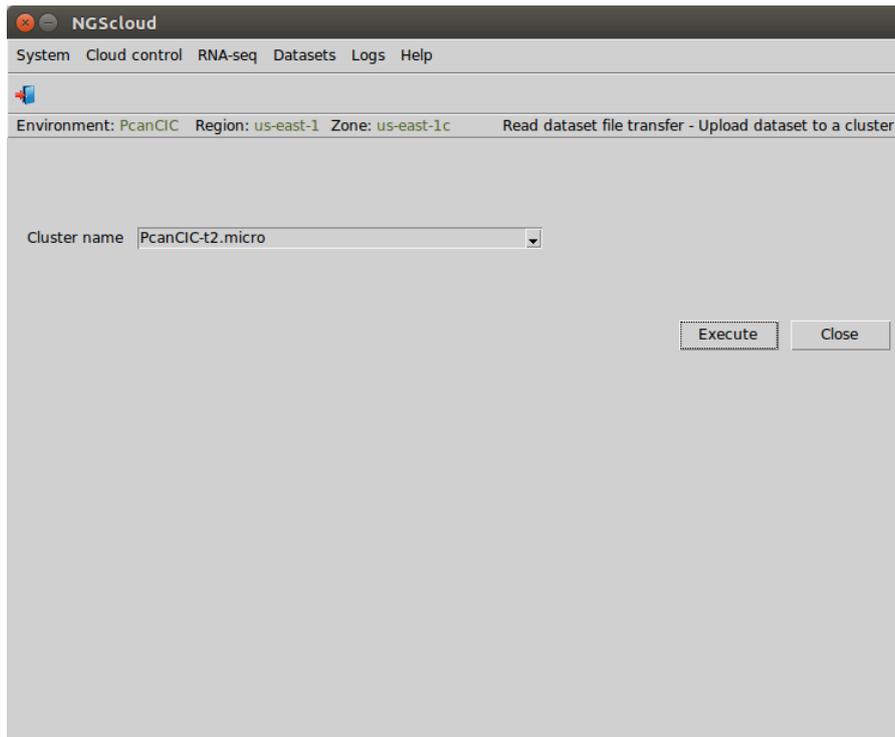
In this example, we can notice that the configure file has three sections: *identification*, with the experiment identification; *file-1*, with the local path of the first read file; and *file-2* with the second one.

It is convenient to perform the file transfer steps when an Internet connection with a large bandwidth is available due to the large size of many of the files necessary to perform full RNA-seq analysis.

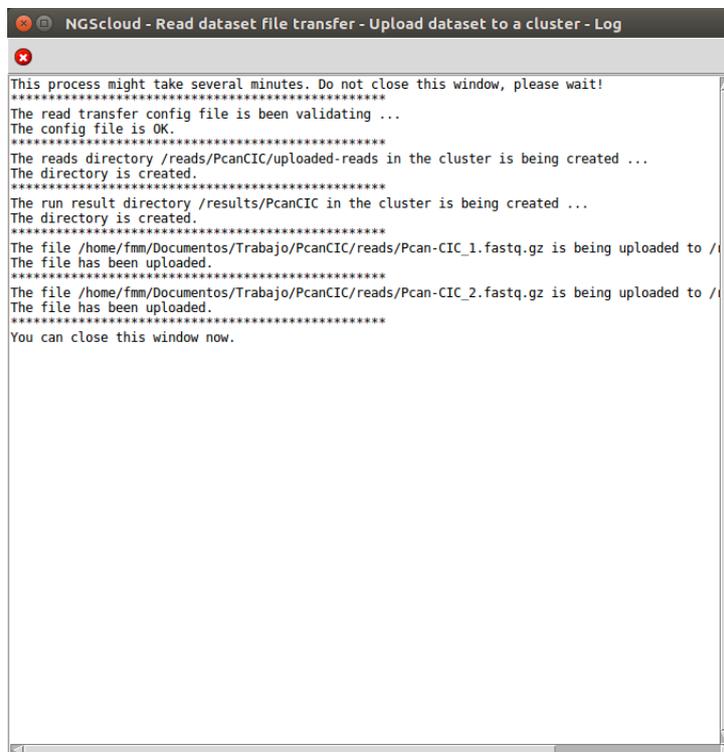
To upload the read files to the cluster, we select the menu item with this path:

Main menu > Dataset > Read dataset file transfer > Upload dataset to a cluster

In the raised window, we select **PcanCIC-t2.micro** in *Cluster name* combo-box; and then we press the *Execute* button:



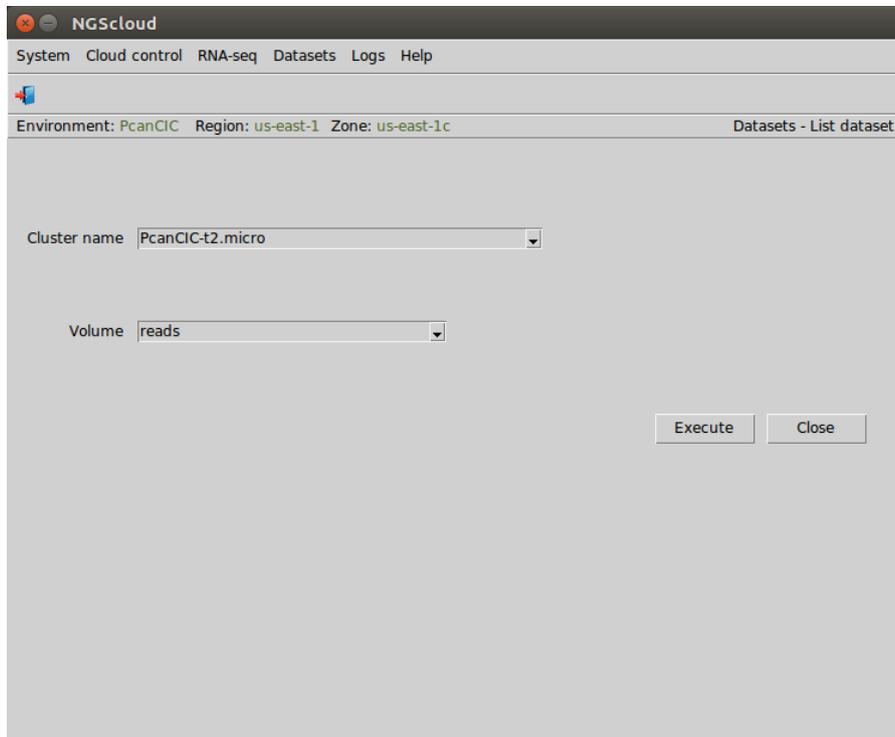
A window is raised with the upload log:



Now, we are going to review the uploaded files. We select the menu item with this path:

Main menu > Dataset > List dataset

In the raised window, we select **PcanCIC-t2.micro** in the *Cluster name* combo-box and **reads** in the *Volume* combo-box. Then we press the *Execute* button:

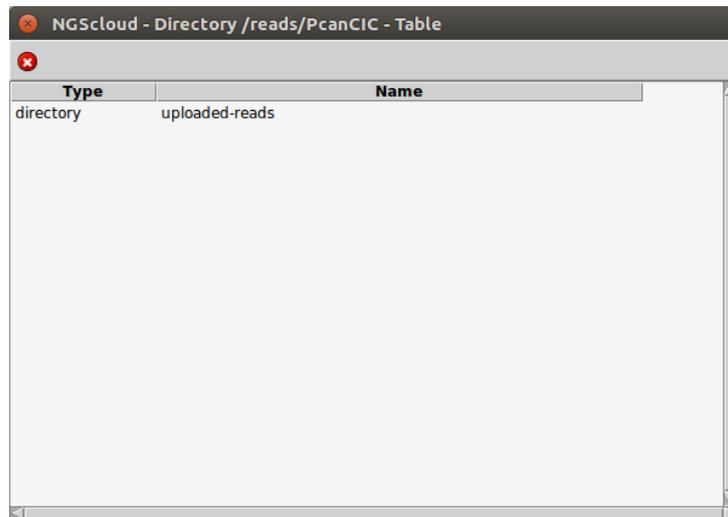


We can check the experiments whose read files have been uploaded in the next window. We click in **PcanCIC** row:

The screenshot shows the NGScloud - Directory /reads - Table window with the following table:

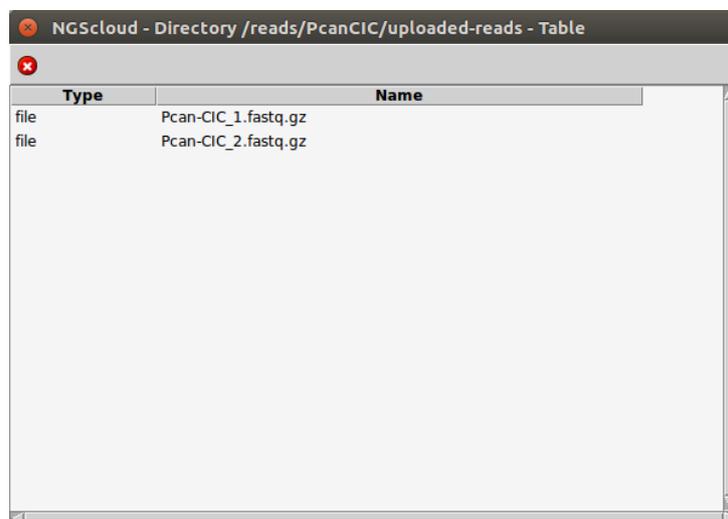
Type	Name
directory	Athaliana01x
directory	PcanCIC
directory	test

Now, a window with the read datasets of the experiment PcanCIC is shown:



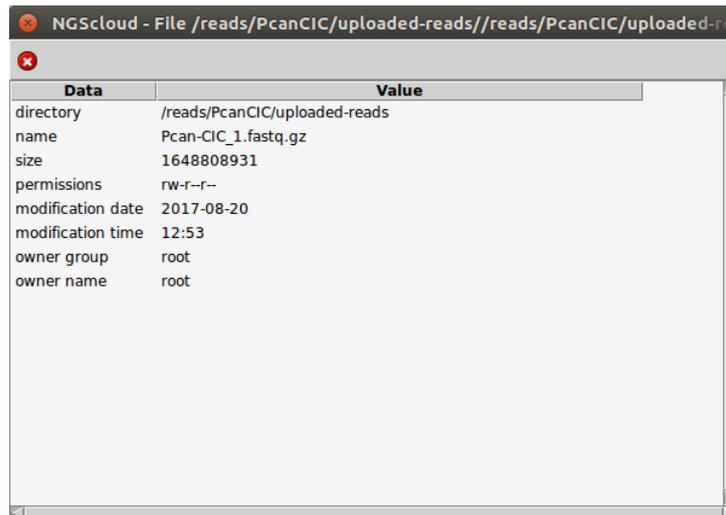
Type	Name
directory	uploaded-reads

So far, we only have one: the dataset corresponding to the uploaded-reads. We click on it and another window appears with the content of the uploaded-reads dataset. In this case, the two files are shown.



Type	Name
file	Pcan-CIC_1.fastq.gz
file	Pcan-CIC_2.fastq.gz

If we click in a file row, e.g. the **Pcan-CIC_1.fastq.gz** one, the characteristics of this file are listed:



Data	Value
directory	/reads/PcanCIC/uploaded-reads
name	Pcan-CIC_1.fastq.gz
size	1648808931
permissions	rw-r--
modification date	2017-08-20
modification time	12:53
owner group	root
owner name	root

Setup bioinformatic applications in the cluster

Bioconda is necessary to setup the bioinformatic applications. To setup Bioconda in the application volume in a cluster, select the menu item with this path:

Main menu > Bioinfo software setup > Miniconda3 (Python & Bioconda environment)

And, in the next windows, type the cluster name.

To setup FastQC in the application volume in a cluster, select the menu item with this path:

Main menu > Bioinfo software setup > FastQC

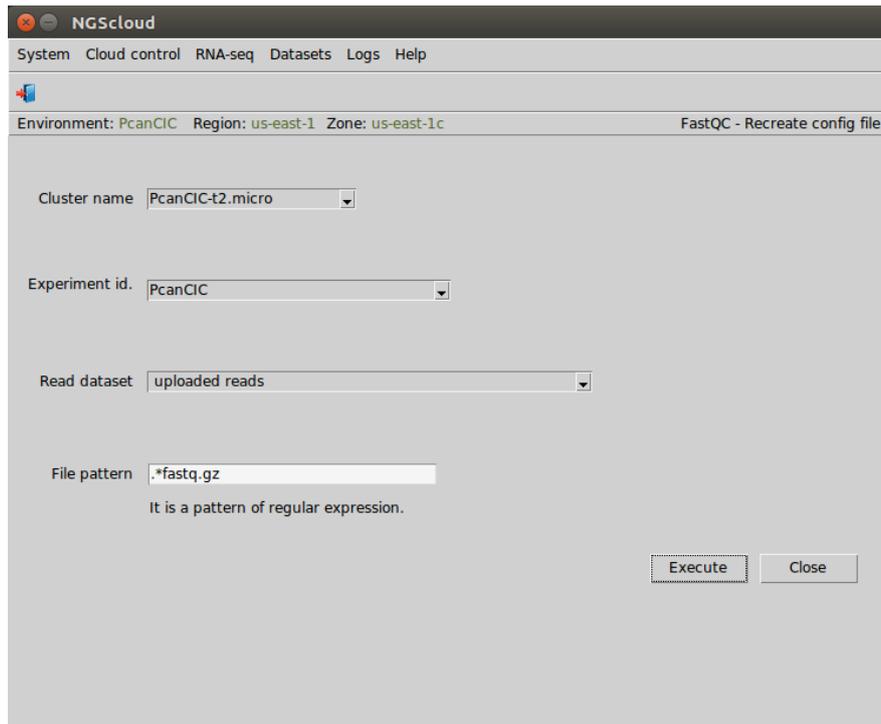
And in the next windows, type the cluster name. Also, install Trimmomatic and Trinity as in the setup FastQC.

Review the quality of reads using FastQC

Now we are going to review the quality of reads using FastQC. First, we create the configuration file, we select the menu item with this path:

Main menu > RNA-seq > Read quality > FastQC > Recreate config file

In the raised window, we select **PcanCIC-t2.micro** in the *Cluster name* combo-box, **PcanCIC** in the *Experiment id* combo-box, **uploaded reads** in the *Read dataset* combo-box, and we type **.*fastq.gz** as the pattern to select the files in the *File pattern* textbox. Then we press the *Execute* button:



In the next window, we can inspect the config file. In this example, there are four sections: *identification*, with the experiment and the read dataset identifications; *FastQC parameters*, with the thread number parameter of FastQC (we modify its value to 1); *file-1*, with the local path of the first read file; and *file-2* with the path of the second one:

```

NGScloud - Edit - ./config/fastqc-config.txt
# You must review the information of this file and update the values with the corresponding ones to the current
#
# The files must be located in the cluster directory /reads/experiment_id/read_dataset_id
# The experiment_id and read_dataset_id names are fixed in the identification section.
#
# You can consult the parameters of FastQC and their meaning in http://www.bioinformatics.babraham.ac.uk/project
# This section has the information identifies the experiment.
[identification]
experiment_id = PcanCIC                # experiment identification
read_dataset_id = uploaded-reads      # read dataset identification

# This section has the information to set the FastQC parameters
[FastQC parameters]
threads = 1                            # number of threads for use

# This section has the information of the first file.
[file-1]
file_name = Pcan-CIC_1.fastq.gz        # read file name

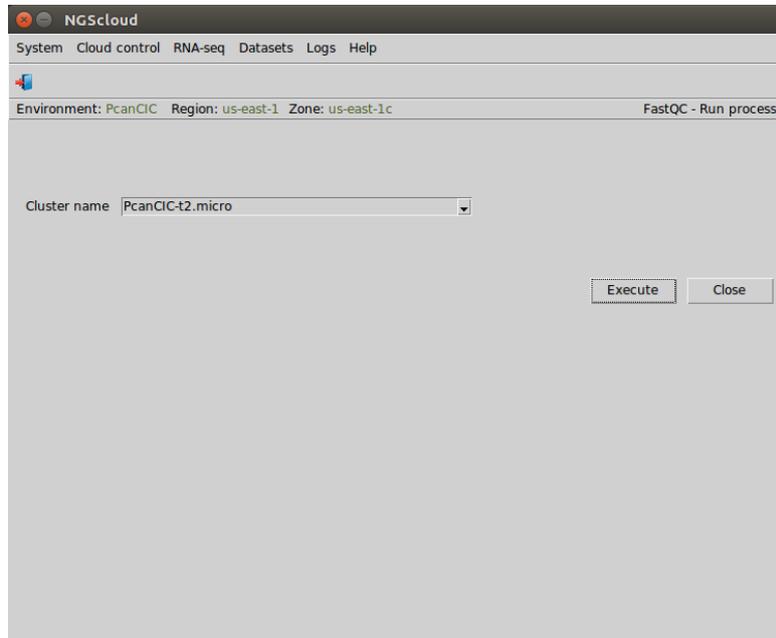
# If there are more files, you must repeat the section file-1 with the data of each file.
# The section identification must be library-n (n is an integer not repeated)
[file-2]
file_name = Pcan-CIC_2.fastq.gz        # read file name

```

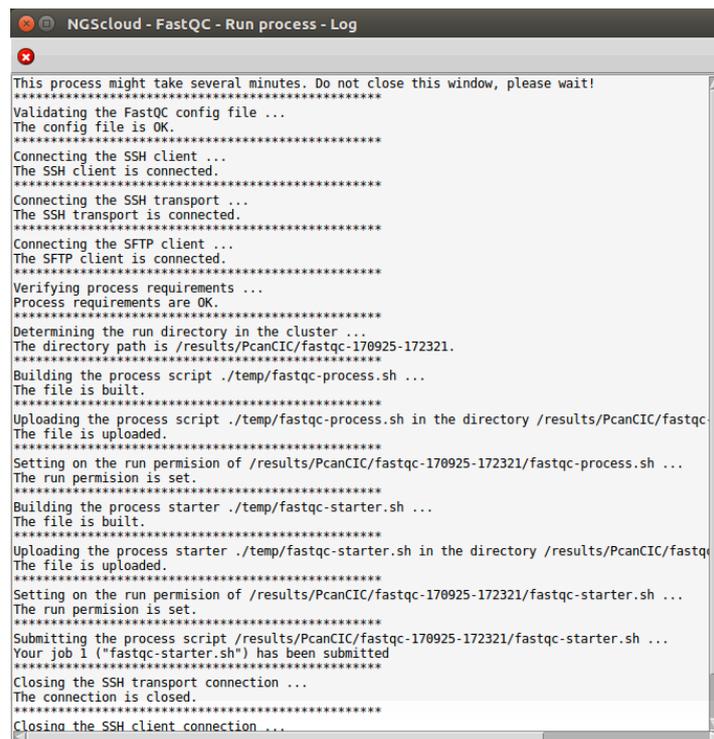
To run the quality process in the cluster, we select the menu item with this path:

Main menu > RNA-seq > Read quality > FastQC > Run read quality process

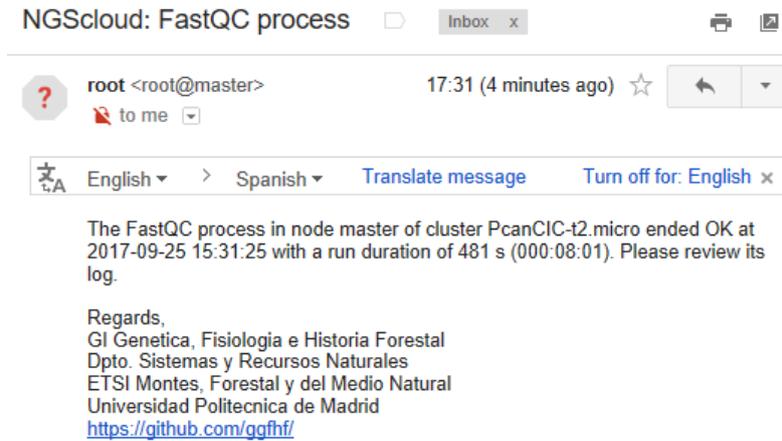
In the raised window, we select **PcanCIC-t2.micro** in the *Cluster name* combo-box; and then we press the *Execute* button:



A window is raised with the submission log:



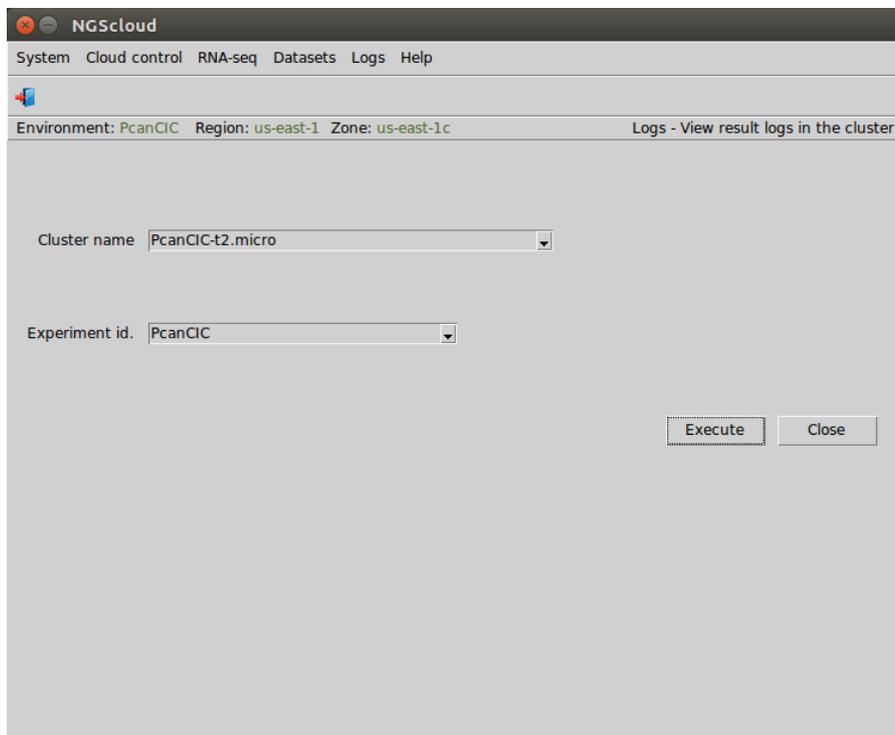
At the end of the run, an email is sent, informing of its completion:



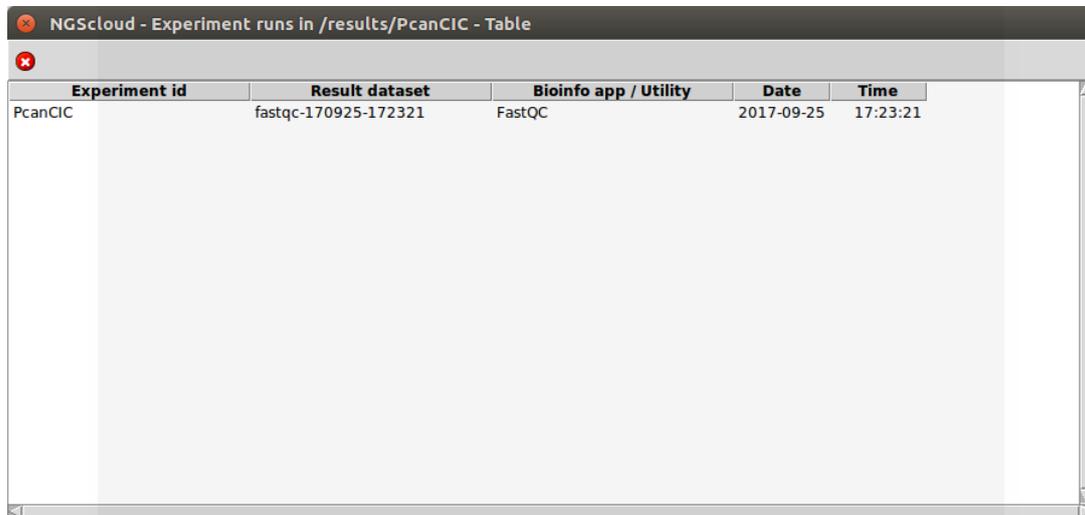
We can view the process log during and after the run. To do so, we select the menu item with this path:

Main menu > Logs > View result logs in the cluster

In the raised window, we select **PcanCIC-t2.micro** in the *Cluster name* combo-box and **PcanCIC** in the *Experiment id* combo-box; and then we press the *Execute* button:

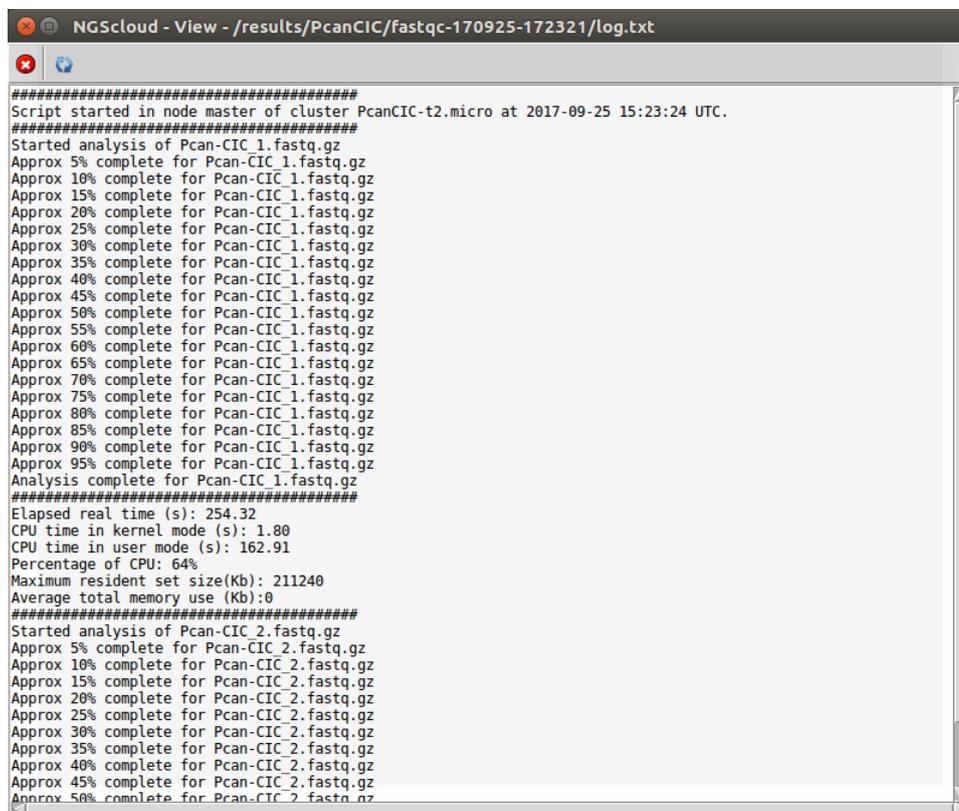


Now, a window with the result datasets of each bioinformatic program run in the experiment PcanCIC is shown:



Experiment id	Result dataset	Bioinfo app / Utility	Date	Time
PcanCIC	fastqc-170925-172321	FastQC	2017-09-25	17:23:21

So far, we have only performed a single run: the dataset **fastqc-170925-172321** corresponding to the last (and unique) FastQC run. Clicking on it, another window appears with its corresponding log.



```

#####
Script started in node master of cluster PcanCIC-t2.micro at 2017-09-25 15:23:24 UTC.
#####
Started analysis of Pcan-CIC_1.fastq.gz
Approx 5% complete for Pcan-CIC_1.fastq.gz
Approx 10% complete for Pcan-CIC_1.fastq.gz
Approx 15% complete for Pcan-CIC_1.fastq.gz
Approx 20% complete for Pcan-CIC_1.fastq.gz
Approx 25% complete for Pcan-CIC_1.fastq.gz
Approx 30% complete for Pcan-CIC_1.fastq.gz
Approx 35% complete for Pcan-CIC_1.fastq.gz
Approx 40% complete for Pcan-CIC_1.fastq.gz
Approx 45% complete for Pcan-CIC_1.fastq.gz
Approx 50% complete for Pcan-CIC_1.fastq.gz
Approx 55% complete for Pcan-CIC_1.fastq.gz
Approx 60% complete for Pcan-CIC_1.fastq.gz
Approx 65% complete for Pcan-CIC_1.fastq.gz
Approx 70% complete for Pcan-CIC_1.fastq.gz
Approx 75% complete for Pcan-CIC_1.fastq.gz
Approx 80% complete for Pcan-CIC_1.fastq.gz
Approx 85% complete for Pcan-CIC_1.fastq.gz
Approx 90% complete for Pcan-CIC_1.fastq.gz
Approx 95% complete for Pcan-CIC_1.fastq.gz
Analysis complete for Pcan-CIC_1.fastq.gz
#####
Elapsed real time (s): 254.32
CPU time in kernel mode (s): 1.80
CPU time in user mode (s): 162.91
Percentage of CPU: 64%
Maximum resident set size(Kb): 211240
Average total memory use (Kb):0
#####
Started analysis of Pcan-CIC_2.fastq.gz
Approx 5% complete for Pcan-CIC_2.fastq.gz
Approx 10% complete for Pcan-CIC_2.fastq.gz
Approx 15% complete for Pcan-CIC_2.fastq.gz
Approx 20% complete for Pcan-CIC_2.fastq.gz
Approx 25% complete for Pcan-CIC_2.fastq.gz
Approx 30% complete for Pcan-CIC_2.fastq.gz
Approx 35% complete for Pcan-CIC_2.fastq.gz
Approx 40% complete for Pcan-CIC_2.fastq.gz
Approx 45% complete for Pcan-CIC_2.fastq.gz
Approx 50% complete for Pcan-CIC_2.fastq.gz

```

In the toolbar, there is a button to refresh the run status. Clicking it, the log will be updated.

All the process logs have:

- A header with the node where the script runs and the time when it started.
- Information about the elapsed time, the CPU usage and the maximum memory is displayed for each run of the bioinformatic program.
- At the bottom, a summary with the status (OK, if all the programs have ended without errors; WRONG, otherwise), the end time, and the duration of the script run.

```

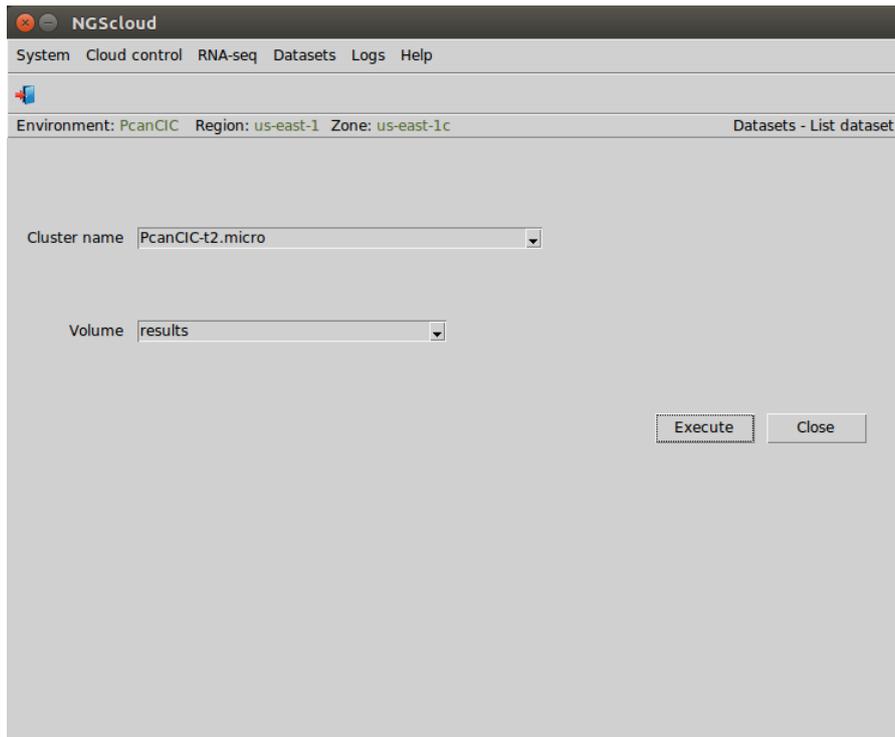
NGScloud - View - /results/PcanCIC/fastqc-170925-172321/log.txt
Approx 65% complete for Pcan-CIC 1.fastq.gz
Approx 90% complete for Pcan-CIC 1.fastq.gz
Approx 95% complete for Pcan-CIC 1.fastq.gz
Analysis complete for Pcan-CIC 1.fastq.gz
#####
Elapsed real time (s): 254.32
CPU time in kernel mode (s): 1.80
CPU time in user mode (s): 162.91
Percentage of CPU: 64%
Maximum resident set size(Kb): 211240
Average total memory use (Kb):0
#####
Started analysis of Pcan-CIC 2.fastq.gz
Approx 5% complete for Pcan-CIC 2.fastq.gz
Approx 10% complete for Pcan-CIC 2.fastq.gz
Approx 15% complete for Pcan-CIC 2.fastq.gz
Approx 20% complete for Pcan-CIC 2.fastq.gz
Approx 25% complete for Pcan-CIC 2.fastq.gz
Approx 30% complete for Pcan-CIC 2.fastq.gz
Approx 35% complete for Pcan-CIC 2.fastq.gz
Approx 40% complete for Pcan-CIC 2.fastq.gz
Approx 45% complete for Pcan-CIC 2.fastq.gz
Approx 50% complete for Pcan-CIC 2.fastq.gz
Approx 55% complete for Pcan-CIC 2.fastq.gz
Approx 60% complete for Pcan-CIC 2.fastq.gz
Approx 65% complete for Pcan-CIC 2.fastq.gz
Approx 70% complete for Pcan-CIC 2.fastq.gz
Approx 75% complete for Pcan-CIC 2.fastq.gz
Approx 80% complete for Pcan-CIC 2.fastq.gz
Approx 85% complete for Pcan-CIC 2.fastq.gz
Approx 90% complete for Pcan-CIC 2.fastq.gz
Approx 95% complete for Pcan-CIC 2.fastq.gz
Analysis complete for Pcan-CIC 2.fastq.gz
#####
Elapsed real time (s): 227.36
CPU time in kernel mode (s): 1.96
CPU time in user mode (s): 162.80
Percentage of CPU: 72%
Maximum resident set size(Kb): 211496
Average total memory use (Kb):0
#####
Script ended OK at 2017-09-25 15:31:25 UTC with a run duration of 481 s (000:08:01).
#####

```

In order to access a list with the output files generated by FastQC, we select the menu item with this path:

Main menu > Dataset > List dataset

In the raised window, we select **PcanCIC-t2.micro** in the *Cluster name* combo-box and **results** in the *Volume* combo-box. Then we press the *Execute* button:

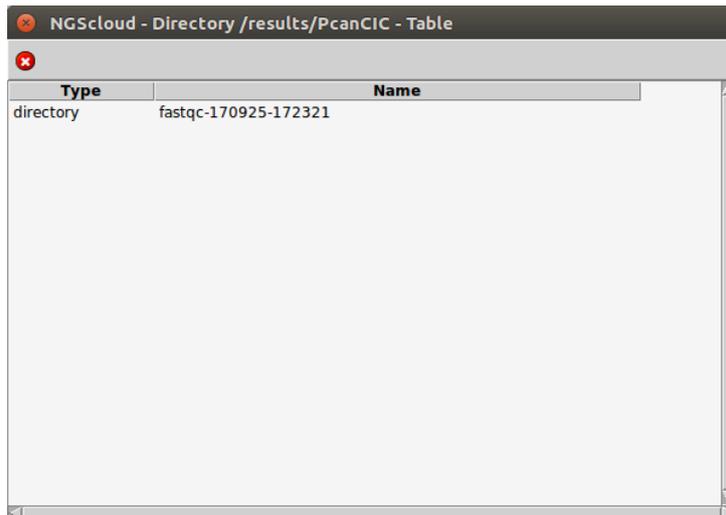


To inspect the experiments that have result datasets, we click in the **PcanCIC** row:

The screenshot shows a window titled 'NGScloud - Directory /results - Table'. It contains a table with two columns: 'Type' and 'Name'. The table lists four entries:

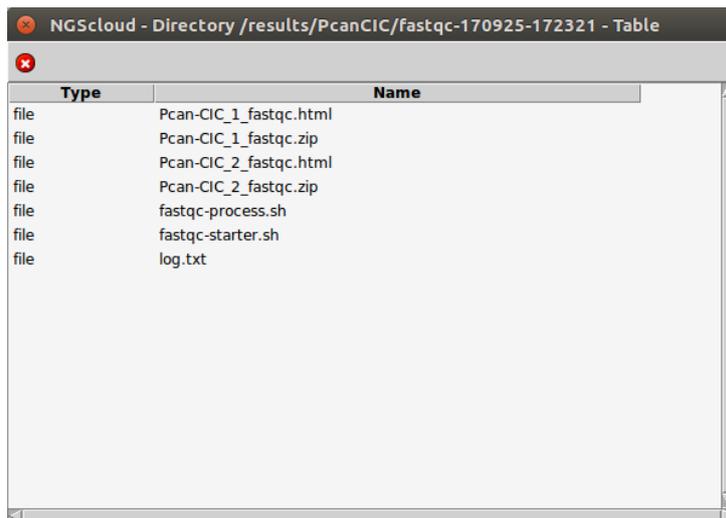
Type	Name
directory	Athaliana01x
directory	PcanCIC
directory	database
directory	test

Next, a window with the result datasets of the experiment PcanCIC is shown:



Type	Name
directory	fastqc-170925-172321

So far, we only have one: the dataset corresponding to the FastQC analysis recently completed. We click on it and another window appears with the content of the files corresponding to this analysis.



Type	Name
file	Pcan-CIC_1_fastqc.html
file	Pcan-CIC_1_fastqc.zip
file	Pcan-CIC_2_fastqc.html
file	Pcan-CIC_2_fastqc.zip
file	fastqc-process.sh
file	fastqc-starter.sh
file	log.txt

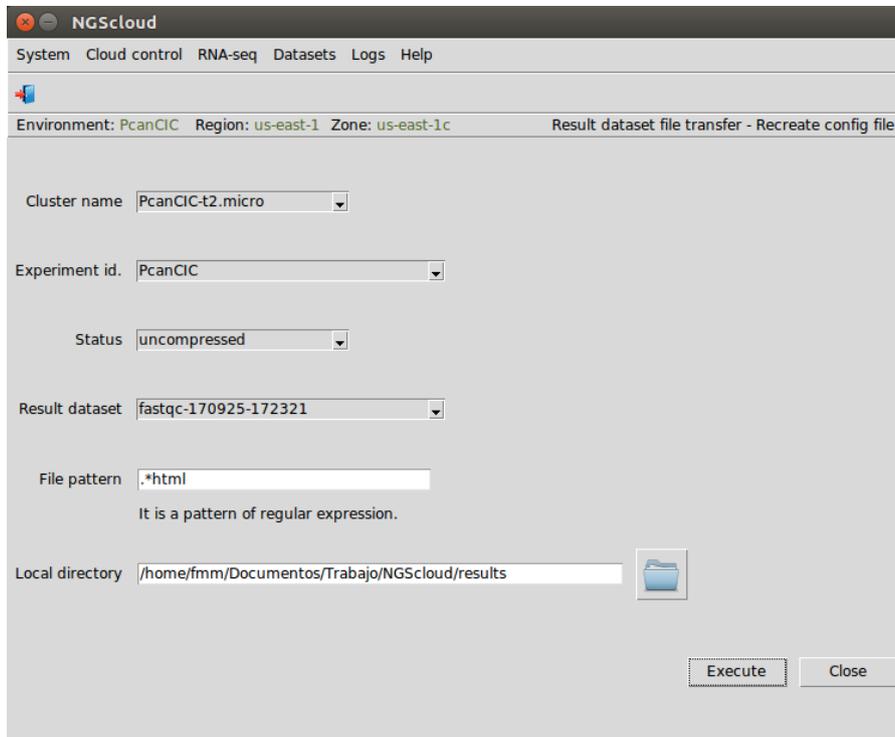
Download the quality analysis results

Next, we are going to review the analysis files generated by FastQC. First, we have to download the ".html" files with the results to a local computer. To do so, we first create the configuration file by selecting the menu item with the following path:

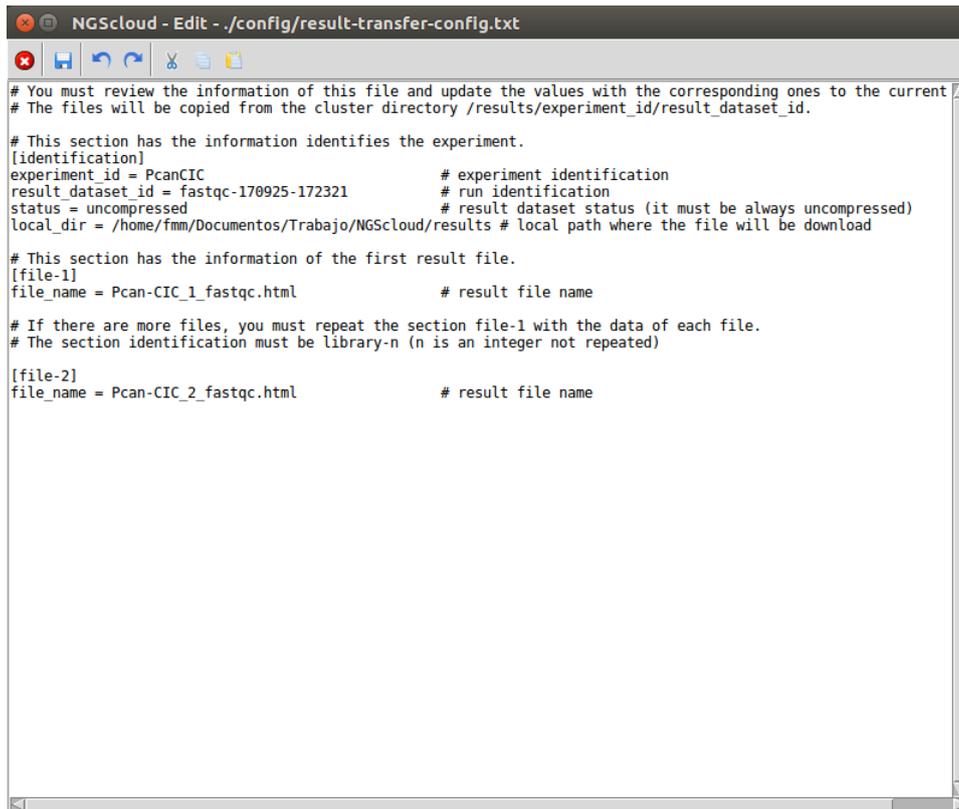
Main menu > Datasets > Result dataset file transfer > Recreate config file

In the raised window, we select **PcanCIC-t2.micro** in the *Cluster name* combo-box, **PcanCIC** in the *Experiment id* combo-box, **uncompressed** in the *Status* combo-box (because the dataset

has not been previously compressed), **fastqc-170925-172321** in the *Result dataset* combo-box, and we type **.*html** as the pattern to select the files in the *File pattern* textbox and the local directory where the files will be downloaded in the *Local directory* textbox (or we select it using the button close to the textbox). Then we press the *Execute* button:



In the next window, we can inspect the config file. It has three sections: *identification*, with the experiment and the result dataset identifications, the compression status of the dataset, and the local directory; *file-1*, with the name of the first result file; and *file-2* with the name of the second one:



```

NGScloud - Edit - ./config/result-transfer-config.txt
# You must review the information of this file and update the values with the corresponding ones to the current
# The files will be copied from the cluster directory /results/experiment_id/result_dataset_id.

# This section has the information identifies the experiment.
[identification]
experiment_id = PcanCIC # experiment identification
result_dataset_id = fastqc-170925-172321 # run identification
status = uncompressed # result dataset status (it must be always uncompressed)
local_dir = /home/fmm/Documentos/Trabajo/NGScloud/results # local path where the file will be download

# This section has the information of the first result file.
[file-1]
file_name = Pcan-CIC_1_fastqc.html # result file name

# If there are more files, you must repeat the section file-1 with the data of each file.
# The section identification must be library-n (n is an integer not repeated)

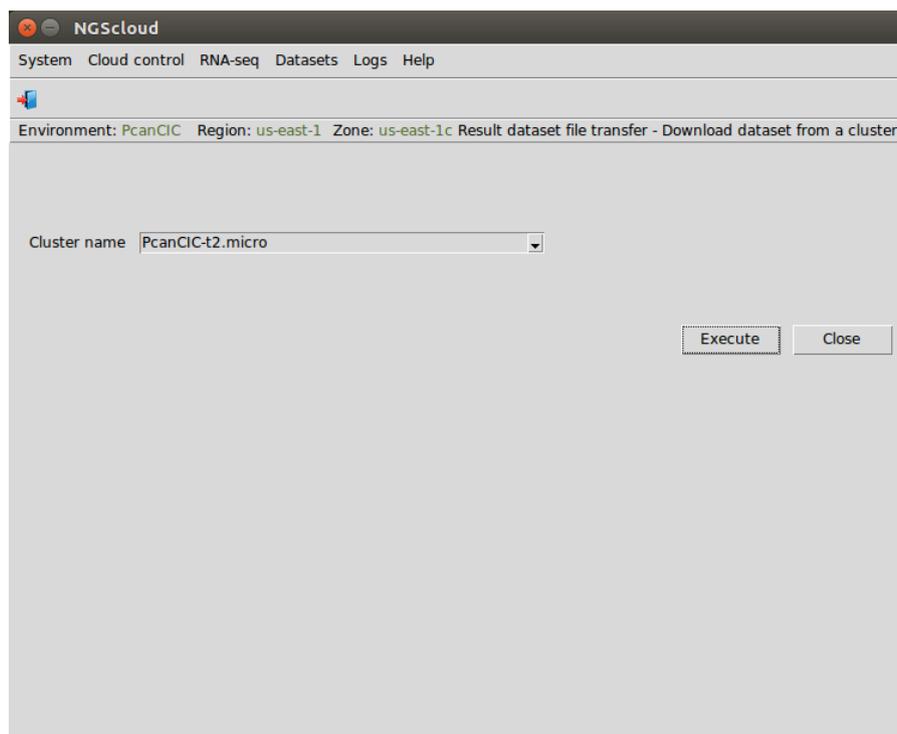
[file-2]
file_name = Pcan-CIC_2_fastqc.html # result file name

```

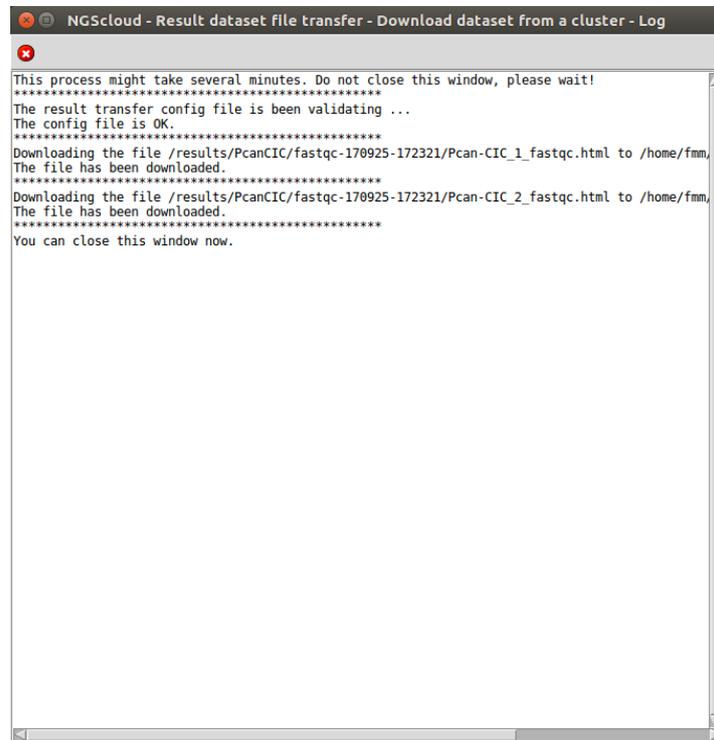
To download the result files from the cluster, we select the menu item with this path:

Main menu > Dataset > Result dataset file transfer > Download dataset from a cluster

In the raised window, we select **PcanCIC-t2.micro** in the *Cluster name* combo-box; and then we press the *Execute* button:



A window is raised with the log corresponding to the download:



```

NGScloud - Result dataset file transfer - Download dataset from a cluster - Log
This process might take several minutes. Do not close this window, please wait!
*****
The result transfer config file is been validating ...
The config file is OK.
*****
Downloading the file /results/PcanCIC/fastqc-170925-172321/Pcan-CIC_1_fastqc.html to /home/fmm,
The file has been downloaded.
*****
Downloading the file /results/PcanCIC/fastqc-170925-172321/Pcan-CIC_2_fastqc.html to /home/fmm,
The file has been downloaded.
*****
You can close this window now.

```

Trim reads using Trimmomatic

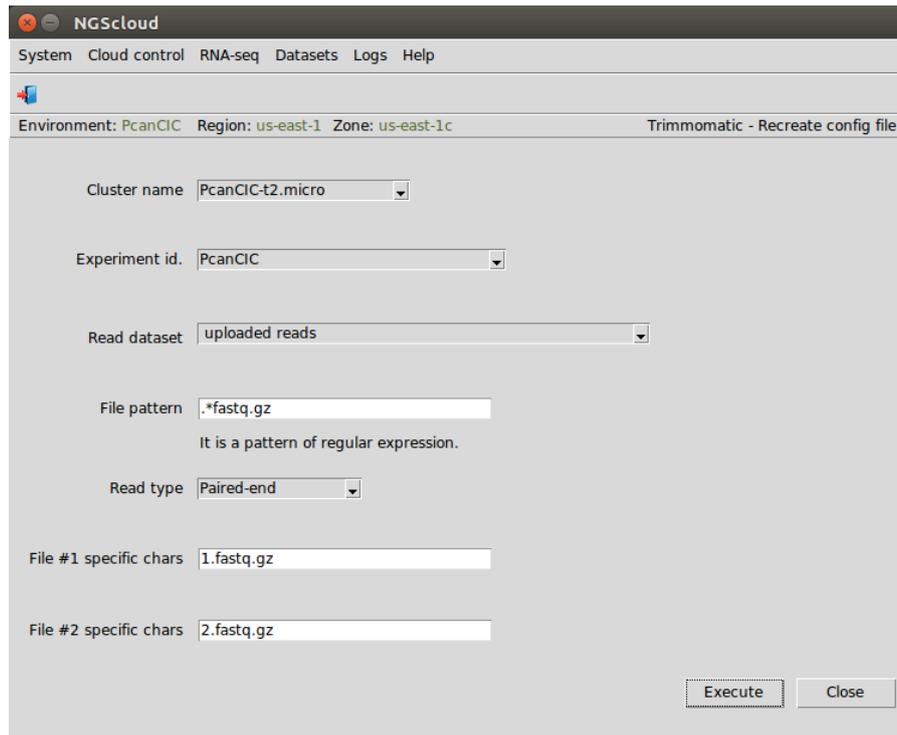
Once we have reviewed the two result files and have decided to cut 12 nucleotides from the start of reads. In this point, we are going to use Trimmomatic to do this step.

First, we create the configuration file by selecting the menu item with this path:

Main menu > RNA-seq > Trimming > Trimmomatic > Recreate config file

In the raised window, we select **PcanCIC-t2.micro** in the *Cluster name* combo-box, **PcanCIC** in the *Experiment id* combo-box, **uploaded reads** in the *Read dataset* combo-box, we type **.*fastq.gz** as the pattern to select the files in *File pattern* textbox and we select **Paired-end** in the *Read type* combo-box; finally, we type **1.fastq.gz** in the *File #1 specific chars* textbox and **2.fastq.gz** in the *File #2 specific chars* textbox. These last two strings are used to distinguish the file of each strand among the selected files by the pattern corresponding to the experiment libraries. In this example, there is only one library.

Then we press the *Execute* button:



The screenshot shows the NGScloud application window with the following configuration for Trimmomatic:

- Cluster name: PcanCIC-t2.micro
- Experiment id: PcanCIC
- Read dataset: uploaded reads
- File pattern: .*fastq.gz (Note: It is a pattern of regular expression.)
- Read type: Paired-end
- File #1 specific chars: 1.fastq.gz
- File #2 specific chars: 2.fastq.gz

Buttons: Execute, Close

In the next window, we visualize the config file. In this example, it has six sections: *identification*, with the experiment and the read dataset identifications; *Trimmomatic parameters*, with the thread number (we modify its value to **1**) and phred quality score; *Trimming step values*, with the step list that Trimmomatic can perform (we modify the headcrop value to **12**); *Trimming step order* with the order in which Trimmomatic must carry out every step indicated in the previous section; *library* with the library type; *library-1* with the two read files for the first library (in this example, we only have one library):

```

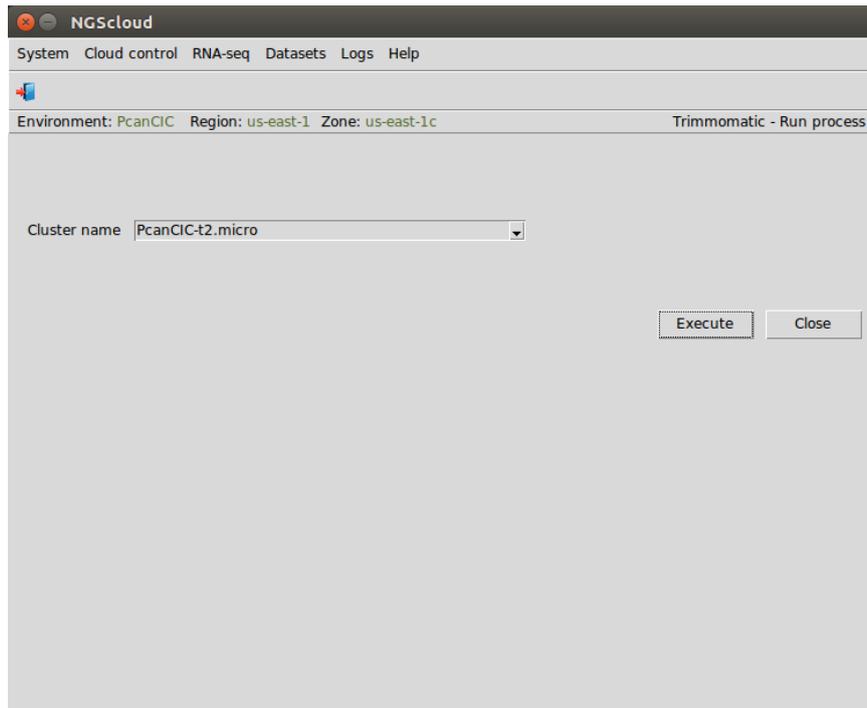
NGScloud - Edit - ./config/trimmo-config.txt
# You must review the information of this file and update the values with the corresponding ones to the current
#
# The files must be located in the cluster directory /reads/experiment_id/read_dataset_id
# The experiment_id and read_dataset_id names are fixed in the identification section.
#
# You can consult the parameters and trimming sets of Trimmomatic and their meaning in http://www.usadellab.org/
# This section has the information identifies the experiment.
[identification]
experiment_id = PcanCIC # experiment identification
read_dataset_id = uploaded-reads # read dataset identification
# This section has the information to set the Trimmomatic parameters.
[Trimmomatic parameters]
threads = 1 # number of threads for use
phred = 64 # 33 (phred 33) or 64 (phred 64)
# This section has the information to set the trimming step values
[Trimming step values]
illuminaclip = NONE # $TRIMMOMATIC_PATH/adapters/fastqWithAdaptersEtc:seedMismatch
slidingwindow = NONE # windowSize:requiredQuality or NONE (do not perform this step)
leading = NONE # quality or NONE (do not perform this step)
trailing = NONE # quality or NONE (do not perform this step)
crop = NONE # length or NONE (do not perform this step)
headcrop = 12 # length or NONE (do not perform this step)
minlen = NONE # length or NONE (do not perform this step)
tophred33 = FALSE # TRUE to performance this step, else FALSE
tophred64 = FALSE # TRUE to performance this step, else FALSE
# This section has the information to set the performance order of every step with value
[Trimming step order]
order = headcrop # if there are more than one step, they must be separated by
# This section has the global information of all libraries.
[library]
read_type = PE # SE (single-end) or PE (paired-end)
# This section has the information of the first library.
[library-1]
read_file_1 = Pcan-CIC_1.fastq.gz # name of the read file in SE read type or the + strand read
read_file_2 = Pcan-CIC_2.fastq.gz # name of the - strand reads file in PE read type or NONE in
# If there are more libraries, you must repeat the section library-1 with the data of each file

```

We run the trimming process in the cluster by selecting the menu item with this path:

Main menu > RNA-seq > Trimming > Trimmomatic > Run trimming process

In the raised window, we select **PcanCIC-t2.micro** in the *Cluster name* combo-box; and then we press the *Execute* button:



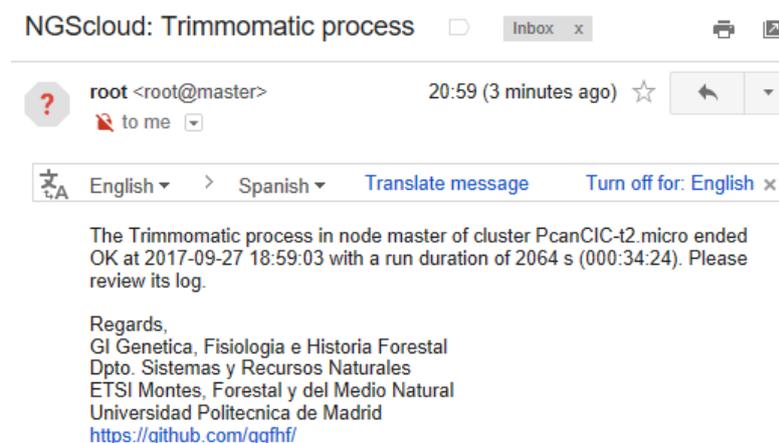
A window is raised with the submission log:

```

NGScloud - Trimmomatic - Run process - Log
This process might take several minutes. Do not close this window, please wait!
*****
Validating the Trimmomatic config file ...
The config file is OK.
*****
Connecting the SSH client ...
The SSH client is connected.
*****
Connecting the SSH transport ...
The SSH transport is connected.
*****
Connecting the SFTP client ...
The SFTP client is connected.
*****
Verifying process requirements ...
Process requirements are OK.
*****
Determining the run directory in the cluster ...
The directory path is /results/PcanCIC/trimmo-170927-202434.
*****
Building the process script ./temp/trimmo-process.sh ...
The file is built.
*****
Uploading the process script ./temp/trimmo-process.sh in the directory /results/PcanCIC/trimmo-170927-202434
The file is uploaded.
*****
Setting on the run permission of /results/PcanCIC/trimmo-170927-202434/trimmo-process.sh ...
The run permission is set.
*****
Building the process starter ./temp/trimmo-starter.sh ...
The file is built.
*****
Uploading the process starter ./temp/trimmo-starter.sh in the directory /results/PcanCIC/trimmo-170927-202434
The file is uploaded.
*****
Setting on the run permission of /results/PcanCIC/trimmo-170927-202434/trimmo-starter.sh ...
The run permission is set.
*****
Submitting the process script /results/PcanCIC/trimmo-170927-202434/trimmo-starter.sh ...
Your job 4 ("trimmo-starter.sh") has been submitted
*****
Closing the SSH transport connection ...
The connection is closed.
*****
Closing the SSH client connection ...

```

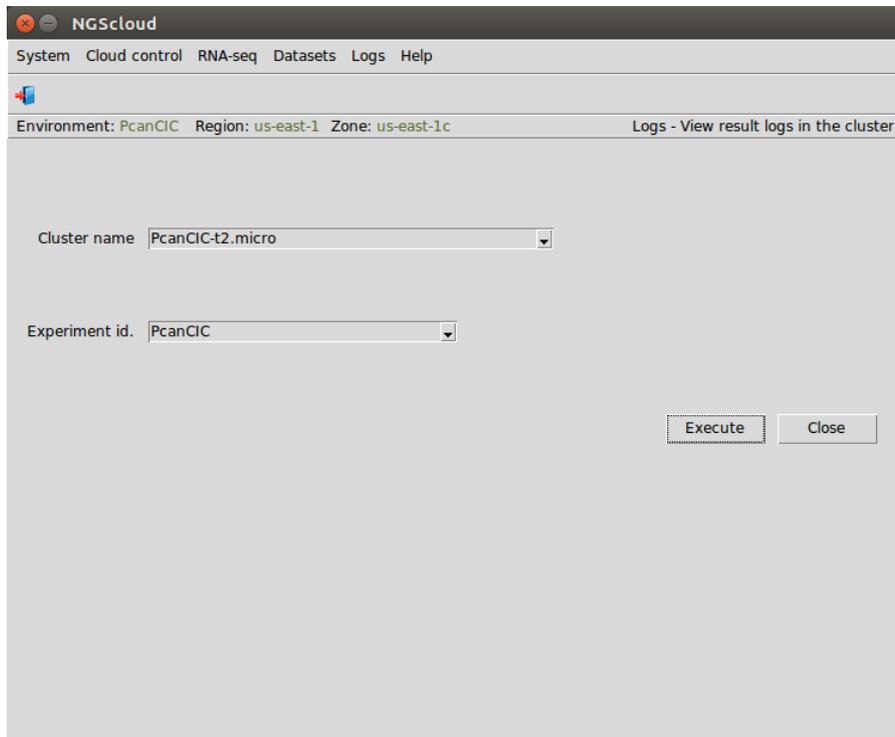
At the end of the run, an email is sent, informing of its completion:



We can view the process log during and after its run. To do so, we select the menu item with this path:

Main menu > Logs > View result logs in the cluster

In the raised window, we select **PcanCIC-t2.micro** in the *Cluster name* combo-box and **PcanCIC** in the *Experiment id* combo-box; and then we press the *Execute* button:



A window with the result datasets for each run of the bioinformatic programs that correspond to the experiment PcanCIC is shown:

The screenshot shows a window titled 'NGScloud - Experiment runs in /results/PcanCIC - Table'. It contains a table with the following data:

Experiment id	Result dataset	Bioinfo app / Utility	Date	Time
PcanCIC	fastqc-170925-172321	FastQC	2017-09-25	17:23:21
PcanCIC	trimmo-170927-202434	Trimmomatic	2017-09-27	20:24:34

At this moment, there are two result datasets: **fastqc-170925-172321**, corresponding to the previous run of FastQC; and **trimmo-170927-202434**, corresponding to the run of

Trimmomatic. We click on this last dataset and another window appears with its corresponding log:

```

#####
Script started in node master of cluster PcanCIC-t2.micro at 2017-09-27 18:24:39 UTC.
#####
TrimmomaticPE: Started with arguments:
-threads 1 -phred64 -trimlog Pcan-CIC_1.fastq.gz.log /reads/PcanCIC/uploaded-reads/Pcan-CIC_1.fastq.gz /reads/P
Input Read Pairs: 21444414 Both Surviving: 21444414 (100.00%) Forward Only Surviving: 0 (0.00%) Reverse Only Sur
TrimmomaticPE: Completed successfully
#####
Elapsed real time (s): 2064.31
CPU time in kernel mode (s): 373.77
CPU time in user mode (s): 1609.26
Percentage of CPU: 96%
Maximum resident set size(Kb): 173584
Average total memory use (Kb):0
#####
Script ended OK at 2017-09-27 18:59:03 UTC with a run duration of 2064 s (000:34:24).
#####

```

[Terminate the cluster with a t2.micro template and create another cluster with a r3.4xlarge template](#)

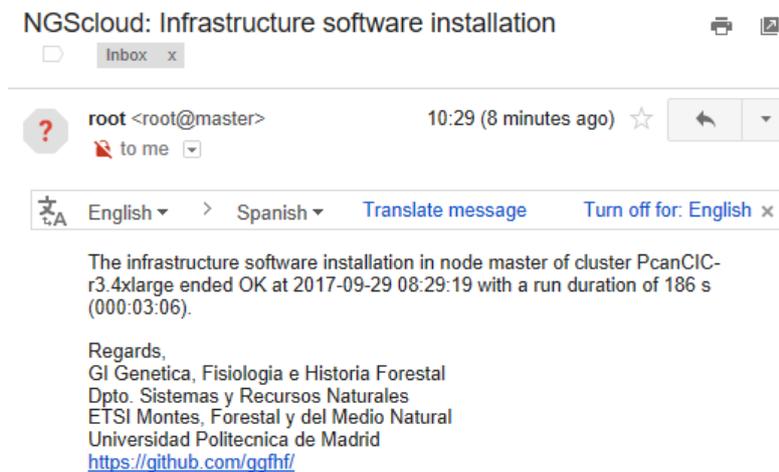
After read trimming, we are going to assembly a preliminary transcriptome using Trinity. Trinity's hardware requirements are very high in terms of CPUs and GiBs of RAM memory. We must terminate the current cluster and create another one fulfilling these requirements. We choose a r3.4xlarge template whose instances have 16 CPUs and 122 GiBs of RAM memory, in order to be able to analyze the large read files of our experiment with Trinity.

To terminate the current cluster, we select the menu item with this path:

Main menu > Cloud control > Cluster operation > Terminate cluster

In the raised window, we select **PcanCIC-t2.micro** in the *Cluster name* combo-box; and then we press the *Execute* button:

When the cluster is started, infrastructure software will be installed. At the end of the installation, an email is sent, informing of its completion:

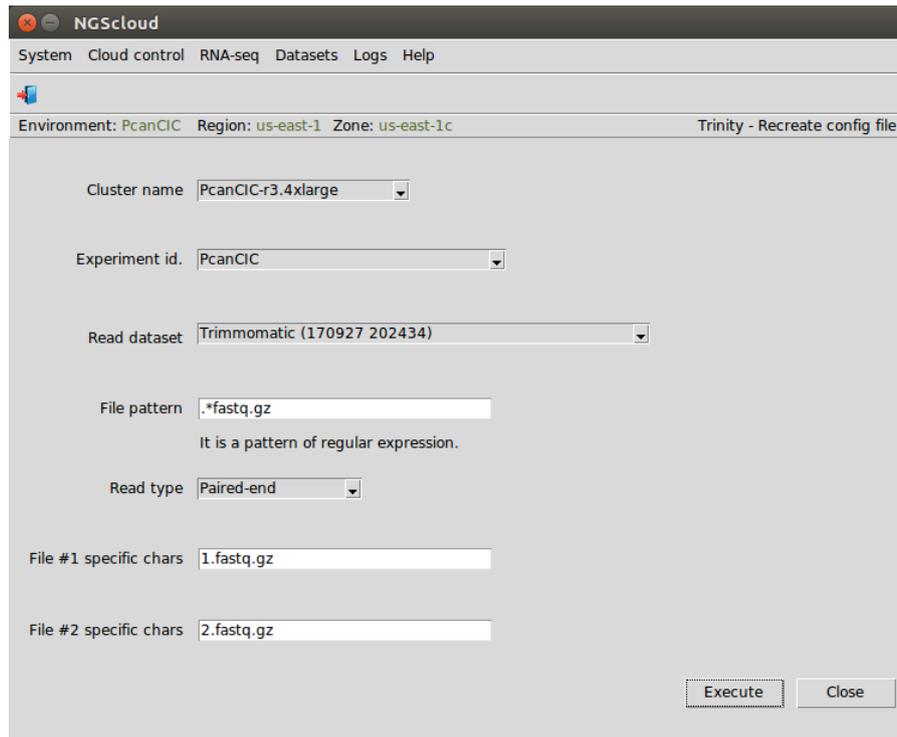


Assembly reads using Trinity

First, we create the config file by selecting the menu item with this path:

Main menu > RNA-seq > De novo assembly > Trinity > Recreate config file

In the raised window, we select **PcanCIC-r3.4xlarge** in the *Cluster name* combo-box, **PcanCIC** in the *Experiment id* combo-box, **Trimmomatic (170927 202434)** in the *Read dataset* combo-box, we type **.*fastq.gz** as the pattern to select the files in *File pattern* textbox, and we select **Paired-end** in the *Read type* combo-box; finally, we type **1.fastq.gz** in the *File #1 specific chars* textbox and **2.fastq.gz** in the *File #2 specific chars* textbox. Then we press the *Execute* button:

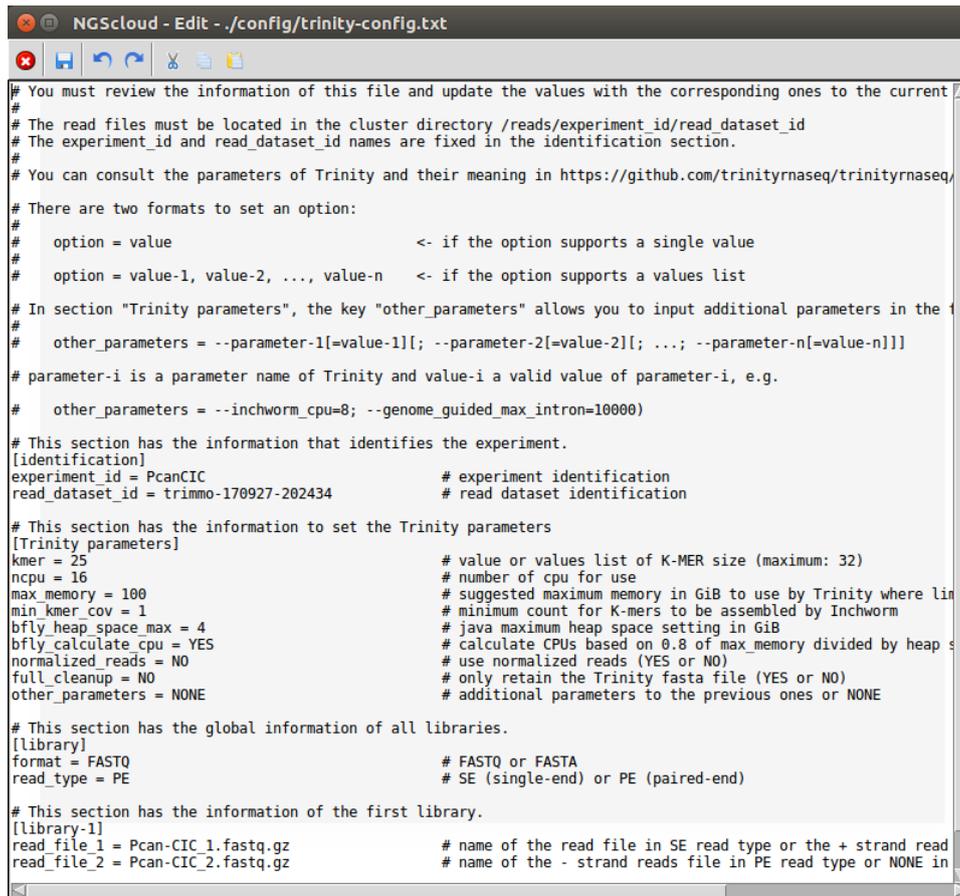


The screenshot shows the NGScloud interface for configuring a Trinity run. The window title is "NGScloud" and it has a menu bar with "System", "Cloud control", "RNA-seq", "Datasets", "Logs", and "Help". The status bar shows "Environment: PcanCIC Region: us-east-1 Zone: us-east-1c" and "Trinity - Recreate config file". The configuration fields are as follows:

Field	Value
Cluster name	PcanCIC-r3.4xlarge
Experiment id.	PcanCIC
Read dataset	Trimmomatic (170927 202434)
File pattern	.*fastq.gz
Read type	Paired-end
File #1 specific chars	1.fastq.gz
File #2 specific chars	2.fastq.gz

Buttons: Execute, Close

In the next window, we can examine the config file. In this example, it has four sections: *identification*, with the experiment and the read dataset identifications; *Trinity parameters*, with several parameters used by Trinity in which we can modify the value of CPUs number to **16**, and the value of suggested maximum memory to **100**; *library* with the format and library type; *library-1* with the two read files for the first library (in this example, we only have one library):



```

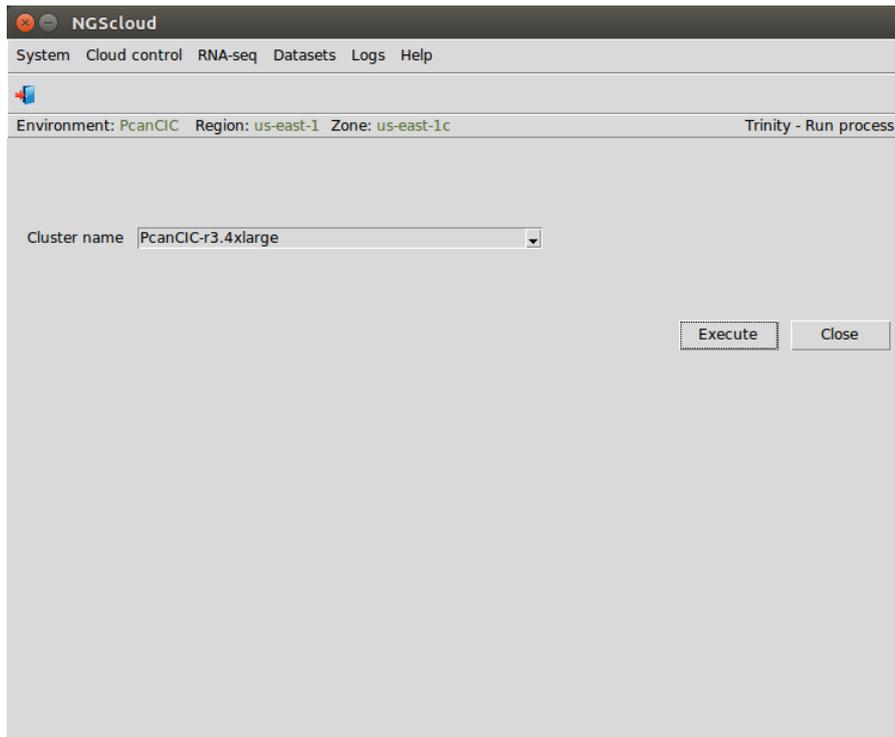
# You must review the information of this file and update the values with the corresponding ones to the current
#
# The read files must be located in the cluster directory /reads/experiment_id/read_dataset_id
# The experiment_id and read_dataset_id names are fixed in the identification section.
#
# You can consult the parameters of Trinity and their meaning in https://github.com/trinityrnaseq/trinityrnaseq/
#
# There are two formats to set an option:
#
# option = value                                <- if the option supports a single value
#
# option = value-1, value-2, ..., value-n       <- if the option supports a values list
#
# In section "Trinity parameters", the key "other_parameters" allows you to input additional parameters in the
#
# other_parameters = --parameter-1[=value-1]; --parameter-2[=value-2]; ...; --parameter-n[=value-n]]
#
# parameter-i is a parameter name of Trinity and value-i a valid value of parameter-i, e.g.
#
# other_parameters = --inchworm_cpu=8; --genome_guided_max_intron=10000)
#
# This section has the information that identifies the experiment.
[identification]
experiment_id = PcanCIC                          # experiment identification
read_dataset_id = trimmo-170927-202434           # read dataset identification
#
# This section has the information to set the Trinity parameters
[Trinity parameters]
kmer = 25                                         # value or values list of K-MER size (maximum: 32)
ncpu = 16                                         # number of cpu for use
max_memory = 100                                 # suggested maximum memory in GiB to use by Trinity where lin
min_kmer_cov = 1                                 # minimum count for K-mers to be assembled by Inchworm
bfly_heap_space_max = 4                          # java maximum heap space setting in GiB
bfly_calculate_cpu = YES                         # calculate CPUs based on 0.8 of max_memory divided by heap s
normalized_reads = NO                            # use normalized reads (YES or NO)
full_cleanup = NO                               # only retain the Trinity fasta file (YES or NO)
other_parameters = NONE                          # additional parameters to the previous ones or NONE
#
# This section has the global information of all libraries.
[library]
format = FASTQ                                   # FASTQ or FASTA
read_type = PE                                  # SE (single-end) or PE (paired-end)
#
# This section has the information of the first library.
[[library-1]
read_file_1 = Pcan-CIC_1.fastq.gz                # name of the read file in SE read type or the + strand read
read_file_2 = Pcan-CIC_2.fastq.gz                # name of the - strand reads file in PE read type or NONE in

```

We run the assembly process in the cluster by selecting the menu item with this path:

Main menu > RNA-seq > De novo assembly > Trinity > Run assembly process

In the raised window, we select **PcanCIC-r3.4xlarge** in the *Cluster name* combo-box; and then we press the *Execute* button:



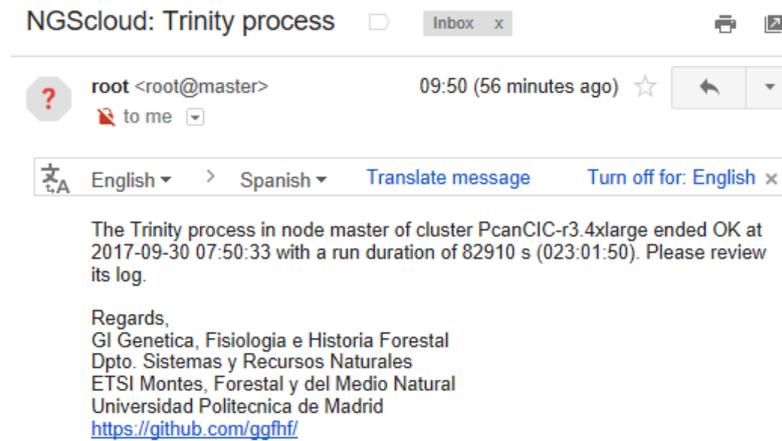
A window is raised showing the submission log:

```

NGScloud - Trinity - Run process - Log
This process might take several minutes. Do not close this window, please wait!
*****
Validating the Trinity config file ...
The config file is OK.
*****
Connecting the SSH client ...
The SSH client is connected.
*****
Connecting the SSH transport ...
The SSH transport is connected.
*****
Connecting the SFTP client ...
The SFTP client is connected.
*****
Verifying process requirements ...
Process requirements are OK.
*****
Determining the run directory for kmer 25 in the cluster ...
The directory path is /results/PcanCIC/trinity-170929-104827.
*****
Building the process script ./temp/trinity-process.sh ...
The file is built.
*****
Uploading the process script ./temp/trinity-process.sh in the directory /results/PcanCIC/trinity-170929-104827
The file is uploaded.
*****
Setting on the run permission of /results/PcanCIC/trinity-170929-104827/trinity-process.sh ...
The run permission is set on.
*****
Building the process starter ./temp/trinity-starter.sh ...
The file is built.
*****
Uploading the process starter ./temp/trinity-starter.sh in the directory /results/PcanCIC/trinity-170929-104827
The file is uploaded.
*****
Setting on the run permission of /results/PcanCIC/trinity-170929-104827/trinity-starter.sh ...
The run permission is set on.
*****
Submitting the process script /results/PcanCIC/trinity-170929-104827/trinity-starter.sh ...
Your job 1 ("trinity-starter.sh") has been submitted
*****
Closing the SSH transport connection ...
The connection is closed.
*****
Closing the SSH client connection ...

```

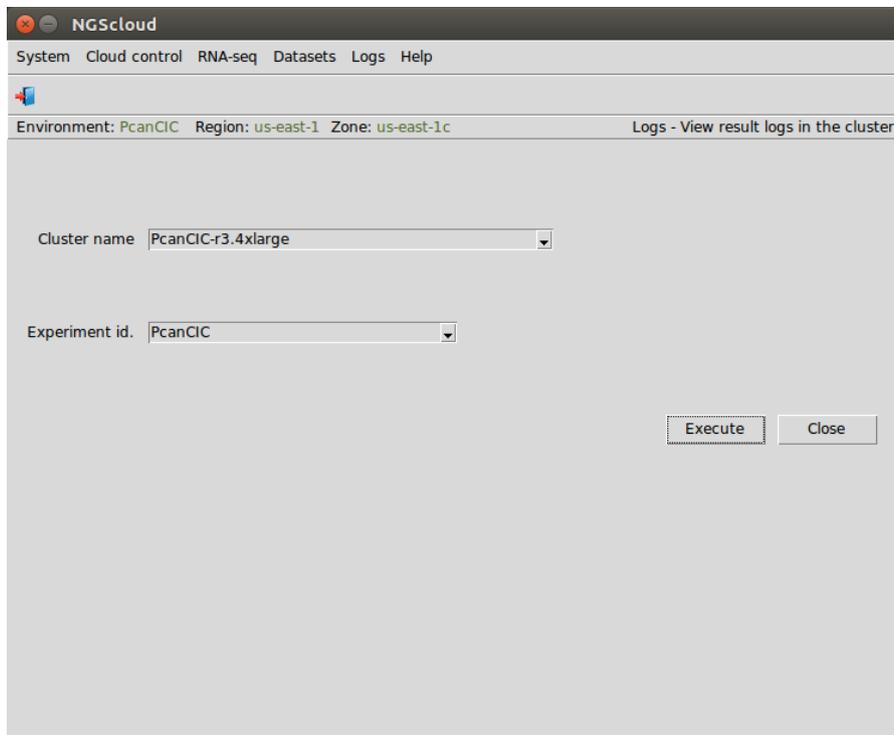
At the end of the run, an email is sent, informing of its completion:



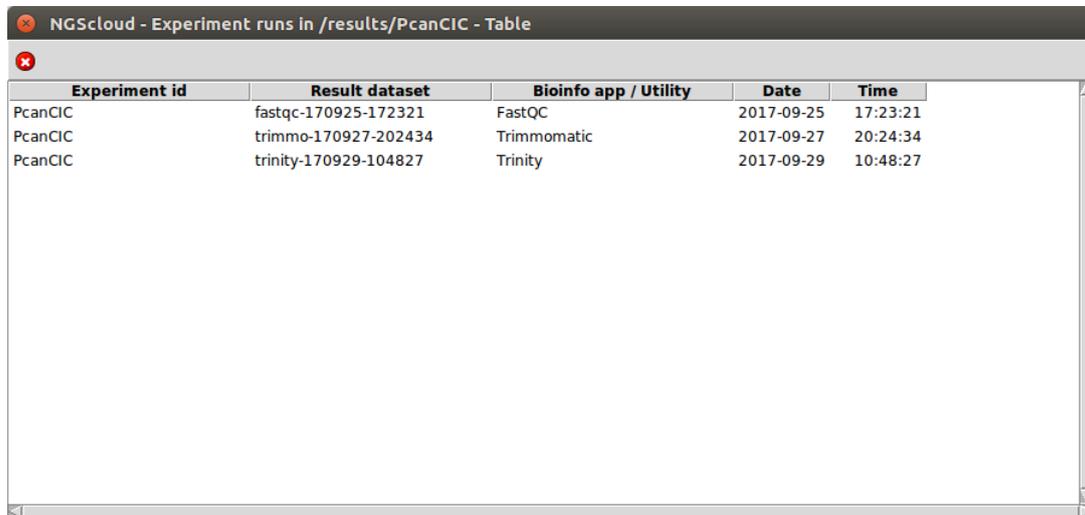
We can view the process log during and after the run. To do so, we select the menu item with this path:

Main menu > Logs > View result logs in the cluster

In the raised window, we select **PcanCIC-r3.4xlarge** in the *Cluster name* combo-box and **PcanCIC** in the *Experiment id* combo-box; and then we press the *Execute* button:

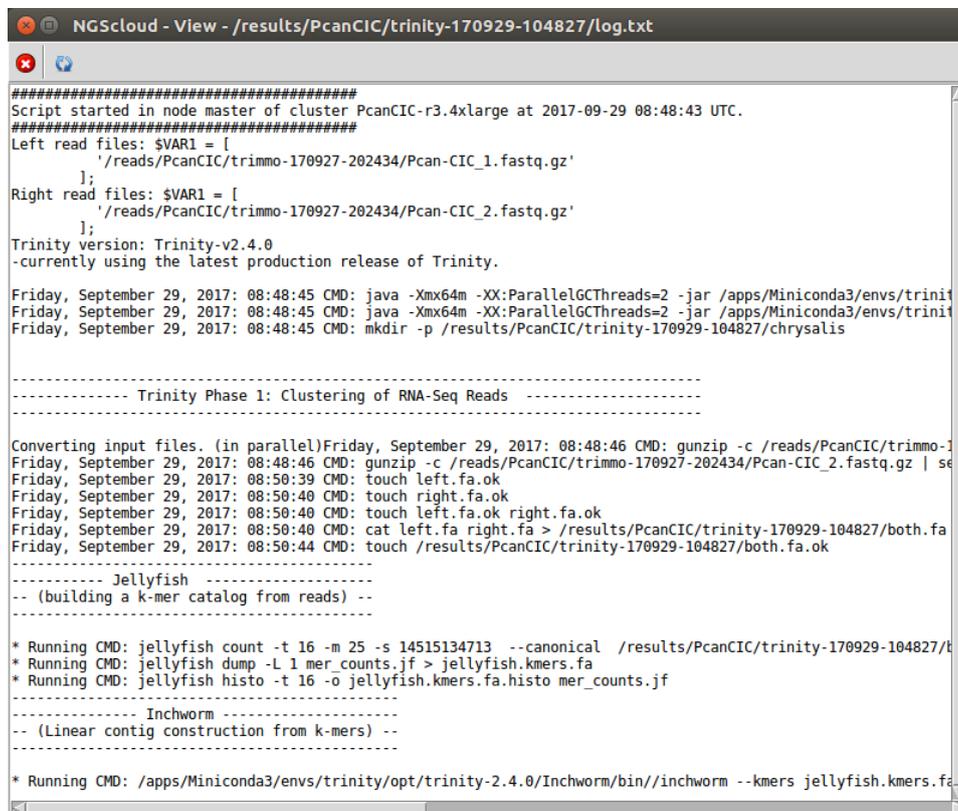


A window with the result datasets for each run of the bioinformatic programs that correspond to the experiment PcanCIC is shown:



Experiment id	Result dataset	Bioinfo app / Utility	Date	Time
PcanCIC	fastqc-170925-172321	FastQC	2017-09-25	17:23:21
PcanCIC	trimmo-170927-202434	Trimmomatic	2017-09-27	20:24:34
PcanCIC	trinity-170929-104827	Trinity	2017-09-29	10:48:27

Three result datasets are shown. We click on the **trinity-170929-104827** row, which is the dataset generated by the Trinity run, and another window appears with its corresponding log:



```

#####
Script started in node master of cluster PcanCIC-r3.4xlarge at 2017-09-29 08:48:43 UTC.
#####
Left read files: $VAR1 = [
    '/reads/PcanCIC/trimmo-170927-202434/Pcan-CIC_1.fastq.gz'
];
Right read files: $VAR1 = [
    '/reads/PcanCIC/trimmo-170927-202434/Pcan-CIC_2.fastq.gz'
];
Trinity version: Trinity-v2.4.0
-currently using the latest production release of Trinity.

Friday, September 29, 2017: 08:48:45 CMD: java -Xmx64m -XX:ParallelGCThreads=2 -jar /apps/Miniconda3/envs/trinity-2.4.0/bin/trinity --single -S /reads/PcanCIC/trimmo-170927-202434/Pcan-CIC_1.fastq.gz --paired -S /reads/PcanCIC/trimmo-170927-202434/Pcan-CIC_2.fastq.gz --output /results/PcanCIC/trinity-170929-104827/
Friday, September 29, 2017: 08:48:45 CMD: java -Xmx64m -XX:ParallelGCThreads=2 -jar /apps/Miniconda3/envs/trinity-2.4.0/bin/trinity --single -S /reads/PcanCIC/trimmo-170927-202434/Pcan-CIC_1.fastq.gz --paired -S /reads/PcanCIC/trimmo-170927-202434/Pcan-CIC_2.fastq.gz --output /results/PcanCIC/trinity-170929-104827/
Friday, September 29, 2017: 08:48:45 CMD: mkdir -p /results/PcanCIC/trinity-170929-104827/chrysalis

----- Trinity Phase 1: Clustering of RNA-Seq Reads -----
-----

Converting input files. (in parallel)Friday, September 29, 2017: 08:48:46 CMD: gunzip -c /reads/PcanCIC/trimmo-170927-202434/Pcan-CIC_1.fastq.gz | sed 's/\.gz$//' > /results/PcanCIC/trinity-170929-104827/both.fa
Friday, September 29, 2017: 08:48:46 CMD: gunzip -c /reads/PcanCIC/trimmo-170927-202434/Pcan-CIC_2.fastq.gz | sed 's/\.gz$//' > /results/PcanCIC/trinity-170929-104827/both.fa
Friday, September 29, 2017: 08:50:39 CMD: touch left.fa.ok
Friday, September 29, 2017: 08:50:40 CMD: touch right.fa.ok
Friday, September 29, 2017: 08:50:40 CMD: touch left.fa.ok right.fa.ok
Friday, September 29, 2017: 08:50:40 CMD: cat left.fa right.fa > /results/PcanCIC/trinity-170929-104827/both.fa
Friday, September 29, 2017: 08:50:44 CMD: touch /results/PcanCIC/trinity-170929-104827/both.fa.ok

----- Jellyfish -----
-- (building a k-mer catalog from reads) --
-----
* Running CMD: jellyfish count -t 16 -m 25 -s 14515134713 --canonical /results/PcanCIC/trinity-170929-104827/both.fa
* Running CMD: jellyfish dump -L 1 mer_counts.jf > jellyfish.kmers.fa
* Running CMD: jellyfish histo -t 16 -o jellyfish.kmers.fa.histo mer_counts.jf

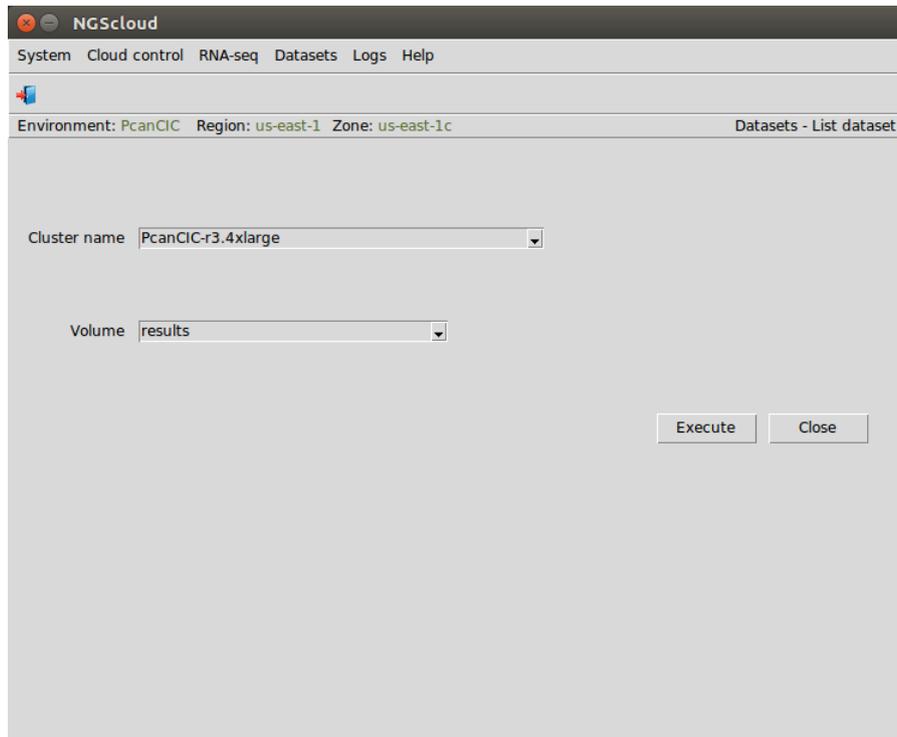
----- Inchworm -----
-- (Linear contig construction from k-mers) --
-----
* Running CMD: /apps/Miniconda3/envs/trinity/opt/trinity-2.4.0/Inchworm/bin/inchworm --kmers jellyfish.kmers.fa

```

Now, we are going to list the assembly generated by Trinity. We select the menu item with this path:

Main menu > Dataset > List dataset

In the raised window, we select **PcanCIC-r3.4xlarge** in the *Cluster name* combo-box and **results** in the *Volume* combo-box. Then we press the *Execute* button:

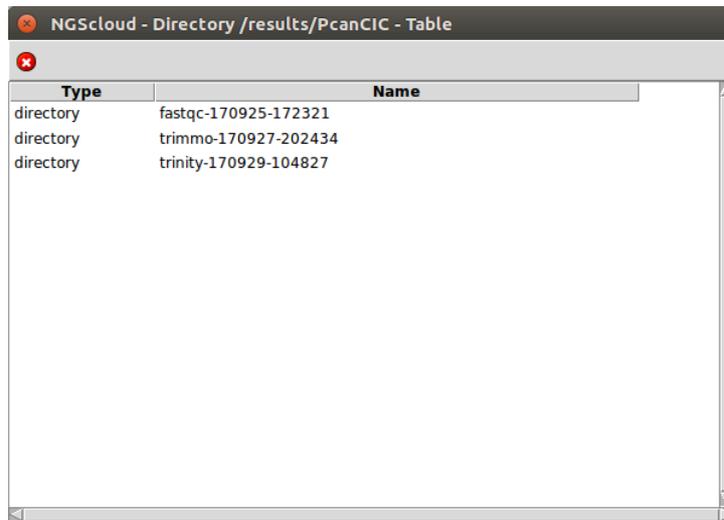


Now, we click in the **PcanCIC** row:

The screenshot shows a window titled 'NGScloud - Directory /results - Table' containing a table with two columns: 'Type' and 'Name'. The table lists four directory entries.

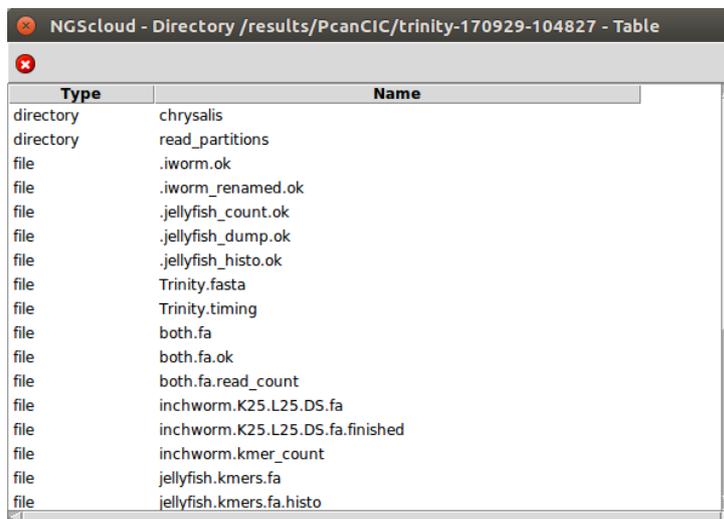
Type	Name
directory	Athaliana01x
directory	PcanCIC
directory	database
directory	test

Next, a window with the result datasets of the experiment PcanCIC is shown:



Type	Name
directory	fastqc-170925-172321
directory	trimmo-170927-202434
directory	trinity-170929-104827

We click on the **trinity-170929-104827** row, and another window appears with the content of the files corresponding to this assembly.



Type	Name
directory	chrysalis
directory	read_partitions
file	.iworm.ok
file	.iworm_renamed.ok
file	.jellyfish_count.ok
file	.jellyfish_dump.ok
file	.jellyfish_histo.ok
file	Trinity.fasta
file	Trinity.timing
file	both.fa
file	both.fa.ok
file	both.fa.read_count
file	inchworm.K25.L25.DS.fa
file	inchworm.K25.L25.DS.fa.finished
file	inchworm.kmer_count
file	jellyfish.kmers.fa
file	jellyfish.kmers.fa.histo

The file **Trinity.fasta** is the one corresponding to the transcriptome. To observe its characteristics, we click on it:

The screenshot shows a window titled "NGScloud - File /results/PcanCIC/trinity-170929-104827//results/PcanCIC/". Inside the window is a table with two columns: "Data" and "Value".

Data	Value
directory	/results/PcanCIC/trinity-170929-104827
name	Trinity.fasta
size	338748408
permissions	rw-r--
modification date	2017-08-30
modification time	07:50
owner group	root
owner name	root

Evaluate the transcriptome using RSEM-EVAL

Next, we are going to evaluate the quality of the transcriptome generated by Trinity with RSEM-EVAL, which is included in the DETONATE package. To do so, we first create the config file by selecting the menu item with the following path:

Main menu > RNA-seq > Assembly quality and transcript quantification > RSEM-EVAL (DETONATE package) > Recreate config file

In the raised window, we select **PcanCIC-r3.4xlarge** in the *Cluster name* combo-box, **PcanCIC** in the *Experiment id* combo-box, **Trimmomatic (170927 202434)** in the *Read dataset* combo-box, we type **.*fastq.gz** as the pattern to select the files in *File pattern* textbox, and we select **Paired-end** in the *Read type* combo-box; we type **1.fastq.gz** in the *File #1 specific chars* textbox and **2.fastq.gz** in the *File #2 specific chars* textbox; finally, we selected **Trinity (170927 104827)** in the *Assembly dataset* combo-box. The *Assembly type* combo-box only is activated when the assembly dataset was generated by SOAPdenovo-Trans; in this case, the combo-box has two items: CONTIGS and SCAFFOLDS. Then we press the *Execute* button:

In the next window, we can examine the config file. In this example, it has four sections: *identification*, with the experiment, read and assembly dataset identifications; *RSEM-EVAL parameters*, with several parameters used by RSEM-EVAL, where we can modify the value of threads number to **16**; *library* with the format and library type; and *library-1* with the two read files for the first library (in this example, we only have one library):

```

# You must review the information of this file and update the values with the corresponding ones to the current
#
# The read files must be located in the cluster directory /reads/experiment_id/read_dataset_id
# The assembly files must be located in the cluster directory /results/experiment_id/assembly_dataset_id
# The experiment_id, read_dataset_id and assembly_dataset_id names are fixed in the identification section.
#
# You can consult the parameters of RSEM-EVAL (DETONATE package) and their meaning in http://deweylab.biostat.wi
#
# This section has the information identifies the assembly result dataset.
[identification]
experiment_id = PcanCIC # experiment identification
read_dataset_id = trimmo-170927-202434 # read dataset identification
assembly_software = trinity # assembly software: sdnt (SOAPdenovo-Trans) or trinity (Trin
assembly_dataset_id = trinity-170929-104827 # assembly dataset identification
assembly_type = NONE # CONTIGS or SCAFFOLDS in SOAPdenovo-Trans; NONE in Trinity,

# This section has the information to set the RSEM-EVAL parameters
[RSEM-EVAL parameters]
num_threads = 12 # number of threads for use
bowtie2_mismatch_rate = 0.1 # maximum mismatch rate allowed (Bowtie 2 parameter)
keep_intermediate_files = NO # keep temporary files generated: YES or NO

# This section has the global information of all libraries.
[library]
format = FASTQ # FASTQ or FASTA
read_type = PE # SE (single-end) or PE (paired-end)
length = 200 # average read length in SE read type or average fragment len

# This section has the information of the first library.
[library-1]
read_file_1 = Pcan-CIC_1.fastq.gz # name of the read file in SE read type or the + strand read
read_file_2 = Pcan-CIC_2.fastq.gz # name of the - strand reads file in PE read type or NONE in

# If there are more libraries, you must repeat the section library-1 with the data of each file.
# The section identification must be library-n (n is an integer not repeated)

```

We run the assembly assessment process in the cluster by selecting the menu item with this path:

Main menu > RNA-seq > Assembly quality and transcript quantification > RSEM-EVAL (DETONATE package) > Run assembly assessment process file

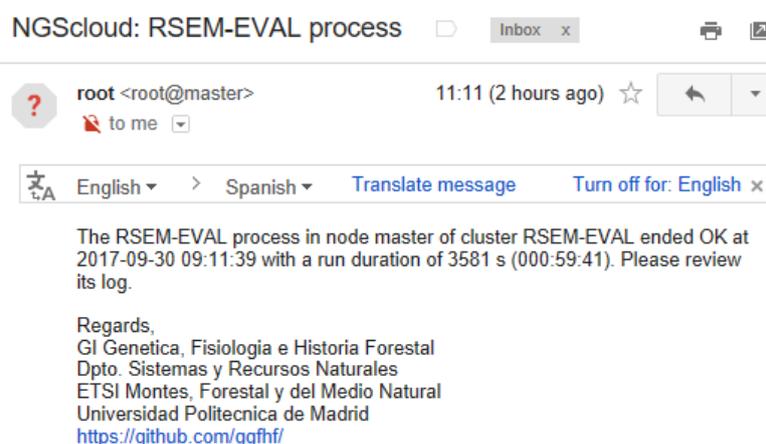
In the raised window, we select **PcanCIC-r3.4xlarge** in the *Cluster name* combo-box; and then we press the *Execute* button:

```

NGScloud - RSEM-EVAL - Run process - Log
This process might take several minutes. Do not close this window, please wait!
*****
Validating the RSEM-EVAL config file ...
The config file is OK.
*****
Connecting the SSH client ...
The SSH client is connected.
*****
Connecting the SSH transport ...
The SSH transport is connected.
*****
Connecting the SFTP client ...
The SFTP client is connected.
*****
Verifying process requirements ...
Process requirements are OK.
*****
Determining the run directory in the cluster ...
The directory path is /results/PcanCIC/rsemeval-170930-101146.
*****
Building the process script ./temp/rsemeval-process.sh ...
The file is built.
*****
Uploading the process script ./temp/rsemeval-process.sh in the directory /results/PcanCIC/rsemeval-170930-101146/rsemeval-process.sh ...
The file is uploaded.
*****
Setting on the run permission of /results/PcanCIC/rsemeval-170930-101146/rsemeval-process.sh ...
The run permission is set.
*****
Building the process starter ./temp/rsemeval-starter.sh ...
The file is built.
*****
Uploading the process starter ./temp/rsemeval-starter.sh in the directory /results/PcanCIC/rsemeval-170930-101146/rsemeval-starter.sh ...
The file is uploaded.
*****
Setting on the run permission of /results/PcanCIC/rsemeval-170930-101146/rsemeval-starter.sh ...
The run permission is set.
*****
Submitting the process script /results/PcanCIC/rsemeval-170930-101146/rsemeval-starter.sh ...
Your job 2 ("rsemeval-starter.sh") has been submitted
*****
Closing the SSH transport connection ...
The connection is closed.
*****
Closing the SSH client connection ...

```

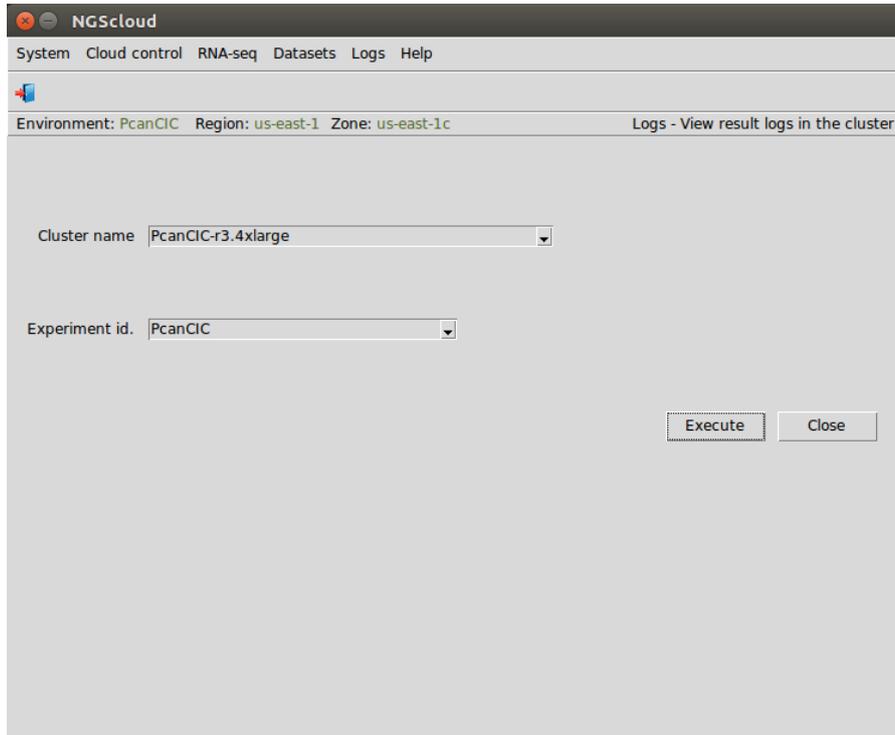
At the end of the run, an email is sent, informing of its completion:



We can view the process log during and after its run. To do so, we select the menu item with this path:

Main menu > Logs > View result logs in the cluster

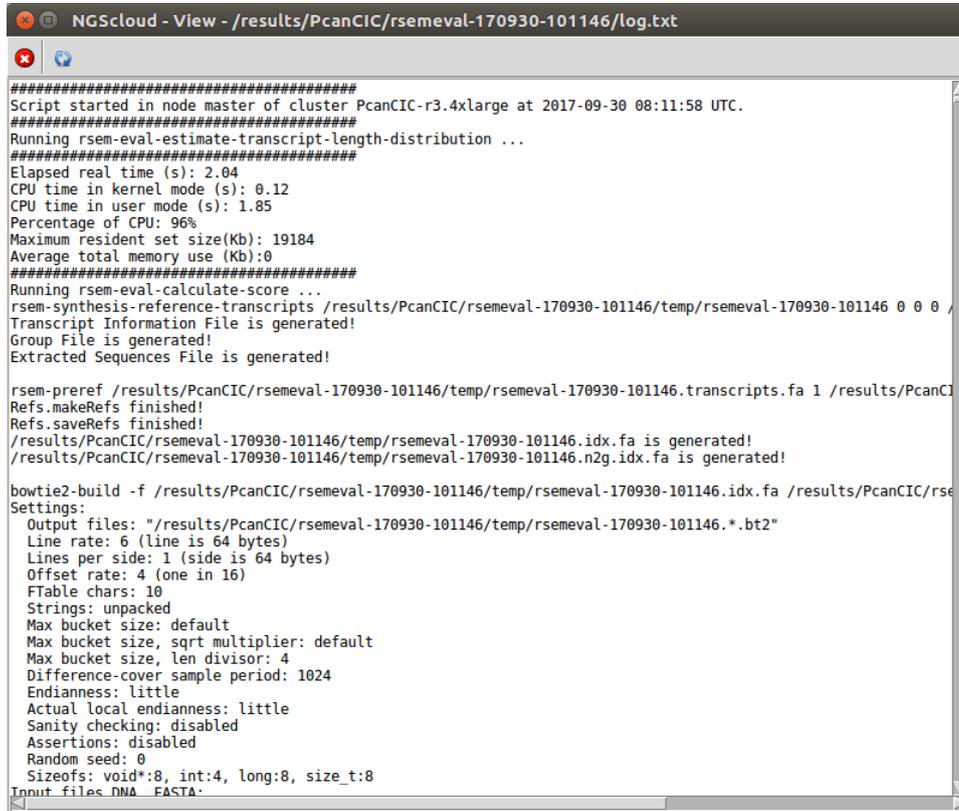
In the raised window, we select **PcanCIC-r3.4xlarge** in the *Cluster name* combo-box and **PcanCIC** in the *Experiment id* combo-box; and then we press the *Execute* button:



A window with the result datasets for each run of the bioinformatic programs that correspond to the experiment PcanCIC is shown:

Experiment id	Result dataset	Bioinfo app / Utility	Date	Time
PcanCIC	fastqc-170925-172321	FastQC	2017-09-25	17:23:21
PcanCIC	rsemeval-170930-101146	RSEM-EVAL	2017-09-30	10:11:46
PcanCIC	trimmo-170927-202434	Trimmomatic	2017-09-27	20:24:34
PcanCIC	trinity-170929-104827	Trinity	2017-09-29	10:48:27

Four result datasets are shown. We click on the **rsemeval-170930-101146** row, the dataset generated by RSEM-EVAL run, and another window appears with its corresponding log:



```

#####
Script started in node master of cluster PcanCIC-r3.4xlarge at 2017-09-30 08:11:58 UTC.
#####
Running rsem-eval-estimate-transcript-length-distribution ...
#####
Elapsed real time (s): 2.04
CPU time in kernel mode (s): 0.12
CPU time in user mode (s): 1.85
Percentage of CPU: 96%
Maximum resident set size(Kb): 19184
Average total memory use (Kb):0
#####
Running rsem-eval-calculate-score ...
rsem-synthesis-reference-transcripts /results/PcanCIC/rsemeval-170930-101146/temp/rsemeval-170930-101146 0 0 0 /
Transcript Information File is generated!
Group File is generated!
Extracted Sequences File is generated!

rsem-preref /results/PcanCIC/rsemeval-170930-101146/temp/rsemeval-170930-101146.transcripts.fa 1 /results/PcanCIC
Refs.makeRefs finished!
Refs.saveRefs finished!
/results/PcanCIC/rsemeval-170930-101146/temp/rsemeval-170930-101146.idx.fa is generated!
/results/PcanCIC/rsemeval-170930-101146/temp/rsemeval-170930-101146.n2g.idx.fa is generated!

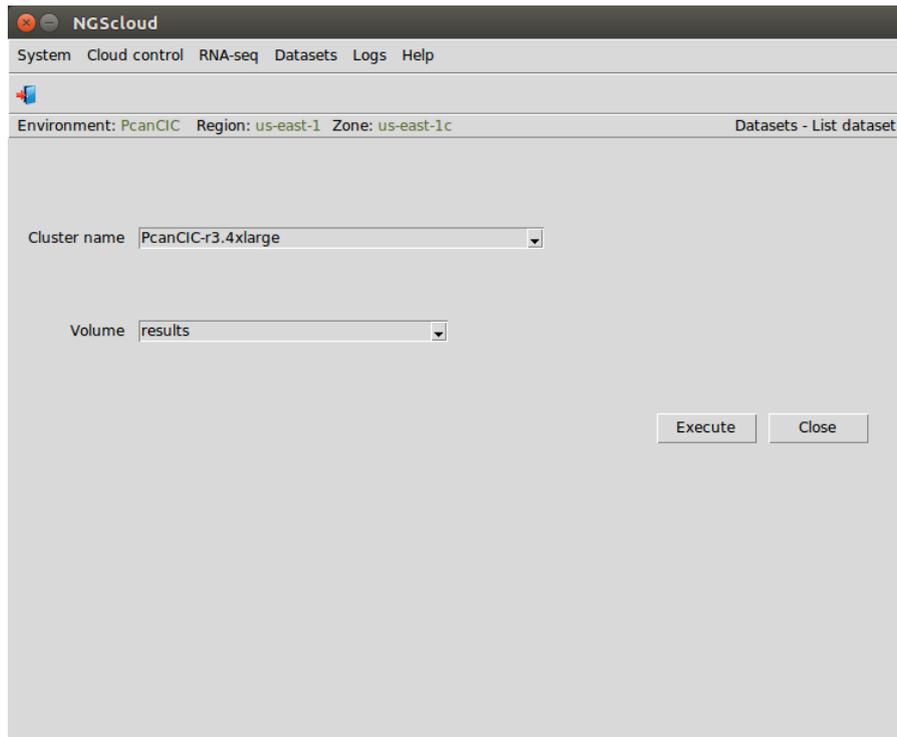
bowtie2-build -f /results/PcanCIC/rsemeval-170930-101146/temp/rsemeval-170930-101146.idx.fa /results/PcanCIC/rse
Settings:
  Output files: "/results/PcanCIC/rsemeval-170930-101146/temp/rsemeval-170930-101146.*.bt2"
  Line rate: 6 (line is 64 bytes)
  Lines per side: 1 (side is 64 bytes)
  Offset rate: 4 (one in 16)
  FTable chars: 10
  Strings: unpacked
  Max bucket size: default
  Max bucket size, sqrt multiplier: default
  Max bucket size, len divisor: 4
  Difference-cover sample period: 1024
  Endianness: little
  Actual local endianness: little
  Sanity checking: disabled
  Assertions: disabled
  Random seed: 0
  Sizeofs: void*:8, int:4, long:8, size_t:8
Input_files DNA FASTA:

```

Now, we are going to view the assessment files generated by RSEM-EVAL. We select the menu item with this path:

Main menu > Dataset > List dataset

In the raised window, we select **PcanCIC-r3.4xlarge** in the *Cluster name* combo-box and **results** in the *Volume* combo-box. Then we press the *Execute* button:

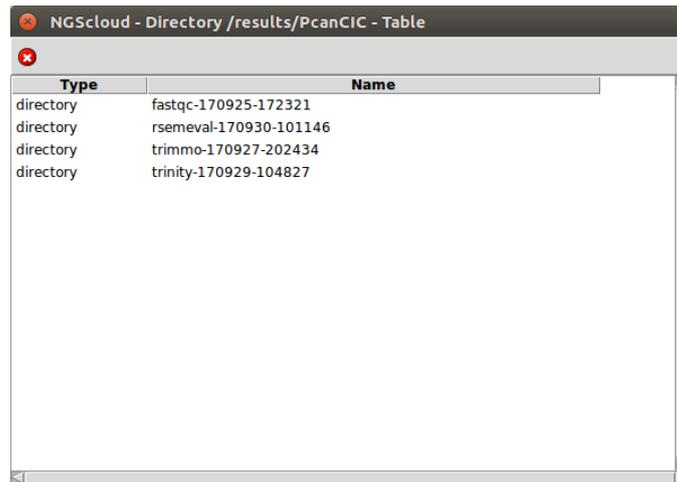


Now, we click in the **PcanCIC** row:

The screenshot shows a window titled 'NGScloud - Directory /results - Table'. The window contains a table with the following data:

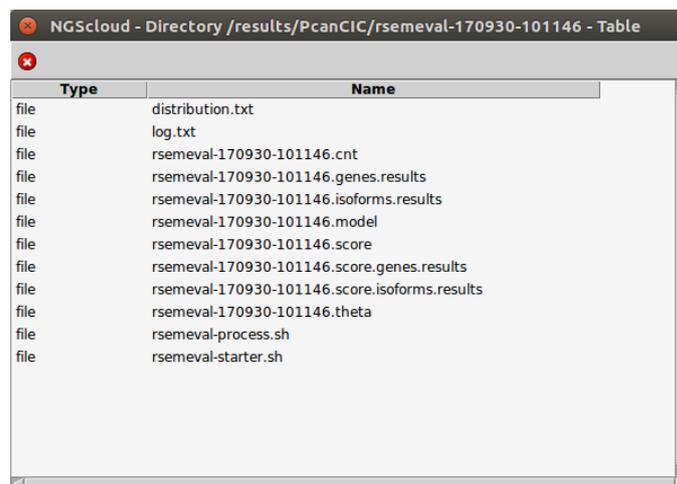
Type	Name
directory	Athaliana01x
directory	PcanCIC
directory	database
directory	test

Next, a window with the result datasets of the experiment PcanCIC is shown:



Type	Name
directory	fastqc-170925-172321
directory	rsemeval-170930-101146
directory	trimmo-170927-202434
directory	trinity-170929-104827

We click on the **rsemeval-170930-101146** row, and another window appears with the content of the files corresponding to this assembly.



Type	Name
file	distribution.txt
file	log.txt
file	rsemeval-170930-101146.cnt
file	rsemeval-170930-101146.genes.results
file	rsemeval-170930-101146.isoforms.results
file	rsemeval-170930-101146.model
file	rsemeval-170930-101146.score
file	rsemeval-170930-101146.score.genes.results
file	rsemeval-170930-101146.score.isoforms.results
file	rsemeval-170930-101146.theta
file	rsemeval-process.sh
file	rsemeval-starter.sh

The files whose name end in ".results" have the information about the assembly assessment.

[Terminate the cluster with r3.4xlarge template and create another cluster with r3.xlarge template](#)

Now we are going to terminate the PcanCIC-r3.4xlarge and to create a cluster with a r3.xlarge template, 4 CPUs and 30.5 GiB of RAM, because it is not necessary to use an instance with many CPUs and large RAM memory in order to .do the task of filtering.

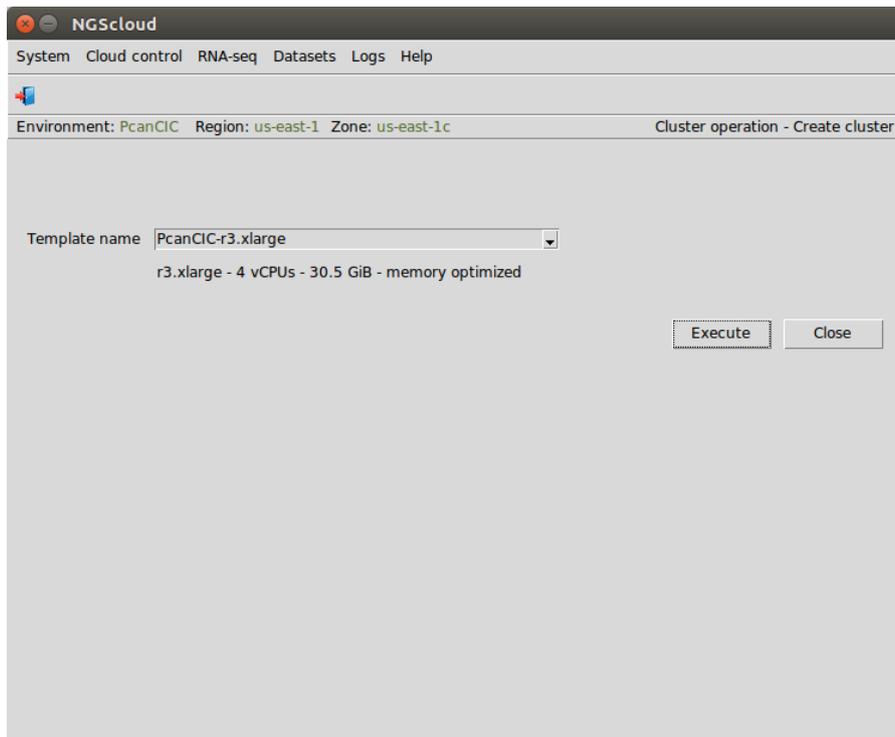
First, we select the menu item with this path:

Main menu > Cloud control > Cluster operation > Terminate cluster

Now we create a cluster with a t2.micro template. We select the menu item with this path:

Main menu > Cloud control > Cluster operation > Create cluster

In the raised window, we select **PcanCIC-r3.xlarge**, the template corresponding to a t2.micro instance type, in *Template name* combo-box; and then we press the *Execute* button:



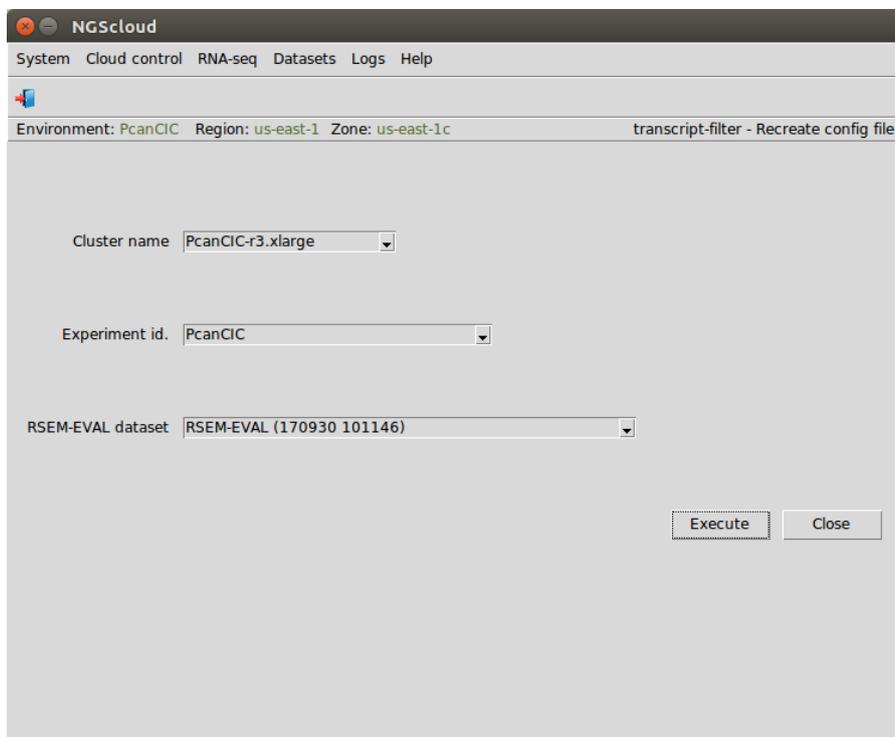
A window is raised displaying the run log:

Transcriptome filtering using transcript-filter

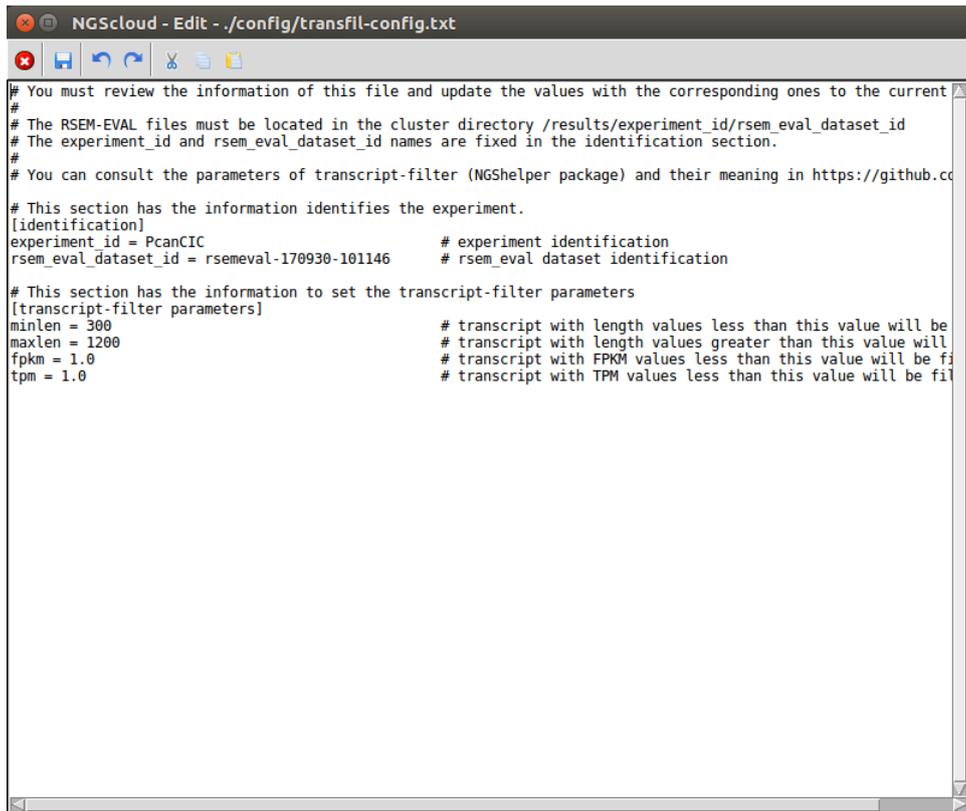
Transcript-filter uses a result of RSEM-EVAL to filter transcripts by length (max or min), or by FPKM or TPM. To do this step, we create the config file by selecting the menu item with this path:

Main menu > RNA-seq > Transcriptome filtering > transcript.filter (NGShelper package) > Recreate config file

In the raised window, we select **PcanCIC-r3.xlarge** in the *Cluster name* combo-box, **PcanCIC** in the *Experiment id* combo-box, and **RSEM-EVAL (170930 101146)** in the *RSEM-EVAL dataset* combo-box. Then we press the *Execute* button:



In the next window, we can examine the config file. There are two sections: *identification*, with the experiment and the RSEM-EVAL dataset identifications; and *transcript-filter parameters*, with the minimum and maximum lengths of transcripts, and the minimum FPKM and TPM values selected by the user:



```

NGScloud - Edit - ./config/transfil-config.txt
# You must review the information of this file and update the values with the corresponding ones to the current
#
# The RSEM-EVAL files must be located in the cluster directory /results/experiment_id/rsem_eval_dataset_id
# The experiment_id and rsem_eval_dataset_id names are fixed in the identification section.
#
# You can consult the parameters of transcript-filter (NGShelper package) and their meaning in https://github.co
# This section has the information identifies the experiment.
[identification]
experiment_id = PcanCIC                # experiment identification
rsem_eval_dataset_id = rsemeval-170930-101146 # rsem_eval dataset identification

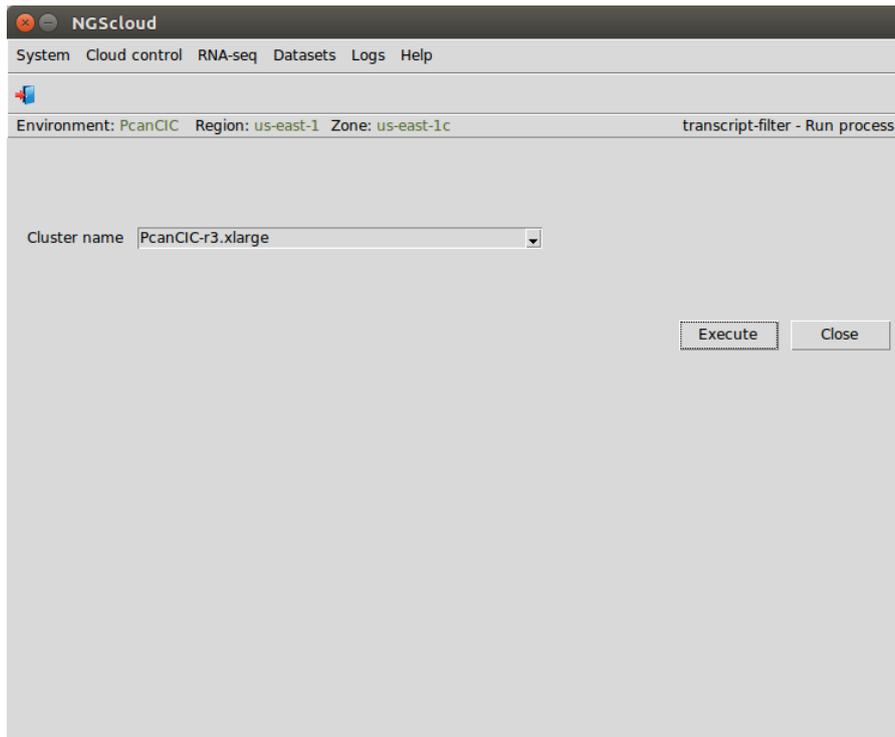
# This section has the information to set the transcript-filter parameters
[transcript-filter parameters]
minlen = 300                          # transcript with length values less than this value will be
maxlen = 1200                          # transcript with length values greater than this value will
fpkm = 1.0                             # transcript with FPKM values less than this value will be fi
tpm = 1.0                              # transcript with TPM values less than this value will be fi

```

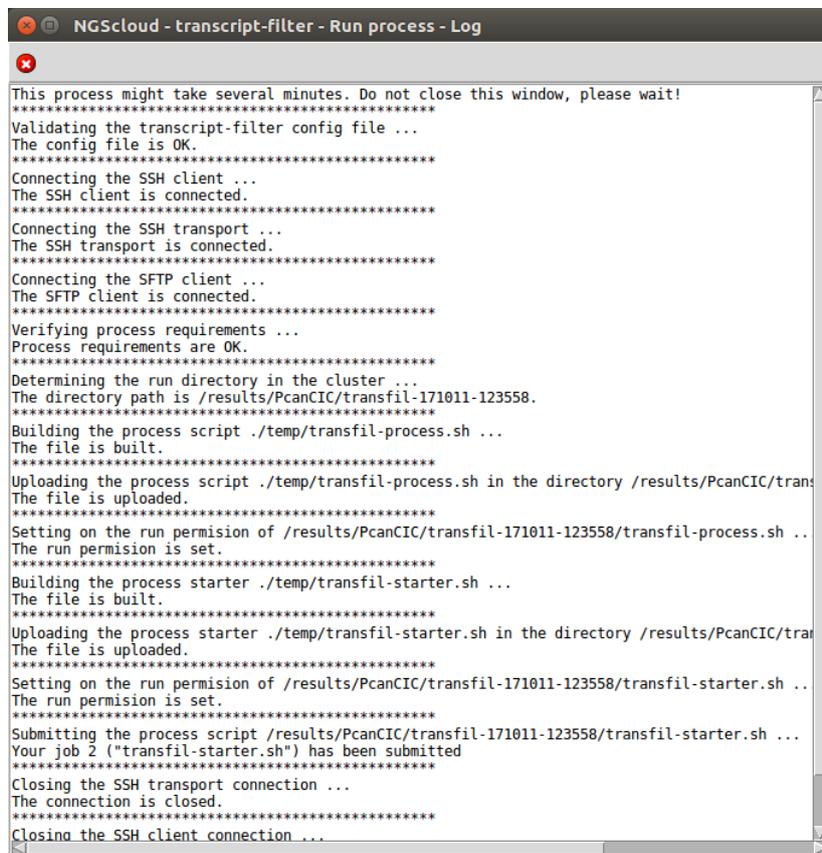
We run the assembly process in the cluster by selecting the menu item with this path:

Main menu > RNA-seq > Transcriptome filtering > transcript.filter (NGShelper package) > Run transcriptome filtering process

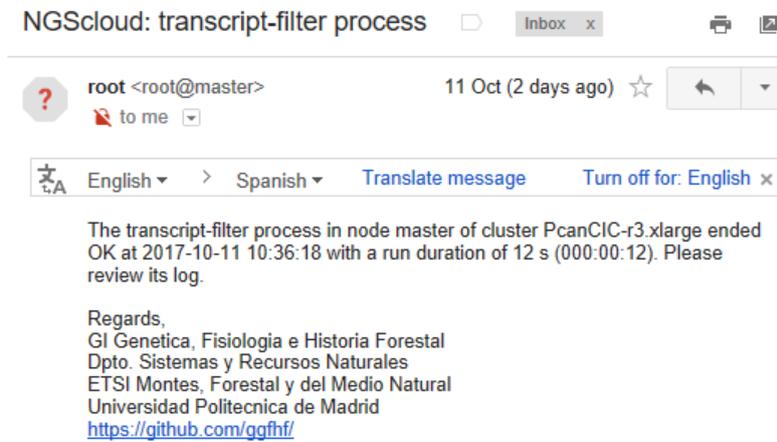
In the raised window, we select **PcanCIC-r.xlarge** in the *Cluster name* combo-box; and then we press the *Execute* button:



A window is raised showing the submission log:



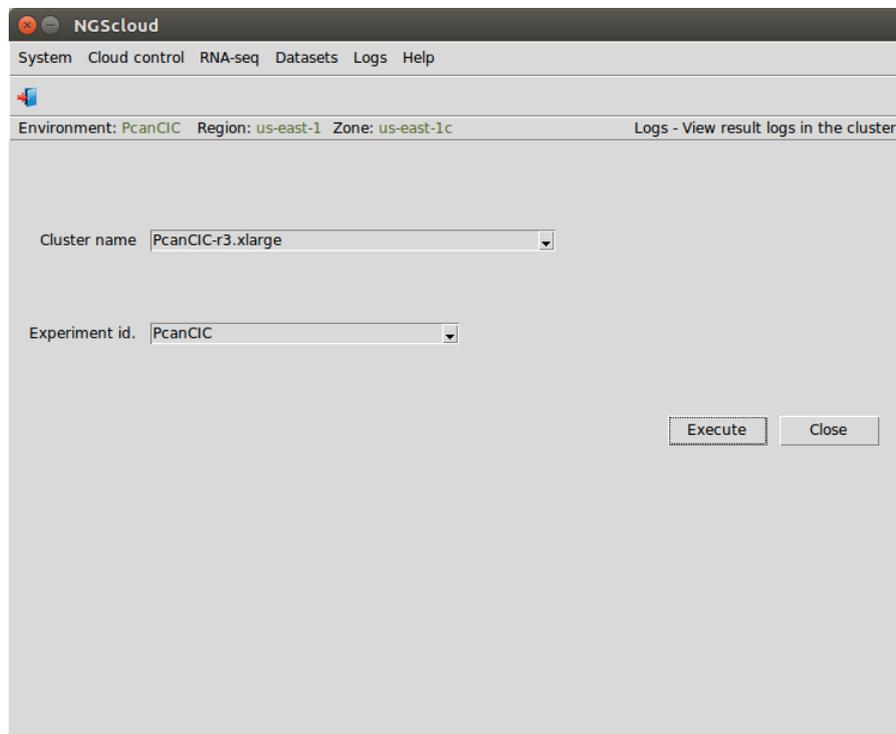
At the end of the run, an email is sent, informing of its completion:



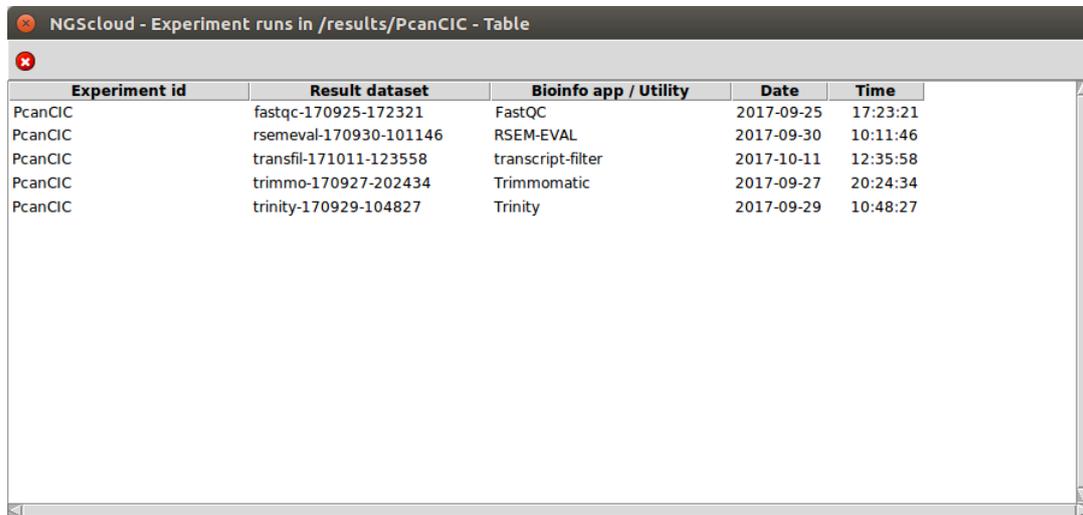
We can view the process log during and after the run. To do so, we select the menu item with this path:

Main menu > Logs > View result logs in the cluster

In the raised window, we select **PcanCIC-r3.xlarge** in the *Cluster name* combo-box and **PcanCIC** in the *Experiment id* combo-box; and then we press the *Execute* button:

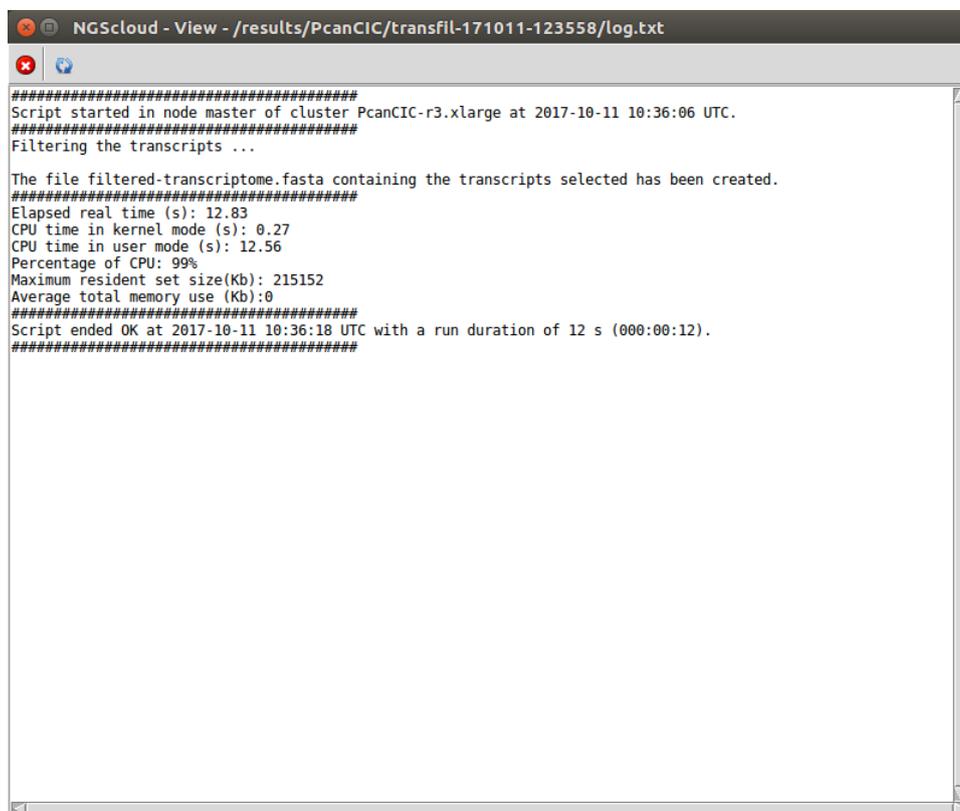


A window with the result datasets for each run of the bioinformatic programs that correspond to the experiment PcanCIC is shown:



Experiment id	Result dataset	Bioinfo app / Utility	Date	Time
PcanCIC	fastqc-170925-172321	FastQC	2017-09-25	17:23:21
PcanCIC	rsemval-170930-101146	RSEM-EVAL	2017-09-30	10:11:46
PcanCIC	transfil-171011-123558	transcript-filter	2017-10-11	12:35:58
PcanCIC	trimmo-170927-202434	Trimmomatic	2017-09-27	20:24:34
PcanCIC	trinity-170929-104827	Trinity	2017-09-29	10:48:27

Five result datasets are shown. We click on the **transfil-171011-123558** row, which is the dataset generated by the Trinity run, and another window appears with its corresponding log:



```

#####
Script started in node master of cluster PcanCIC-r3.xlarge at 2017-10-11 10:36:06 UTC.
#####
Filtering the transcripts ...

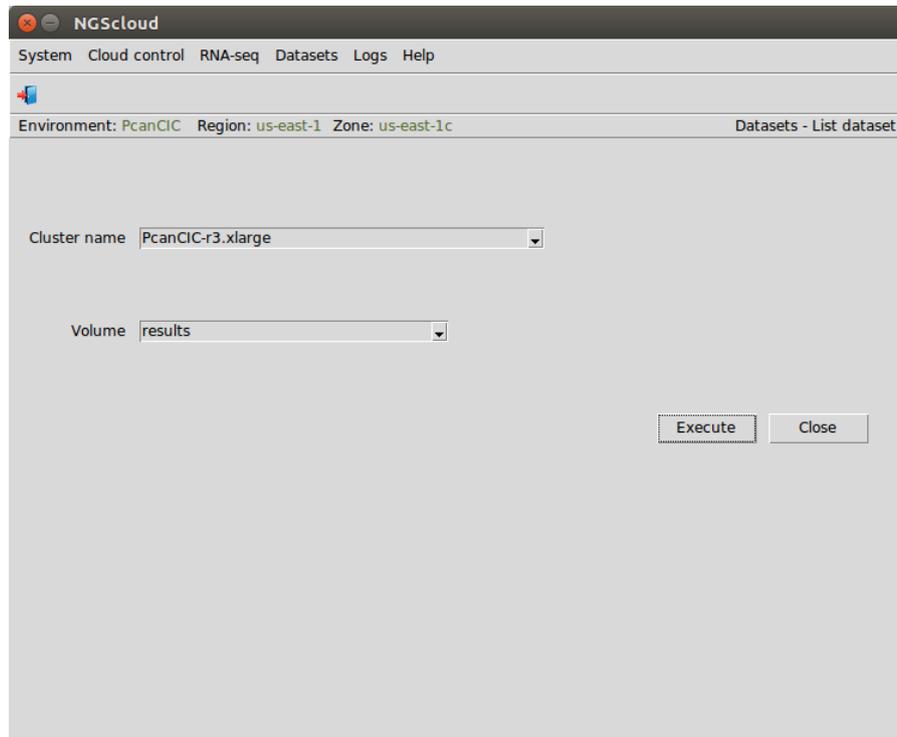
The file filtered-transcriptome.fasta containing the transcripts selected has been created.
#####
Elapsed real time (s): 12.83
CPU time in kernel mode (s): 0.27
CPU time in user mode (s): 12.56
Percentage of CPU: 99%
Maximum resident set size(Kb): 215152
Average total memory use (Kb):0
#####
Script ended OK at 2017-10-11 10:36:18 UTC with a run duration of 12 s (000:00:12).
#####

```

Now, we are going to list the assembly generated by transcript-filter. We select the menu item with this path:

Main menu > Dataset > List dataset

In the raised window, we select **PcanCIC-r3.xlarge** in the *Cluster name* combo-box and **results** in the *Volume* combo-box. Then we press the *Execute* button:

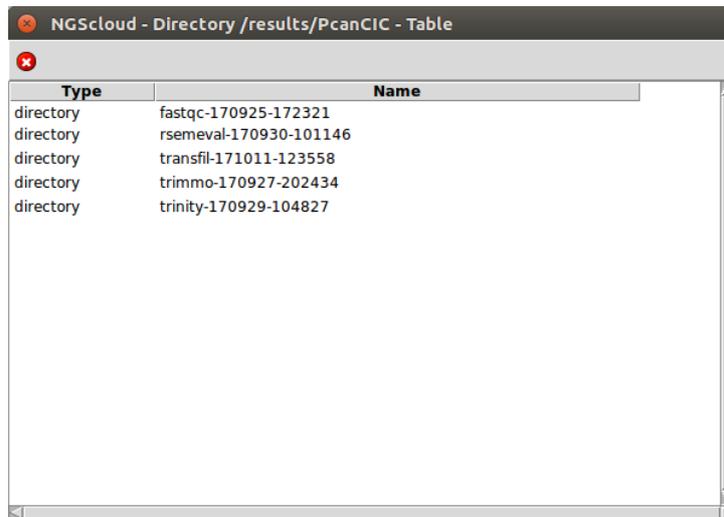


Now, we click in the **PcanCIC** row:

The screenshot shows a table window titled 'NGScloud - Directory /results - Table'. The table has two columns: 'Type' and 'Name'. The data rows are as follows:

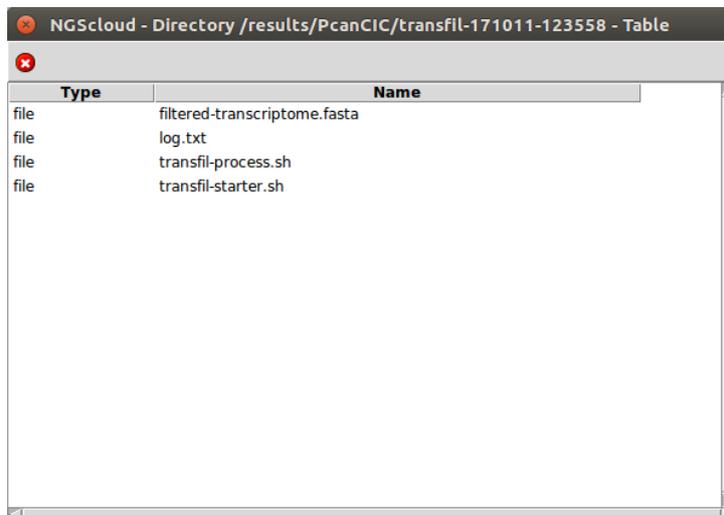
Type	Name
directory	Athaliana01x
directory	PcanCIC
directory	database
directory	test

Next, a window with the result datasets of the experiment PcanCIC is shown:



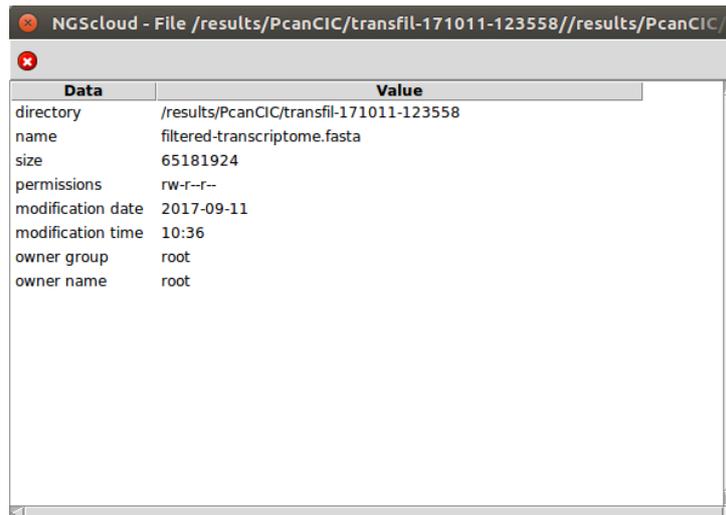
Type	Name
directory	fastqc-170925-172321
directory	rsemeval-170930-101146
directory	transfil-171011-123558
directory	trimmo-170927-202434
directory	trinity-170929-104827

We click on the **transfil-171011-123558** row, and another window appears with the content of the files corresponding to this assembly.



Type	Name
file	filtered-transcriptome.fasta
file	log.txt
file	transfil-process.sh
file	transfil-starter.sh

The file **filtered-transcriptome.fasta** is the one corresponding to the transcriptome generated by transcript.filter. To inspect its characteristics, we click on it:



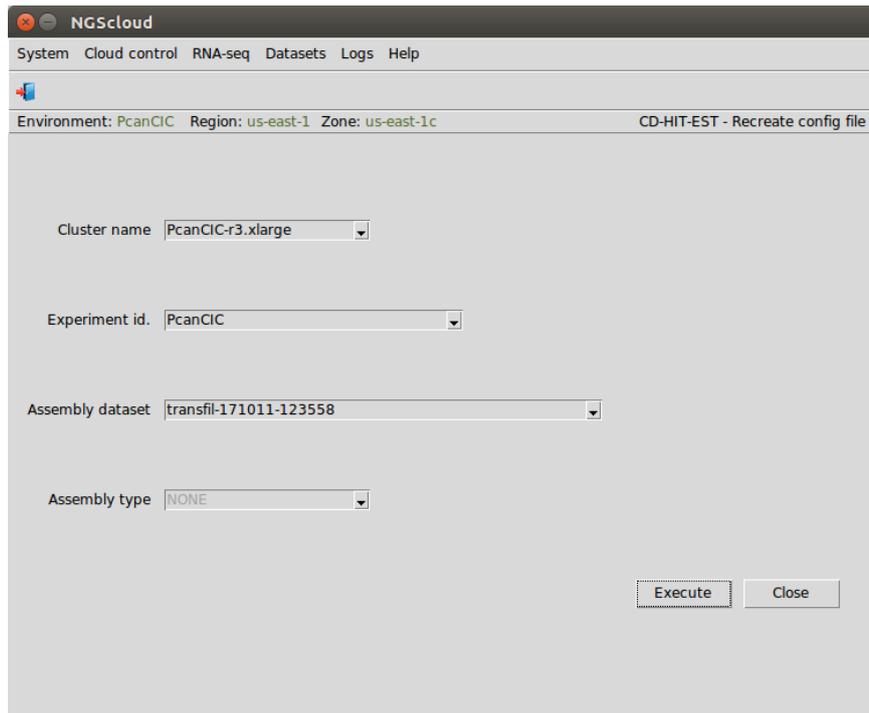
Data	Value
directory	/results/PcanCIC/transfil-171011-123558
name	filtered-transcriptome.fasta
size	65181924
permissions	rw-r--
modification date	2017-09-11
modification time	10:36
owner group	root
owner name	root

Transcriptome clustering using CD-HIT-EST

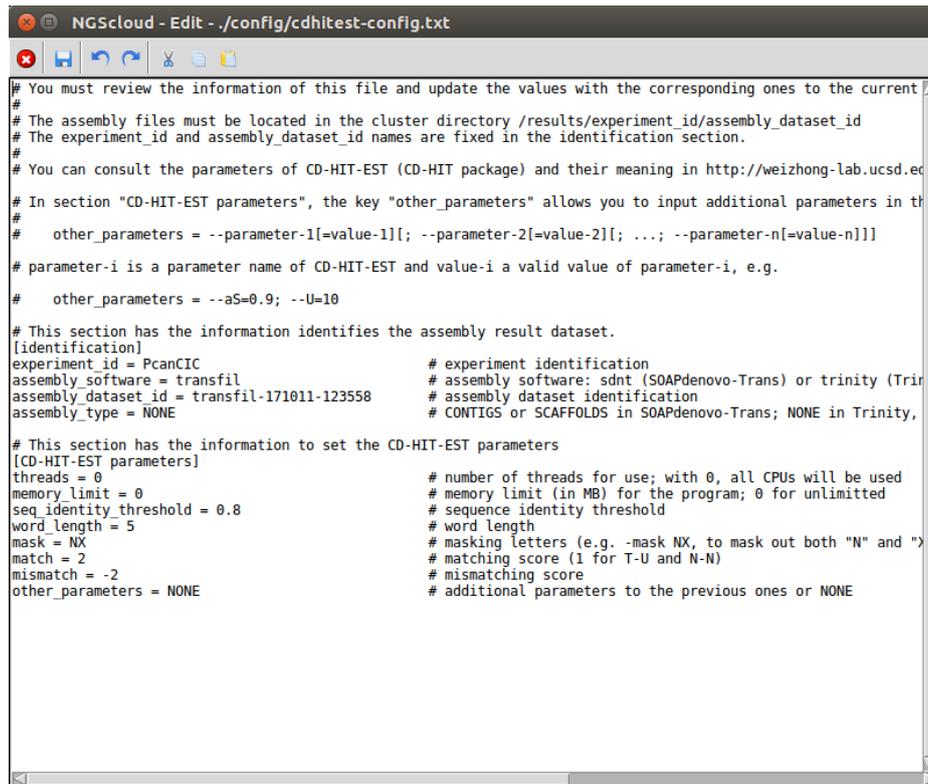
First, we create the config file by selecting the menu item with this path:

Main menu > RNA-seq > Transcriptome filtering > CD-HIT-EST (CD-HIT package) > Recreate config file

In the raised window, we select **PcanCIC-r3.xlarge** in the *Cluster name* combo-box, **PcanCIC** in the *Experiment id* combo-box, and **transcript-filter (171011 123538)** in the *Assembly dataset* combo-box. The *Assembly type* combo-box is only activated when the assembly dataset was generated by SOAPdenovo-Trans; in this case, the combo-box has two items: CONTIGS and SCAFFOLDS. Then we press the *Execute* button:



In the next window, we can examine the config file. In this example, it has two sections: *identification*, with the experiment and the assembly dataset identifications; *CD-HIT-EST parameters*, with several parameters used by CD-HIT-EST, where we can modify the value of threads number to **0** (this value indicates that all CPUs will be used), the value of memory_limit to **0** (this value indicates unlimited value), and the value sequence identity threshold to **0.8** (or any desired value above 0.8):



```

# You must review the information of this file and update the values with the corresponding ones to the current
#
# The assembly files must be located in the cluster directory /results/experiment_id/assembly_dataset_id
# The experiment_id and assembly_dataset_id names are fixed in the identification section.
#
# You can consult the parameters of CD-HIT-EST (CD-HIT package) and their meaning in http://weizhong-lab.ucsd.edu
#
# In section "CD-HIT-EST parameters", the key "other_parameters" allows you to input additional parameters in the
#
#   other_parameters = --parameter-1[=value-1]; --parameter-2[=value-2]; ...; --parameter-n[=value-n]]
#
# parameter-i is a parameter name of CD-HIT-EST and value-i a valid value of parameter-i, e.g.
#
#   other_parameters = --aS=0.9; --U=10
#
# This section has the information identifies the assembly result dataset.
[identification]
experiment_id = PcanCIC                # experiment identification
assembly_software = transfil           # assembly software: sdnt (SOAPdenovo-Trans) or trinity (Trinity)
assembly_dataset_id = transfil-171011-123558 # assembly dataset identification
assembly_type = NONE                  # CONTIGS or SCAFFOLDS in SOAPdenovo-Trans; NONE in Trinity,

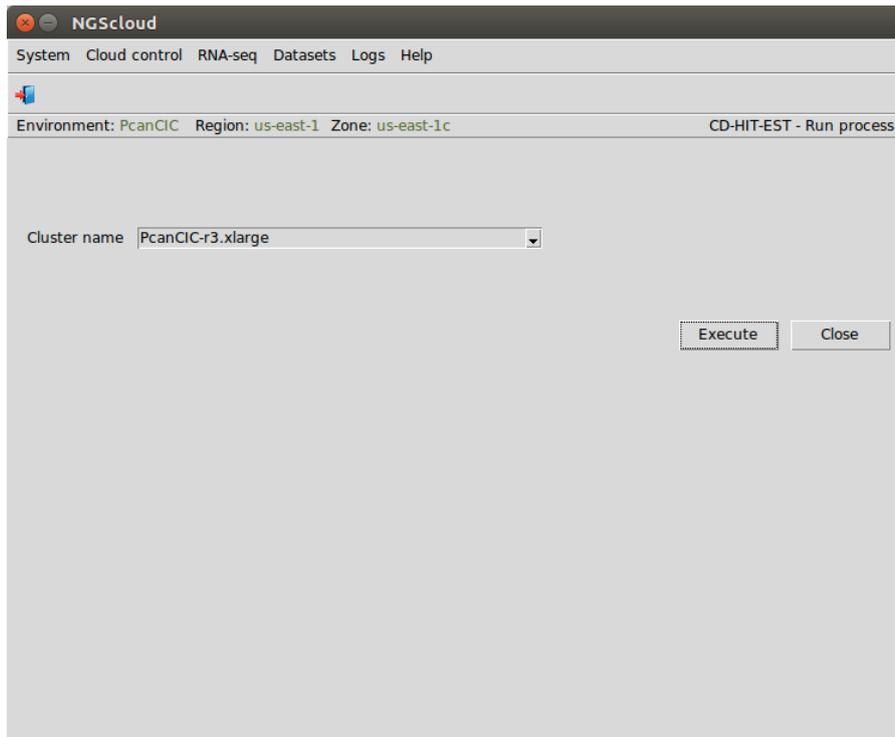
# This section has the information to set the CD-HIT-EST parameters
[CD-HIT-EST parameters]
threads = 0                            # number of threads for use; with 0, all CPUs will be used
memory_limit = 0                        # memory limit (in MB) for the program; 0 for unlimited
seq_identity_threshold = 0.8            # sequence identity threshold
word_length = 5                         # word length
mask = NX                               # masking letters (e.g. -mask NX, to mask out both "N" and "X")
match = 2                               # matching score (1 for T-U and N-N)
mismatch = -2                           # mismatching score
other_parameters = NONE                 # additional parameters to the previous ones or NONE

```

We run the assembly process in the cluster by selecting the menu item with this path:

Main menu > RNA-seq > Transcriptome filtering > CD-HIT-EST (CD-HIT package) > Run transcriptome filtering process

In the raised window, we select **PcanCIC-r3.xlarge** in the *Cluster name* combo-box; and then we press the *Execute* button:



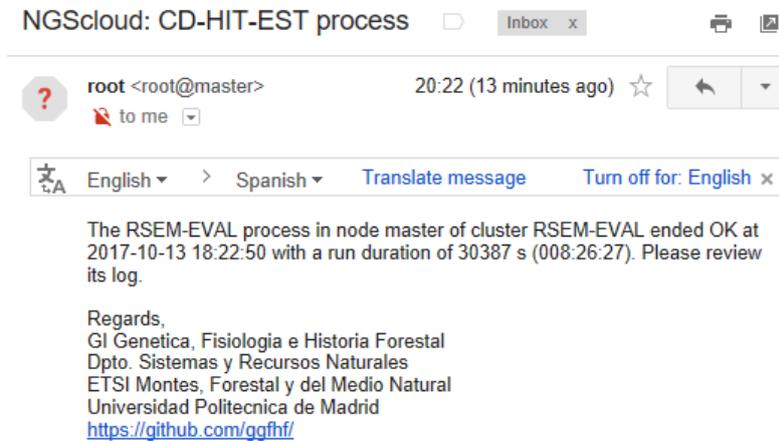
A window is raised showing the submission log:

```

NGScloud - CD-HIT-EST - Run process - Log
This process might take several minutes. Do not close this window, please wait!
*****
Validating the CD-HIT-EST config file ...
The config file is OK.
*****
Connecting the SSH client ...
The SSH client is connected.
*****
Connecting the SSH transport ...
The SSH transport is connected.
*****
Connecting the SFTP client ...
The SFTP client is connected.
*****
Verifying process requirements ...
Process requirements are OK.
*****
Determining the run directory in the cluster ...
The directory path is /results/PcanCIC/cdhitest-171013-115617.
*****
Building the process script ./temp/cdhitest-process.sh ...
The file is built.
*****
Uploading the process script ./temp/cdhitest-process.sh in the directory /results/PcanCIC/cdhitest-171013-115617/cdhitest-process.sh ...
The file is uploaded.
*****
Setting on the run permission of /results/PcanCIC/cdhitest-171013-115617/cdhitest-process.sh ...
The run permission is set.
*****
Building the process starter ./temp/cdhitest-starter.sh ...
The file is built.
*****
Uploading the process starter ./temp/cdhitest-starter.sh in the directory /results/PcanCIC/cdhitest-171013-115617/cdhitest-starter.sh ...
The file is uploaded.
*****
Setting on the run permission of /results/PcanCIC/cdhitest-171013-115617/cdhitest-starter.sh ...
The run permission is set.
*****
Submitting the process script /results/PcanCIC/cdhitest-171013-115617/cdhitest-starter.sh ...
Your job 1 ("cdhitest-starter.sh") has been submitted
*****
Closing the SSH transport connection ...
The connection is closed.
*****
Closing the SSH client connection ...

```

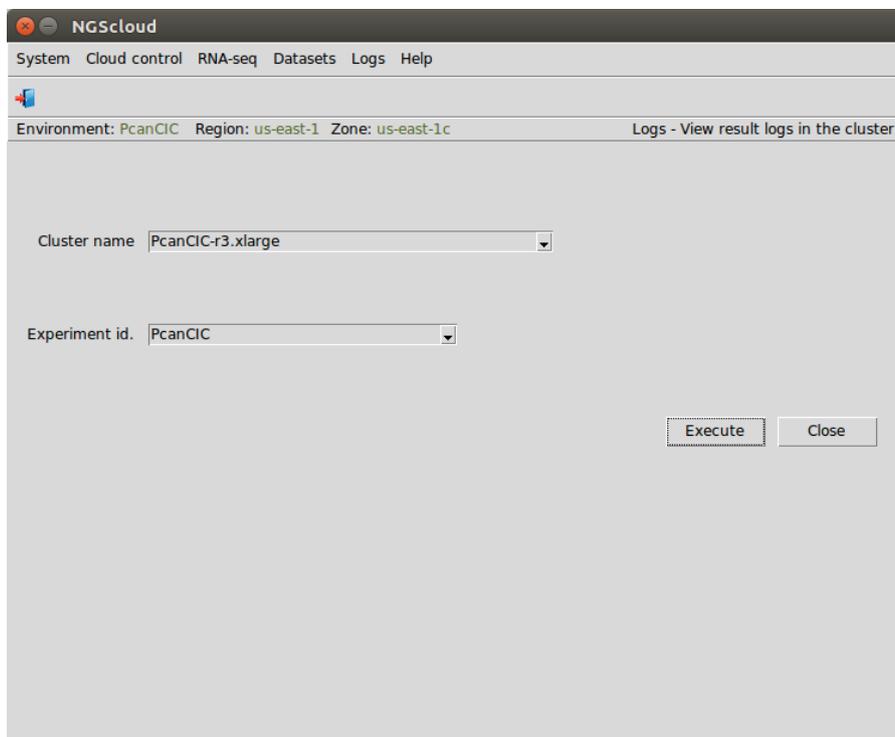
At the end of the run, an email is sent, informing of its completion:



We can view the process log during and after the run. To do so, we select the menu item with this path:

Main menu > Logs > View result logs in the cluster

In the raised window, we select **PcanCIC-r3.xlarge** in the *Cluster name* combo-box and **PcanCIC** in the *Experiment id* combo-box; and then we press the *Execute* button:



A window with the result datasets for each run of the bioinformatic programs that correspond to the experiment PcanCIC is shown:

NGScloud - Experiment runs in /results/PcanCIC - Table

Experiment id	Result dataset	Bioinfo app / Utility	Date	Time
PcanCIC	cdhitest-171013-115617	CD-HIT-EST	2017-10-13	11:56:17
PcanCIC	fastqc-170925-172321	FastQC	2017-09-25	17:23:21
PcanCIC	rsemeval-170930-101146	RSEM-EVAL	2017-09-30	10:11:46
PcanCIC	transfil-171011-123558	transcript-filter	2017-10-11	12:35:58
PcanCIC	trimmo-170927-202434	Trimmomatic	2017-09-27	20:24:34
PcanCIC	trinity-170929-104827	Trinity	2017-09-29	10:48:27

Six result datasets are shown. We click on the **cdhitest-171013-115617** row, which is the dataset generated by the CD-HIT-EST run, and another window appears with its corresponding log:

NGScloud - View - /results/PcanCIC/cdhitest-171013-115617/log.txt

```
#####
Script started in node master of cluster PcanCIC-r3.xlarge at 2017-10-13 09:56:23 UTC.
#####
Running CD-HIT-EST process ...
=====
Program: CD-HIT, V4.6 (+OpenMP), Sep 18 2016, 07:20:34
Command: cd-hit-est -T 0 -M 0 -i
         /results/PcanCIC/transfil-171011-123558/filtered-transcriptome.fasta
         -c 0.8 -n 5 -mask NX -match 2 -mismatch -2 -o
         /results/PcanCIC/cdhitest-171013-115617/clustered-transcriptome.fasta

Started: Fri Oct 13 09:56:23 2017
=====
Output
-----
total number of CPUs in the system is 4
Actual number of CPUs to be used: 4

total seq: 85629
longest and shortest : 1200 and 300
Total letters: 42874226
Sequences have been sorted

Approximated minimal memory consumption:
Sequence      : 53M
Buffer        : 4 X 13M = 53M
Table         : 2 X 1M = 2M
Miscellaneous : 1M
Total         : 110M

Table limit with the given memory limit:
Max number of representatives: 1500000
Max number of word counting entries: 114412500

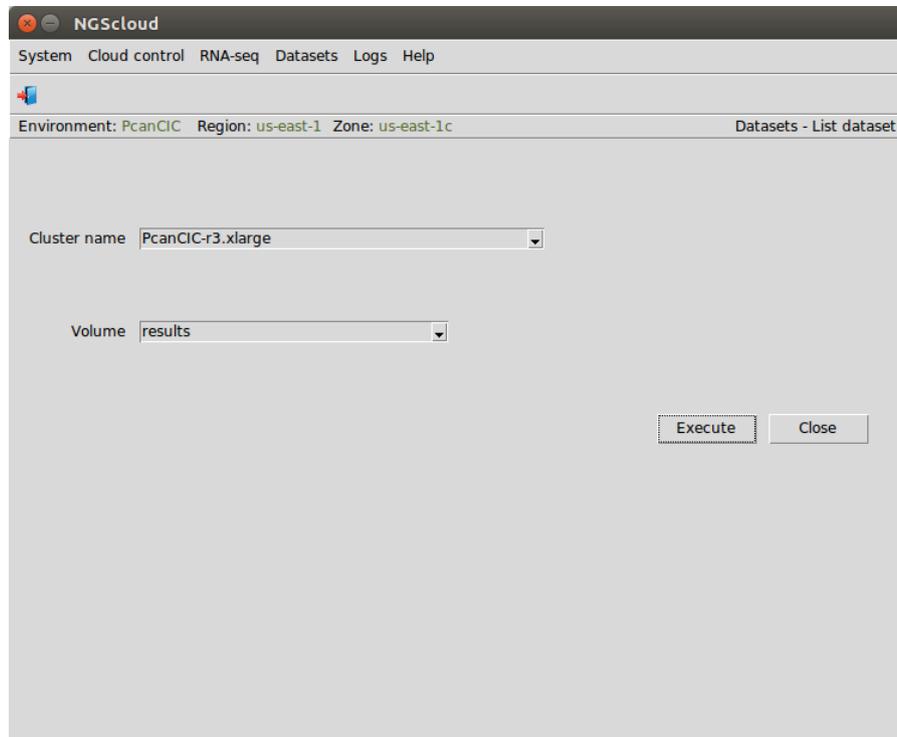
# comparing sequences from      0 to 14271
.....
..... 10000 finished      8777 clusters
..... new table with 12338 representatives

# comparing sequences from 14271 to 26164
.....
..... 20000 finished      17165 clusters
```

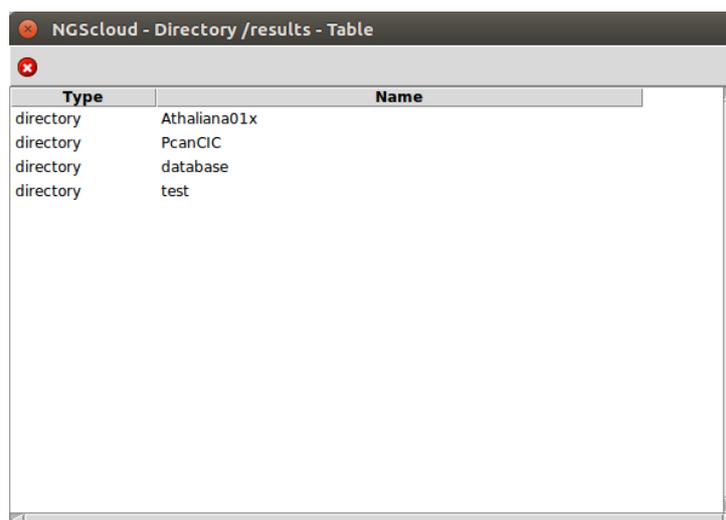
Now, we are going to list the assembly generated by CD-HIT-EST. We select the menu item with this path:

Main menu > Dataset > List dataset

In the raised window, we select **PcanCIC-r4.xlarge** in the *Cluster name* combo-box and **results** in the *Volume* combo-box. Then we press the *Execute* button:



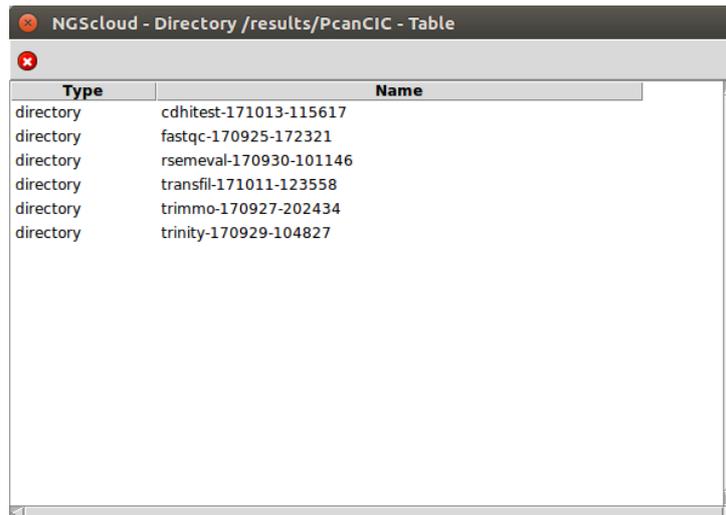
Now, we click in the **PcanCIC** row:



The screenshot shows a window titled "NGScloud - Directory /results - Table". It contains a table with two columns: "Type" and "Name". The table has four rows of data.

Type	Name
directory	Athaliana01x
directory	PcanCIC
directory	database
directory	test

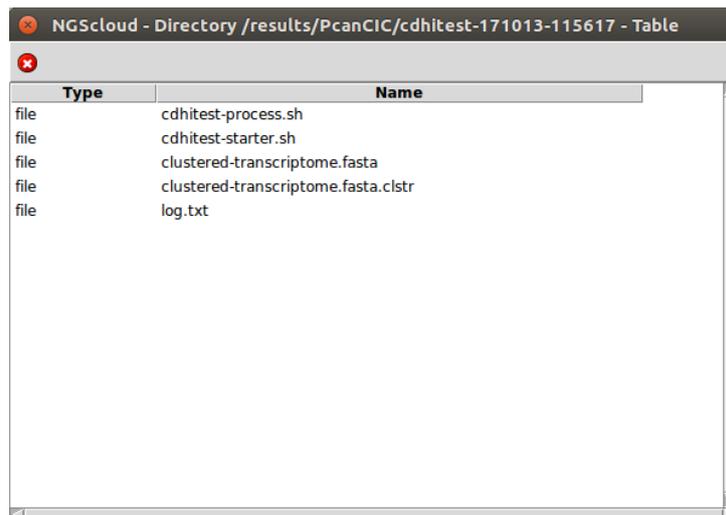
Next, a window with the result datasets of the experiment PcanCIC is shown:



NGScloud - Directory /results/PcanCIC - Table

Type	Name
directory	cdhitest-171013-115617
directory	fastqc-170925-172321
directory	rsemeval-170930-101146
directory	transfil-171011-123558
directory	trimmo-170927-202434
directory	trinity-170929-104827

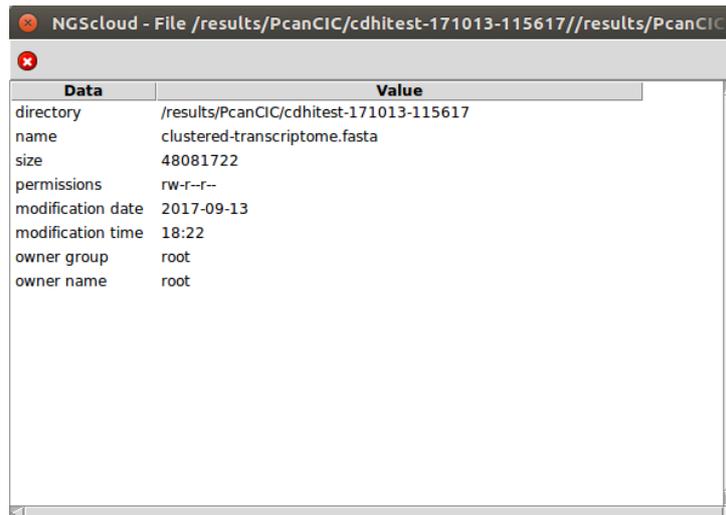
We click on the **cdhitest-171013-115617** row, and another window appears with the content of the files corresponding to this assembly.



NGScloud - Directory /results/PcanCIC/cdhitest-171013-115617 - Table

Type	Name
file	cdhitest-process.sh
file	cdhitest-starter.sh
file	clustered-transcriptome.fasta
file	clustered-transcriptome.fasta.clstr
file	log.txt

The file **Trinity.fasta** is the one corresponding to the transcriptome. To observe its characteristics, we click on it:



Data	Value
directory	/results/PcanCIC/cdhitest-171013-115617
name	clustered-transcriptome.fasta
size	48081722
permissions	rw-r--
modification date	2017-09-13
modification time	18:22
owner group	root
owner name	root

Terminate the cluster with r3.xlarge template and create another cluster with c3.xlarge template

Now we are going to terminate the PcanCIC-r3.xlarge and to create a cluster with a c3.xlarge template with the same CPU number but with less memory amount because it is not necessary in order to annotate the transcriptome and so we will save money.

First, we select the menu item with this path:

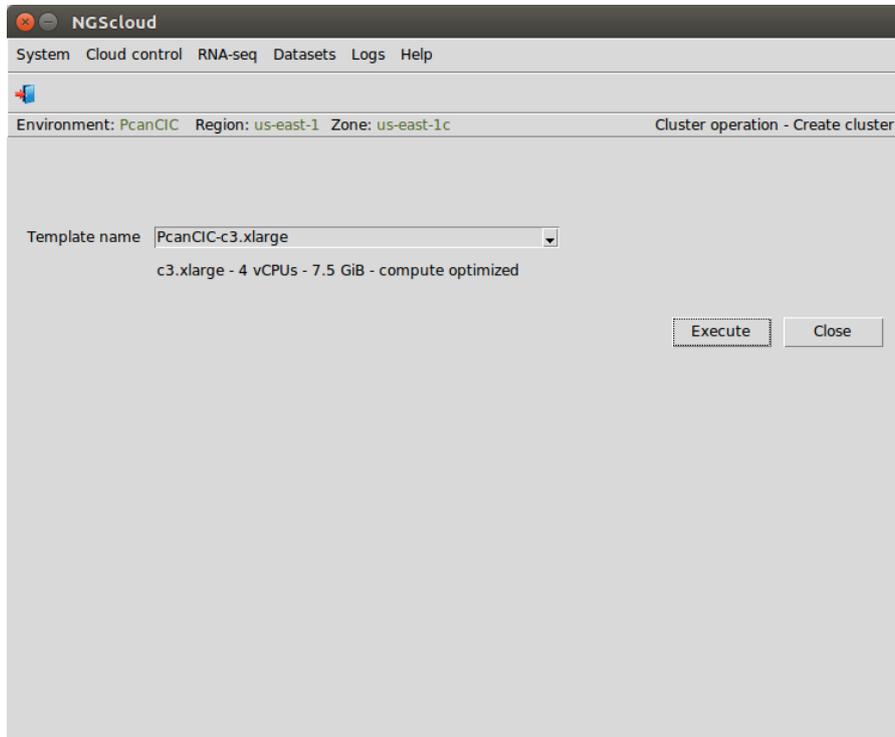
Main menu > Cloud control > Cluster operation > Terminate cluster

In the raised window, we select **PcanCIC-r3.xlarge** in the *Cluster name* combo-box; and then we press the *Execute* button:

Now we create a cluster with a c3.xlarge template. We select the menu item with this path:

Main menu > Cloud control > Cluster operation > Create cluster

In the raised window, we select **PcanCIC-c3.xlarge**, the template corresponding to a c3.xlarge instance type, in *Template name* combo-box; and then we press the *Execute* button:



A window is raised displaying the run log:

Upload the protein database to the cluster

Previously we had downloaded the FASTA protein files from TAIR (The Arabidopsis Information Resource) whose URL is

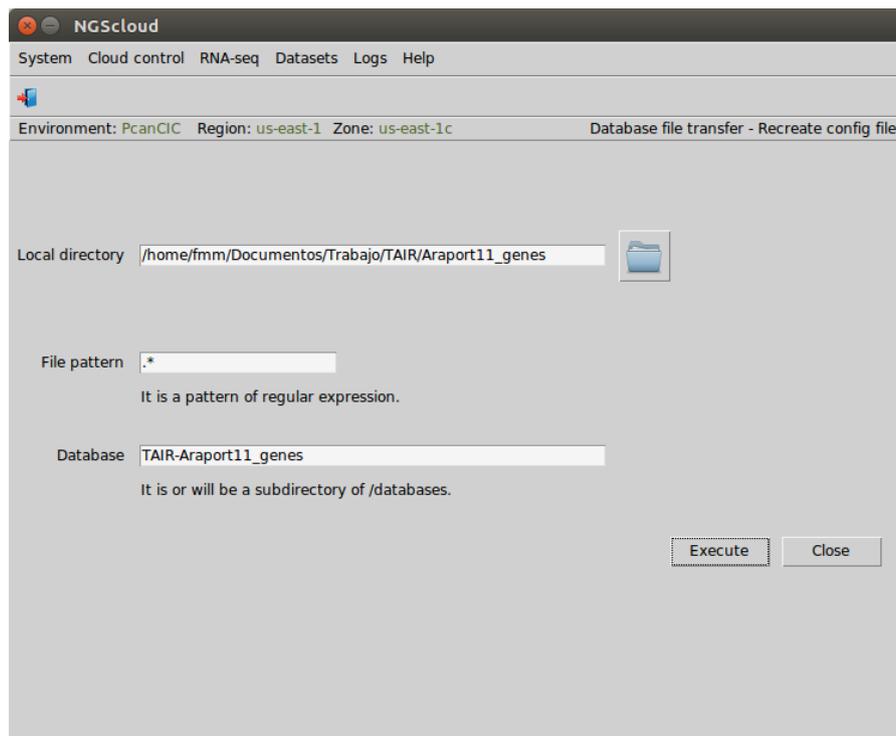
https://www.arabidopsis.org/download_files/Proteins/Araport11_protein_lists/Araport11_genes.201606.pep.fasta.gz

to the local computer and we had built the database with the program makeblastdb. The name of database file passed to makeblastdb is **Araport11_genes**.

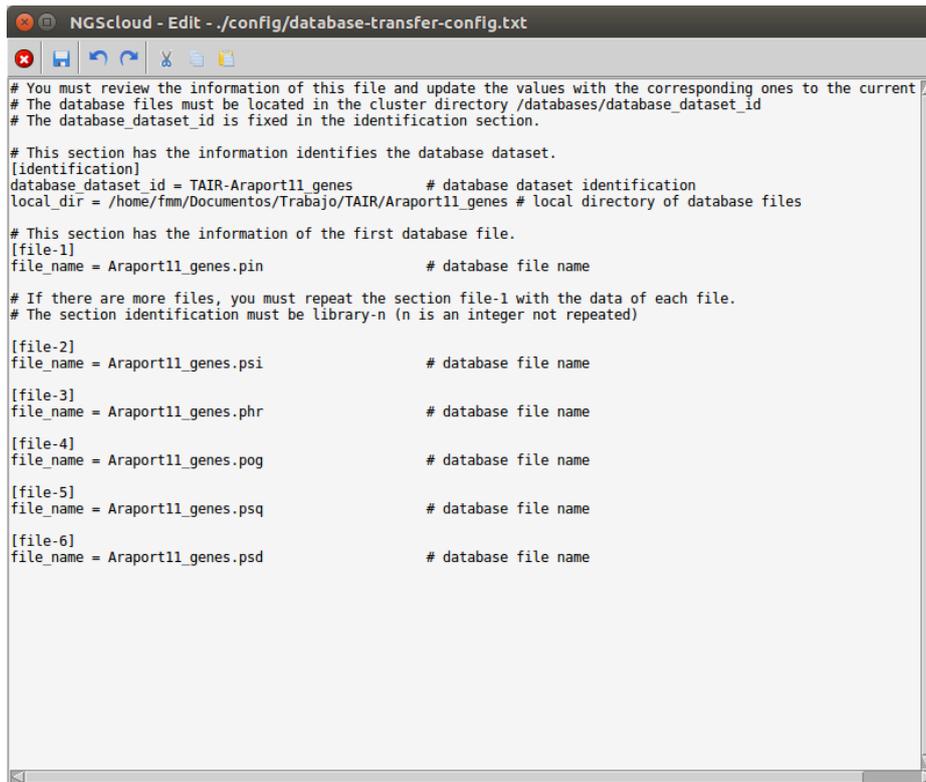
To create a config file to upload the database files to a cluster, select the menu item with this path:

Main menu > Dataset > Database file transfer > Recreate config file

In the raised window, we type the local directory where the database is in the *Local directory* textbox (or we select it using the next button), *.** as the pattern to select the files in the *File pattern* textbox, and **TAIR-Araport11_genes** in the *Database* textbox. Then we press the *Execute* button:



In the next window, we can edit the config file created. In this example, we can notice that the configure file has a section *identification*, with the database identification and the local directory where the database is; and several section *file-i*, with the name of each file of the local directory:



```

# You must review the information of this file and update the values with the corresponding ones to the current
# The database files must be located in the cluster directory /databases/database_dataset_id
# The database_dataset_id is fixed in the identification section.

# This section has the information identifies the database dataset.
[identification]
database_dataset_id = TAIR-Araport11_genes # database dataset identification
local_dir = /home/fmm/Documentos/Trabajo/TAIR/Araport11_genes # local directory of database files

# This section has the information of the first database file.
[file-1]
file_name = Araport11_genes.pin # database file name

# If there are more files, you must repeat the section file-1 with the data of each file.
# The section identification must be library-n (n is an integer not repeated)

[file-2]
file_name = Araport11_genes.psi # database file name

[file-3]
file_name = Araport11_genes.phr # database file name

[file-4]
file_name = Araport11_genes.pog # database file name

[file-5]
file_name = Araport11_genes.psq # database file name

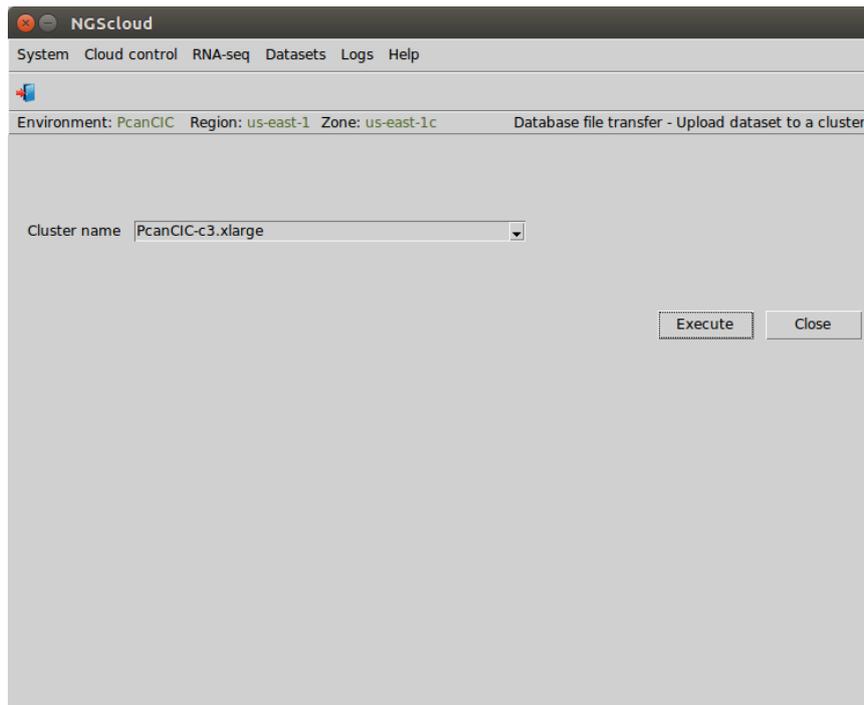
[file-6]
file_name = Araport11_genes.psd # database file name

```

To upload the database files to the cluster, we select the menu item with this path:

Main menu > Dataset > Database file transfer > Upload dataset to a cluster

In the raised window, we select **PcanCIC-c3.large** in *Cluster name* combo-box; and then we press the *Execute* button:



A window is raised with the upload log:

```

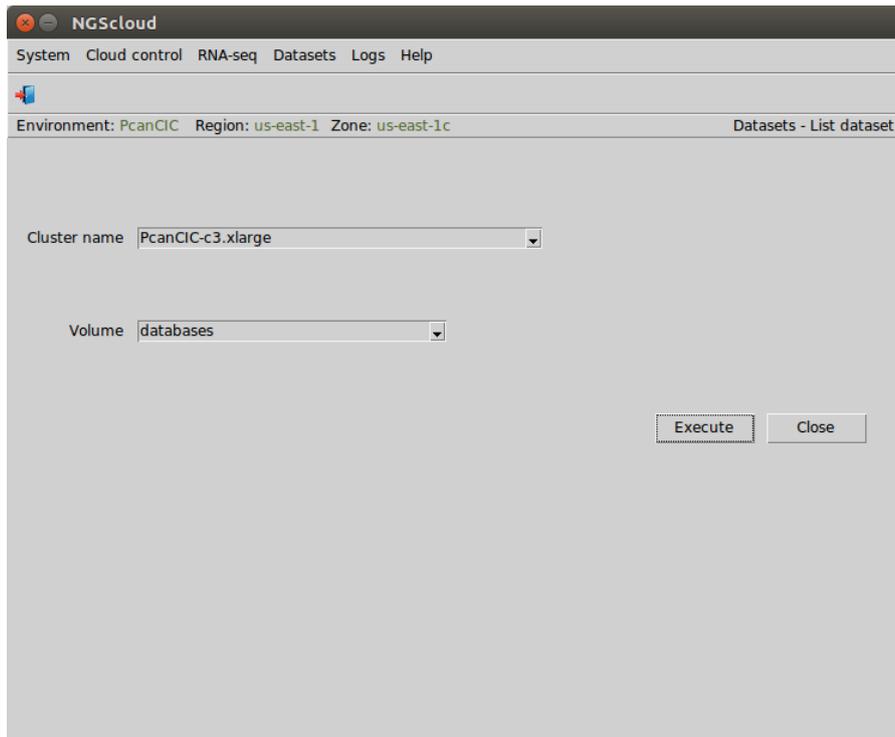
NGScloud - Database file transfer - Upload dataset to a cluster - Log
This process might take several minutes. Do not close this window, please wait!
*****
The database transfer config file is been validating ...
The config file is OK.
*****
The database directory /databases/TAIR-Araport11_genes in the cluster is being created ...
The directory is created.
*****
The file Araport11_genes.pin is being uploaded to /databases/TAIR-Araport11_genes ...
The file has been uploaded.
*****
The file Araport11_genes.psi is being uploaded to /databases/TAIR-Araport11_genes ...
The file has been uploaded.
*****
The file Araport11_genes.phr is being uploaded to /databases/TAIR-Araport11_genes ...
The file has been uploaded.
*****
The file Araport11_genes.pog is being uploaded to /databases/TAIR-Araport11_genes ...
The file has been uploaded.
*****
The file Araport11_genes.psq is being uploaded to /databases/TAIR-Araport11_genes ...
The file has been uploaded.
*****
The file Araport11_genes.psd is being uploaded to /databases/TAIR-Araport11_genes ...
The file has been uploaded.
*****
You can close this window now.

```

Now, we are going to review the uploaded files. We select the menu item with this path:

Main menu > Dataset > List dataset

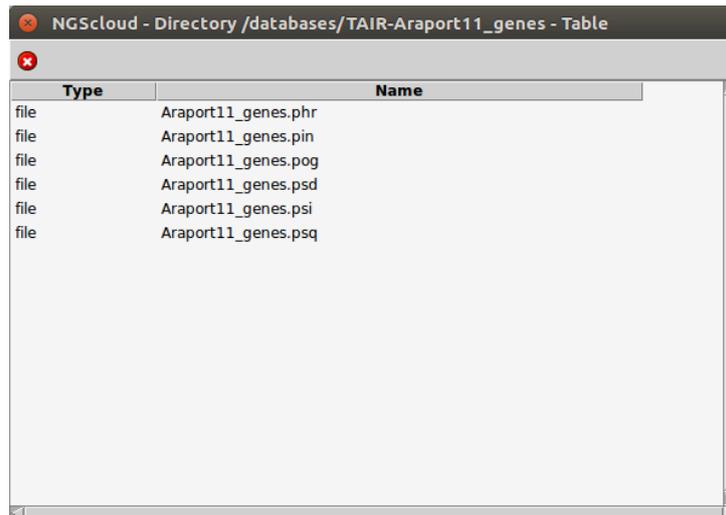
In the raised window, we select **PcanCIC-c3.xlarge** in the *Cluster name* combo-box and **databases** in the *Volume* combo-box. Then we press the *Execute* button:



We can check the experiments whose read files have been uploaded in the next window. We click in **TAIR-Araport11_genes** row:

Type	Name
directory	TAIR-Araport11_genes

So far, we only have one: the dataset corresponding to the uploaded-database. We click on it and another window appears with the content of the uploaded-database:



The screenshot shows a window titled "NGScloud - Directory /databases/TAIR-Araport11_genes - Table". Inside the window is a table with two columns: "Type" and "Name". The table contains six rows of file entries.

Type	Name
file	Araport11_genes.phr
file	Araport11_genes.pin
file	Araport11_genes.pog
file	Araport11_genes.psd
file	Araport11_genes.psi
file	Araport11_genes.psq

Add nodes to the cluster with a c3.xlarge template

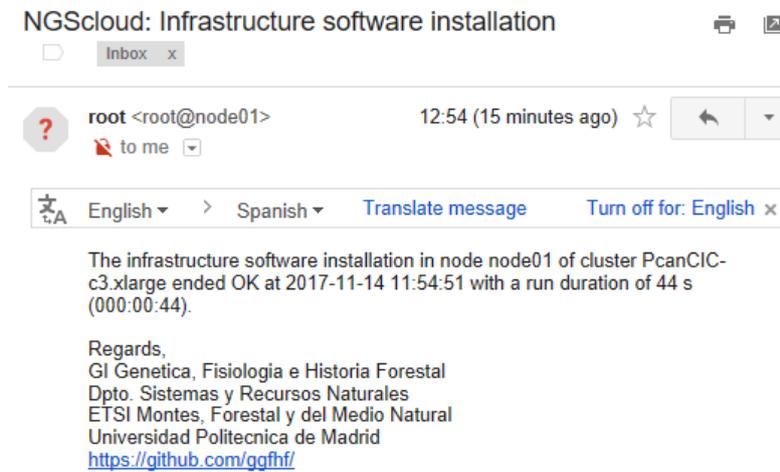
We are going to use transcriptome-blastx to annotate the transcriptome. This program supports parallelization, so it can use several nodes to increase the run speed. In this example, we are going to add 4 nodes to the cluster **PcanCIC-c3.xlarge**. So, we will have 5 nodes running, the master and the 4 subsidiary nodes, in such a way one node distributing the work to the other 4.

To add a node to a cluster, we select the menu item with this path:

Main menu > Cloud control > Node operation > Add node in a cluster

In the raised window, we select **PcanCIC-c3.xlarge** in the *Cluster name* combo-box and **node01** in the *Volume* combo-box. Then we press the *Execute* button:

At the end of the run, an email is sent, informing of its completion:

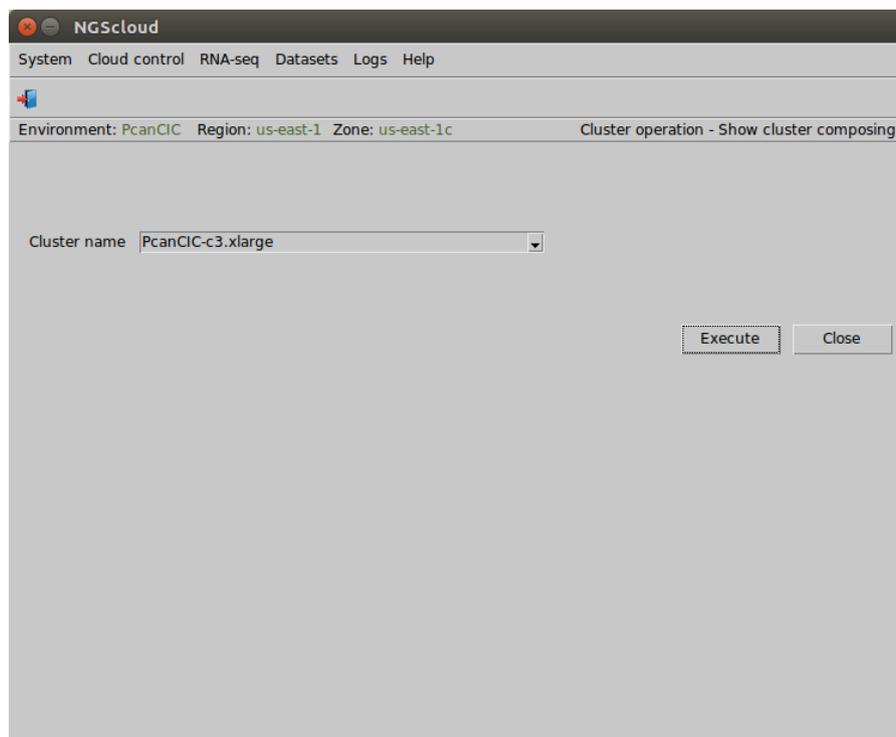


We repeat these actions to add the nodes **node2**, **node3** and **node04**.

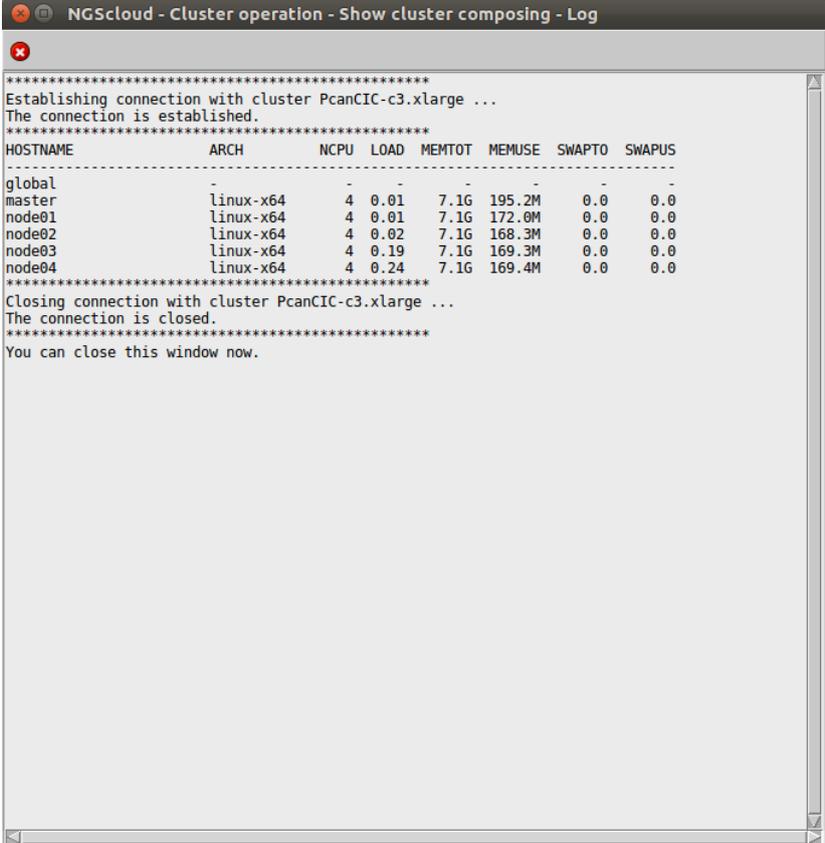
Now we are going to inspect the cluster composition selecting the menu item with this path:

Main menu > Cloud control > Cluster operation > Show cluster composing

In the raised window, we select **PcanCIC-c3.xlarge** in the *Cluster name* combo-box and press the *Execute* button:



A window is raised with the cluster composition and the characteristics of the nodes:



```

*****
Establishing connection with cluster PcanCIC-c3.xlarge ...
The connection is established.
*****
HOSTNAME          ARCH          NCPU  LOAD  MEMTOT  MEMUSE  SWAPTO  SWAPUS
-----
global            -             -     -     -        -        -        -
master            linux-x64     4  0.01  7.1G    195.2M  0.0     0.0
node01            linux-x64     4  0.01  7.1G    172.0M  0.0     0.0
node02            linux-x64     4  0.02  7.1G    168.3M  0.0     0.0
node03            linux-x64     4  0.19  7.1G    169.3M  0.0     0.0
node04            linux-x64     4  0.24  7.1G    169.4M  0.0     0.0
*****
Closing connection with cluster PcanCIC-c3.xlarge ...
The connection is closed.
*****
You can close this window now.

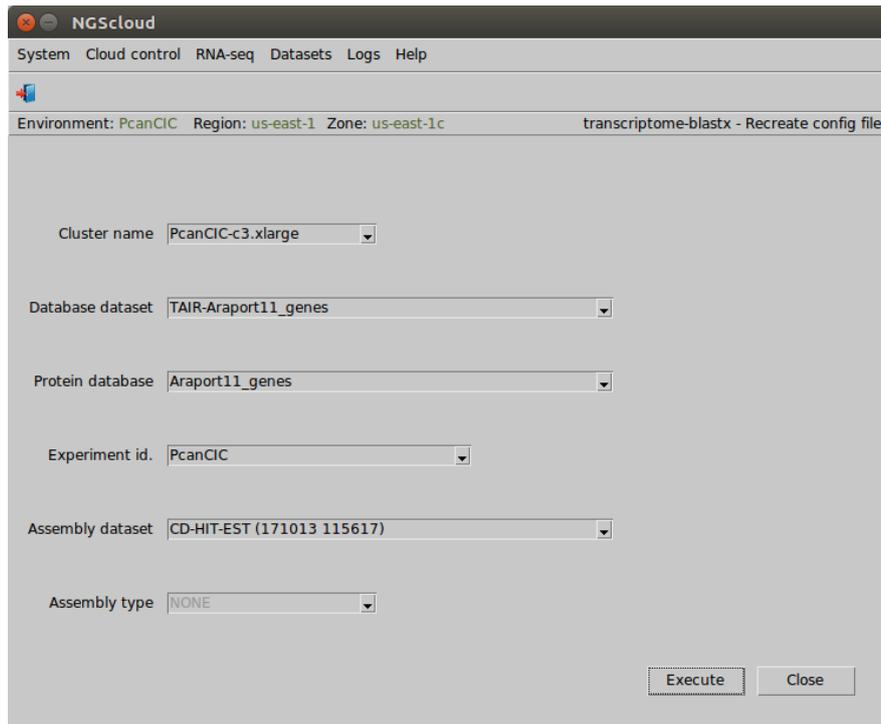
```

Annotate the filtered and clustered transcriptome using transcriptome-blastx

First, we create the config file by selecting the menu item with this path:

Main menu > RNA-seq > Annotation > transcriptome-blastx (NGShelper package) > Recreate config file

In the raised window, we select **PcanCIC-c3.xlarge** in the *Cluster name* combo-box, **TAIR-Araport11_genes** in the *Database dataset* combo-box, **Araport11_genes** in the *Database file* combo-box, **PcanCIC** in the *Experiment id* combo-box, and **CD-HIT-EST (171013 115617)** in the *Assembly dataset* combo-box. The *Assembly type* combo-box is only activated when the assembly dataset was generated by SOAPdenovo-Trans; in this case, the combo-box has two items: CONTIGS and SCAFFOLDS. Then we press the *Execute* button:



In the next window, we can examine the config file. There are two sections: *identification*, with the database, experiment and assembly identifications; and *transcriptome-blastx parameters*, with the parameters used by transcriptome-blastx. We modify the node number to **4** and the threads number by node to **4** (every node has 4 CPUs):

```

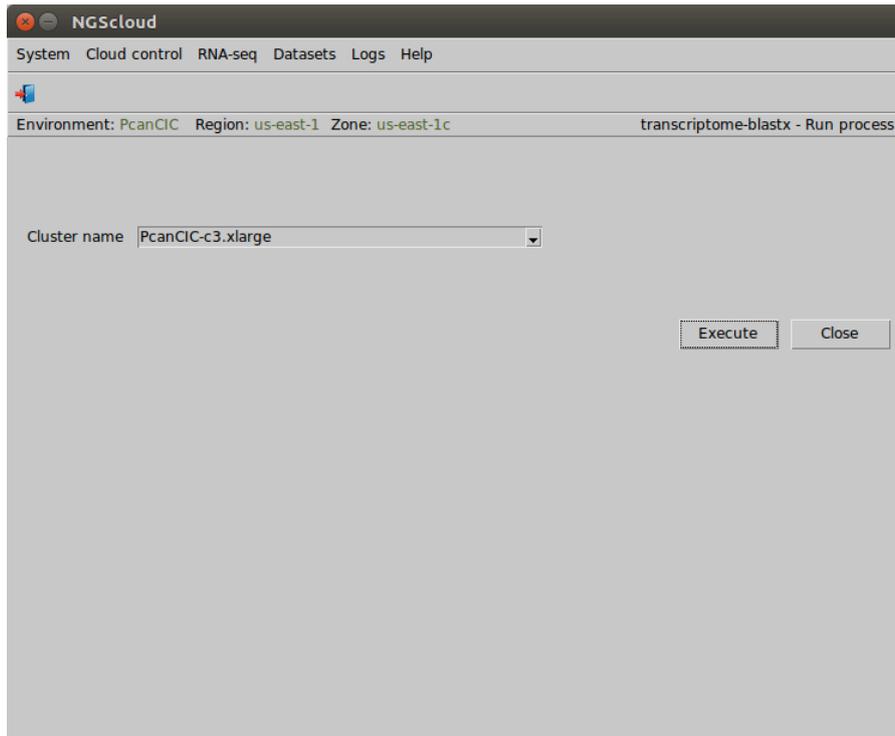
NGScloud - Edit - ./config/transbastx-config.txt
# You must review the information of this file and update the values with the corresponding ones to the current
#
# The database files must be located in the cluster directory /databases/database_dataset_id
# The assembly files must be located in the cluster directory /results/experiment_id/assembly_dataset_id
# The experiment_id, database_dataset_id and assembly_dataset_id names are fixed in the identification section.
#
# You can consult the parameters of transcriptome-blastx (NGShelper package) and their meaning in https://github.com
#
# This section has the information identifies the experiment.
[identification]
database_dataset_id = TAIR-Araport11_genes      # database dataset identification
protein_database_name = Araport11_genes        # protein database name
experiment_id = PcanCIC                        # experiment identification
assembly_software = cdhitest                   # assembly software: sdnt (SOAPdenovo-Trans) or trinity (Trinity)
assembly_dataset_id = cdhitest-171013-115617   # assembly dataset identification
assembly_type = NONE                          # CONTIGS or SCAFFOLDS in SOAPdenovo-Trans; NONE in Trinity
#
# This section has the information to set the transcriptome-blastx parameters
[transcriptome-blastx parameters]
node_number = 4                               # node number (they must be started)
blastx_thread_number = 4                      # threads number using by blastx in every node
e_value = 1E-6                                # expectation value (E-value) threshold for saving hits
max_target_seqs = 10                          # maximum number of aligned sequences to keep
max_hsp = 999999                              # maximum number of HSPs per subject sequence to save for each
qcov_hsp_perc = 0.0                          # alignments below the specified query coverage per HSPs are

```

We run the annotation process in the cluster by selecting the menu item with this path:

Main menu > RNA-seq > Annotation > transcriptome-blastx (NGShelper package) > Run annotation process

In the raised window, we select **PcanCIC-c3.xlarge** in the *Cluster name* combo-box; and then we press the *Execute* button:



A window is raised showing the submission log:

```

NGScloud - transcriptome-blastx - Run process - Log
This process might take several minutes. Do not close this window, please wait!
*****
Validating the transcriptome-blastx config file ...
The config file is OK.
*****
Connecting the SSH client ...
The SSH client is connected.
*****
Connecting the SSH transport ...
The SSH transport is connected.
*****
Connecting the SFTP client ...
The SFTP client is connected.
*****
Verifying process requirements ...
Process requirements are OK.
*****
Determining the run directory in the cluster ...
The directory path is /results/PcanCIC/transbastx-171114-133352.
*****
Building the process script ./temp/transbastx-process.sh ...
The file is built.
*****
Uploading the process script ./temp/transbastx-process.sh in the directory /results/PcanCIC/transbastx-171114-133352/transbastx-process.sh ...
The file is uploaded.
*****
Setting on the run permission of /results/PcanCIC/transbastx-171114-133352/transbastx-process.sh ...
The run permission is set.
*****
Building the process starter ./temp/transbastx-starter.sh ...
The file is built.
*****
Uploading the process starter ./temp/transbastx-starter.sh in the directory /results/PcanCIC/transbastx-171114-133352/transbastx-starter.sh ...
The file is uploaded.
*****
Setting on the run permission of /results/PcanCIC/transbastx-171114-133352/transbastx-starter.sh ...
The run permission is set.
*****
Submitting the process script /results/PcanCIC/transbastx-171114-133352/transbastx-starter.sh ...
Your job 1 ("transbastx-starter.sh") has been submitted
*****
Closing the SSH transport connection ...
The connection is closed.
*****
Closing the SSH client connection ...

```

At the end of the run, an email is sent, informing of its completion:

```

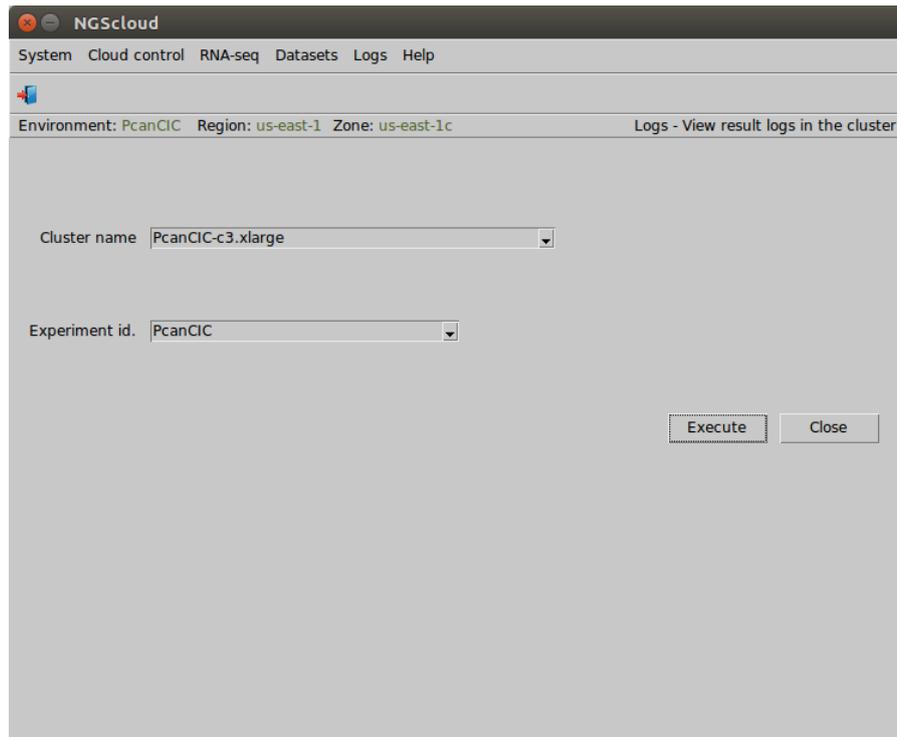
NGScloud: transcriptome-blastx process
root <root@master> 14:02 (37 minutes ago)
to me
English Spanish Translate message Turn off for: English
The transcript-filter process in node master of cluster PcanCIC-c3.xlarge ended
OK at 2017-11-14 13:02:28 with a run duration of 1711 s (000:28:31). Please
review its log.
Regards,
GI Genetica, Fisiologia e Historia Forestal
Dpto. Sistemas y Recursos Naturales
ETSI Montes, Forestal y del Medio Natural
Universidad Politecnica de Madrid
https://github.com/ggfhf/

```

We can view the process log during and after the run. To do so, we select the menu item with this path:

Main menu > Logs > View result logs in the cluster

In the raised window, we select **PcanCIC-c3.xlarge** in the *Cluster name* combo-box and **PcanCIC** in the *Experiment id* combo-box; and then we press the *Execute* button:



A window with the result datasets for each run of the bioinformatic programs that correspond to the experiment PcanCIC is shown:

Experiment id	Result dataset	Bioinfo app / Utility	Date	Time
PcanCIC	cdhitest-171013-115617	CD-HIT-EST	2017-10-13	11:56:17
PcanCIC	fastqc-170925-172321	FastQC	2017-09-25	17:23:21
PcanCIC	rsemval-170930-101146	RSEM-EVAL	2017-09-30	10:11:46
PcanCIC	transbastx-171114-133352	transcriptome-blastx	2017-11-14	13:33:52
PcanCIC	transfil-171011-123558	transcript-filter	2017-10-11	12:35:58
PcanCIC	trimmo-170927-202434	Trimmomatic	2017-09-27	20:24:34
PcanCIC	trinity-170929-104827	Trinity	2017-09-29	10:48:27

Seven result datasets are shown. We click on the **transbastx-171114-133352** row, which is the dataset generated by the transcriptome-blastx run, and another window appears with its corresponding log:

```

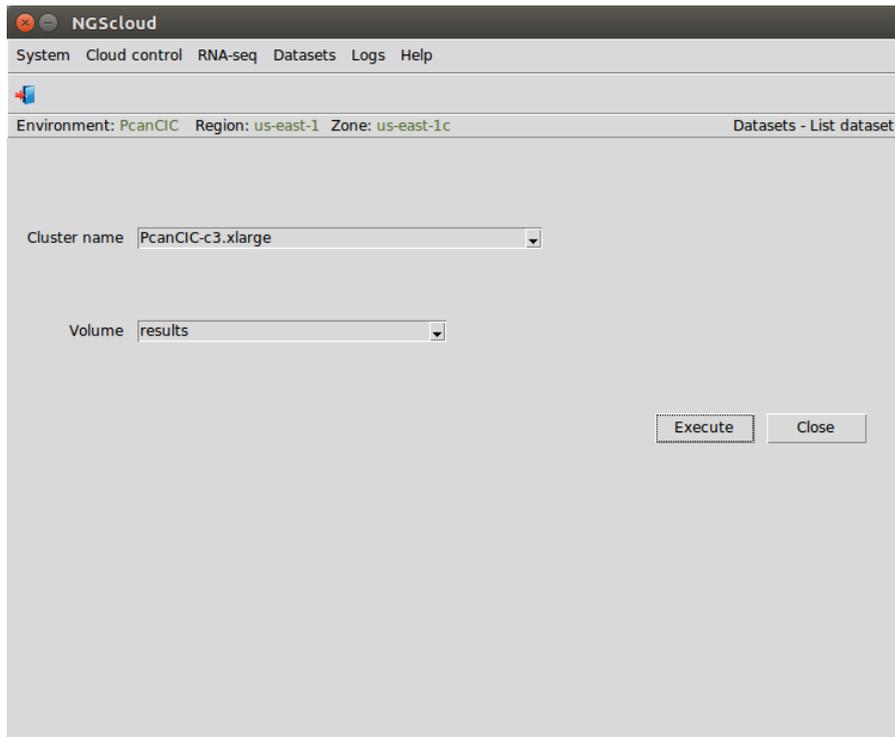
#####
Script started in node master of cluster PcanCIC-c3.xlarge at 2017-11-14 12:33:57 UTC.
#####
Running the transcriptome blastx process ...
Creating the node transcript files ...
There are 68297 transcripts in the transcriptome file.
The transcripts files are created.
Creating the control files ...
The control files are created.
Building the blastx process and watcher scripts ...
The blastx process and watcher scripts are built.
Submitting the blastx process scripts ...
Your job 2 ("blastx-00-process.sh") has been submitted
Your job 3 ("blastx-01-process.sh") has been submitted
Your job 4 ("blastx-02-process.sh") has been submitted
Your job 5 ("blastx-03-process.sh") has been submitted
The blastx process scripts are submitted.
Verifying blastx processes ...
2017-11-14 12:34:07 UTC ... STARTING: 04 - RUNNING: 00 - OK: 00 - WRONG: 00 - OTHERS: 00
2017-11-14 12:35:07 UTC ... STARTING: 00 - RUNNING: 04 - OK: 00 - WRONG: 00 - OTHERS: 00
2017-11-14 12:36:07 UTC ... STARTING: 00 - RUNNING: 04 - OK: 00 - WRONG: 00 - OTHERS: 00
2017-11-14 12:37:07 UTC ... STARTING: 00 - RUNNING: 04 - OK: 00 - WRONG: 00 - OTHERS: 00
2017-11-14 12:38:07 UTC ... STARTING: 00 - RUNNING: 04 - OK: 00 - WRONG: 00 - OTHERS: 00
2017-11-14 12:39:08 UTC ... STARTING: 00 - RUNNING: 04 - OK: 00 - WRONG: 00 - OTHERS: 00
2017-11-14 12:40:08 UTC ... STARTING: 00 - RUNNING: 04 - OK: 00 - WRONG: 00 - OTHERS: 00
2017-11-14 12:41:08 UTC ... STARTING: 00 - RUNNING: 04 - OK: 00 - WRONG: 00 - OTHERS: 00
2017-11-14 12:42:08 UTC ... STARTING: 00 - RUNNING: 04 - OK: 00 - WRONG: 00 - OTHERS: 00
2017-11-14 12:43:08 UTC ... STARTING: 00 - RUNNING: 04 - OK: 00 - WRONG: 00 - OTHERS: 00
2017-11-14 12:44:08 UTC ... STARTING: 00 - RUNNING: 04 - OK: 00 - WRONG: 00 - OTHERS: 00
2017-11-14 12:45:08 UTC ... STARTING: 00 - RUNNING: 04 - OK: 00 - WRONG: 00 - OTHERS: 00
2017-11-14 12:46:08 UTC ... STARTING: 00 - RUNNING: 04 - OK: 00 - WRONG: 00 - OTHERS: 00
2017-11-14 12:47:08 UTC ... STARTING: 00 - RUNNING: 04 - OK: 00 - WRONG: 00 - OTHERS: 00
2017-11-14 12:48:08 UTC ... STARTING: 00 - RUNNING: 04 - OK: 00 - WRONG: 00 - OTHERS: 00
2017-11-14 12:49:08 UTC ... STARTING: 00 - RUNNING: 04 - OK: 00 - WRONG: 00 - OTHERS: 00
2017-11-14 12:50:08 UTC ... STARTING: 00 - RUNNING: 04 - OK: 00 - WRONG: 00 - OTHERS: 00
2017-11-14 12:51:08 UTC ... STARTING: 00 - RUNNING: 04 - OK: 00 - WRONG: 00 - OTHERS: 00
2017-11-14 12:52:08 UTC ... STARTING: 00 - RUNNING: 04 - OK: 00 - WRONG: 00 - OTHERS: 00
2017-11-14 12:53:08 UTC ... STARTING: 00 - RUNNING: 04 - OK: 00 - WRONG: 00 - OTHERS: 00
2017-11-14 12:54:08 UTC ... STARTING: 00 - RUNNING: 04 - OK: 00 - WRONG: 00 - OTHERS: 00
2017-11-14 12:55:08 UTC ... STARTING: 00 - RUNNING: 04 - OK: 00 - WRONG: 00 - OTHERS: 00
2017-11-14 12:56:08 UTC ... STARTING: 00 - RUNNING: 04 - OK: 00 - WRONG: 00 - OTHERS: 00
2017-11-14 12:57:08 UTC ... STARTING: 00 - RUNNING: 04 - OK: 00 - WRONG: 00 - OTHERS: 00
2017-11-14 12:58:08 UTC ... STARTING: 00 - RUNNING: 04 - OK: 00 - WRONG: 00 - OTHERS: 00

```

Now, we are going to list the annotation generated by transcriptome-blastx. We select the menu item with this path:

Main menu > Dataset > List dataset

In the raised window, we select **PcanCIC-c3.xlarge** in the *Cluster name* combo-box and **results** in the *Volume* combo-box. Then we press the *Execute* button:

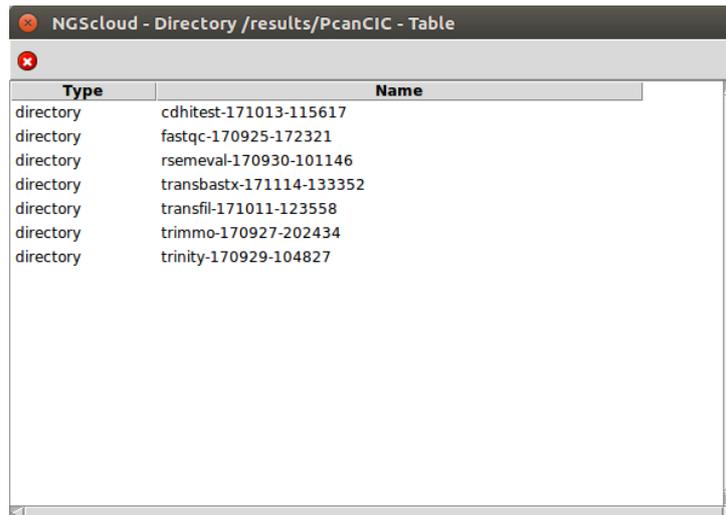


Now, we click in the **PcanCIC** row:

The screenshot shows a window titled 'NGScloud - Directory /results - Table'. It contains a table with two columns: 'Type' and 'Name'. The table lists five entries, all of which are 'directory' type.

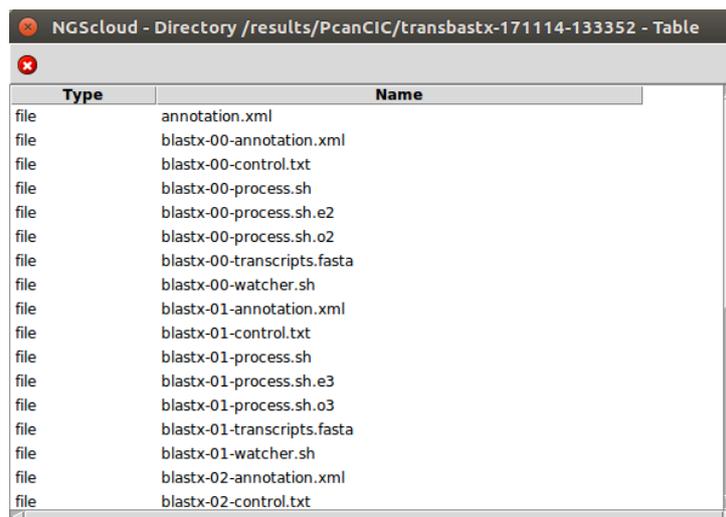
Type	Name
directory	Athaliana01x
directory	PcanCIC
directory	database
directory	reference
directory	test

Next, a window with the result datasets of the experiment PcanCIC is shown:



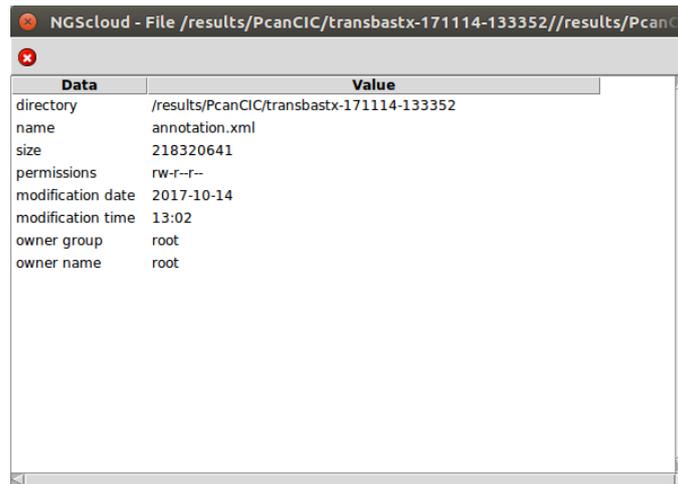
Type	Name
directory	cdhitest-171013-115617
directory	fastqc-170925-172321
directory	rsemeval-170930-101146
directory	transbastx-171114-133352
directory	transfil-171011-123558
directory	trimmo-170927-202434
directory	trinity-170929-104827

We click on the **transbastx-171114-133352** row, and another window appears with the content of the files corresponding to this assembly.



Type	Name
file	annotation.xml
file	blastx-00-annotation.xml
file	blastx-00-control.txt
file	blastx-00-process.sh
file	blastx-00-process.sh.e2
file	blastx-00-process.sh.o2
file	blastx-00-transcripts.fasta
file	blastx-00-watcher.sh
file	blastx-01-annotation.xml
file	blastx-01-control.txt
file	blastx-01-process.sh
file	blastx-01-process.sh.e3
file	blastx-01-process.sh.o3
file	blastx-01-transcripts.fasta
file	blastx-01-watcher.sh
file	blastx-02-annotation.xml
file	blastx-02-control.txt

The file **annotation.xml** is the one corresponding to the complete annotation, after concatenating the annotation files of all nodes. To observe its characteristics, we click on it:



Data	Value
directory	/results/PcanCIC/transbastx-171114-133352
name	annotation.xml
size	218320641
permissions	rw-r--
modification date	2017-10-14
modification time	13:02
owner group	root
owner name	root

Terminate the cluster with c3.xlarge template and create another cluster with t2.micro template

Now we are going to terminate the PcanCIC-c3.xlarge and to create a cluster with a t2.micro template, because it is not necessary to use an instance with many CPUs and large RAM memory in order to download them to our local computer.

First, we select the menu item with this path:

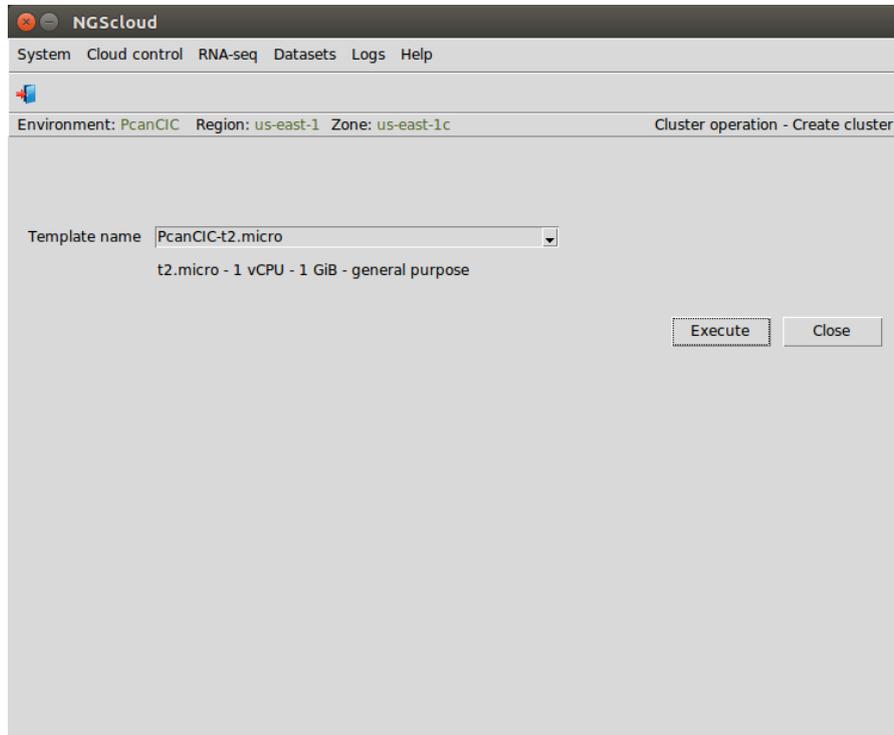
Main menu > Cloud control > Cluster operation > Terminate cluster

In the raised window, we select **PcanCIC-c3.xlarge** in the *Cluster name* combo-box; and then we press the *Execute* button:

Now we create a cluster with a t2.micro template. We select the menu item with this path:

Main menu > Cloud control > Cluster operation > Create cluster

In the raised window, we select **PcanCIC-t2.micro**, the template corresponding to a t2.micro instance type, in *Template name* combo-box; and then we press the *Execute* button:



A window is raised displaying the run log:

```

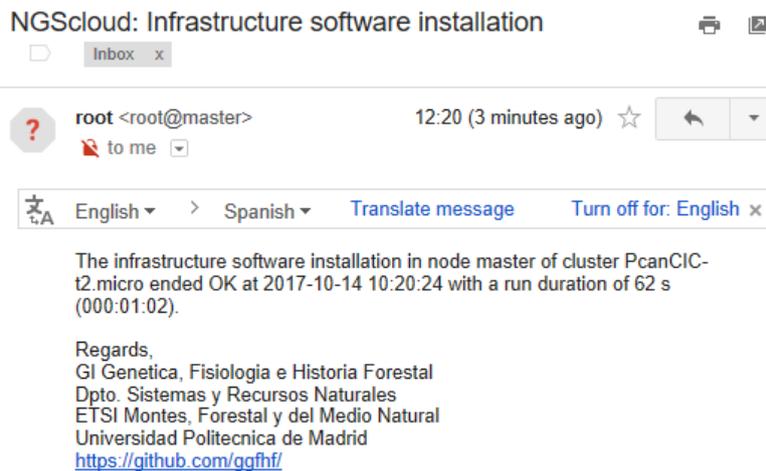
NGScloud - Cluster operation - Create cluster - Log
This process might take several minutes. Do not close this window, please wait!
*****
Verifying process requirements ...
Process requirements are OK.
*****
Creating the cluster PcanCIC-t2.micro using StarCluster ...

StarCluster - (http://star.mit.edu/cluster) (v. 0.95.6)
Software Tools for Academics and Researchers (STAR)
Please submit bug reports to starcluster@mit.edu

>>> Validating cluster template settings...
>>> Cluster template settings are valid
>>> Starting cluster...
>>> Launching a 1-node cluster...
>>> Creating security group @sc-PcanCIC-t2.micro...
Reservation:r-0fb43e9cea93167a5
>>> Waiting for instances to propagate...
0/1 | | 0%
1/1 | | 100%
>>> Waiting for cluster to come up... (updating every 30s)
>>> Waiting for all nodes to be in a 'running' state...
0/1 | | 0%
0/1 | | 0%
1/1 | | 100%
>>> Waiting for SSH to come up on all nodes...
0/1 | | 0%
0/1 | | 0%
0/1 | | 0%
0/1 | | 0%
0/1 | | 0%
0/1 | | 0%
0/1 | | 0%
0/1 | | 0%
0/1 | | 0%
0/1 | | 0%
0/1 | | 0%
1/1 | | 100%
>>> Waiting for cluster to come up took 1.106 mins
>>> The master node is ec2-54-152-166-158.compute-1.amazonaws.com
>>> Configuring cluster...
>>> Attaching volume vol-00eee59f776a0fd83 to master node on /dev/sdz ...
>>> Attaching volume vol-06calb476b0410hf1 to master node on /dev/sdy

```

When the cluster is started, infrastructure software will be installed. At the end of the installation, an email is sent, informing of its completion:



Download the transcriptome, evaluation and annotation files

In this step, we are going to download the transcriptome generated by Trinity and the filtered and clustered transcriptome, the complete annotation file, and the result files yielded by RSEM-EVAL. Due to the size of the transcriptome file, we are going to compress it previously.

To compress the assembly file, we first create the configuration file by selecting the menu item with the following path:

Main menu > Datasets > Result dataset file compression/decompression > Recreate config file

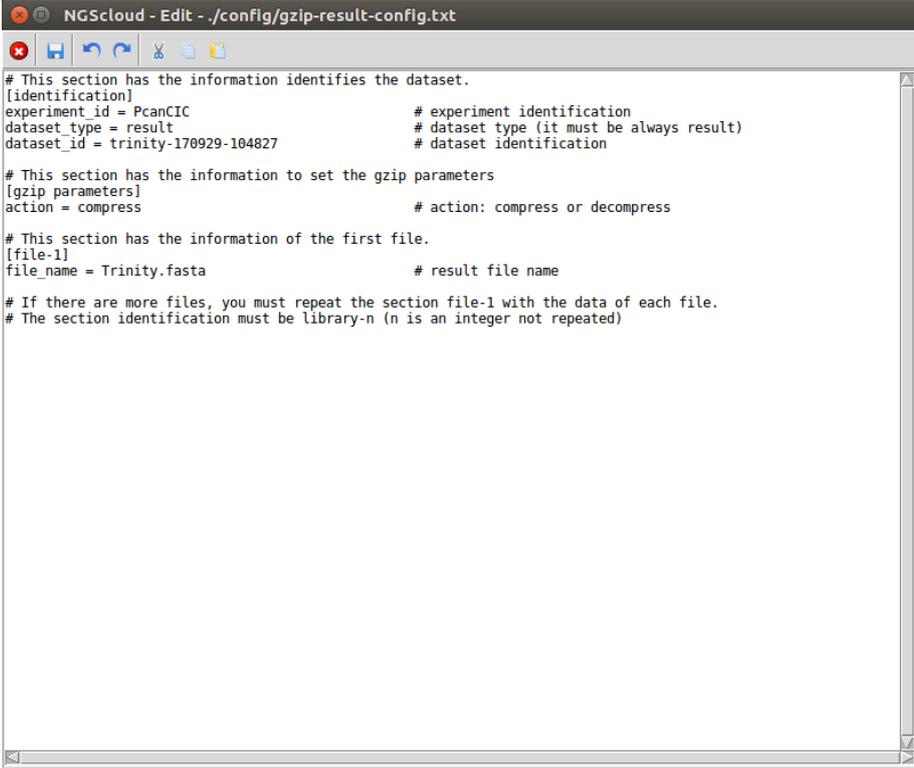
In the raised window, we select **PcanCIC-t2.micro** in the *Cluster name* combo-box, **PcanCIC** in the *Experiment id* combo-box, **uncompressed** in the *Status* combo-box (because the dataset has not been previously compressed), **trinity-170929-104827** in the *Result dataset* combo-box, and we type **Trinity.fasta.gz** as the pattern to select the files in the *File pattern* textbox and the local directory where the files will be downloaded in the *Local directory* textbox (or we select it using the button close to the textbox). Then we press the *Execute* button:

The screenshot shows a web browser window titled "NGScloud" with a menu bar containing "System", "Cloud control", "RNA-seq", "Datasets", "Logs", and "Help". Below the menu bar, the environment information is displayed: "Environment: PcanCIC Region: us-east-1 Zone: us-east-1c Result dataset file compression/decompression - Recreate conf". The main content area contains several form fields:

- Cluster name:** A dropdown menu with "PcanCIC-t2.micro" selected.
- Whole dataset?:** A dropdown menu with "selected files" selected.
- Action:** A dropdown menu with "compress" selected.
- Experiment id:** A dropdown menu with "PcanCIC" selected.
- Result dataset:** A dropdown menu with "Trinity (170929 104827)" selected.
- File pattern:** A text input field containing "Trinity.fasta". Below it, a note reads "It is a pattern of regular expression."

At the bottom right of the form, there are two buttons: "Execute" and "Close".

In the next window, we can examine the config file. In this example, it has three sections: *identification*, with the dataset type and the experiment and dataset identifications, the action to do; and *file-1* with the file name of the Trinity assembly (in this example, we have only selected this file):



```

# This section has the information identifies the dataset.
[identification]
experiment_id = PcanCIC                # experiment identification
dataset_type = result                  # dataset type (it must be always result)
dataset_id = trinity-170929-104827    # dataset identification

# This section has the information to set the gzip parameters
[gzip parameters]
action = compress                      # action: compress or decompress

# This section has the information of the first file.
[file-1]
file_name = Trinity.fasta              # result file name

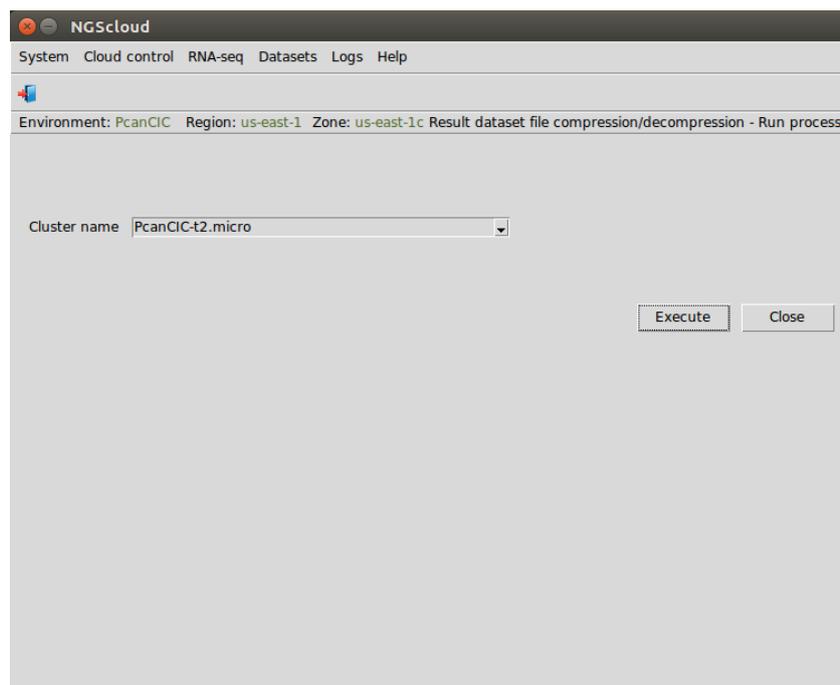
# If there are more files, you must repeat the section file-1 with the data of each file.
# The section identification must be library-n (n is an integer not repeated)

```

To compress the assembly file, we select the menu item with this path:

Main menu > Dataset > Result dataset compression/decompression > Run compression/decompression process

In the raised window, we select **PcanCIC-t2.micro** in the *Cluster name* combo-box; and then we press the *Execute* button:



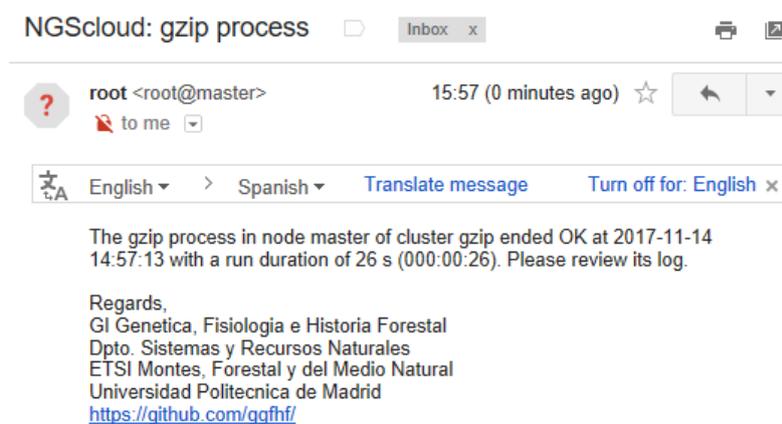
A window is raised with the submission log:

```

NGScloud - Result dataset file compression/decompression - Run process - Log
The config file is OK.
*****
Connecting the SSH client ...
The SSH client is connected.
*****
Connecting the SSH transport ...
The SSH transport is connected.
*****
Connecting the SFTP client ...
The SFTP client is connected.
*****
Verifying process requirements ...
Process requirements are OK.
*****
Determining the run directory in the cluster ...
The directory path is /results/PcanCIC/gzip-171114-155633.
*****
Building the process script ./temp/gzip-result-process.sh ...
The file is built.
*****
Uploading the process script ./temp/gzip-result-process.sh in the directory /results/PcanCIC/g
The file is uploaded.
*****
Setting on the run permission of /results/PcanCIC/gzip-171114-155633/gzip-result-process.sh ...
The run permission is set.
*****
Building the process starter ./temp/gzip-result-starter.sh ...
The file is built.
*****
Uploading the process starter ./temp/gzip-result-starter.sh in the directory /results/PcanCIC/g
The file is uploaded.
*****
Setting on the run permission of /results/PcanCIC/gzip-171114-155633/gzip-result-starter.sh ...
The run permission is set.
*****
Submitting the process script /results/PcanCIC/gzip-171114-155633/gzip-result-starter.sh ...
Your job 1 ("gzip-result-starter.sh") has been submitted
*****
Closing the SSH transport connection ...
The connection is closed.
*****
Closing the SSH client connection ...
The connection is closed.
*****
You can close this window now.

```

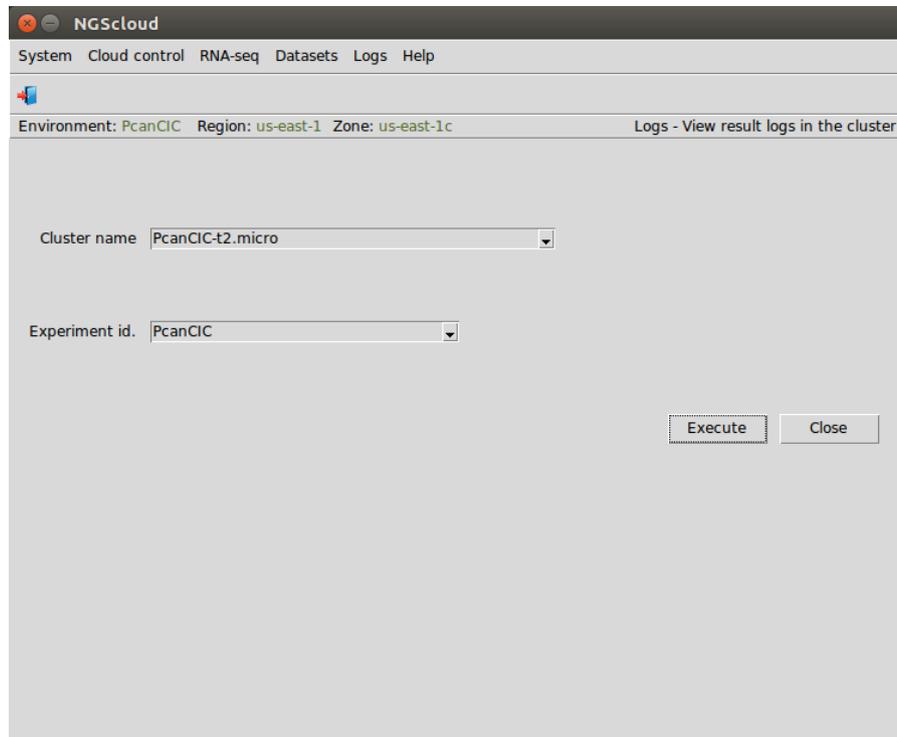
At the end of the run, an email is sent, informing of its completion:



We can view the process log during and after its run. To do so, we select the menu item with this path:

Main menu > Logs > View result logs in the cluster

In the raised window, we select **PcanCIC-t2.micro** in the *Cluster name* combo-box and **PcanCIC** in the *Experiment id* combo-box; and then we press the *Execute* button:



A window with the result datasets for each run of the bioinformatic programs that correspond to the experiment PcanCIC is shown:

The screenshot shows a table window titled 'NGScloud - Experiment runs in /results/PcanCIC - Table'. The table has five columns: 'Experiment id', 'Result dataset', 'Bioinfo app / Utility', 'Date', and 'Time'. The data rows are as follows:

Experiment id	Result dataset	Bioinfo app / Utility	Date	Time
PcanCIC	cdhitest-171013-115617	CD-HIT-EST	2017-10-13	11:56:17
PcanCIC	fastqc-170925-172321	FastQC	2017-09-25	17:23:21
PcanCIC	gzip-171114-155633	gzip	2017-11-14	15:56:33
PcanCIC	rsemeval-170930-101146	RSEM-EVAL	2017-09-30	10:11:46
PcanCIC	transbastx-171114-133352	transcriptome-blastx	2017-11-14	13:33:52
PcanCIC	transfil-171011-123558	transcript-filter	2017-10-11	12:35:58
PcanCIC	trimmo-170927-202434	Trimmomatic	2017-09-27	20:24:34
PcanCIC	trinity-170929-104827	Trinity	2017-09-29	10:48:27

Five result datasets are shown. We click on the **gzip-171114-155633** row, the dataset generated by gzip, the compression program, and another window appears with its corresponding log:

```

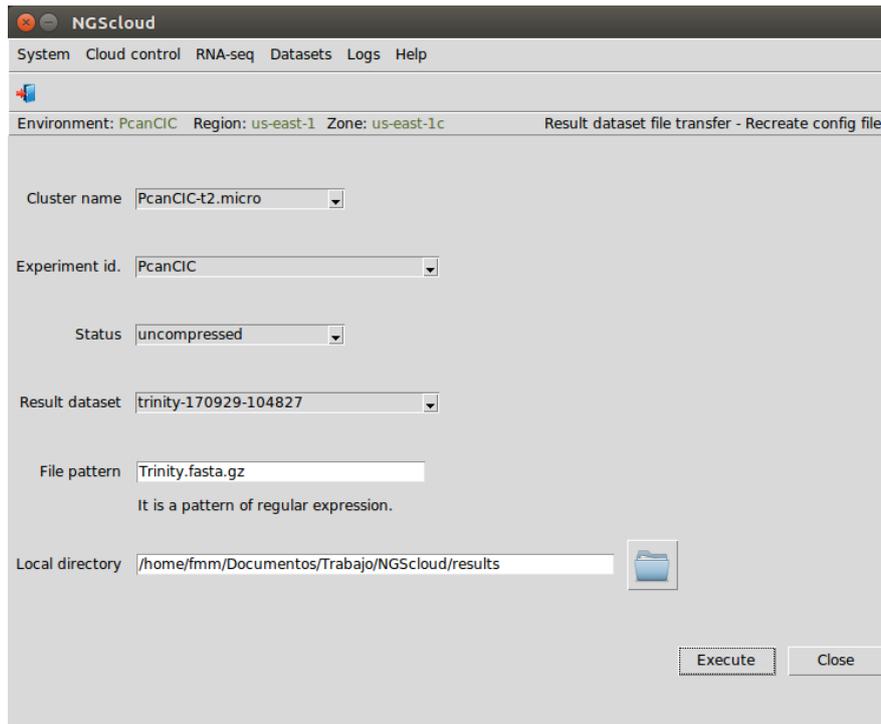
#####
Script started in node master of cluster PcanCIC-t2.micro at 2017-11-14 14:56:47 UTC.
#####
Compressing/decompressing /results/PcanCIC/trinity-170929-104827/Trinity.fasta ...
Elapsed real time (s): 26.35
CPU time in kernel mode (s): 0.16
CPU time in user mode (s): 26.13
Percentage of CPU: 99%
Maximum resident set size(Kb): 1016
Average total memory use (Kb):0
#####
Script ended OK at 2017-11-14 14:57:13 UTC with a run duration of 26 s (000:00:26).
#####

```

Next, we are going to download the compressed assembly file. To do so, we first create the configuration file by selecting the menu item with the following path:

Main menu > Datasets > Result dataset file transfer > Recreate config file

In the raised window, we select **PcanCIC-t2.micro** in the *Cluster name* combo-box, **PcanCIC** in the *Experiment id* combo-box, **uncompressed** in the *Status* combo-box (because the dataset has not been previously compressed), **trinity-170929-104827** in the *Result dataset* combo-box, and we type **Trinity.fasta.gz** as the pattern to select the files in the *File pattern* textbox and the local directory where the files will be downloaded in the *Local directory* textbox (or we select it using the button close to the textbox). Then we press the *Execute* button:



In the next window, we can examine the config file. In this example, it has two sections: *identification*, with the experiment and result dataset identifications, the status of the dataset and the local path where files will be download; and *file-1* with the file name of the compressed Trinity assembly (in this example, we have only selected this file):

```

NGScloud - Edit - ./config/result-transfer-config.txt
# You must review the information of this file and update the values with the corresponding ones to the current
# The files will be copied from the cluster directory /results/experiment_id/result_dataset_id.

# This section has the information identifies the experiment.
[identification]
experiment_id = PcanCIC # experiment identification
result_dataset_id = trinity-170929-104827 # run identification
status = uncompressed # result dataset status (it must be always uncompressed)
local_dir = /home/fmm/Documentos/Trabajo/NGScloud/results # local path where the file will be download

# This section has the information of the first result file.
[file-1]
file_name = Trinity.fasta.gz # result file name

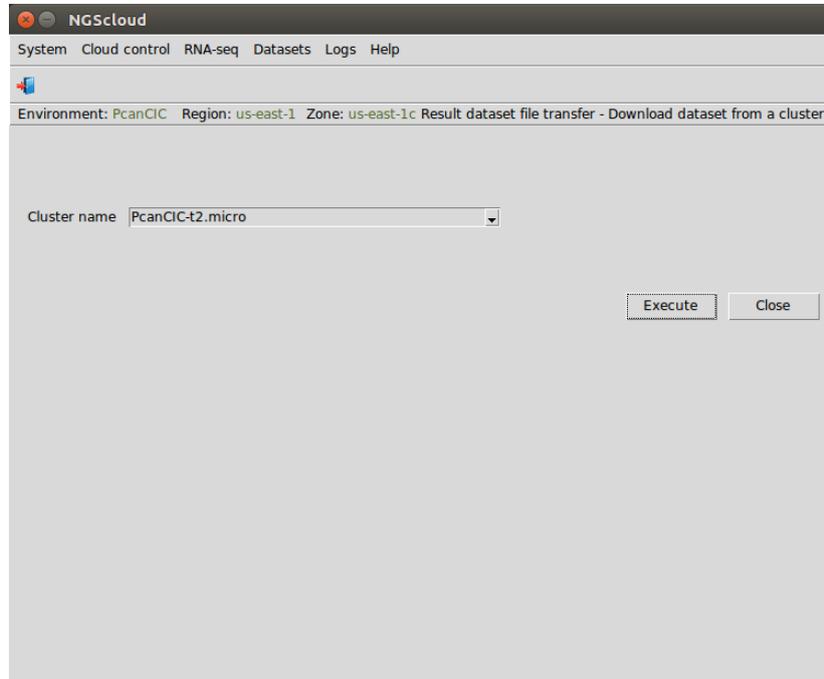
# If there are more files, you must repeat the section file-1 with the data of each file.
# The section identification must be library-n (n is an integer not repeated)

```

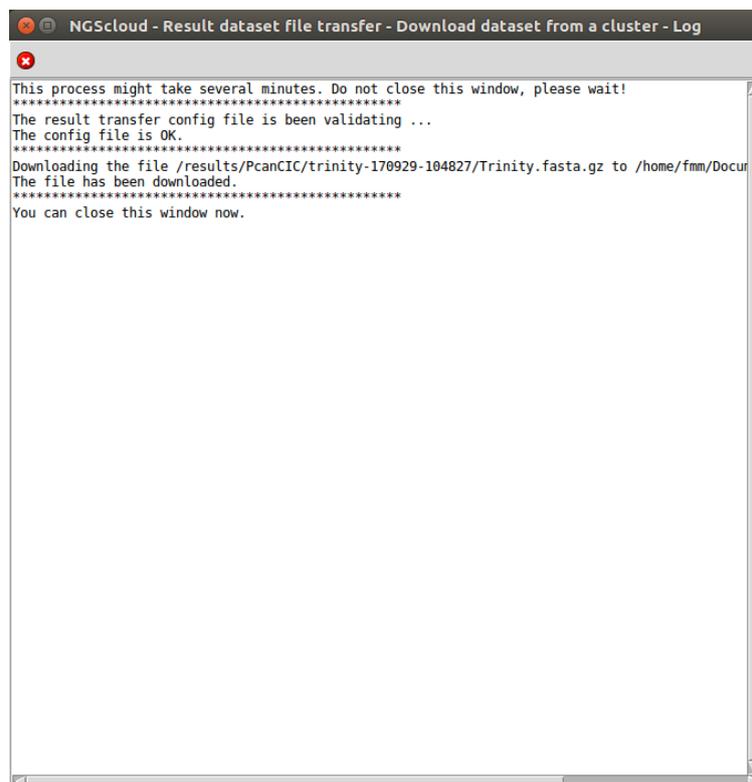
To download the compressed Trinity assembly from the cluster, we select the menu item with this path:

Main menu > Dataset > Result dataset file transfer > Download dataset from a cluster

In the raised window, we select **PcanCIC-t2.micro** in the *Cluster name* combo-box; and then we press the *Execute* button:



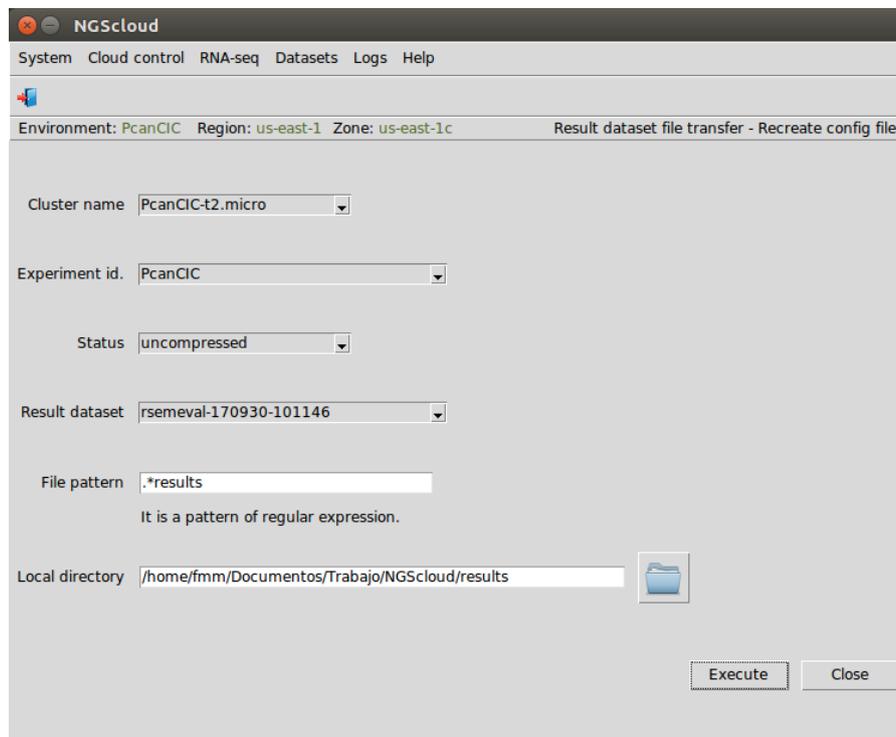
A window is raised with the log corresponding to the download:



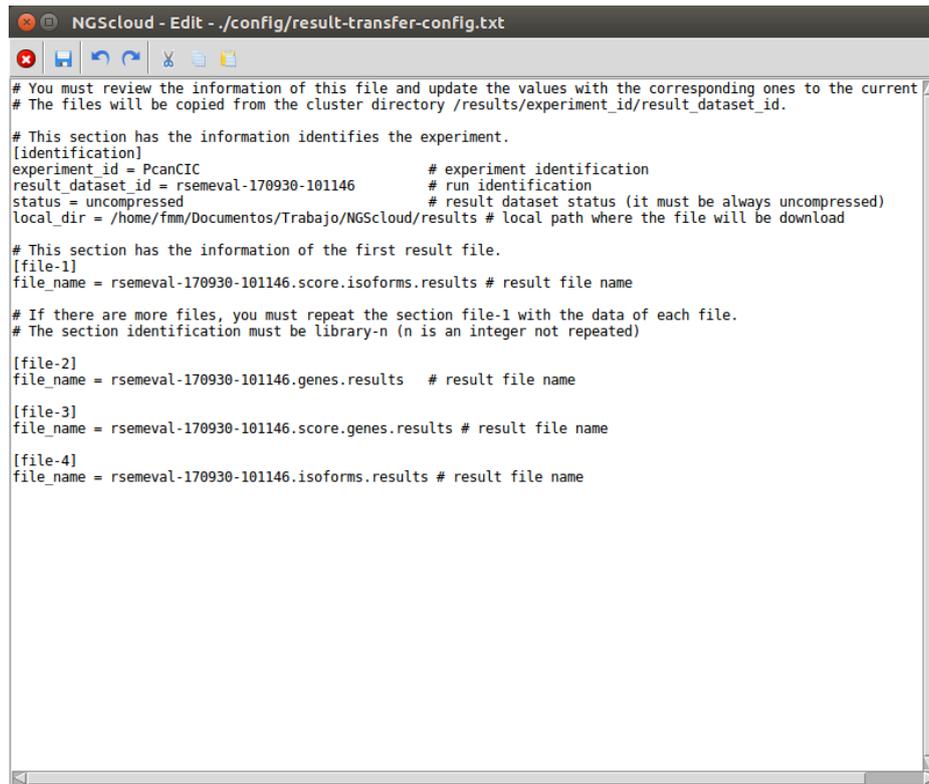
Next, we are going to download the assembly assessment files generated by RSEM-EVAL. We have to download the ".results" files generated by this program. To do so, we first create the configuration file by selecting the menu item with the following path:

Main menu > Datasets > Result dataset file transfer > Recreate config file

In the raised window, we select **PcanCIC-t2.micro** in the *Cluster name* combo-box, **PcanCIC** in the *Experiment id* combo-box, **uncompressed** in the *Status* combo-box (because the dataset has not been previously compressed), **rsemeval-170930-101146** in the *Result dataset* combo-box, and we type **.*results** as the pattern to select the files in the *File pattern* textbox and the local directory where the files will be downloaded in the *Local directory* textbox (or we select it using the button close to the textbox). Then we press the *Execute* button:



In the next window, we can examine the config file. In this example, it has five sections: *identification*, with the experiment and result dataset identifications, the status of the dataset and the local path where files will be downloaded; and *file-1* to *file-4* with the file names of four result files:



```

# You must review the information of this file and update the values with the corresponding ones to the current
# The files will be copied from the cluster directory /results/experiment_id/result_dataset_id.

# This section has the information identifies the experiment.
[identification]
experiment_id = PcanCIC # experiment identification
result_dataset_id = rsemeval-170930-101146 # run identification
status = uncompressed # result dataset status (it must be always uncompressed)
local_dir = /home/fmm/Documentos/Trabajo/NGScloud/results # local path where the file will be download

# This section has the information of the first result file.
[file-1]
file_name = rsemeval-170930-101146.score.isoforms.results # result file name

# If there are more files, you must repeat the section file-1 with the data of each file.
# The section identification must be library-n (n is an integer not repeated)
[file-2]
file_name = rsemeval-170930-101146.genes.results # result file name

[file-3]
file_name = rsemeval-170930-101146.score.genes.results # result file name

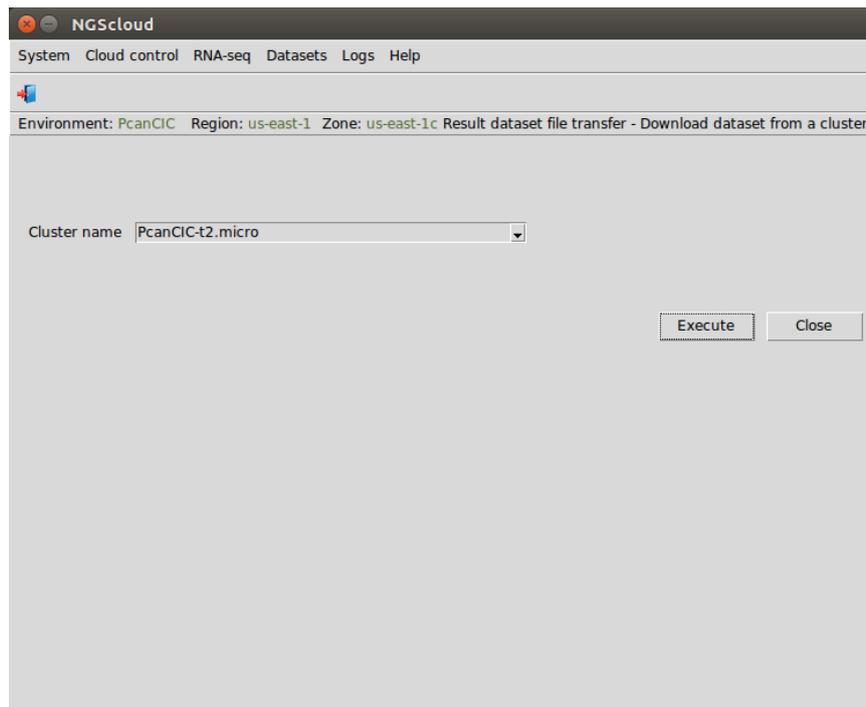
[file-4]
file_name = rsemeval-170930-101146.isoforms.results # result file name

```

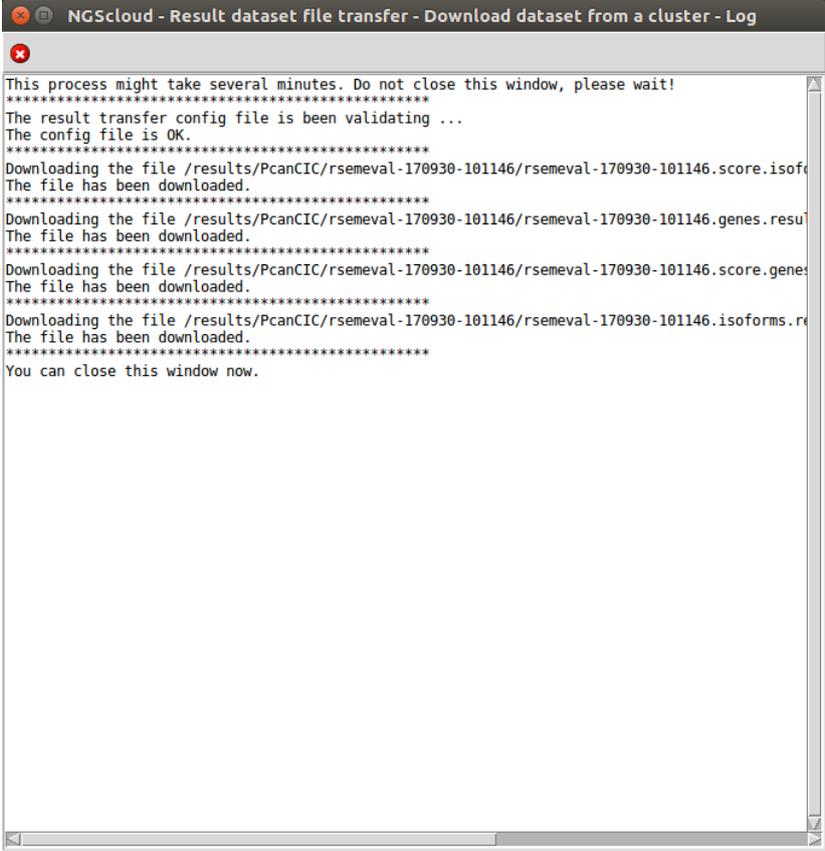
To download the result files from the cluster, we select the menu item with this path:

Main menu > Dataset > Result dataset file transfer > Download dataset from a cluster

In the raised window, we select **PcanCIC-t2.micro** in the *Cluster name* combo-box; and then we press the *Execute* button:



A window is raised with the log corresponding to the download:



```

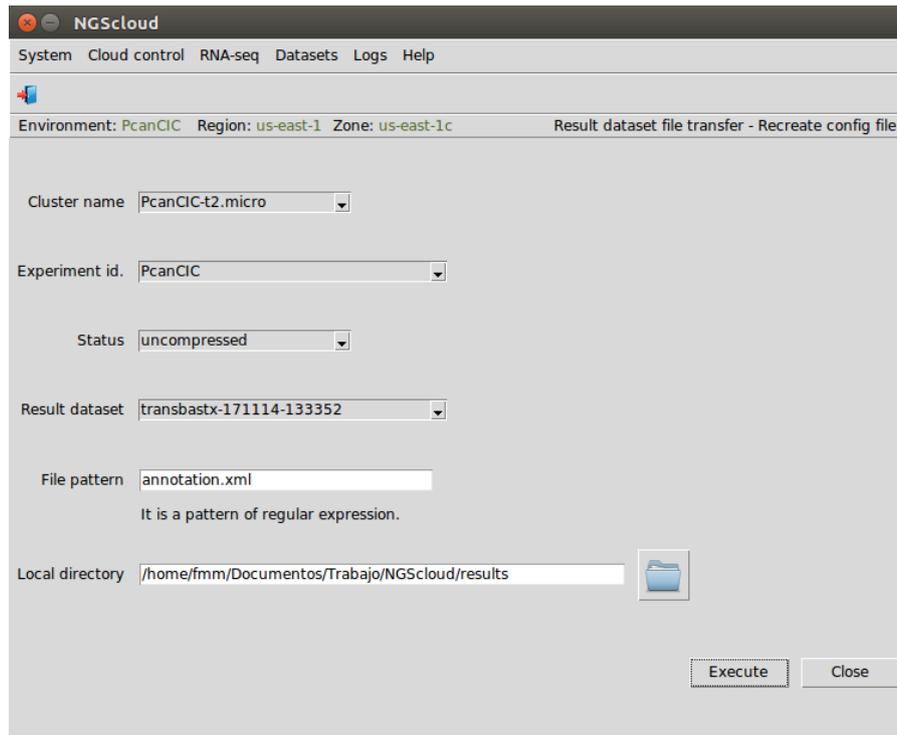
NGScloud - Result dataset file transfer - Download dataset from a cluster - Log
This process might take several minutes. Do not close this window, please wait!
*****
The result transfer config file is been validating ...
The config file is OK.
*****
Downloading the file /results/PcanCIC/rsemeval-170930-101146/rsemeval-170930-101146.score.isofo
The file has been downloaded.
*****
Downloading the file /results/PcanCIC/rsemeval-170930-101146/rsemeval-170930-101146.genes.resu
The file has been downloaded.
*****
Downloading the file /results/PcanCIC/rsemeval-170930-101146/rsemeval-170930-101146.score.gene
The file has been downloaded.
*****
Downloading the file /results/PcanCIC/rsemeval-170930-101146/rsemeval-170930-101146.isoforms.r
The file has been downloaded.
*****
You can close this window now.

```

And finally, we are going to download the annotation file generated by transcriptome-blastx. We have to download the file `annotation.xml` generated by this program. To do so, we first create the configuration file by selecting the menu item with the following path:

Main menu > Datasets > Result dataset file transfer > Recreate config file

In the raised window, we select **PcanCIC-t2.micro** in the *Cluster name* combo-box, **PcanCIC** in the *Experiment id* combo-box, **uncompressed** in the *Status* combo-box (because the dataset has not been previously compressed), **transbastx-171114-133353** in the *Result dataset* combo-box, and we type **annotation.xml** as the pattern to select the files in the *File pattern* textbox and the local directory where the files will be downloaded in the *Local directory* textbox (or we select it using the button close to the textbox). Then we press the *Execute* button:



In the next window, we can examine the config file. In this example, it has two sections: *identification*, with the experiment and result dataset identifications, the status of the dataset and the local path where files will be downloaded; and *file-1* with the file name of the annotation file:

```

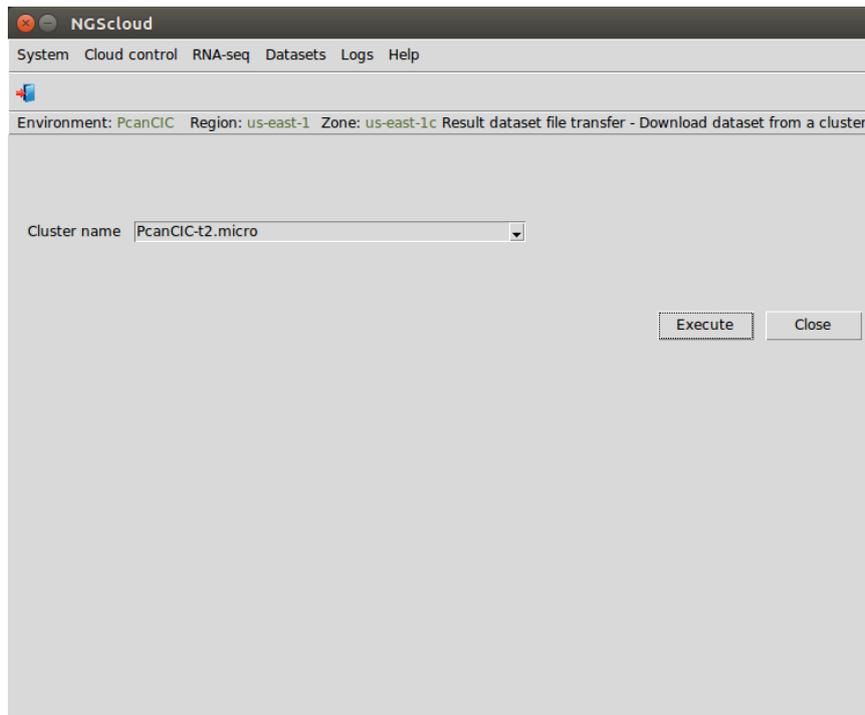
NGScloud - Edit - ./config/result-transfer-config.txt
# You must review the information of this file and update the values with the corresponding ones to the current
# The files will be copied from the cluster directory /results/experiment_id/result_dataset_id.
# This section has the information identifies the experiment.
[identification]
experiment_id = PcanCIC # experiment identification
result_dataset_id = transbastx-171114-133352 # run identification
status = uncompressed # result dataset status (it must be always uncompressed)
local_dir = /home/fmm/Documentos/Trabajo/NGScloud/results # local path where the file will be download
# This section has the information of the first result file.
[file-1]
file_name = annotation.xml # result file name
# If there are more files, you must repeat the section file-1 with the data of each file.
# The section identification must be library-n (n is an integer not repeated)

```

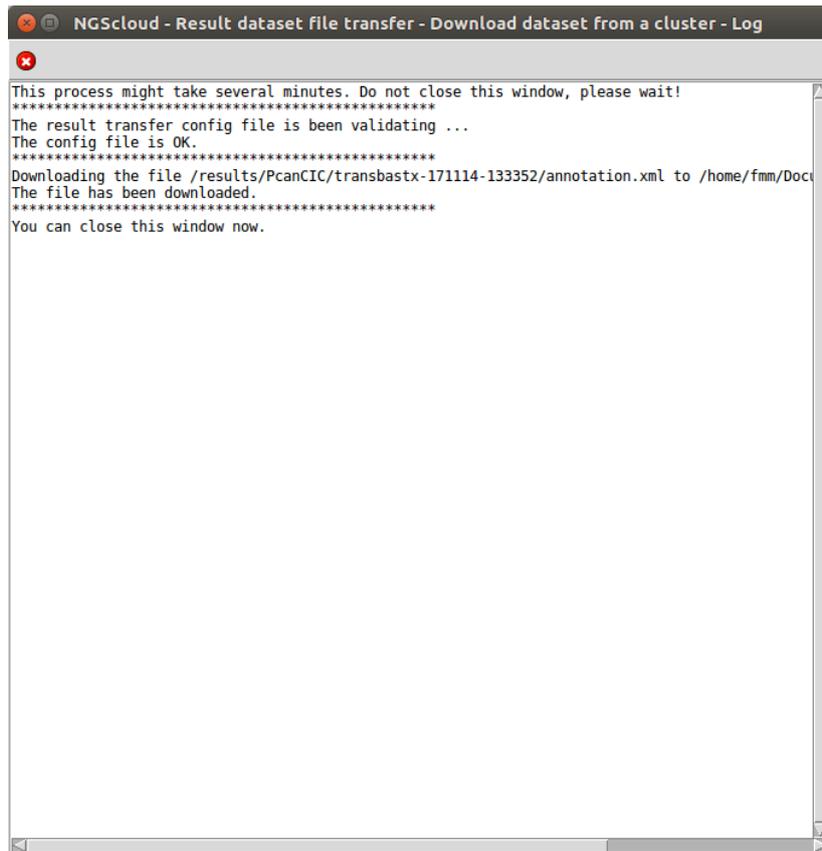
To download the result files from the cluster, we select the menu item with this path:

Main menu > Dataset > Result dataset file transfer > Download dataset from a cluster

In the raised window, we select **PcanCIC-t2.micro** in the *Cluster name* combo-box; and then we press the *Execute* button:



A window is raised with the log corresponding to the download:

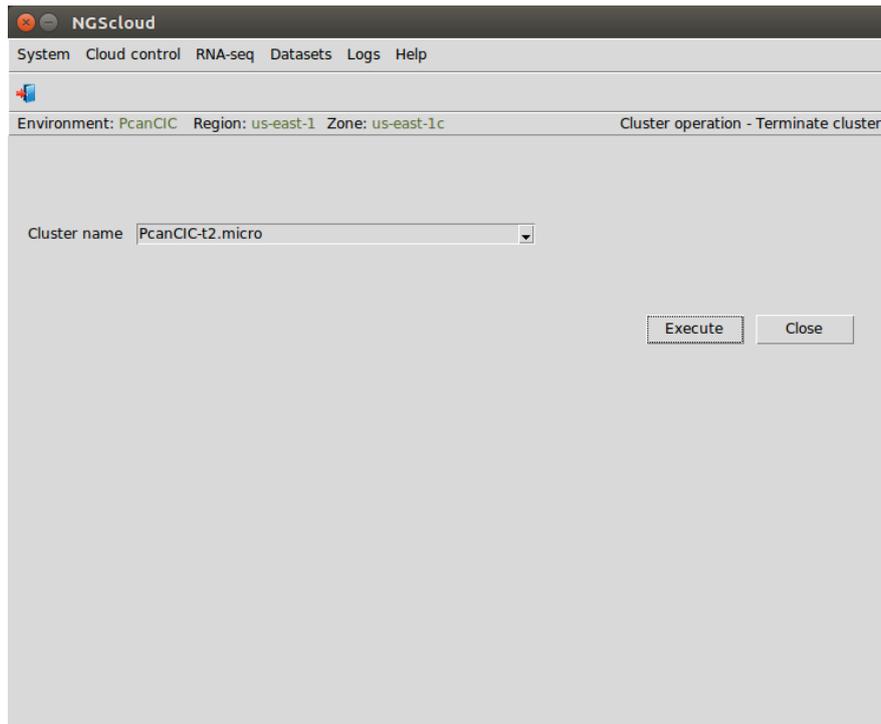


Terminate the cluster with the t2.micro template

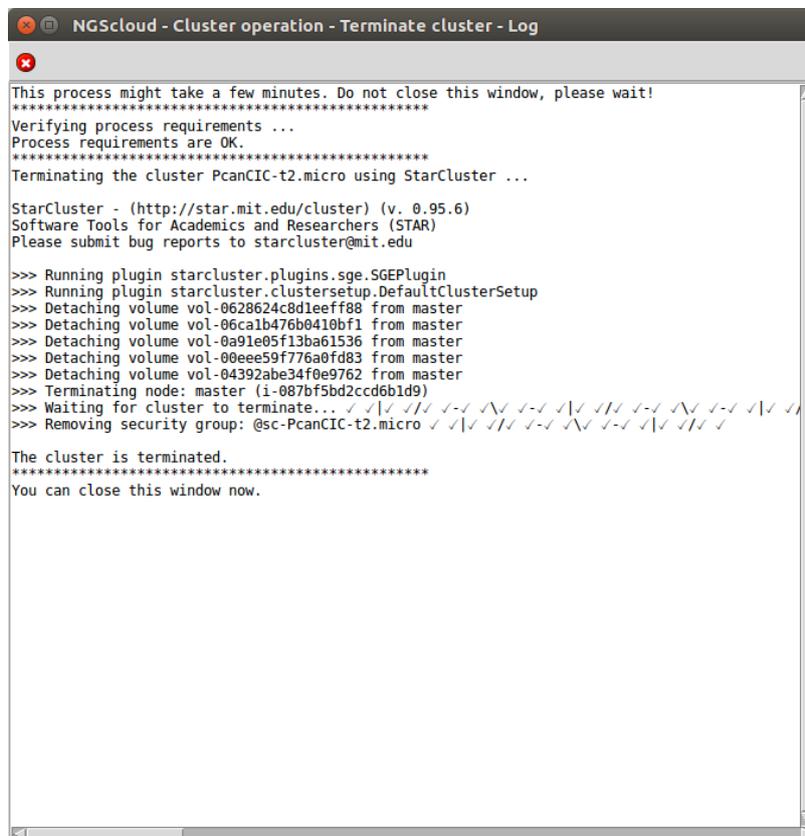
Once the analysis is complete, we terminate the cluster PcanCIC-t2.micro. To do so, we select the menu item with this path:

Main menu > Cloud control > Cluster operation > Terminate cluster

In the raised window, we select **PcanCIC-t2.micro** in the *Cluster name* combo-box; and then we press the *Execute* button:



A window is raised displaying the run log:



How-to

Below is the menu path of the tasks the most common tasks:

[How to display this manual](#)

Main menu > Help > View help ... or pressing F1 key

[How to recreate the NGScloud config file](#)

Main menu > Cloud control > Configuration > Recreate NGScloud config file

[How to create a new environment](#)

Main menu > Cloud control > Set environment

[How to change to another environment](#)

Main menu > Cloud control > Set environment

[How to view characteristics of a cluster template](#)

Main menu > Cloud control > Configuration > List cluster templates

[How to create a cluster](#)

Main menu > Cloud control > Cluster operation > Create cluster

[How to terminate a cluster](#)

Main menu > Cloud control > Cluster operation > Terminate cluster

[How to list the running clusters](#)

Main menu > Cloud control > Cluster operation > List clusters

[How to create a volume](#)

Main menu > Cloud control > Volume Operation > Create volume

[How to remove a volume](#)

Main menu > Cloud control > Volume Operation > Remove volume

[How to list the created volumes](#)

Main menu > Cloud control > Volume Operation > List volumes

[How to link a volume in cluster templates](#)

Main menu > Cloud control > Configuration > Link volume in a cluster templates

[How to add a node in a cluster](#)

Main menu > Cloud control > Node operation > Add node in a cluster

[How to remove a node in a cluster](#)

Main menu > Cloud control > Node operation > Remove node in a cluster

[How to open a terminal of a cluster](#)

Main menu > Cloud control > Open a terminal

How to set up a bioinformatic software in a cluster

Main menu > Cloud control > Bioinfo software setup > bioinformatic software to use

How to run a RNA-seq bioinformatic software in a cluster

Main menu > RNA-seq > "Task of RNA-seq workflow" > "Bioinformatic software" > Recreate config file

Main menu > RNA-seq > "Task of RNA-seq workflow" > "Bioinformatic software" > "Run process"

Task of RNA-seq workflow	Bioinformatic software
Read quality	FastQC
Trimming	Trimmomatic
Digital normalization	Insilico_read_normalization (Trinity package)
De novo assembly	SOAPdenovo-Trans
	Trinity
Reference-based assembly	STAR
Assembly quality and transcript quantification	QUAST
	rnaQUAST
	RSEM-EVAL (DETONATE package)
Transcriptome filtering	CD-HIT-EST (CD-HIT package)
	transcript-filter (NGShelper package)
Annotation	Transcriptome-blast (NGShelper package)

How to display datasets of a volume

Main menu > Dataset > List dataset

How to display the contents of a dataset

Main menu > Dataset > List dataset

How to upload reference files to a cluster

Main menu > Datasets > Reference dataset file transfer > Recreate config file

Main menu > Datasets > Reference dataset file transfer > Upload dataset to a cluster

How to compress/decompress reference files in a cluster

Main menu > Datasets > Reference dataset file compression/decompression > Recreate config file

Main menu > Datasets > Reference dataset file compression/decompression > Run compression/decompression process

How to upload database files to a cluster

Main menu > Datasets > Database dataset file transfer > Recreate config file

Main menu > Datasets > Database dataset file transfer > Upload dataset to a cluster

How to compress/decompress database files in a cluster

Main menu > Datasets > Database dataset file compression/decompression > Recreate config file

Main menu > Datasets > Database dataset file compression/decompression > Run compression/decompression process

[How to upload read files to a cluster](#)

Main menu > Datasets > Read dataset file transfer > Recreate config file

Main menu > Datasets > Read dataset file transfer > Upload dataset to a cluster

[How to compress/decompress read files in a cluster](#)

Main menu > Datasets > Read dataset file compression/decompression > Recreate config file

Main menu > Datasets > Read dataset file compression/decompression > Run compression/decompression process

[How to download results files from a cluster](#)

Main menu > Datasets > Result dataset file transfer > Recreate config file

Main menu > Datasets > Result dataset file transfer > Download dataset from a cluster

[How to compress/decompress results files in a cluster](#)

Main menu > Datasets > Result dataset file compression/decompression > Recreate config file

Main menu > Datasets > Result dataset file compression/decompression > Run compression/decompression process

[How to view submission logs in the local computer](#)

Main menu > Logs > View submission logs in the local computer

[How to view result logs in the cluster](#)

Main menu > Logs > View result logs in the cluster