
PM4NGS Documentation

Release 0.0.18.dev4+gd7bbf86

Roberto Vera Alvarez

Feb 04, 2022

Contents

1	Installation	3
1.1	Requirements	3
1.2	PM4NGS with Conda/BioConda (Recommended)	3
1.2.1	Conda installation	3
1.2.2	PM4NGS conda installation	3
1.2.3	PM4NGS conda env activation	3
1.3	PM4NGS with Python virtual environment	4
1.3.1	PM4NGS python installation	4
1.3.2	PM4NGS python env activation	4
2	PM4NGS on Ubuntu	5
2.1	Installing PM4NGS	5
2.2	Using PM4NGS	5
2.3	Running the ChIP-exo demo	6
2.4	Running the Jupyter Server	7
3	PM4NGS on Windows Subsystem for Linux	9
3.1	Activating WSL	9
3.2	Installing Docker	10
3.3	Installing Ubuntu on WSL	13
3.4	Installing PM4NGS	14
3.5	Using PM4NGS on WSL	14
3.6	Running the ChIP-exo demo	15
3.7	Running the Jupyter Server	17
4	PM4NGS on GCP instance with Ubuntu	19
4.1	Creating a GCP instance with Ubuntu 20.02 LTS	19
4.2	Installing PM4NGS on the GCP instance with Ubuntu	22
4.3	Installing PM4NGS	22
4.4	Using PM4NGS	23
4.5	Running the ChIP-exo demo	23
4.6	Running the Jupyter Server	26
4.7	Open a firewall rule for accessing the Jupyter Server	27
5	Sample sheet	45
5.1	Single-end example	45
5.2	Paired-end example	45

5.3	Processing data from the NCBI SRA	46
5.4	Sample sheet column names and description	46
6	Differential Gene expression and GO enrichment from RNA-Seq data	47
6.1	Input requirements	48
6.2	Pipeline command line	49
6.3	Creating the DGA and GO enrichment from RNA-Seq data project	49
6.4	Jupyter server	52
6.5	Data processing	53
6.6	CWL workflows	53
6.6.1	SRA download and QC workflow	53
6.6.2	Samples Trimming	55
6.6.3	STAR based alignment and sorting	55
6.6.4	RNA-Seq quantification and QC workflow using TPMCalculator	56
6.6.5	Differential Gene Expression analysis from RNA-Seq data	58
6.6.6	GO enrichment from RNA-Seq data	61
6.7	Demo	62
7	Differential Binding detection from ChIP-Seq data	65
7.1	Input requirements	66
7.2	Pipeline command line	66
7.3	Creating the Differential Binding detection from ChIP-Seq data project	67
7.4	Jupyter server	70
7.5	Data processing	70
7.6	CWL workflows	71
7.6.1	SRA download and QC workflow	71
7.6.2	Samples Trimming	72
7.6.3	BWA based alignment and quality control workflow	72
7.6.4	Peak caller workflow using MACS2	74
7.6.5	MEME Motif detection workflow	76
7.6.6	MEME databases	77
7.6.7	Differential binding analysis with Diffbind	78
7.7	Demo	79
8	Detection of binding events from ChIP-exo data	81
8.1	Input requirements	82
8.2	Pipeline command line	82
8.3	Creating the Detection of binding events from ChIP-exo data project	83
8.4	Jupyter server	86
8.5	Data processing	86
8.6	CWL workflows	87
8.6.1	SRA download and QC workflow	87
8.6.2	Samples Trimming	88
8.6.3	BWA based alignment and quality control workflow	88
8.6.4	Peak caller workflow using MACE	90
8.6.5	MEME Motif detection workflow	91
8.6.6	MEME databases	92
8.7	Demo	93
9	Transcriptome Annotation pipeline for non-model organisms	95
9.1	Annotation Workflow	95
9.2	GCP configuration	96
9.3	Input requirements	96
9.4	Pipeline command line	96
9.5	Creating the annotation project	96

9.6	Jupyter server	98
9.7	Data processing	98
9.8	Demo	98
10	Create a new PM4NGS pipeline	101
10.1	PM4NGS based pipeline folder structure	101
10.2	CWL tools and workflows specifications	103
11	Utils	105
11.1	STAR Genomic Indexes	105
11.2	BWA Genomic Indexes	106
11.3	Create BED from GTF	106
12	PM4NGS available genomes	107
13	Introduction	109
14	Features	111
15	Citation	113
16	Help and Support	115
17	Public Domain notice	117



1.1 Requirements

1. Poppler (<https://poppler.freedesktop.org/>)

1.2 PM4NGS with Conda/BioConda (Recommended)

1.2.1 Conda installation

Conda should be already installed and configured using these commands:

```
localhost:~> wget https://repo.anaconda.com/miniconda/Miniconda3-latest-Linux-x86_64.sh
↪ sh
localhost:~> sh Miniconda3-latest-Linux-x86_64.sh
localhost:~> conda config --add channels defaults
localhost:~> conda config --add channels bioconda
localhost:~> conda config --add channels conda-forge
```

1.2.2 PM4NGS conda installation

PM4NGS should be installed in a Conda virtual environment named *pm4ngs*:

```
localhost:~> conda create -n pm4ngs pm4ngs
```

1.2.3 PM4NGS conda env activation

For activating the conda env:

```
localhost:~> conda activate pm4ngs
localhost:~> pm4ngs-create -v
PM4NGS version: 0.1.dev0+d20200819
```

1.3 PM4NGS with Python virtual environment

1.3.1 PM4NGS python installation

Python 3.6 or above should be installed.

```
localhost:~> python3 -m venv pm4ngs_venv
localhost:~> source pm4ngs_venv/bin/activate
(pm4ngs_venv) localhost:~> pip install wheel
(pm4ngs_venv) localhost:~> pip install pm4ngs
```

1.3.2 PM4NGS python env activation

For activating the virtual env:

```
localhost:~> source pm4ngs_venv/bin/activate
(pm4ngs_venv) localhost:~> pm4ngs-create -v
PM4NGS version: 0.0.4
```

CHAPTER 2

PM4NGS on Ubuntu

Runs these commands on a terminal to to prepare the instance to run PM4NGS

```
veraalva@perseo:~$ sudo apt-get update
veraalva@perseo:~$ sudo apt-get install docker.io python3 python3-pip python3-venv
→python3-dev poppler-utils gcc nodejs tree
veraalva@perseo:~$ sudo usermod -aG docker $USER
```

Close and reopen the terminal to set the docker group in the user.

2.1 Installing PM4NGS

Creates a Python virtual environment named: **pm4ngs_venv** for installing PM4NGS

```
veraalva@perseo:~$ python3 -m venv pm4ngs_venv
veraalva@perseo:~$ source pm4ngs_venv/bin/activate
(pm4ngs_venv) veraalva@perseo:~$ pip install wheel
(pm4ngs_venv) veraalva@perseo:~$ pip install pm4ngs
```

2.2 Using PM4NGS

Open a terminal and activate the **pm4ngs_venv** virtual environment

```
veraalva@perseo:~$ source pm4ngs_venv/bin/activate
(pm4ngs_venv) veraalva@perseo:~$ pm4ngs-chipexo --version
PM4NGS version: 0.0.4
(pm4ngs_venv) veraalva@perseo:~$
```

2.3 Running the ChIP-exo demo

Open a terminal and activate the **pm4ngs_venv** virtual environment

```
veraalva@perseo:~$ source pm4ngs_venv/bin/activate
(pm4ngs_venv) veraalva@perseo:~$ pm4ngs-chipexo-demo
Generating demo for ChIP-exo data analysis project
Downloading file: pm4ngs_chipexo_demo_config.yaml
Downloading file: pm4ngs_chipexo_demo_sample_data.csv
Using config file: pm4ngs_chipexo_demo_config.yaml
{
  "author_name": "Roberto Vera Alvarez",
  "user_email": "veraalva@ncbi.nlm.nih.gov",
  "project_name": "pm4ngs-chipexo",
  "dataset_name": "PRJNA338159",
  "is_data_in_SRA": "y",
  "sequencing_technology": "single-end",
  "create_demo": "n",
  "number_spots": "1000000",
  "organism": "Escherichia coli",
  "genome_name": "NC_000913.3",
  "genome_dir": "{{ cookiecutter.genome_name }}",
  "aligner_index_dir": "{{ cookiecutter.genome_dir }}/BWA/",
  "genome_fasta": "{{ cookiecutter.genome_dir }}/NC_000913.3.fa",
  "genome_gtf": "{{ cookiecutter.genome_dir }}/NC_000913.3.gtf",
  "genome_chromsizes": "{{ cookiecutter.genome_dir }}/NC_000913.3.sizes",
  "use_docker": "y",
  "max_number_threads": "32"
}
Cloning Git repo: https://github.com/ncbi/cwl-ngs-workflows-cbb to /home/veraalva/
↳pm4ngs-chipexo/bin/cwl
Updating CWLs dockerPull and SoftwareRequirement from: /home/veraalva/pm4ngs-chipexo/
↳requirements/conda-env-dependencies.yaml
bamscale with version 0.0.3 update image to: quay.io/biocontainers/bamscale:0.0.3--
↳ha85820d_0
/Users/veraalva/my_ngs_project/bin/cwl/tools/bamscale/bamscale-docker.yml with
↳old image replaced: quay.io/biocontainers/bamscale:0.0.5--h18f8b1d_1
bedtools with version 2.29.2 update image to: quay.io/biocontainers/bedtools:2.29.2--
↳hc088bd4_0
/Users/veraalva/my_ngs_project/bin/cwl/tools/bedtools/bedtools-docker.yml with
↳old image replaced: quay.io/biocontainers/bedtools:2.28.0--hdf88d34_0
bioconductor-diffbind with version 2.16.0 update image to: quay.io/biocontainers/
↳bioconductor-diffbind:2.16.0--r40h5f743cb_0
/Users/veraalva/my_ngs_project/bin/cwl/tools/R/deseq2-pca.cwl with old image
↳replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
/Users/veraalva/my_ngs_project/bin/cwl/tools/R/macscutoff.cwl with old image
↳replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
/Users/veraalva/my_ngs_project/bin/cwl/tools/R/dga_heatmaps.cwl with old image
↳replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
/Users/veraalva/my_ngs_project/bin/cwl/tools/R/diffbind.cwl with old image
↳replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
/Users/veraalva/my_ngs_project/bin/cwl/tools/R/edgeR-2conditions.cwl with old
↳image replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
/Users/veraalva/my_ngs_project/bin/cwl/tools/R/volcano_plot.cwl with old image
↳replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
/Users/veraalva/my_ngs_project/bin/cwl/tools/R/readQC.cwl with old image
↳replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
```

(continues on next page)

(continued from previous page)

```

/Users/veraalva/my_nginx_project/bin/cwl/tools/R/deseq2-2conditions.cwl with old_
→image replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
bwa with version 0.7.17 update image to: quay.io/biocontainers/bwa:0.7.17--hed695b0_7
/Users/veraalva/my_nginx_project/bin/cwl/tools/bwa/bwa-docker.yml with old image_
→replaced: quay.io/biocontainers/bwa:0.7.17--h84994c4_5
There is not biocontainer image for gffread version 0.12.1
homer with version 4.11 update image to: quay.io/biocontainers/homer:4.11--
→pl526h9a982cc_2
/Users/veraalva/my_nginx_project/bin/cwl/tools/homer/homer-docker.yml with old_
→image replaced: quay.io/biocontainers/homer:4.11--pl526h2bce143_2
mace with version 1.2 update image to: quay.io/biocontainers/mace:1.2--py27h99da42f_0
/Users/veraalva/my_nginx_project/bin/cwl/tools/mace/mace-docker.yml with old image_
→replaced: quay.io/biocontainers/mace:1.2--py27h99da42f_1
meme with version 5.1.1 update image to: quay.io/biocontainers/meme:5.1.1--
→py37pl526h072abfd_3
/Users/veraalva/my_nginx_project/bin/cwl/tools/meme/meme-docker.yml with old image_
→replaced: quay.io/biocontainers/meme:5.1.1--py27pl526h53063a7_3
Copying file /Users/veraalva/Work/Developer/Python/pm4ngs/pm4ngs-chipexo/example/
→pm4ngs_chipexo_demo_sample_data.csv to /Users/veraalva/my_nginx_project/data/my_
→dataset_name/sample_table.csv
6 files loaded
Using table:
  sample_name file                condition  replicate
0  SRR4011416   Exp_O2_growth_no_rifampicin      1
1  SRR4011417   Exp_O2_growth_no_rifampicin      2
2  SRR4011421   Exp_O2_growth_rifampicin         1
3  SRR4011425   Exp_O2_growth_rifampicin         2
4  SRR4011418   Stat_O2_growth_no_rifampicin      1
5  SRR4011419   Stat_O2_growth_no_rifampicin      2
Done

```

2.4 Running the Jupyter Server

Open a terminal and activate the **pm4ngs_venv** virtual environment

```

veraalva@perseo:~$ source pm4ngs_venv/bin/activate
(pm4ngs_venv) veraalva@perseo:~$ jupyter notebook --no-browser
[I 17:04:45.633 NotebookApp] Serving notebooks from local directory: /home/veraalva
[I 17:04:45.633 NotebookApp] Jupyter Notebook 6.1.4 is running at:
[I 17:04:45.634 NotebookApp] http://localhost:8888/?
→token=90bcbcd8a87e5421cf451e6a58d88bfa212355b36f0ed7f1a
[I 17:04:45.634 NotebookApp] or http://127.0.0.1:8888/?
→token=90bcbcd8a87e5421cf451e6a58d88bfa212355b36f0ed7f1a
[I 17:04:45.634 NotebookApp] Use Control-C to stop this server and shut down all_
→kernels (twice to skip confirmation).
[C 17:04:45.637 NotebookApp]

To access the notebook, open this file in a browser:
  file:///home/veraalva/.local/share/jupyter/runtime/nbserver-522-open.html
Or copy and paste one of these URLs:
  http://localhost:8888/?token=90bcbcd8a87e5421cf451e6a58d88bfa212355b36f0ed7f1a
  or http://127.0.0.1:8888/?token=90bcbcd8a87e5421cf451e6a58d88bfa212355b36f0ed7f1a

```

Copy the URL with localhost in a browser.

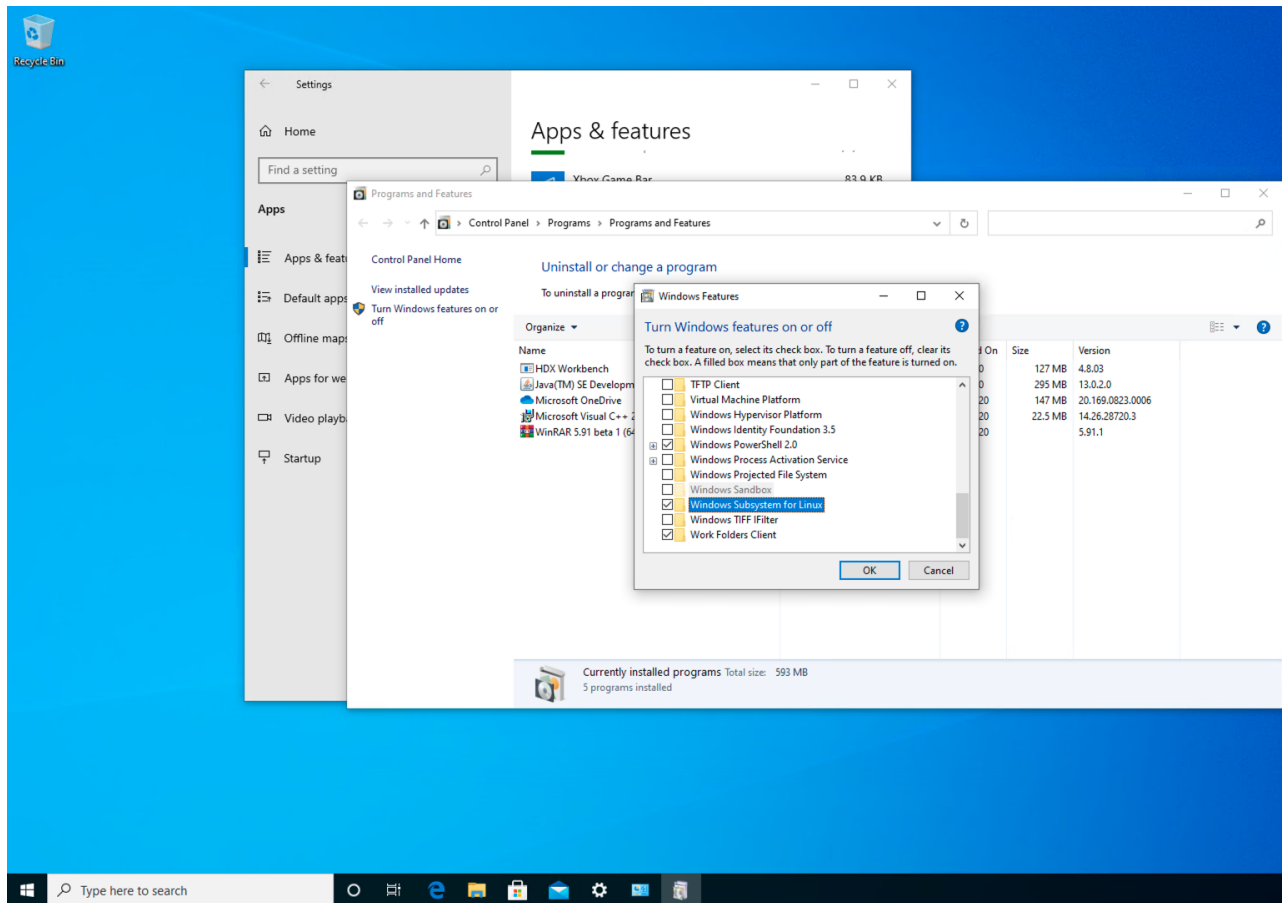
PM4NGS on Windows Subsystem for Linux

PM4NGS can be installed on Microsoft Windows 10 using the Windows Subsystem for Linux (WSL) and Docker.

3.1 Activating WSL

Follow the steps for activating WSL

1. Open the Settings app
2. Go to Apps -> Apps & Features
3. Scroll down to the Programs and Feature link
4. Click the link. The Programs and Features dialog will be opened.
5. On the left, click the link Turn Windows Features on or off.
6. The dialog Windows Features will appear on the screen. Scroll down to the option named Windows Subsystem for Linux and enable it
7. Click OK to apply the changes you made. Windows will install WSL.



8. Reboot the operating system when prompted.

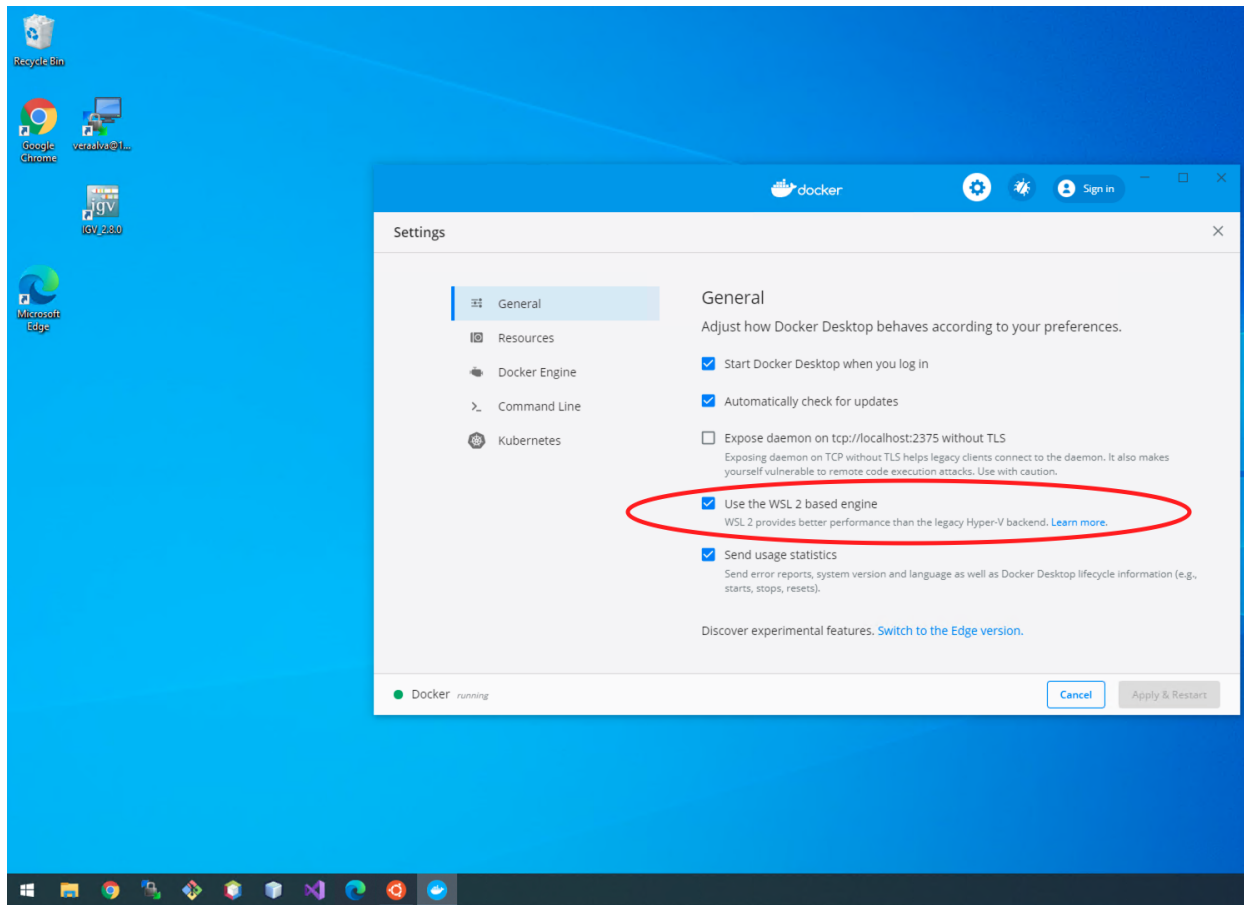
3.2 Installing Docker

Before you install the Docker Desktop WSL 2 backend, you must complete the following steps:

1. Install Windows 10, version 2004 or higher. The Docker Desktop Edge release also supports Windows 10, version 1903 or higher.
2. Enable WSL 2 feature on Windows. For detailed instructions, refer to the [Microsoft documentation](#).

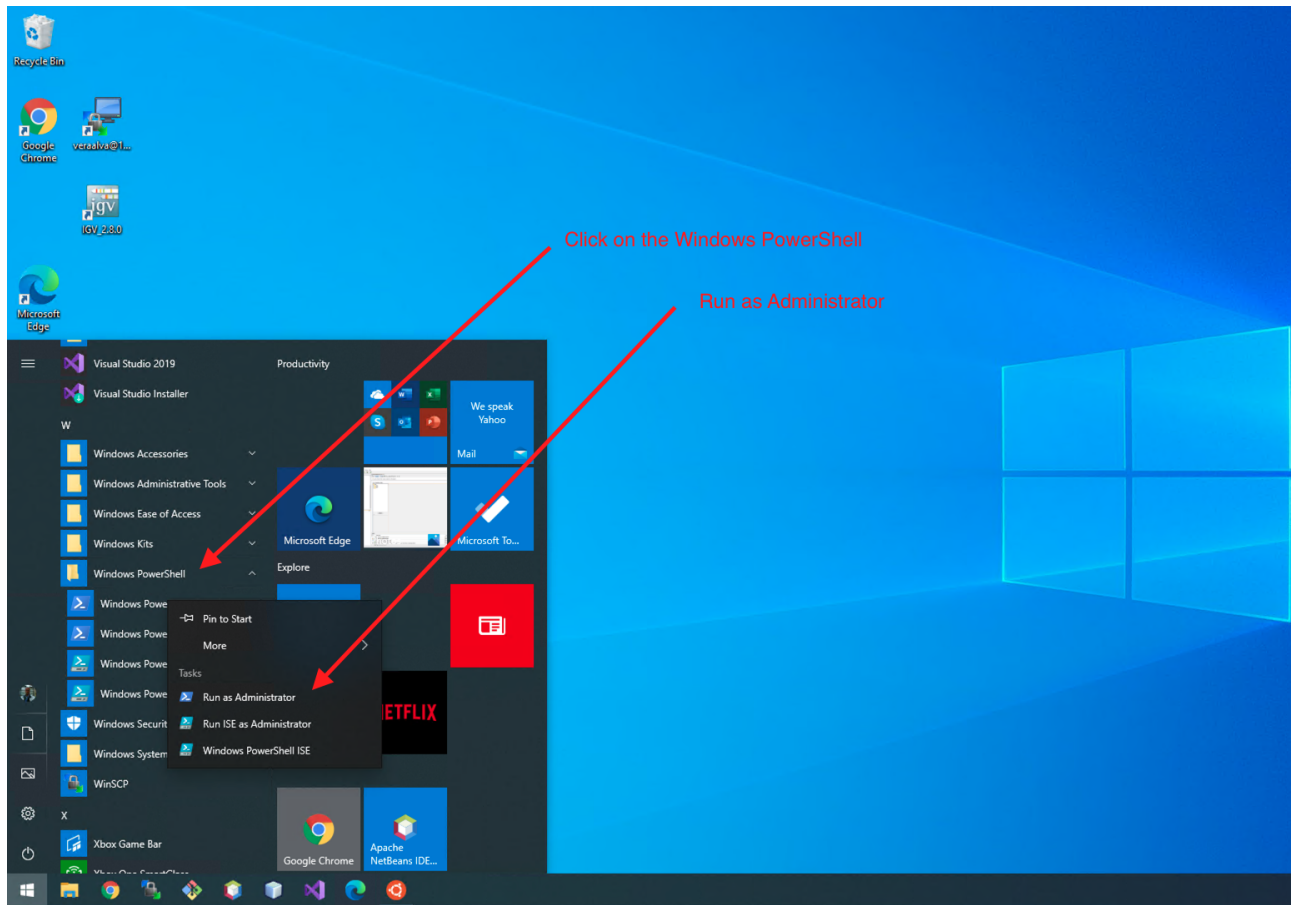
Download and install the [Linux kernel update package](#).

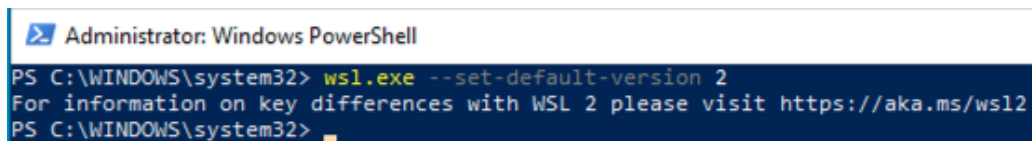
Download [Docker Desktop Stable 2.3.0.2](#) or a later release. After installed, open the Docker Setting and activate the WSL 2. Then, click on Apply and Restart.



Open a PowerShell as Administrator for activating WSL 2 as default option runnign the command

```
wsl.exe --set-default-version 2
```

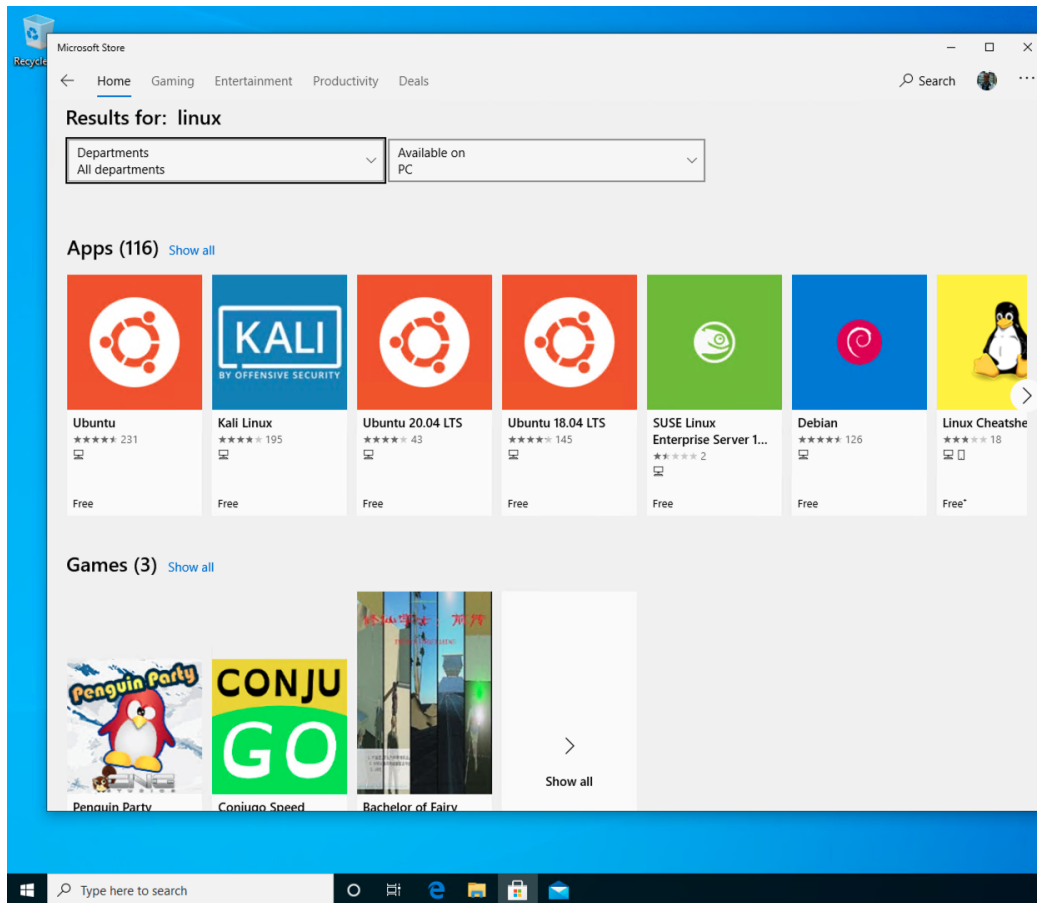


A screenshot of a Windows PowerShell terminal window titled "Administrator: Windows PowerShell". The terminal has a dark blue background. The command prompt shows the user at the C:\WINDOWS\system32 directory. They have entered the command `wsl.exe --set-default-version 2`. The terminal displays the output: "For information on key differences with WSL 2 please visit <https://aka.ms/wsl2>". The prompt then returns to `PS C:\WINDOWS\system32>` with a cursor.

```
Administrator: Windows PowerShell
PS C:\WINDOWS\system32> wsl.exe --set-default-version 2
For information on key differences with WSL 2 please visit https://aka.ms/wsl2
PS C:\WINDOWS\system32>
```

3.3 Installing Ubuntu on WSL

Open Microsoft Store and search for Linux, then install the Linux distribution that you like. We recommend Ubuntu 20.04 LTS.



Lunch Ubuntu

3.4 Installing PM4NGS

Runs these commands on the Ubuntu terminal to to prepare the instance to run PM4NGS

```
veraalva@perseo:~$ sudo apt-get update
veraalva@perseo:~$ sudo apt-get install python3 python3-pip python3-venv python3-dev_
↳poppler-utils gcc nodejs docker.io
veraalva@perseo:~$ sudo usermod -aG docker $USER
```

Close and reopen the terminal to set the docker group in the user.

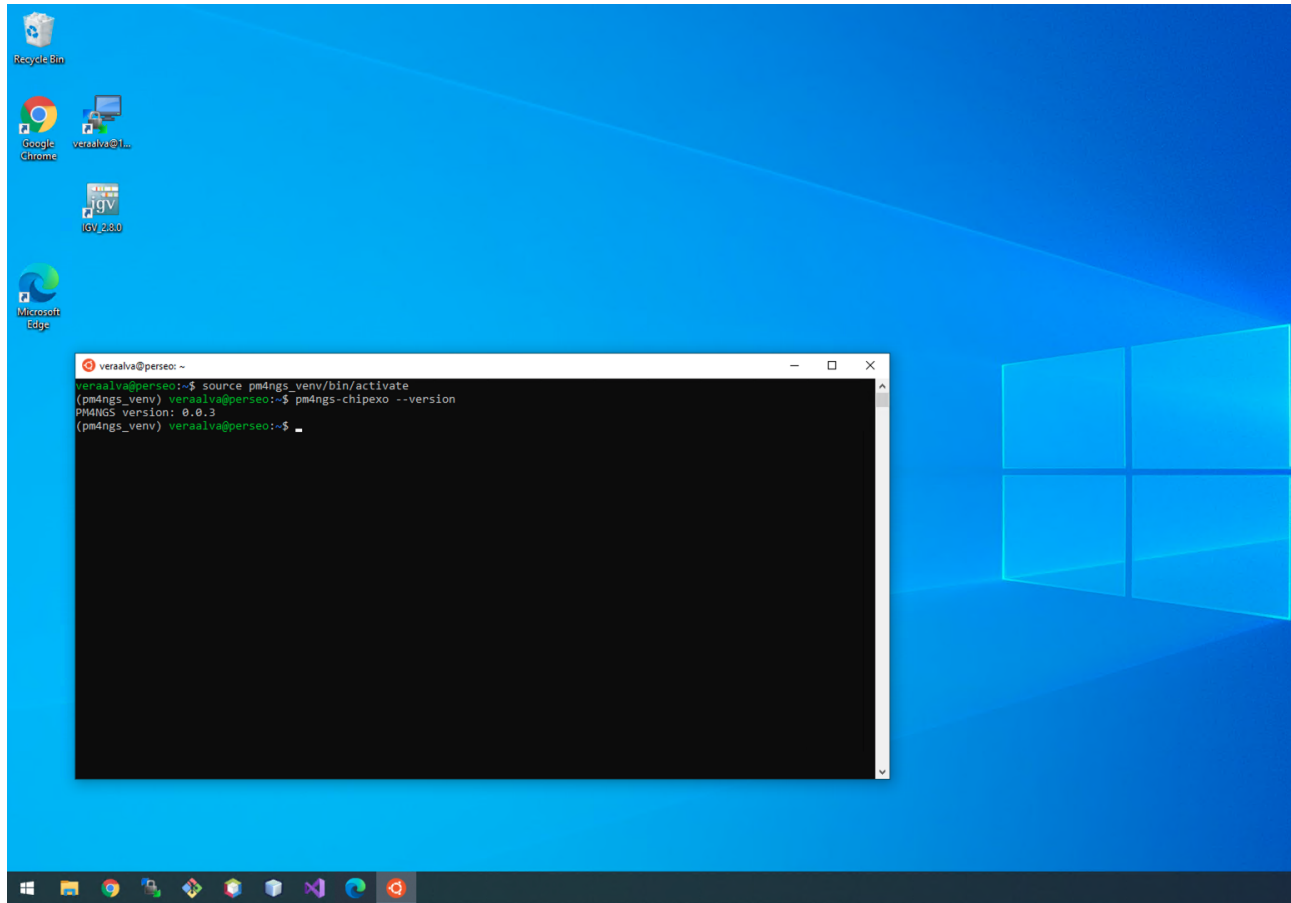
Creates a Python virtual environment named: **pm4ngs_venv** for installing PM4NGS

```
veraalva@perseo:~$ python3 -m venv pm4ngs_venv
veraalva@perseo:~$ source pm4ngs_venv/bin/activate
(pm4ngs_venv) veraalva@perseo:~$ pip install wheel
(pm4ngs_venv) veraalva@perseo:~$ pip install pm4ngs
```

3.5 Using PM4NGS on WSL

Open the Ubuntu app in Windows and activate the **pm4ngs_venv** virtual environment


```
veraalva@perseo:~$ source pm4ngs_venv/bin/activate
(pm4ngs_venv) veraalva@perseo:~$ pm4ngs-chipexo --version
PM4NGS version: 0.0.4
(pm4ngs_venv) veraalva@perseo:~$
```



3.6 Running the ChIP-exo demo

Open the Ubuntu app in Windows and activate the **pm4ngs_venv** virtual environment

```
veraalva@perseo:~$ source pm4ngs_venv/bin/activate
(pm4ngs_venv) veraalva@perseo:~$ pm4ngs-chipexo-demo
Generating demo for ChIP-exo data analysis project
Downloading file: pm4ngs_chipexo_demo_config.yaml
Downloading file: pm4ngs_chipexo_demo_sample_data.csv
Using config file: pm4ngs_chipexo_demo_config.yaml
{
  "author_name": "Roberto Vera Alvarez",
  "user_email": "veraalva@ncbi.nlm.nih.gov",
  "project_name": "pm4ngs-chipexo",
  "dataset_name": "PRJNA338159",
  "is_data_in_SRA": "y",
  "sequencing_technology": "single-end",
  "create_demo": "n",
```

(continues on next page)

(continued from previous page)

```

    "number_spots": "1000000",
    "organism": "Escherichia coli",
    "genome_name": "NC_000913.3",
    "genome_dir": "{{ cookiecutter.genome_name }}",
    "aligner_index_dir": "{{ cookiecutter.genome_dir }}/BWA/",
    "genome_fasta": "{{ cookiecutter.genome_dir }}/NC_000913.3.fa",
    "genome_gtf": "{{ cookiecutter.genome_dir }}/NC_000913.3.gtf",
    "genome_chromsizes": "{{ cookiecutter.genome_dir }}/NC_000913.3.sizes",
    "use_docker": "y",
    "max_number_threads": "32"
}

```

Cloning Git repo: <https://github.com/ncbi/cwl-ngs-workflows-cbb> to /home/veraalva/
 ↪ pm4ngs-chipexo/bin/cwl

Updating CWLs dockerPull and SoftwareRequirement from: /home/veraalva/pm4ngs-chipexo/
 ↪ requirements/conda-env-dependencies.yaml

bamscale with version 0.0.3 update image to: quay.io/biocontainers/bamscale:0.0.3--
 ↪ ha85820d_0

/Users/veraalva/my_ngs_project/bin/cwl/tools/bamscale/bamscale-docker.yml with
 ↪ old image replaced: quay.io/biocontainers/bamscale:0.0.5--h18f8b1d_1

bedtools with version 2.29.2 update image to: quay.io/biocontainers/bedtools:2.29.2--
 ↪ hc088bd4_0

/Users/veraalva/my_ngs_project/bin/cwl/tools/bedtools/bedtools-docker.yml with
 ↪ old image replaced: quay.io/biocontainers/bedtools:2.28.0--hdf88d34_0

bioconductor-diffbind with version 2.16.0 update image to: quay.io/biocontainers/
 ↪ bioconductor-diffbind:2.16.0--r40h5f743cb_0

/Users/veraalva/my_ngs_project/bin/cwl/tools/R/deseq2-pca.cwl with old image
 ↪ replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2

/Users/veraalva/my_ngs_project/bin/cwl/tools/R/macscutoff.cwl with old image
 ↪ replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2

/Users/veraalva/my_ngs_project/bin/cwl/tools/R/dga_heatmaps.cwl with old image
 ↪ replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2

/Users/veraalva/my_ngs_project/bin/cwl/tools/R/diffbind.cwl with old image
 ↪ replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2

/Users/veraalva/my_ngs_project/bin/cwl/tools/R/edgeR-2conditions.cwl with old
 ↪ image replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2

/Users/veraalva/my_ngs_project/bin/cwl/tools/R/volcano_plot.cwl with old image
 ↪ replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2

/Users/veraalva/my_ngs_project/bin/cwl/tools/R/readQC.cwl with old image
 ↪ replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2

/Users/veraalva/my_ngs_project/bin/cwl/tools/R/deseq2-2conditions.cwl with old
 ↪ image replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2

bwa with version 0.7.17 update image to: quay.io/biocontainers/bwa:0.7.17--hed695b0_7

/Users/veraalva/my_ngs_project/bin/cwl/tools/bwa/bwa-docker.yml with old image
 ↪ replaced: quay.io/biocontainers/bwa:0.7.17--h84994c4_5

There is not biocontainer image **for** gffread version 0.12.1

homer with version 4.11 update image to: quay.io/biocontainers/homer:4.11--
 ↪ pl526h9a982cc_2

/Users/veraalva/my_ngs_project/bin/cwl/tools/homer/homer-docker.yml with old
 ↪ image replaced: quay.io/biocontainers/homer:4.11--pl526h2bce143_2

mace with version 1.2 update image to: quay.io/biocontainers/mace:1.2--py27h99da42f_0

/Users/veraalva/my_ngs_project/bin/cwl/tools/mace/mace-docker.yml with old image
 ↪ replaced: quay.io/biocontainers/mace:1.2--py27h99da42f_1

meme with version 5.1.1 update image to: quay.io/biocontainers/meme:5.1.1--
 ↪ py37pl526h072abfd_3

/Users/veraalva/my_ngs_project/bin/cwl/tools/meme/meme-docker.yml with old image
 ↪ replaced: quay.io/biocontainers/meme:5.1.1--py27pl526h53063a7_3

Copying file /Users/veraalva/Work/Developer/Python/pm4ngs/pm4ngs-chipexo/example/
 ↪ pm4ngs_chipexo_demo_sample_data.csv to /Users/veraalva/my_ngs_project/data/my_
 ↪ dataset_name/sample_table.csv

(continues on next page)

(continued from previous page)

```

6 files loaded
Using table:
  sample_name file                condition  replicate
0  SRR4011416      Exp_O2_growth_no_rifampicin      1
1  SRR4011417      Exp_O2_growth_no_rifampicin      2
2  SRR4011421      Exp_O2_growth_rifampicin         1
3  SRR4011425      Exp_O2_growth_rifampicin         2
4  SRR4011418      Stat_O2_growth_no_rifampicin     1
5  SRR4011419      Stat_O2_growth_no_rifampicin     2
Done

```

3.7 Running the Jupyter Server

Open the Ubuntu app in Windows and activate the **pm4ngs_venv** virtual environment

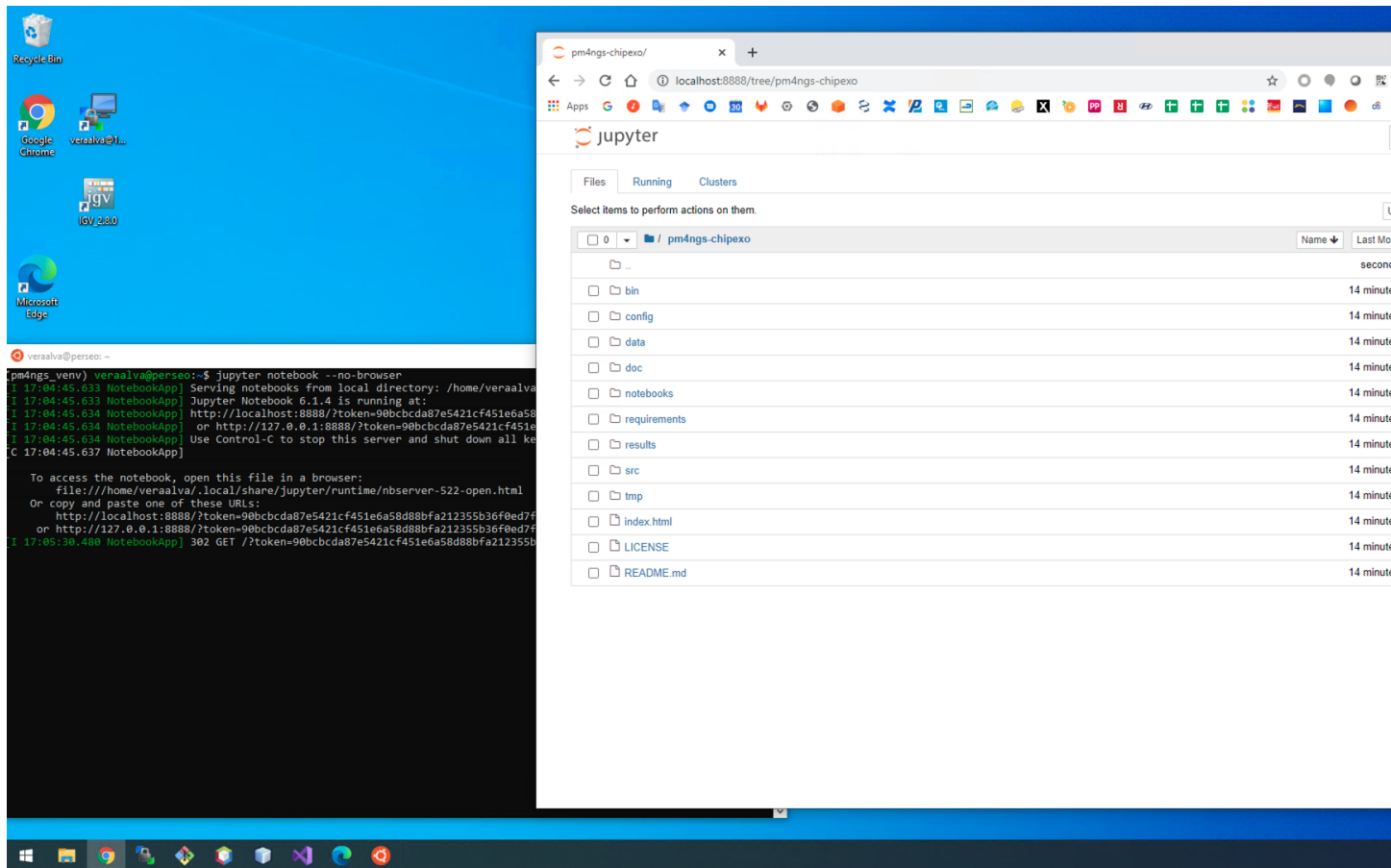
```

veraalva@perseo:~$ source pm4ngs_venv/bin/activate
(pm4ngs_venv) veraalva@perseo:~$ jupyter notebook --no-browser
[I 17:04:45.633 NotebookApp] Serving notebooks from local directory: /home/veraalva
[I 17:04:45.633 NotebookApp] Jupyter Notebook 6.1.4 is running at:
[I 17:04:45.634 NotebookApp] http://localhost:8888/?
↪token=90bcbcd87e5421cf451e6a58d88bfa212355b36f0ed7f1a
[I 17:04:45.634 NotebookApp] or http://127.0.0.1:8888/?
↪token=90bcbcd87e5421cf451e6a58d88bfa212355b36f0ed7f1a
[I 17:04:45.634 NotebookApp] Use Control-C to stop this server and shut down all
↪kernels (twice to skip confirmation).
[C 17:04:45.637 NotebookApp]

To access the notebook, open this file in a browser:
    file:///home/veraalva/.local/share/jupyter/runtime/nbserver-522-open.html
Or copy and paste one of these URLs:
    http://localhost:8888/?token=90bcbcd87e5421cf451e6a58d88bfa212355b36f0ed7f1a
    or http://127.0.0.1:8888/?token=90bcbcd87e5421cf451e6a58d88bfa212355b36f0ed7f1a

```

Copy the URL with localhost in a browser:



PM4NGS on GCP instance with Ubuntu

4.1 Creating a GCP instance with Ubuntu 20.02 LTS

Creates a VM instance running Ubuntu

i You have a draft that wasn't submitted, click Restore to keep working on it Restore

Name ?
Name is permanent
instance-1

Labels ? (Optional)
+ Add label

Region ?
Region is permanent
us-central1 (Iowa)


Zone ?
Zone is permanent
us-central1-c

Machine configuration

Machine family
General-purpose Memory-optimized Compute-optimized
Machine types for common workloads, optimized for cost and flexibility

Series
E2
CPU platform selection based on availability

Machine type
e2-standard-16 (16 vCPU, 64 GB memory)


	vCPU	Memory	GPUs
	16	64 GB	-

⌵ CPU platform and GPU

Confidential VM service ?
☐ Enable the Confidential Computing service on this VM instance.

Container ?
☐ Deploy a container image to this VM instance. [Learn more](#)

Boot disk ?



New 500 GB standard persistent disk
Image
Ubuntu 20.04 LTS Change

\$411.35 monthly estimate
That's about \$0.563 hourly
Pay for what you use: No upfront costs and per second bi

Item	Estimated costs
16 vCPUs + 64 GB memory	\$391.35/month
500 GB standard persistent disk	\$20.00/month
Sustained use discount ?	- \$0.00/month
Total	\$411.35/month

[Compute Engine pricing](#) ↗

⌵ Less

Select a 16 vCPU machine type

Select Ubuntu 20.04 LTS with 500 GB disk

Select boot disk with Ubuntu 20.04 LTS with 500 GB of standard persistent disk.

Boot disk

Select an image or snapshot to create a boot disk; or attach an existing disk. Can't find what you need?

Public images Custom images Snapshots Existing disks

Operating system
Ubuntu

Version
Ubuntu 20.04 LTS

amd64 focal image built on 2020-09-17, supports Shielded VM features ?

Boot disk type ? **Size (GB)** ?

Standard persistent disk 500

Click on the **SSH** button for accessing the instance

VM instances

Click SSH for accessing the VM terminal

Filter VM instances

	Name ^	Zone	Recommendation	In use by	Internal IP	External IP	Connect
<input type="checkbox"/>	instance-1	us-central1-c			10.128.0.21 (nic0)	34.123.78.105	SSH

A terminal is available after accessing through SSH

```

veraalva@instance-1: ~
ssh.cloud.google.com/projects/cbb-research-dl/zones/us-central1-c/instances/instance-1?nonAdminProxySessionReason=1&aut...
System information as of Wed Oct  7 13:33:55 UTC 2020

System load:  0.08          Processes:           248
Usage of /:   0.3% of 484.41GB Users logged in:       0
Memory usage: 0%          IPv4 address for ens4: 10.128.0.21
Swap usage:   0%

1 update can be installed immediately.
0 of these updates are security updates.
To see these additional updates run: apt list --upgradable

The list of available updates is more than a week old.
To check for new updates run: sudo apt update

The programs included with the Ubuntu system are free software;
the exact distribution terms for each program are described in the
individual files in /usr/share/doc/*/copyright.

Ubuntu comes with ABSOLUTELY NO WARRANTY, to the extent permitted by
applicable law.

veraalva@instance-1:~$ 

```

4.2 Installing PM4NGS on the GCP instance with Ubuntu

Runs these commands on a terminal to prepare the instance to run PM4NGS

```

veraalva@instance-1:~$ sudo apt-get update
veraalva@instance-1:~$ sudo apt-get install docker.io python3 python3-pip python3-
↳venv python3-dev poppler-utils gcc nodejs tree
veraalva@instance-1:~$ sudo usermod -aG docker $USER
veraalva@instance-1:~$ logout

```

Close and reopen the terminal to set the docker group in the user. Then, click on the SSH button again to re-launch the terminal.

4.3 Installing PM4NGS

Creates a Python virtual environment named: **pm4ngs_venv** for installing PM4NGS

```

veraalva@instance-1:~$ python3 -m venv pm4ngs_venv
veraalva@instance-1:~$ source pm4ngs_venv/bin/activate
(pm4ngs_venv) veraalva@instance-1:~$ pip install wheel
(pm4ngs_venv) veraalva@instance-1:~$ pip install pm4ngs

```


4.4 Using PM4NGS

Open a terminal and activate the **pm4ngs_venv** virtual environment

```
veraalva@instance-1:~$ source pm4ngs_venv/bin/activate
(pm4ngs_venv) veraalva@instance-1:~$ pm4ngs-chipexo --version
PM4NGS version: 0.0.4
(pm4ngs_venv) veraalva@instance-1:~$
```

4.5 Running the ChIP-exo demo

Open a terminal and activate the **pm4ngs_venv** virtual environment

```
veraalva@instance-1:~$ source pm4ngs_venv/bin/activate
(pm4ngs_venv) veraalva@perseo:~$ pm4ngs-chipexo-demo
Generating demo for ChIP-exo data analysis project
Downloading file: pm4ngs_chipexo_demo_config.yaml
Downloading file: pm4ngs_chipexo_demo_sample_data.csv
Using config file: pm4ngs_chipexo_demo_config.yaml
{
  "author_name": "Roberto Vera Alvarez",
  "user_email": "veraalva@ncbi.nlm.nih.gov",
  "project_name": "pm4ngs-chipexo",
  "dataset_name": "PRJNA338159",
  "is_data_in_SRA": "y",
  "sequencing_technology": "single-end",
  "create_demo": "n",
  "number_spots": "1000000",
  "organism": "Escherichia coli",
  "genome_name": "NC_000913.3",
  "genome_dir": "{{ cookiecutter.genome_name }}",
  "aligner_index_dir": "{{ cookiecutter.genome_dir }}/BWA/",
  "genome_fasta": "{{ cookiecutter.genome_dir }}/NC_000913.3.fa",
  "genome_gtf": "{{ cookiecutter.genome_dir }}/NC_000913.3.gtf",
  "genome_chromsizes": "{{ cookiecutter.genome_dir }}/NC_000913.3.sizes",
  "use_docker": "y",
  "max_number_threads": "32"
}
Cloning Git repo: https://github.com/ncbi/cwl-ngs-workflows-cbb to /home/veraalva/
↳ pm4ngs-chipexo/bin/cwl
Updating CWLs dockerPull and SoftwareRequirement from: /home/veraalva/pm4ngs-chipexo/
↳ requirements/conda-env-dependencies.yaml
bamscale with version 0.0.3 update image to: quay.io/biocontainers/bamscale:0.0.3--
↳ ha85820d_0
/Users/veraalva/my_ngs_project/bin/cwl/tools/bamscale/bamscale-docker.yml with
↳ old image replaced: quay.io/biocontainers/bamscale:0.0.5--h18f8b1d_1
bedtools with version 2.29.2 update image to: quay.io/biocontainers/bedtools:2.29.2--
↳ hc088bd4_0
/Users/veraalva/my_ngs_project/bin/cwl/tools/bedtools/bedtools-docker.yml with
↳ old image replaced: quay.io/biocontainers/bedtools:2.28.0--hdf88d34_0
bioconductor-diffbind with version 2.16.0 update image to: quay.io/biocontainers/
↳ bioconductor-diffbind:2.16.0--r40h5f743cb_0
/Users/veraalva/my_ngs_project/bin/cwl/tools/R/deseq2-pca.cwl with old image
↳ replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
/Users/veraalva/my_ngs_project/bin/cwl/tools/R/macs-cutoff.cwl with old image
↳ replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2 (continues on next page)
```

(continued from previous page)

```

/Users/veraalva/my_nginx_project/bin/cwl/tools/R/dga_heatmaps.cwl with old image_
→replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
/Users/veraalva/my_nginx_project/bin/cwl/tools/R/diffbind.cwl with old image_
→replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
/Users/veraalva/my_nginx_project/bin/cwl/tools/R/edgeR-2conditions.cwl with old_
→image replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
/Users/veraalva/my_nginx_project/bin/cwl/tools/R/volcano_plot.cwl with old image_
→replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
/Users/veraalva/my_nginx_project/bin/cwl/tools/R/readQC.cwl with old image_
→replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
/Users/veraalva/my_nginx_project/bin/cwl/tools/R/deseq2-2conditions.cwl with old_
→image replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
bwa with version 0.7.17 update image to: quay.io/biocontainers/bwa:0.7.17--hed695b0_7
/Users/veraalva/my_nginx_project/bin/cwl/tools/bwa/bwa-docker.yml with old image_
→replaced: quay.io/biocontainers/bwa:0.7.17--h84994c4_5
There is not biocontainer image for gffread version 0.12.1
homer with version 4.11 update image to: quay.io/biocontainers/homer:4.11--
→pl526h9a982cc_2
/Users/veraalva/my_nginx_project/bin/cwl/tools/homer/homer-docker.yml with old_
→image replaced: quay.io/biocontainers/homer:4.11--pl526h2bce143_2
mace with version 1.2 update image to: quay.io/biocontainers/mace:1.2--py27h99da42f_0
/Users/veraalva/my_nginx_project/bin/cwl/tools/mace/mace-docker.yml with old image_
→replaced: quay.io/biocontainers/mace:1.2--py27h99da42f_1
meme with version 5.1.1 update image to: quay.io/biocontainers/meme:5.1.1--
→py37pl526h072abfd_3
/Users/veraalva/my_nginx_project/bin/cwl/tools/meme/meme-docker.yml with old image_
→replaced: quay.io/biocontainers/meme:5.1.1--py27pl526h53063a7_3
Copying file /Users/veraalva/Work/Developer/Python/pm4ngs/pm4ngs-chipexo/example/
→pm4ngs_chipexo_demo_sample_data.csv to /Users/veraalva/my_nginx_project/data/my_
→dataset_name/sample_table.csv
6 files loaded
Using table:
  sample_name file                condition replicate
0  SRR4011416   Exp_O2_growth_no_rifampicin      1
1  SRR4011417   Exp_O2_growth_no_rifampicin      2
2  SRR4011421   Exp_O2_growth_rifampicin         1
3  SRR4011425   Exp_O2_growth_rifampicin         2
4  SRR4011418   Stat_O2_growth_no_rifampicin      1
5  SRR4011419   Stat_O2_growth_no_rifampicin      2
Done

```

The terminal will look like the next image.

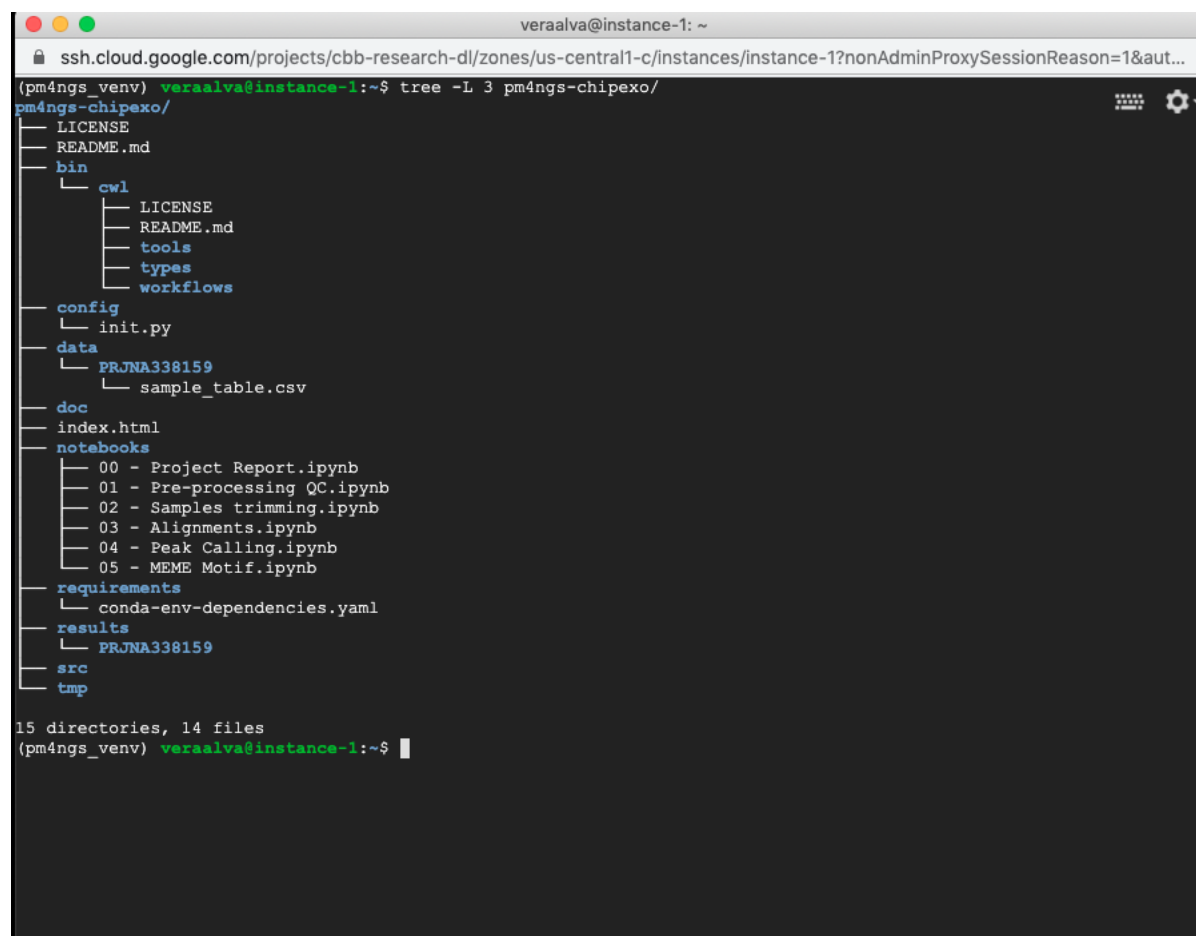
```

veraalva@instance-1: ~
ssh.cloud.google.com/projects/cbb-research-dl/zones/us-central1-c/instances/instance-1?nonAdminProxySessionReason=1&aut...
5f743cb_0
/home/veraalva/pm4ngs-chipexo/bin/cwl/tools/R/readQC.cwl with old image replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
/home/veraalva/pm4ngs-chipexo/bin/cwl/tools/R/deseq2-pca.cwl with old image replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
/home/veraalva/pm4ngs-chipexo/bin/cwl/tools/R/diffbind.cwl with old image replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
/home/veraalva/pm4ngs-chipexo/bin/cwl/tools/R/macscutoff.cwl with old image replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
/home/veraalva/pm4ngs-chipexo/bin/cwl/tools/R/deseq2-2conditions.cwl with old image replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
/home/veraalva/pm4ngs-chipexo/bin/cwl/tools/R/dga_heatmaps.cwl with old image replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
/home/veraalva/pm4ngs-chipexo/bin/cwl/tools/R/edgeR-2conditions.cwl with old image replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
/home/veraalva/pm4ngs-chipexo/bin/cwl/tools/R/volcano_plot.cwl with old image replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
bwa with version 0.7.17 update image to: quay.io/biocontainers/bwa:0.7.17--hed695b0_7
/home/veraalva/pm4ngs-chipexo/bin/cwl/tools/bwa/bwa-docker.yml with old image replaced: quay.io/biocontainers/bwa:0.7.17--hed695b0_7
There is not biocontainer image for gffread version 0.12.1
homer with version 4.11 update image to: quay.io/biocontainers/homer:4.11--pl526h9a982cc_2
/home/veraalva/pm4ngs-chipexo/bin/cwl/tools/homer/homer-docker.yml with old image replaced: quay.io/biocontainers/homer:4.11--pl526h9a982cc_2
mace with version 1.2 update image to: quay.io/biocontainers/mace:1.2--py27h99da42f_0
/home/veraalva/pm4ngs-chipexo/bin/cwl/tools/mace/mace-docker.yml with old image replaced: quay.io/biocontainers/mace:1.2--py27h99da42f_0
meme with version 5.1.1 update image to: quay.io/biocontainers/meme:5.1.1--py37pl526h072abfd_3
/home/veraalva/pm4ngs-chipexo/bin/cwl/tools/meme/meme-docker.yml with old image replaced: quay.io/biocontainers/meme:5.1.1--py37pl526h072abfd_3
Copying file /home/veraalva/pm4ngs-chipexo_demo_sample_data.csv to /home/veraalva/pm4ngs-chipexo/data/PRJNA338159/sample_table.csv
6 files loaded
Using table:
  sample_name  file              condition  replicate
0  SRR4011416    Exp_02_growth_no_rifampicin  1
1  SRR4011417    Exp_02_growth_no_rifampicin  2
2  SRR4011421    Exp_02_growth_rifampicin     1
3  SRR4011425    Exp_02_growth_rifampicin     2
4  SRR4011418    Stat_02_growth_no_rifampicin  1
5  SRR4011419    Stat_02_growth_no_rifampicin  2
Done
(pm4ngs_venv) veraalva@instance-1:~$

```

Running the command tree to show the project structure

```
(pm4ngs_venv) veraalva@instance-1:~$ tree -L 3 pm4ngs-chipexo/
```



```

veraalva@instance-1: ~
ssh.cloud.google.com/projects/cbb-research-dl/zones/us-central1-c/instances/instance-1?nonAdminProxySessionReason=1&aut...
(pm4ngs_venv) veraalva@instance-1:~$ tree -L 3 pm4ngs-chipexo/
pm4ngs-chipexo/
├── LICENSE
├── README.md
├── bin
│   └── cwl
│       ├── LICENSE
│       ├── README.md
│       ├── tools
│       ├── types
│       └── workflows
├── config
│   └── init.py
├── data
│   └── PRJNA338159
│       └── sample_table.csv
├── doc
├── index.html
├── notebooks
│   ├── 00 - Project Report.ipynb
│   ├── 01 - Pre-processing QC.ipynb
│   ├── 02 - Samples trimming.ipynb
│   ├── 03 - Alignments.ipynb
│   ├── 04 - Peak Calling.ipynb
│   └── 05 - MEME Motif.ipynb
├── requirements
│   └── conda-env-dependencies.yaml
├── results
│   └── PRJNA338159
├── src
└── tmp

15 directories, 14 files
(pm4ngs_venv) veraalva@instance-1:~$

```

4.6 Running the Jupyter Server

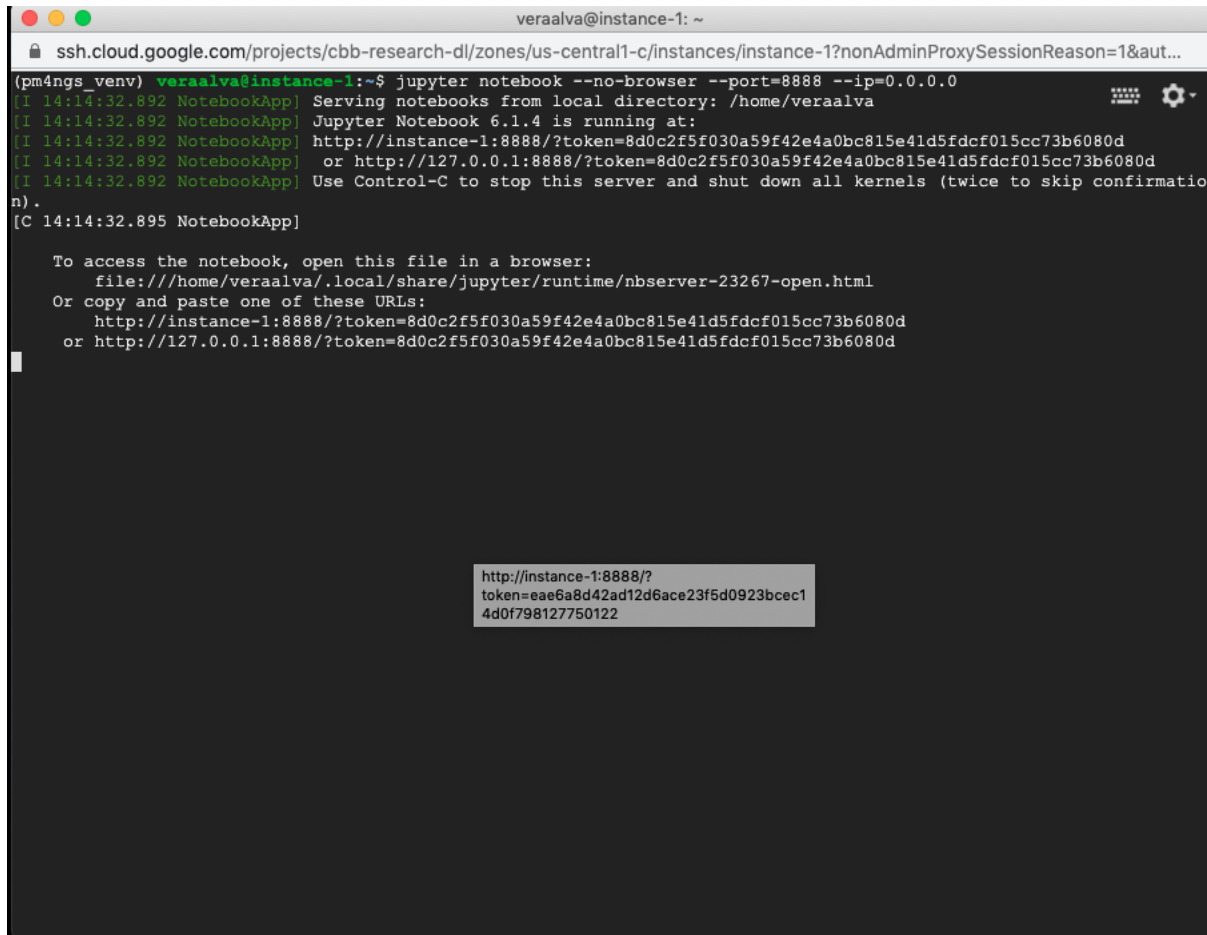
Open a terminal and activate the **pm4ngs_venv** virtual environment and run the jupyter server. As the GCP instance is a remote computer, we need to run the jupyter server with the **-port** and **-ip** options.

```

(pm4ngs_venv) veraalva@instance-1:~$ jupyter notebook --no-browser --port=8888 --ip=0.
↪0.0.0
[I 14:12:52.956 NotebookApp] Serving notebooks from local directory: /home/veraalva
[I 14:12:52.956 NotebookApp] Jupyter Notebook 6.1.4 is running at:
[I 14:12:52.956 NotebookApp] http://instance-1:8888/?
↪token=eae6a8d42ad12d6ace23f5d0923bcec14d0f798127750122
[I 14:12:52.956 NotebookApp] or http://127.0.0.1:8888/?
↪token=eae6a8d42ad12d6ace23f5d0923bcec14d0f798127750122
[I 14:12:52.956 NotebookApp] Use Control-C to stop this server and shut down all
↪kernels (twice to skip confirmatio
n).
[C 14:12:52.959 NotebookApp]

To access the notebook, open this file in a browser:
    file:///home/veraalva/.local/share/jupyter/runtime/nbserver-23251-open.html
Or copy and paste one of these URLs:
    http://instance-1:8888/?token=eae6a8d42ad12d6ace23f5d0923bcec14d0f798127750122
    or http://127.0.0.1:8888/?token=eae6a8d42ad12d6ace23f5d0923bcec14d0f798127750122

```



```

veraalva@instance-1: ~
ssh.cloud.google.com/projects/cbb-research-dl/zones/us-central1-c/instances/instance-1?nonAdminProxySessionReason=1&aut...
(pm4ngs_venv) veraalva@instance-1:~$ jupyter notebook --no-browser --port=8888 --ip=0.0.0.0
[I 14:14:32.892 NotebookApp] Serving notebooks from local directory: /home/veraalva
[I 14:14:32.892 NotebookApp] Jupyter Notebook 6.1.4 is running at:
[I 14:14:32.892 NotebookApp] http://instance-1:8888/?token=8d0c2f5f030a59f42e4a0bc815e41d5fdcf015cc73b6080d
[I 14:14:32.892 NotebookApp] or http://127.0.0.1:8888/?token=8d0c2f5f030a59f42e4a0bc815e41d5fdcf015cc73b6080d
[I 14:14:32.892 NotebookApp] Use Control-C to stop this server and shut down all kernels (twice to skip confirmation).
[C 14:14:32.895 NotebookApp]

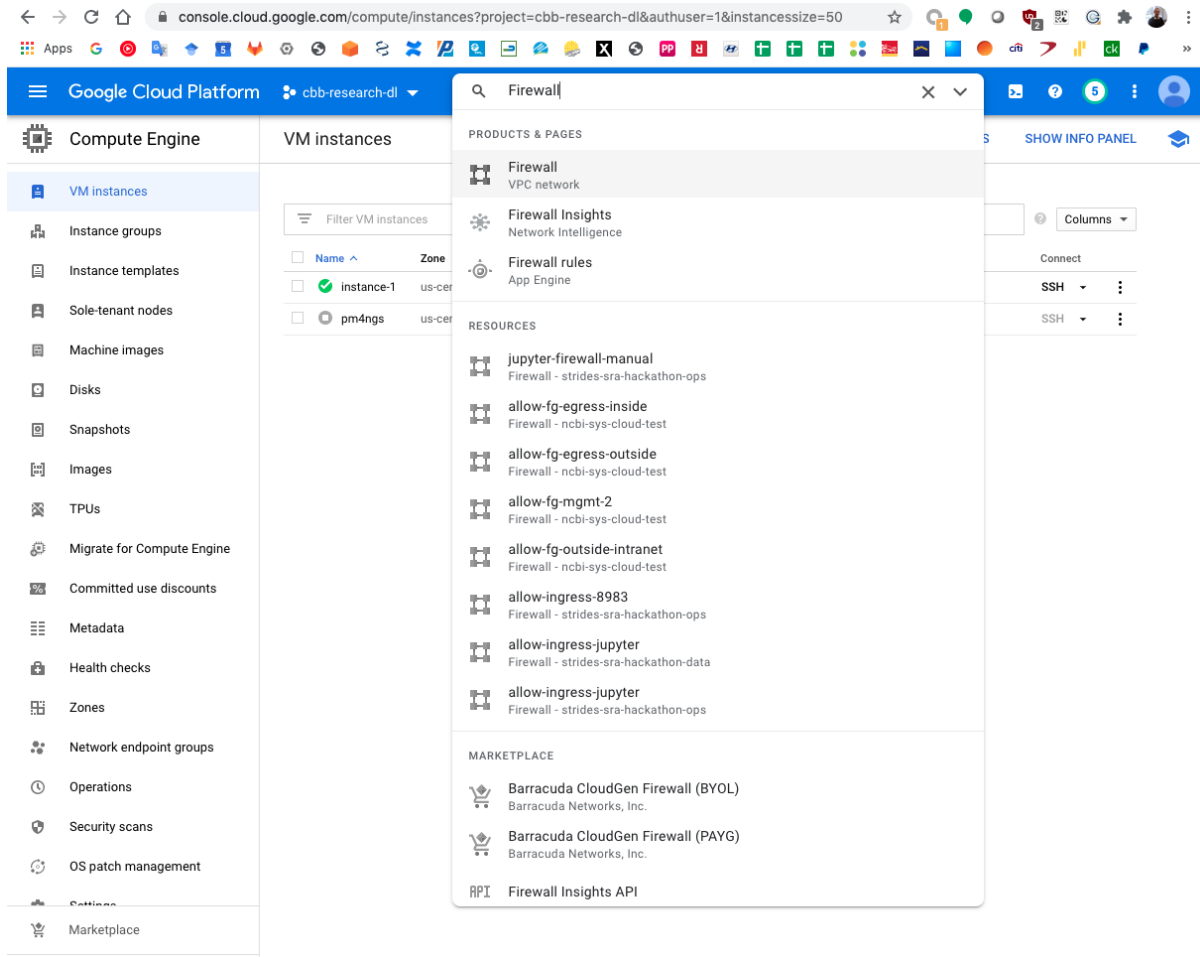
To access the notebook, open this file in a browser:
    file:///home/veraalva/.local/share/jupyter/runtime/nbserver-23267-open.html
Or copy and paste one of these URLs:
    http://instance-1:8888/?token=8d0c2f5f030a59f42e4a0bc815e41d5fdcf015cc73b6080d
    or http://127.0.0.1:8888/?token=8d0c2f5f030a59f42e4a0bc815e41d5fdcf015cc73b6080d

```

<http://instance-1:8888/?token=eae6a8d42ad12d6ace23f5d0923bcec14d0f798127750122>

4.7 Open a firewall rule for accessing the Jupyter Server

A GCP firewall rule should be created to access the Jupyter server remotely. From your desktop computer. Search for **Firewall** in the GCP search bar.



Click on the **Create Firewall Rule** button.

console.cloud.google.com/networking/firewalls/list?authuser=1&project=cbb-research-dl

Google Cloud Platform cbb-research-dl Search products and resources

VPC network

- VPC networks
- External IP addresses
- Firewall**
- Routes
- VPC network peering
- Shared VPC
- Serverless VPC access
- Packet mirroring

Firewall

+ CREATE FIREWALL RULE REFRESH CONFIGURE LOGS DELETE

Firewall policies inherited by this project BETA

Filter table

Enforcement order	Policy name	Firewall rules	Inherited from	Located at
No rows to display				

Click to create a new rule

Firewall rules in this project

Firewall rules control incoming or outgoing traffic to an instance. By default, incoming traffic from outside your network is blocked. [Learn more](#)

Note: App Engine firewalls are managed [here](#).

Filter table

Name	Type	Targets	Filters	Protocols / ports	Action	Priority	Network	↑
<input type="checkbox"/> default-allow-http	Ingress	http-server	IP ranges: 0.0.0.0/0	tcp:80	Allow	1000	default	▼
<input type="checkbox"/> default-allow-icmp	Ingress	Apply to all	IP ranges: 0.0.0.0/0	icmp	Allow	65534	default	▼
<input type="checkbox"/> default-allow-internal	Ingress	Apply to all	IP ranges: 10.0.0.0/8	tcp:0-65535 udp:0-65535 icmp	Allow	65534	default	▼
<input type="checkbox"/> default-allow-rdp	Ingress	Apply to all	IP ranges: 0.0.0.0/0	tcp:3389	Allow	65534	default	▼
<input type="checkbox"/> default-allow-ssh	Ingress	Apply to all	IP ranges: 0.0.0.0/0	tcp:22	Allow	65534	default	▼

Add the rules name

console.cloud.google.com/networking/firewalls/add?authuser=1&project=cbb-research-dl

Google Cloud Platform cbb-research-dl Search products and resources

VPC network

- VPC networks
- External IP addresses
- Firewall
- Routes
- VPC network peering
- Shared VPC
- Serverless VPC access
- Packet mirroring

Create a firewall rule

Firewall rules control incoming or outgoing traffic to an instance. By default, incoming traffic from outside your network is blocked. [Learn more](#)

Name *
jupyter
Lowercase letters, numbers, hyphens allowed

Description

Logs
Turning on firewall logs can generate a large number of logs which can increase costs in Stackdriver. [Learn more](#)
☐ On
☒ Off

Network *
default

Priority *
1000
Priority can be 0 - 65535 [Check priority of other firewall rules](#)

Direction of traffic ?
☒ Ingress
☐ Egress

Action on match ?
☒ Allow
☐ Deny

Targets
Specified target tags

Target tags *

Add the Jupyter port used that is **8888** and click on create.

console.cloud.google.com/networking/firewalls/add?authuser=1&project=cbb-research-dl

Google Cloud Platform cbb-research-dl Search products and resources

VPC network

- VPC networks
- External IP addresses
- Firewall
- Routes
- VPC network peering
- Shared VPC
- Serverless VPC access
- Packet mirroring

Create a firewall rule

☒ Ingress
☐ Egress

Action on match ?
☒ Allow
☐ Deny

Targets
All instances in the network

Source filter
IP ranges

Source IP ranges *
0.0.0.0/0 for example, 0.0.0.0/0, 192.168.2.0/24

Second source filter
None

Protocols and ports ?
☐ Allow all
☒ Specified protocols and ports

☒ tcp : 8888
☐ udp : all
☐ Other protocols
protocols, comma separated, e.g. ah, sctp

DISABLE RULE

CREATE CANCEL

Select All instances in the Network

Type IP range 0.0.0.0/0

Type Jupyter port: 8888

Click to create the rule

The new rule is created and available

console.cloud.google.com/networking/firewalls/list?authuser=1&project=cbb-research-dl

Google Cloud Platform cbb-research-dl Search products and resources

VPC network

- VPC networks
- External IP addresses
- Firewall**
- Routes
- VPC network peering
- Shared VPC
- Serverless VPC access
- Packet mirroring

Firewall

+ CREATE FIREWALL RULE REFRESH CONFIGURE LOGS DELETE

Firewall policies inherited by this project **BETA**

Filter table

Enforcement order	Policy name	Firewall rules	Description	Inherited from	Located at
No rows to display					

Firewall rules in this project

Firewall rules control incoming or outgoing traffic to an instance. By default, incoming traffic from outside your network is blocked. [Learn more](#)

Note: App Engine firewalls are managed [here](#).

Filter table

Name	Type	Targets	Filters	Protocols / ports	Action	Priority	Network	↑
<input type="checkbox"/> default-allow-http	Ingress	http-server	IP ranges: 0.0.0.0/0	tcp:80	Allow	1000	default	▼
<input type="checkbox"/> jupyter	Ingress	Apply to all	IP ranges: 0.0.0.0/0	tcp:8888	Allow	1000	default	▼
<input type="checkbox"/> default-allow-icmp	Ingress	Apply to all	IP ranges: 0.0.0.0/0	icmp	Allow	65534	default	▼
<input type="checkbox"/> default-allow-internal	Ingress	Apply to all	IP ranges: 10.0.0.0/8	tcp:0-65535 udp:0-65535 icmp	Allow	65534	default	▼
<input type="checkbox"/> default-allow-rdp	Ingress	Apply to all	IP ranges: 0.0.0.0/0	tcp:3389	Allow	65534	default	▼
<input type="checkbox"/> default-allow-ssh	Ingress	Apply to all	IP ranges: 0.0.0.0/0	tcp:22	Allow	65534	default	▼

Go back to the VM instances to copy the instance public IP

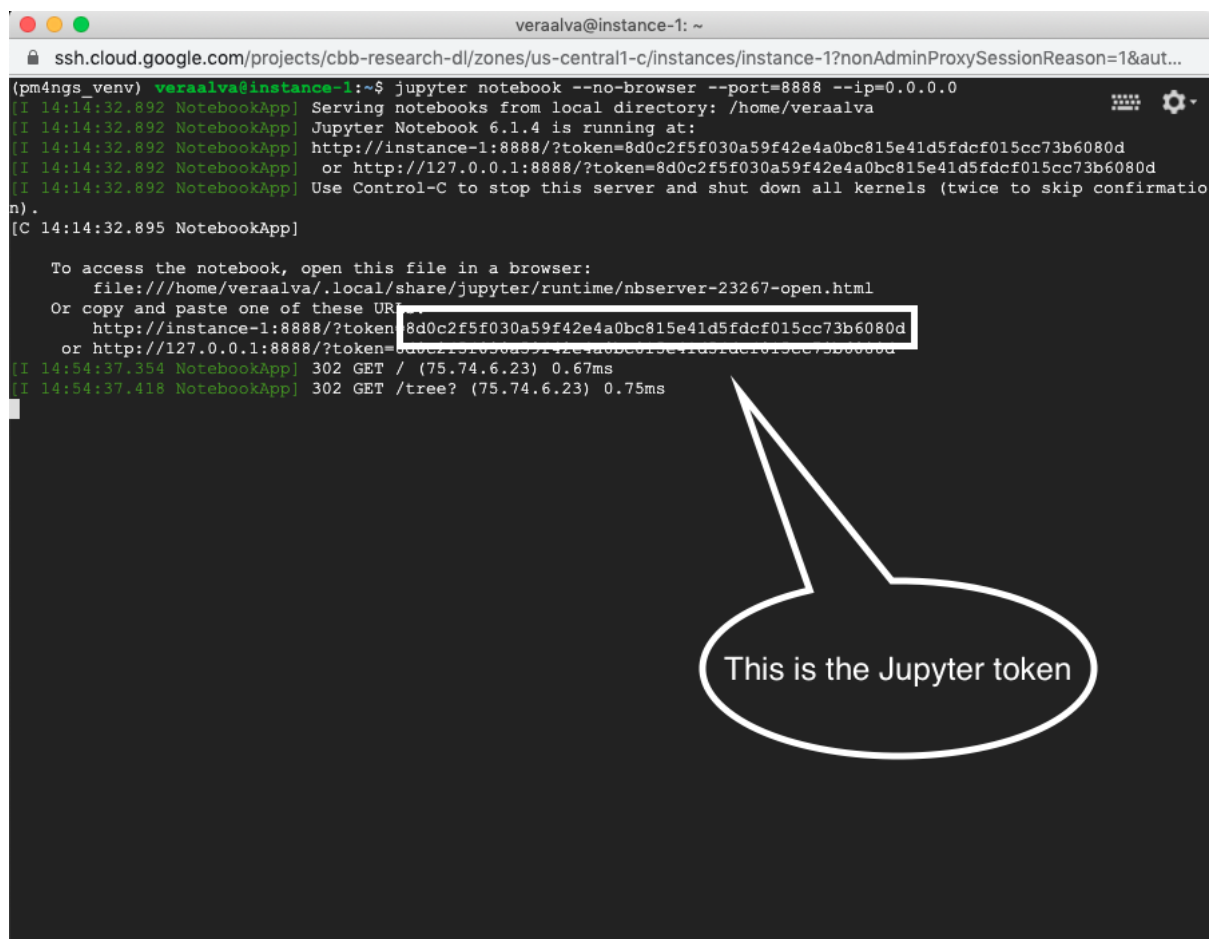
The screenshot shows the Google Cloud Platform console. On the left, the 'VPC network' sidebar is visible. The main area displays 'Firewall policies in this project' and 'Firewall rules in this project'. A search bar at the top right is open, showing 'VM instances' as the selected product. Below the search bar, a list of VM instances is displayed, including 'codeathon-vm-1', 'domains-taxonomy-vm', 'fortigate-1-vm', 'fortigate-2-vm', 'fortigate-5-vm', 'fortigate-6-vm', 'log-z-tensorflow-1-vm', and 'lukas-test-vm'. The 'fortigate-5-vm' instance is highlighted.

Copy the instance public IP to the clipboard

The screenshot shows the Google Cloud Platform console with the 'VM instances' page selected. A table of VM instances is displayed, including 'instance-1' with an internal IP of 10.128.0.21 and an external IP of 34.123.78.105. A red callout bubble points to the 'External IP' column with the text 'Click to copy public IP'. A 'Copy to clipboard' button is visible next to the external IP address.

Copy the URL with localhost in a browser adding :8888 which is the Jupyter server port.

The screenshot shows a web browser with the address bar containing the URL '34.123.78.105:8888'. A red callout bubble points to the address bar with the text 'Paste the IP in a browser, add :8888'. The browser's search bar also shows the URL '34.123.78.105:8888 - Google Search'.

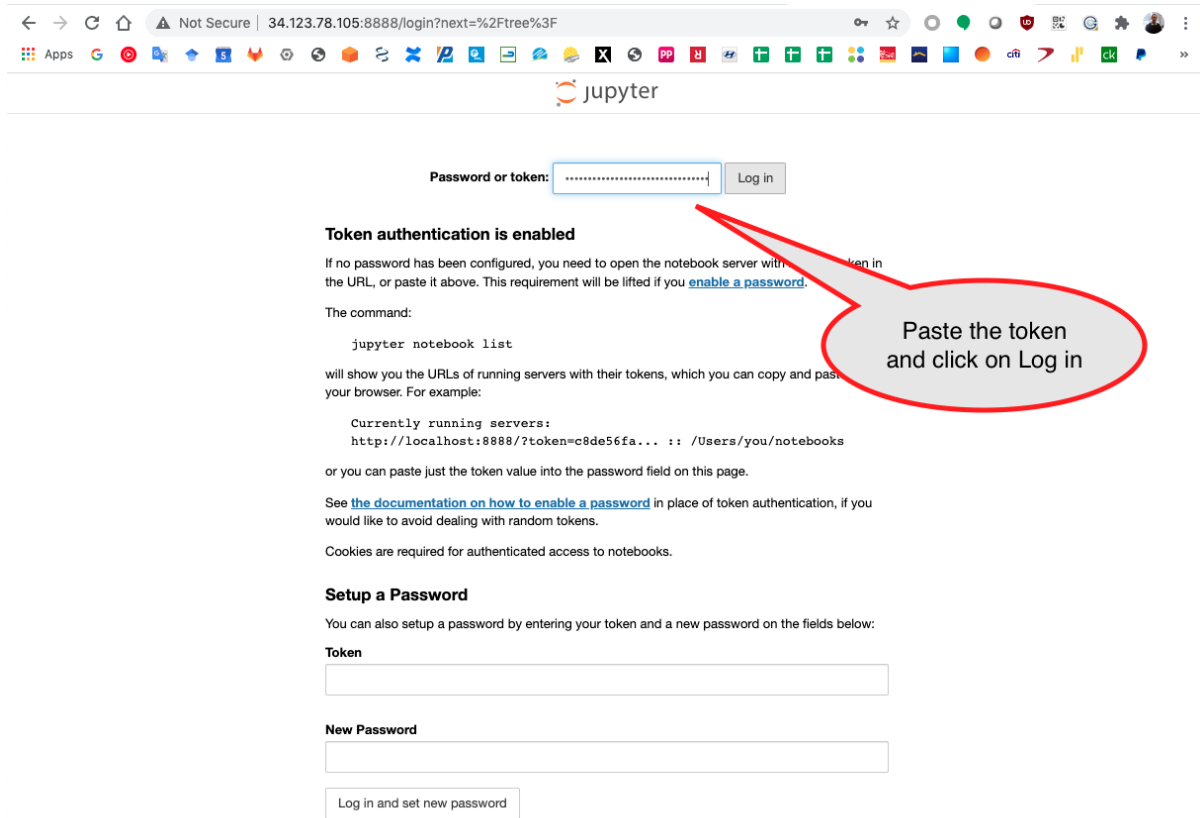
A terminal window titled 'veraalva@instance-1: ~' showing the execution of 'jupyter notebook --no-browser --port=8888 --ip=0.0.0.0'. The output indicates that Jupyter Notebook 6.1.4 is running at 'http://instance-1:8888/?token=8d0c2f5f030a59f42e4a0bc815e41d5fddf015cc73b6080d'. A callout bubble points to the token value, stating 'This is the Jupyter token'.

```
(pm4ngs_venv) veraalva@instance-1:~$ jupyter notebook --no-browser --port=8888 --ip=0.0.0.0
[I 14:14:32.892 NotebookApp] Serving notebooks from local directory: /home/veraalva
[I 14:14:32.892 NotebookApp] Jupyter Notebook 6.1.4 is running at:
[I 14:14:32.892 NotebookApp] http://instance-1:8888/?token=8d0c2f5f030a59f42e4a0bc815e41d5fddf015cc73b6080d
[I 14:14:32.892 NotebookApp] or http://127.0.0.1:8888/?token=8d0c2f5f030a59f42e4a0bc815e41d5fddf015cc73b6080d
[I 14:14:32.892 NotebookApp] Use Control-C to stop this server and shut down all kernels (twice to skip confirmation).

[C 14:14:32.895 NotebookApp]

To access the notebook, open this file in a browser:
    file:///home/veraalva/.local/share/jupyter/runtime/nbserver-23267-open.html
Or copy and paste one of these URLs:
    http://instance-1:8888/?token=8d0c2f5f030a59f42e4a0bc815e41d5fddf015cc73b6080d
    or http://127.0.0.1:8888/?token=8d0c2f5f030a59f42e4a0bc815e41d5fddf015cc73b6080d
[I 14:54:37.354 NotebookApp] 302 GET / (75.74.6.23) 0.67ms
[I 14:54:37.418 NotebookApp] 302 GET /tree? (75.74.6.23) 0.75ms
```

Paste the token in the input bar and click Log in



The screenshot shows a web browser window with the address bar displaying "34.123.78.105:8888/login?next=%2Ftree%3F". The page title is "Jupyter". Below the title, there is a "Password or token:" label followed by a text input field containing a masked password (dots) and a "Log in" button. A red callout bubble with a pointer to the input field contains the text "Paste the token and click on Log in". Below the login form, the text "Token authentication is enabled" is displayed. This is followed by instructions: "If no password has been configured, you need to open the notebook server with the token in the URL, or paste it above. This requirement will be lifted if you [enable a password](#)." The command "jupyter notebook list" is shown, with a note that it will show URLs of running servers with their tokens. An example URL is provided: "http://localhost:8888/?token=c8de56fa... :: /Users/you/notebooks". It is noted that the token value can be pasted into the password field. A link to "the documentation on how to enable a password" is provided. A note states "Cookies are required for authenticated access to notebooks." Below this, the "Setup a Password" section is shown, with instructions to enter a token and a new password. There are two input fields labeled "Token" and "New Password", and a "Log in and set new password" button.

Not Secure | 34.123.78.105:8888/login?next=%2Ftree%3F

Jupyter

Password or token: Log in

Token authentication is enabled

If no password has been configured, you need to open the notebook server with the token in the URL, or paste it above. This requirement will be lifted if you [enable a password](#).

The command:

```
jupyter notebook list
```

will show you the URLs of running servers with their tokens, which you can copy and paste into your browser. For example:

```
Currently running servers:
http://localhost:8888/?token=c8de56fa... :: /Users/you/notebooks
```

or you can paste just the token value into the password field on this page.

See [the documentation on how to enable a password](#) in place of token authentication, if you would like to avoid dealing with random tokens.

Cookies are required for authenticated access to notebooks.

Setup a Password

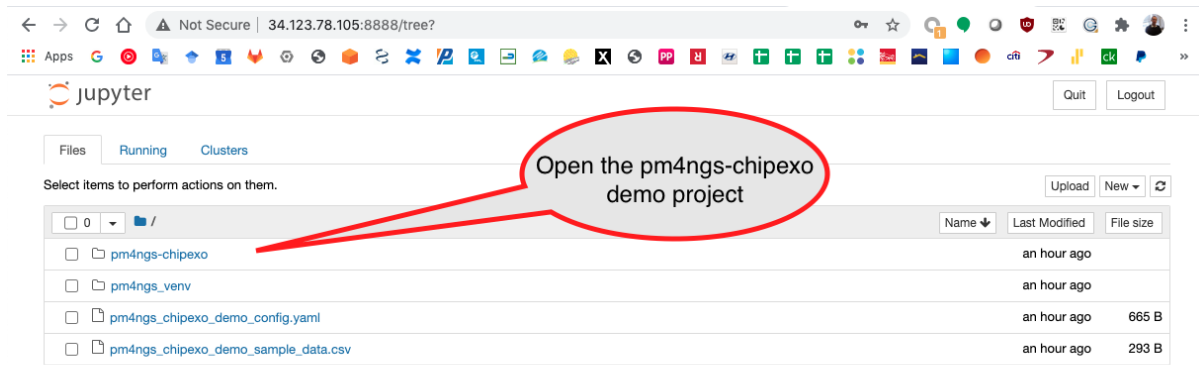
You can also setup a password by entering your token and a new password on the fields below:

Token

New Password

Log in and set new password

Open the **pm4ngs-chipexo** directory

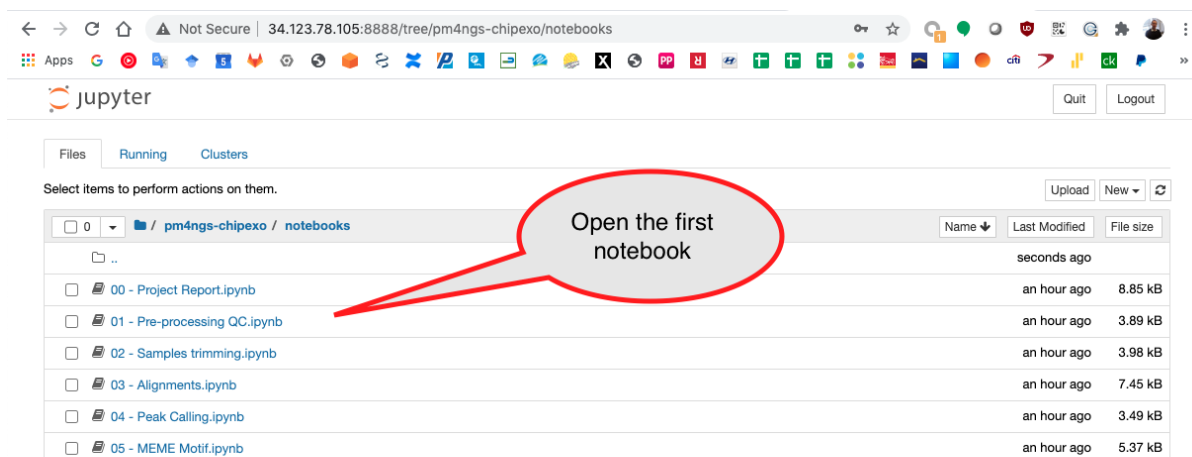


Then, open the **notebooks** directory

The screenshot shows a JupyterLab web interface. At the top, the browser address bar shows '34.123.78.105:8888/tree/pm4ngs-chipexo'. The JupyterLab header includes the 'jupyter' logo and a 'Quit' button. Below the header, the 'Files' tab is active, showing a file explorer for the 'pm4ngs-chipexo' directory. A red box highlights the title 'ChIP-exo PM4NGS organizational structure'. A red speech bubble points to the 'notebooks' directory with the text 'Open the notebooks directory'. The file explorer table lists the following files and directories:

Name	Last Modified	File size
..	seconds ago	
bin	an hour ago	
config	an hour ago	
data	an hour ago	
doc	an hour ago	
notebooks	an hour ago	
requirements	an hour ago	
results	an hour ago	
src	an hour ago	
tmp	an hour ago	
index.html	an hour ago	87 B
LICENSE	an hour ago	961 B
README.md	an hour ago	1.15 kB

Start running the notebook **01 - Pre-processing QC.ipynb**



Open a different VM terminal to run the command **htop** to see the process running. In this case we are seeing multiple **fastq-dump** command being executed.

```

veraalva@instance-1: ~
ssh.cloud.google.com/projects/cbb-research-dl/zones/us-central1-c/instances/instance-1?nonAdminProxySessionReason=1&aut...

1  [ ] 4.6% 5 [ ] 78.1% 9 [ ] 80.8% 13 [ ] 2.6%
2  [ ] 0.7% 6 [ ] 87.4% 10 [ ] 86.8% 14 [ ] 0.0%
3  [ ] 73.5% 7 [ ] 82.9% 11 [ ] 5.3% 15 [ ] 0.0%
4  [ ] 1.3% 8 [ ] 0.0% 12 [ ] 6.0% 16 [ ] 6.0%
Mem [ ] 1.53G/62.8G Tasks: 72, 285 thr; 6 running
Swp [ ] 0K/0K Load average: 3.59 1.16 0.41
Uptime: 01:41:41

  PID USER     PRI  NI  VIRT   RES   SHR  S  CPU% MEM%   TIME+  Command
24998 veraalva  20    0 96196 91728 2900 R 90.5  0.1  1:07.20 /usr/local/ncbi/sra-tools/bin/fastq-dump.2.10.8 --
24999 veraalva  20    0 106M 102M 3044 R 90.5  0.2  1:03.23 /usr/local/ncbi/sra-tools/bin/fastq-dump.2.10.8 --
24997 veraalva  20    0 98484 94312 3056 R 87.9  0.1  1:06.22 /usr/local/ncbi/sra-tools/bin/fastq-dump.2.10.8 --
24995 veraalva  20    0 107M 103M 2948 R 86.6  0.2  1:03.48 /usr/local/ncbi/sra-tools/bin/fastq-dump.2.10.8 --
24996 veraalva  20    0 104M 100M 3096 R 83.9  0.2  0:59.84 /usr/local/ncbi/sra-tools/bin/fastq-dump.2.10.8 --
24994 veraalva  20    0 104M 100M 2956 S 76.7  0.2  1:01.33 /usr/local/ncbi/sra-tools/bin/fastq-dump.2.10.8 --
25010 veraalva  20    0 107M 103M 2948 S 5.3  0.2  0:04.02 /usr/local/ncbi/sra-tools/bin/fastq-dump.2.10.8 --
25015 veraalva  20    0 104M 100M 3096 S 5.3  0.2  0:03.90 /usr/local/ncbi/sra-tools/bin/fastq-dump.2.10.8 --
25011 veraalva  20    0 96196 91728 2900 S 4.6  0.1  0:03.83 /usr/local/ncbi/sra-tools/bin/fastq-dump.2.10.8 --
25014 veraalva  20    0 104M 100M 2956 S 4.0  0.2  0:03.85 /usr/local/ncbi/sra-tools/bin/fastq-dump.2.10.8 --
25008 veraalva  20    0 98484 94312 3056 S 4.0  0.1  0:03.94 /usr/local/ncbi/sra-tools/bin/fastq-dump.2.10.8 --
25009 veraalva  20    0 106M 102M 3044 S 2.0  0.2  0:03.87 /usr/local/ncbi/sra-tools/bin/fastq-dump.2.10.8 --
25188 veraalva  20    0 8652 4324 3168 R 0.7  0.0  0:00.41 htop
4768 root      20    0 2172M 49596 27428 S 0.7  0.1  0:25.74 /usr/bin/containerd
25028 veraalva  20    0 8712 4412 3168 S 0.0  0.0  0:00.79 htop
25054 veraalva  20    0 96196 91728 2900 S 0.0  0.1  0:00.02 /usr/local/ncbi/sra-tools/bin/fastq-dump.2.10.8 --
24359 veraalva  20    0 1253M 60524 35860 S 0.0  0.1  0:00.07 docker run -i --mount=type=bind,source=/home/veraa
4770 root      20    0 2172M 49596 27428 S 0.0  0.1  0:05.53 /usr/bin/containerd
22150 root      20    0 2172M 49596 27428 S 0.0  0.1  0:01.11 /usr/bin/containerd
22355 root      20    0 2172M 49596 27428 S 0.0  0.1  0:00.04 /usr/bin/containerd
23267 veraalva  20    0 236M 60492 15348 S 0.0  0.1  0:01.56 /home/veraalva/pm4ngs_venv/bin/python3 /home/veraa
4777 root      20    0 2172M 49596 27428 S 0.0  0.1  0:01.68 /usr/bin/containerd
22680 root      20    0 2172M 49596 27428 S 0.0  0.1  0:01.16 /usr/bin/containerd
22679 root      20    0 2172M 49596 27428 S 0.0  0.1  0:01.61 /usr/bin/containerd
25187 veraalva  20    0 98484 94312 3056 S 0.0  0.1  0:00.01 /usr/local/ncbi/sra-tools/bin/fastq-dump.2.10.8 --
1334 root      20    0 30068 22720 8828 S 0.0  0.0  0:00.46 /usr/bin/python3 /usr/bin/google_network_daemon
25049 veraalva  20    0 106M 102M 3044 S 0.0  0.2  0:00.01 /usr/local/ncbi/sra-tools/bin/fastq-dump.2.10.8 --
4775 root      20    0 2172M 49596 27428 S 0.0  0.1  0:01.26 /usr/bin/containerd
22526 root      20    0 2499M 98632 49300 S 0.0  0.1  0:00.35 /usr/bin/dockerd -H fd:// --containerd=/run/contai
24734 root      20    0 106M 5296 4684 S 0.0  0.0  0:00.02 containerd-shim -namespace moby -workdir /var/lib/
25175 veraalva  20    0 13916 6020 4544 S 0.0  0.0  0:00.01 sshd: veraalva@pts/3
1 root      20    0 164M 13120 8616 S 0.0  0.0  0:04.62 /sbin/init
F1Help F2Setup F3Search F4Filter F5Tree F6SortBy F7Nice F8Nice F9Kill F10Quit

```

Wait for that process to finish. The log can be checked running the **Checking command output** cell

Checking command output

Execute next cell until it prints: **Run completed**

```
In [5]: check_cwl_command_log(log_file)
```

```
Process no completed.
Please, do not proceed to next cells
Check log file: download.log
```

In the VM terminal you can use the command **tail** to see the process log

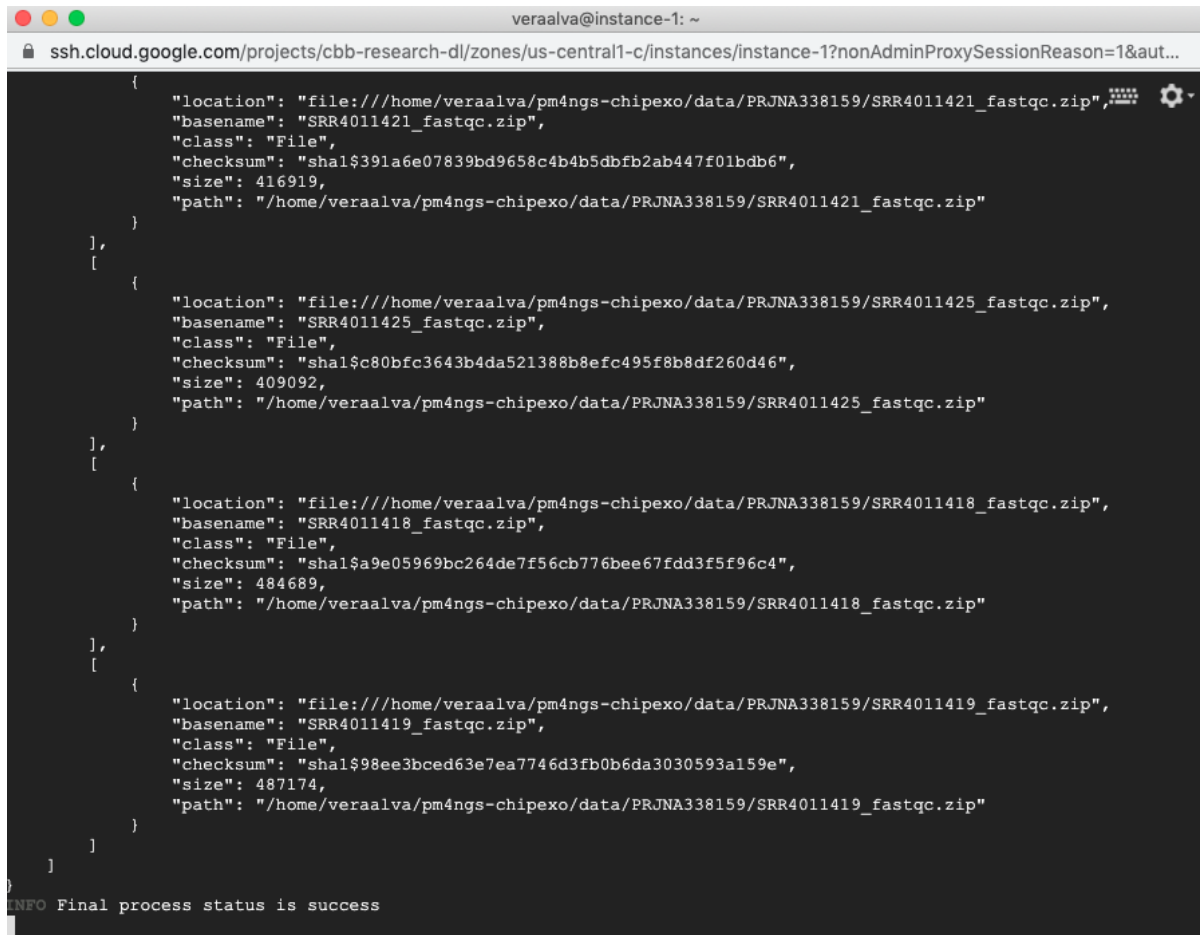
```
(pm4ngs_venv) veraalva@instance-1:~$ tail -f -n 40 pm4ngs-chipexo/data/PRJNA338159/
↪download.log
```

```

veraalva@instance-1: ~
ssh.cloud.google.com/projects/cbb-research-dl/zones/us-central1-c/instances/instance-1?nonAdminProxySessionReason=1&aut...
veraalva@instance-1:~$ tail -f -n 40 pm4ngs-chipexo/data/PRJNA338159/download.log
INFO [job fastq_dump_4] /home/veraalva/pm4ngs-chipexo/tmp/09vsirtc$ docker \
run \
-i \
--mount=type=bind,source=/home/veraalva/pm4ngs-chipexo/tmp/09vsirtc,target=/PnUjLa \
--mount=type=bind,source=/home/veraalva/pm4ngs-chipexo/tmp/871rjshq,target=/tmp \
--mount=type=bind,source=/home/veraalva/.ncbi,target=/PnUjLa/.ncbi,readonly \
--workdir=/PnUjLa \
--user=1001:1002 \
--rm \
--env=TMPDIR=/tmp \
--env=HOME=/PnUjLa \
--cidfile=/home/veraalva/pm4ngs-chipexo/tmp/6e8sdz2/20201007151206-550393.cid \
ncbi/sra-tools:latest \
fastq-dump \
--gzip \
SRR4011425
INFO [job fastq_dump_2] /home/veraalva/pm4ngs-chipexo/tmp/cfbxww59$ docker \
run \
-i \
--mount=type=bind,source=/home/veraalva/pm4ngs-chipexo/tmp/cfbxww59,target=/PnUjLa \
--mount=type=bind,source=/home/veraalva/pm4ngs-chipexo/tmp/9w4635q3,target=/tmp \
--mount=type=bind,source=/home/veraalva/.ncbi,target=/PnUjLa/.ncbi,readonly \
--workdir=/PnUjLa \
--user=1001:1002 \
--rm \
--env=TMPDIR=/tmp \
--env=HOME=/PnUjLa \
--cidfile=/home/veraalva/pm4ngs-chipexo/tmp/0n2rww2z/20201007151206-555879.cid \
ncbi/sra-tools:latest \
fastq-dump \
--gzip \
SRR4011417
Read 2115858 spots for SRR4011421
Written 2115858 spots for SRR4011421
INFO [job fastq_dump_3] Max memory used: 0MiB
INFO [job fastq_dump_3] completed success
Read 13231701 spots for SRR4011425
Written 13231701 spots for SRR4011425
INFO [job fastq_dump_4] Max memory used: 0MiB
INFO [job fastq_dump_4] completed success

```

The process will finish with a message: **Final process status is success**



```

veraalva@instance-1: ~
ssh.cloud.google.com/projects/cbb-research-dl/zones/us-central1-c/instances/instance-1?nonAdminProxySessionReason=1&aut...

{
  "location": "file:///home/veraalva/pm4ngs-chipexo/data/PRJNA338159/SRR4011421_fastqc.zip",
  "basename": "SRR4011421_fastqc.zip",
  "class": "File",
  "checksum": "sha1$391a6e07839bd9658c4b4b5dbfb2ab447f01bdb6",
  "size": 416919,
  "path": "/home/veraalva/pm4ngs-chipexo/data/PRJNA338159/SRR4011421_fastqc.zip"
},
[
  {
    "location": "file:///home/veraalva/pm4ngs-chipexo/data/PRJNA338159/SRR4011425_fastqc.zip",
    "basename": "SRR4011425_fastqc.zip",
    "class": "File",
    "checksum": "sha1$c80bfc3643b4da521388b8efc495f8b8df260d46",
    "size": 409092,
    "path": "/home/veraalva/pm4ngs-chipexo/data/PRJNA338159/SRR4011425_fastqc.zip"
  },
  {
    "location": "file:///home/veraalva/pm4ngs-chipexo/data/PRJNA338159/SRR4011418_fastqc.zip",
    "basename": "SRR4011418_fastqc.zip",
    "class": "File",
    "checksum": "sha1$a9e05969bc264de7f56cb776bee67fdd3f5f96c4",
    "size": 484689,
    "path": "/home/veraalva/pm4ngs-chipexo/data/PRJNA338159/SRR4011418_fastqc.zip"
  },
  {
    "location": "file:///home/veraalva/pm4ngs-chipexo/data/PRJNA338159/SRR4011419_fastqc.zip",
    "basename": "SRR4011419_fastqc.zip",
    "class": "File",
    "checksum": "sha1$98ee3bcd63e7ea7746d3fb0b6da3030593a159e",
    "size": 487174,
    "path": "/home/veraalva/pm4ngs-chipexo/data/PRJNA338159/SRR4011419_fastqc.zip"
  }
]
]
INFO Final process status is success

```

Running the **Checking command output** cell again

Checking command output
 Execute next cell until it prints: **Run completed**

```

In [6]: check_cwl_command_log(log_file)
Process completed.
Proceed to next cells

```

Finish the **01 - Pre-processing QC.ipynb** notebook and go to the project report **00 - Project Report.ipynb**. Execute the first and second cell to visualize the Pre-processing report table.

← → ↻ ⌂ Not Secure | 34.123.78.105:8888/notebooks/pm4ngs-chipexo/notebooks/00%20-%20Project%20R... ☆

Apps

Jupyter 00 - Project Report (unsaved changes) Logout

File Edit View Insert Cell Kernel Widgets Help Not Trusted Python 3

Run

Out[2]: The raw code for this IPython notebook is by default hidden for easier reading. To toggle on/off the raw code, click [here](#).

1. Pre-processing QC

Info

Type	Link
Notebook	01 - Pre-processing QC
Results	PRJNA338159

FastQC report

Sample	Fastq	FastQC Report	No of Reads in fastq	Seq Len	%GC	Poor Quality	Fail Tests
SRR4011416	1295.96MB	0.63MB	16,846,671	101	46	0	Per base sequence content Sequence Duplication Levels Adapter Content
SRR4011417	1362.60MB	0.63MB	17,268,284	101	45	0	Sequence Duplication Levels Adapter Content
SRR4011418	1469.33MB	0.63MB	18,868,657	101	46	0	Per base sequence content Sequence Duplication Levels Adapter Content
SRR4011419	1288.91MB	0.63MB	16,423,038	101	46	0	Per base sequence content Sequence Duplication Levels Adapter Content
SRR4011421	99.40MB	0.59MB	2,115,858	51	45	0	Per base sequence content Sequence Duplication Levels
SRR4011425	619.82MB	0.59MB	13,231,701	51	45	0	Per base sequence content Sequence Duplication Levels

Follow this procedure for each notebook in the project.

CHAPTER 5

Sample sheet

A CSV file describing samples, conditions and replicate number is required during the PM4NGS project creation.

PM4NGS will copy the sample sheet file to the folder **data/{{dataset_name}}** with the standard name **sample_table.csv**.

5.1 Single-end example

This is an sample sheet example for **single-end** sequencing technology data:

sample_name	file	condition	replicate
SRR4011416	net/rawdata/SRR4011416.fastq.gz	Exp_O2_growth_no_rifampicin	1
SRR4011417	net/rawdata/SRR4011417.fastq.gz	Exp_O2_growth_no_rifampicin	2
SRR4011421	net/rawdata/SRR4011421.fastq.gz	Exp_O2_growth_rifampicin	1
SRR4011425	net/rawdata/SRR4011425.fastq.gz	Exp_O2_growth_rifampicin	2
SRR4011418	net/rawdata/SRR4011418.fastq.gz	Stat_O2_growth_no_rifampicin	1
SRR4011419	net/rawdata/SRR4011419.fastq.gz	Stat_O2_growth_no_rifampicin	2

Table source: [sample_sheet_single_end.csv](#)

5.2 Paired-end example

For **paired-end** sequencing technology data, the **l** should be used to separate forward and reverse fastq files:

sample_name	file	condition	replicate
SRR2126784	http://myserver.net/SRR2126784_1.fastq.gz http://myserver.net/SRR2126784_2.fastq.gz	PAIRED NACT	1
SRR2126785	http://myserver.net/SRR2126785_1.fastq.gz http://myserver.net/SRR2126785_2.fastq.gz	PAIRED NACT	1
SRR2126786	http://myserver.net/SRR2126786_1.fastq.gz http://myserver.net/SRR2126786_2.fastq.gz	PAIRED NACT	1
SRR2126787	http://myserver.net/SRR2126787_1.fastq.gz http://myserver.net/SRR2126787_2.fastq.gz	PAIRED NACT	1
SRR3383790	http://myserver.net/SRR3383790_1.fastq.gz http://myserver.net/SRR3383790_2.fastq.gz	PAIRED NACT	1

Source: `sample_sheet_paired_end.csv`

PM4NGS will copy or download the raw fastq files to the `data/{{dataset_name}}/` directory during the project creation if the `[-copy-rawdata]` is used.

5.3 Processing data from the NCBI SRA

For data in the NCBI SRA database, the `file` column should be empty. PM4NGS will download the files during the pre-processing quality control step.

sample_name	file	condition	replicate
SRR7549105		ES	1
SRR7549106		ES	2
SRR7549109		MEF	1
SRR7549110		MEF	2
SRR7549114		rB	1
SRR7549113		rB	2

Source: `sample_sheet_sra.csv`

5.4 Sample sheet column names and description

Note: Columns names are required and are case sensitive.

Columns

- **sample_name:** Sample names. It can be different of sample file name.
- **file:** This is the absolute path or URL to the raw fastq file.

For paired-end data the files should be separated using the unix pipe `|` as **SRR4053795_1.fastq.gz|SRR4053795_2.fastq.gz** must exist.

The data files will be copied to the folder `data/{{dataset_name}}/`.

- **condition:** Conditions to group the samples. Use only alphanumeric characters.

For RNASeq projects the differential gene expression will be generated comparing these conditions. If there are multiple conditions all comparisons will be generated. It must be at least two conditions.

For ChIPSeq projects differential binding events will be detected comparing these conditions. If there are multiple conditions all comparisons will be generated. It must be at least two conditions.

For ChIPexo projects the samples of the same condition will be grouped for the peak calling with MACE.

- **replicate:** Replicate number for samples.

Differential Gene expression and GO enrichment from RNA-Seq data

Differential expression (DE) analysis allows the comparison of RNA expression levels between multiple conditions.

The differential gene expression and GO enrichment pipeline is comprised of five steps, shown in next Table. The first step involves downloading samples from the NCBI SRA database, if necessary, or executing the pre-processing quality control tools on local samples. Subsequently, sample trimming, alignment, and quantification processes are executed. Once all samples are processed, groups of differentially expressed genes are identified per condition, using DESeq and EdgeR. Over- and under-expressed genes are reported by each program, and the interception of their results is computed.

Finally, once differentially expressed genes are identified, a GO enrichment analysis is executed to provide key biological processes, molecular functions, and cellular components for identified genes.

Table 1: DGA and Go enrichment, RNASeq pipeline

Step	Jupyter Notebook	Workflow CWL	Tool	Input	Output	CWL Tool
Sample Download and Quality Control	01 - Pre-processing QC	download_quality_SRA-Tools.cwl	SRA-Tools	SRA accession	Fastq	fastq-dump.cwl
			FastQC	Fastq	FastQC HTML and Zip	fastqc.cwl
Trimming	02 - Samples trimming		Trimmomatic	Fastq	Fastq	trimmomatic-PE.cwl trimmomatic-SE.cwl
Alignments and Quantification	03 - Alignments and Quantification	rnaseq-alignment-quantification.cwl	STAR	Fastq	BAM	star.cwl
			Samtools	SAM	BAM Sorted BAM BAM index BAM stats BAM flagstats	samtools-view.cwl samtools-sort.cwl samtools-index.cwl samtools-stats.cwl samtools-flagstat.cwl
			IGVtools	Sorted BAM	TDF	igvtools-count.cw
			RSeQC	Sorted BAM	Alignment quality control (TXT and PDF)	rseqc-bam_stat.cwl rseqc-infer_experiment.cwl rseqc-junction_annotation.cwl rseqc-junction_saturation.cwl rseqc-read_distribution.cwl rseqc-read_quality.cwl
			TPMCalculator	Sorted BAM	Read counts (TSV) TPM values (TSV)	tpmcalculator.cwl
Differential Gene Analysis	04 - DGA		DeSeq2	Read Matrix	TSV	deseq2-2conditions.cwl
			EdgeR	Read Matrix	TSV	edgeR-2conditions.cwl
Go enrichment	05 - GO enrichment		goenrichment	gene IDs	TSV	Python code in the notebook

6.1 Input requirements

The input requirement for the DGA pipeline is the [Sample sheet](#) file.

6.2 Pipeline command line

The RNASeq based project can be created using the following command line:

```
localhost:~> pm4ngs-rnaseq
usage: Generate a PM4NGS project for RNA-Seq data analysis [-h] [-v]
                                     --sample-sheet
                                     SAMPLE_SHEET
                                     [--config-file CONFIG_FILE]
                                     [--copy-rawdata]
```

Options:

- **sample-sheet:** Sample sheet with the samples metadata
- **config-file:** YAML file with configuration for project creation
- **copy-rawdata:** Copy the raw data defined in the sample sheet to the project structure. (The data can be hosted locally or in an http server)

6.3 Creating the DGA and GO enrichment from RNA-Seq data project

The **pm4ngs-rnaseq** command line executed with the **--sample-sheet** option will let you type the different variables required for creating and configuring the project. The default value for each variable is shown in the brackets. After all questions are answered, the CWL workflow files will be cloned from the github repo [ncbi/cwl-ngs-workflows-cbb](https://github.com/ncbi/cwl-ngs-workflows-cbb) to the folder **bin/cwl**.

```
localhost:~> pm4ngs-rnaseq --sample-sheet my-sample-sheet.tsv
Generating RNA-Seq data analysis project
author_name [Roberto Vera Alvarez]:
email [veraalva@ncbi.nlm.nih.gov]:
project_name [my_ngs_project]:
dataset_name [my_dataset_name]:
is_data_in_SRA [y]:
Select sequencing_technology:
1 - single-end
2 - paired-end
Choose from 1, 2 [1]: 2
create_demo [y]: y
number_spots [1000000]:
organism [human]:
genome_name [hg38]:
genome_dir [hg38]:
aligner_index_dir [hg38/STAR]:
genome_fasta [hg38/genome.fa]:
genome_gtf [hg38/genes.gtf]:
genome_bed [hg38/genes.bed]:
fold_change [2.0]:
fdr [0.05]:
use_docker [y]:
max_number_threads [16]:
Cloning Git repo: https://github.com/ncbi/cwl-ngs-workflows-cbb to /home/veraalva/tmp/
↪ cookiecutter/my_ngs_project/bin/cwl
```

(continues on next page)

(continued from previous page)

```

Updating CWLs dockerPull and SoftwareRequirement from: /home/veraalva/my_ngs_project/
↪requirements/conda-env-dependencies.yaml
bioconductor-diffbind with version 2.16.0 update image to: quay.io/biocontainers/
↪bioconductor-diffbind:2.16.0--r40h5f743cb_0
  /home/veraalva/my_ngs_project/bin/cwl/tools/R/deseq2-pca.cwl with old image_
↪replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
  /home/veraalva/my_ngs_project/bin/cwl/tools/R/macs-cutoff.cwl with old image_
↪replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
  /home/veraalva/my_ngs_project/bin/cwl/tools/R/dga_heatmaps.cwl with old image_
↪replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
  /home/veraalva/my_ngs_project/bin/cwl/tools/R/diffbind.cwl with old image_
↪replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
  /home/veraalva/my_ngs_project/bin/cwl/tools/R/edgeR-2conditions.cwl with old_
↪image replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
  /home/veraalva/my_ngs_project/bin/cwl/tools/R/volcano_plot.cwl with old image_
↪replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
  /home/veraalva/my_ngs_project/bin/cwl/tools/R/readQC.cwl with old image replaced:_
↪quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
  /home/veraalva/my_ngs_project/bin/cwl/tools/R/deseq2-2conditions.cwl with old_
↪image replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
Copying file /home/veraalva/Work/Developer/Python/pm4ngs/pm4ngs-rnaseq/example/pm4ngs_
↪rnaseq_demo_sample_data.csv to /Users/veraalva/my_ngs_project/data/my_dataset_name/
↪sample_table.csv
20 files loaded
Using table:
  sample_name file          condition replicate
0   SRR2126784          PRE_NACT             1
1   SRR2126785          PRE_NACT             1
2   SRR2126786          PRE_NACT             1
3   SRR2126787          PRE_NACT             1
4   SRR3383790          PRE_NACT             1
5   SRR3383791          PRE_NACT             1
6   SRR3383792          PRE_NACT             1
7   SRR3383794          PRE_NACT             1
8   SRR3383796          PRE_NACT             1
9   SRR3383797          PRE_NACT             1
10  SRR3383798          PRE_NACT             1
11  SRR2126788      POST_NACT_CRIS3             1
12  SRR2126789      POST_NACT_CRIS3             1
13  SRR2126790      POST_NACT_CRIS3             1
14  SRR2126791      POST_NACT_CRIS3             1
15  SRR2126792      POST_NACT_CRIS3             1
16  SRR2126793      POST_NACT_CRIS3             1
17  SRR2126794      POST_NACT_CRIS3             1
18  SRR2126795      POST_NACT_CRIS3             1
19  SRR2126796      POST_NACT_CRIS3             1
Done

```

The **pm4ngs-rnaseq** command line will create a project structure as:

```

.
├── LICENSE
├── README.md
├── bin
│   └── cwl
├── config
│   └── init.py

```

(continues on next page)

(continued from previous page)

```

├── data
│   └── my_dataset_name
│       └── sample_table.csv
├── doc
├── index.html
├── notebooks
│   ├── 00 - Project Report.ipynb
│   ├── 01 - Pre-processing QC.ipynb
│   ├── 02 - Samples trimming.ipynb
│   ├── 03 - Alignments and Quantification.ipynb
│   ├── 04 - DGA.ipynb
│   └── 05 - GO enrichment.ipynb
├── requirements
│   └── conda-env-dependencies.yaml
├── results
│   └── my_dataset_name
├── src
└── tmp

61 directories, 239 files

```

Note: RNASeq based project variables

- **author_name:** Default: [Roberto Vera Alvarez]
- **email:** Default: [veraalva@ncbi.nlm.nih.gov]
- **project_name:** Name of the project with no space nor especial characters. This will be used as project folder's name.
Default: [my_ngs_project]
- **dataset_name:** Dataset to process name with no space nor especial characters. This will be used as folder name to group the data. This folder will be created under the **data/{{dataset_name}}** and **results/{{dataset_name}}**.
Default: [my_dataset_name]
- **is_data_in_SRA:** If the data is in the SRA set this to y. A CWL workflow to download the data from the SRA database to the folder **data/{{dataset_name}}** and execute FastQC on it will be included in the **01 - Pre-processing QC.ipynb** notebook.
Set this option to **n**, if the fastq files names and location are included in the sample sheet.
Default: [y]
- **Select sequencing technology:** Select one of the available sequencing technologies in your data.
Values: 1 - single-end, 2 - paired-end
- **create_demo:** If the data is downloaded from the SRA and this option is set to y, only the number of spots specified in the next variable will be downloaded. Useful to test the workflow.
Default: [y]: y
- **number_spots:** Number of sport to download from the SRA database. It is ignored is the **create_demo** is set to **n**.
Default: [1000000]
- **organism:** Organism to process, e.g. human. This is used to link the selected genes to the NCBI gene database.

Default: [human]

- **genome_name:** Genome name , e.g. hg38 or mm10.

Default: [hg38]

- **genome_dir:** Absolute path to the directory with the genome annotation (genome.fa, genes.gtf) to be used by the workflow or the name of the genome.

If the name of the genome is used, PM4NGS will include a cell in the **03 - Alignments and Quantification.ipynb** notebook to download the genome files. The genome data will be at **data/{{dataset_name}}/{{genome_name}}/**

Default: [hg38]

- **aligner_index_dir:** Absolute path to the directory with the genome indexes for STAR.

If **{{genome_name}}/STAR** is used, PM4NGS will include a cell in the **03 - Alignments and Quantification.ipynb** notebook to create the genome indexes for STAR.

Default: [hg38/STAR]

- **genome_fasta:** Absolute path to the genome fasta file

If **{{genome_name}}/genome.fa** is used, PM4NGS will use the downloaded fasta file.

Default: [hg38/genome.fa]

- **genome_gtf:** Absolute path to the genome GTF file

If **{{genome_name}}/genes.gtf** is used, PM4NGS will use the downloaded GTF file.

Default: [hg38/gene.gtf]

- **genome_bed:** Absolute path to the genome BED file

If **{{genome_name}}/genes.bed** is used, PM4NGS will use the downloaded BED file.

Default: [hg38/genes.bed]

- **fold_change:** A real number used as fold change cutoff value for the DG analysis, e.g. 2.0.

Default: [2.0]

- **fdr:** Adjusted P-Value to be used as cutoff in the DG analysis, e.g. 0.05.

Default: [0.05]

- **use_docker:** Set this to y if you will be using Docker. Otherwise Conda needs to be installed in the computer.

Default: [y]

- **max_number_threads:** Number of threads available in the computer.

Default: [16]

6.4 Jupyter server

PM4NGS uses Jupyter as interface for users. After project creation the jupyter server should be started as shown below. The server will open a browser windows showing the project's structure just created.

```
localhost:~> jupyter notebook
```

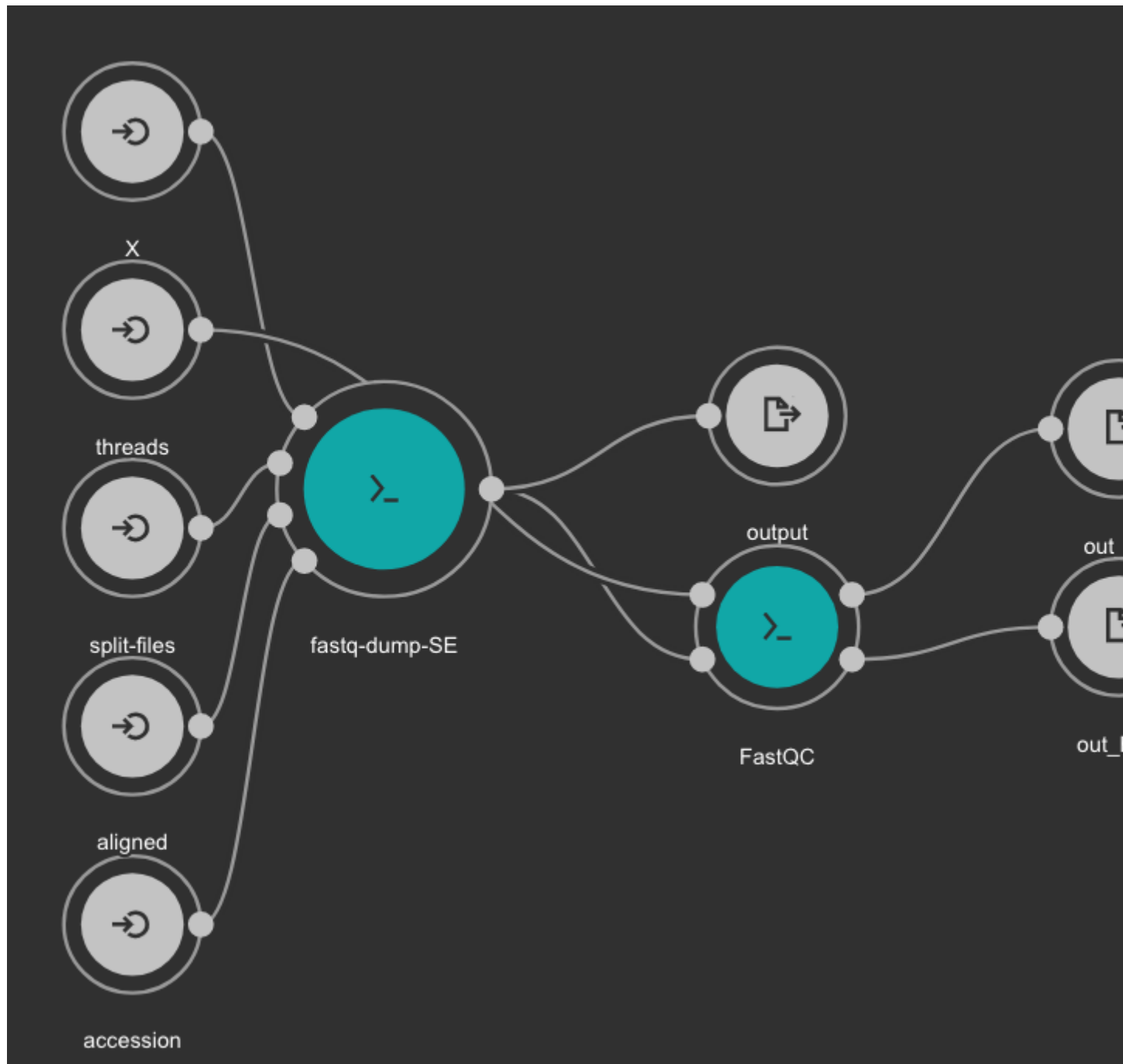
6.5 Data processing

Start executing the notebooks from 01 to 05 waiting for each step completion. The **00 - Project Report.ipynb** notebook can be executed after each notebooks to see the progress in the analysis.

6.6 CWL workflows

6.6.1 SRA download and QC workflow

This CWL workflow is designed to download FASTQ files from the NCBI SRA database using `fastq-dump` and then, execute `fastqc` generating a quality control report of the sample.



Inputs

- **accession:** SRA accession ID. Type: string. Required.
- **aligned:** Used it to download only aligned reads. Type: boolean. Optional.
- **split-files:** Dump each read into separate file. Files will receive suffix corresponding to read number. Type: boolean. Optional.
- **threads:** Number of threads. Type: int. Default: 1. Optional.

- **X**: Maximum spot id. Optional.

Outputs

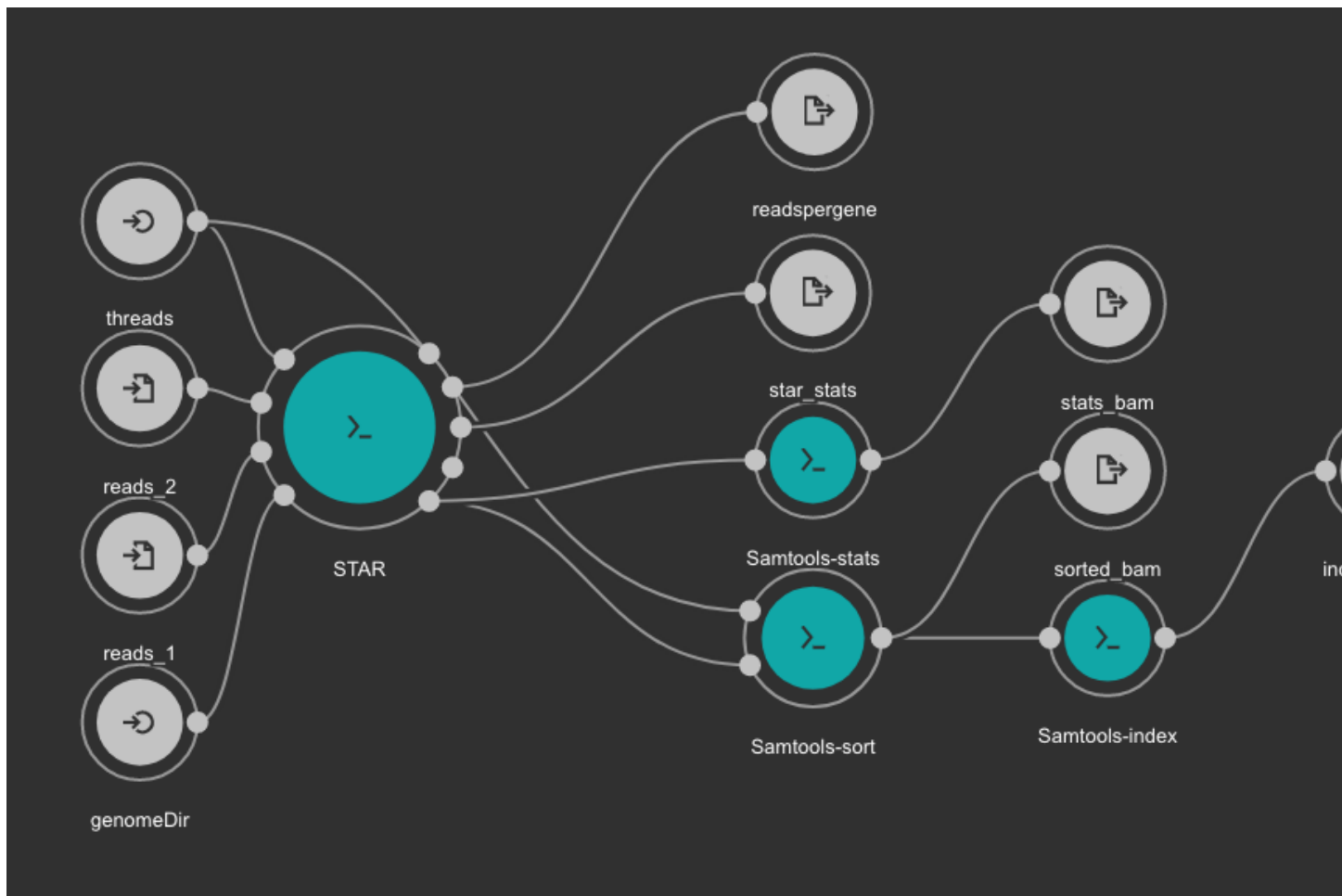
- **output**: Fastq files downloaded. Type: File[]
- **out_zip**: FastQC report ZIP file. Type: File[]
- **out_html**: FastQC report HTML. Type: File[]

6.6.2 Samples Trimming

Our workflows uses [Trimmomatic](#) for read trimming. The Jupyter notebooks uses some basic Trimmomatic options that need to be modified depending on the FastQC quality control report generated for the sample.

6.6.3 STAR based alignment and sorting

This workflows use [STAR](#) for alignning RNA-Seq reads to a genome. The obtained BAM file is sorted using [SAMtools](#). Statistics outputs from STAR and SAMtools are returned as well.



Inputs

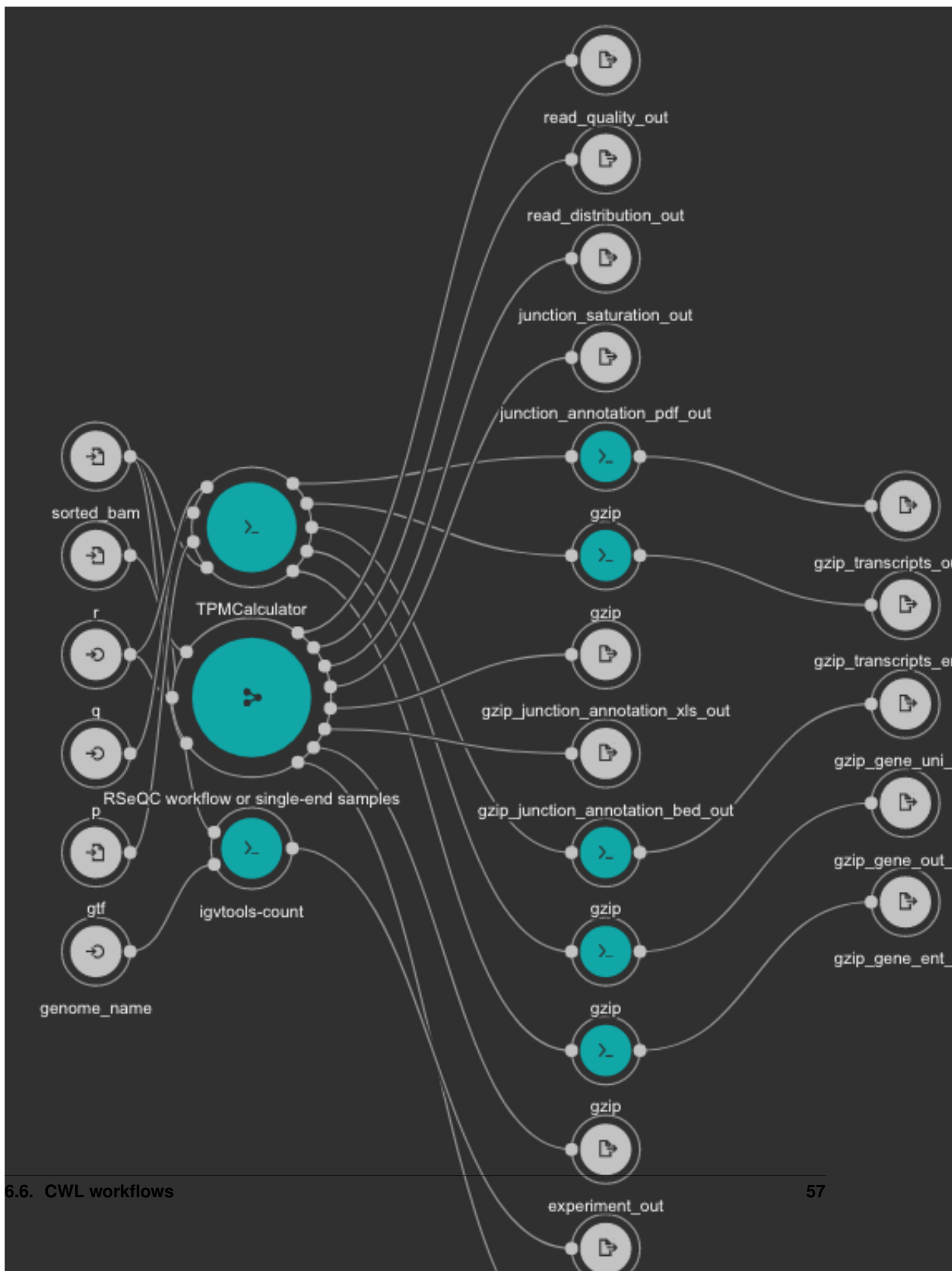
- **genomeDir**: Aligner indexes directory. Type: Directory. Required. Variable ALIGNER_INDEX in the Jupyter Notebooks.
- **threads**: Number of threads. Type: int. Default: 1. Optional.
- **reads_1**: FastQ file to be processed for paired-end reads _1. Type: File. Required.
- **reads_2**: FastQ file to be processed for paired-end reads _2. Type: File. Required.

Outputs

- **sorted_bam**: Final BAM file filtered and sorted. Type: File.
- **indexed_bam**: BAM index file. Type: File.
- **star_stats**: STAR alignment statistics. Type: File.
- **readspergene**: STAR reads per gene output. Type: File.
- **stats_bam**: SAMtools stats output. Type: File.

6.6.4 RNA-Seq quantification and QC workflow using TPMCalculator

This workflow uses [TPMCalculator](#) to quantify the abundance of genes and transcripts from the sorted BAM file. Additionally, [RSeQC](#) is executed to generate multiple quality control outputs from the sorted BAM file. At the end, a TDF file is generated using [igvtools](#) from the BAM file for a quick visualization.



Inputs

- **gtf**: Genome GTF file. Variable GENOME_GTF in the Jupyter Notebooks. Type: File. Required.
- **genome_name**: Genome name as defined in IGV for TDF conversion. Type: string. Required.
- **q**: Minimum MAPQ value to use reads. We recommend 255. Type: int. Required.
- **r**: Reference Genome in BED format used by RSeQC. Variable GENOME_BED in the Jupyter Notebooks. Type: File. Required.
- **sorted_bam**: Sorted BAM file to quantify. Type: File. Required.

Outputs

- **bam_to_tdf_out**: TDF file created with igvtools from the BAM file for quick visualization. Type: File.
- **gzip_gene_ent_out**: TPMCalculator gene ENT output gzipped. Type: File.
- **gzip_gene_out_out**: TPMCalculator gene OUT output gzipped. Type: File.
- **gzip_gene_uni_out**: TPMCalculator gene UNI output gzipped. Type: File.
- **gzip_transcripts_ent_out**: TPMCalculator transcript ENT output gzipped. Type: File.
- **gzip_transcripts_out_out**: TPMCalculator transcript OUT output gzipped. Type: File.
- **bam_stat_out**: RSeQC BAM stats output. Type: File.
- **experiment_out**: RSeQC experiment output. Type: File.
- **gzip_junction_annotation_bed_out**: RSeQC junction annotation bed. Type: File.
- **gzip_junction_annotation_xls_out**: RSeQC junction annotation xls. Type: File.
- **junction_annotation_pdf_out**: RSeQC junction annotation PDF figure. Type: File.
- **junction_saturation_out**: RSeQC junction saturation output. Type: File.
- **read_distribution_out**: RSeQC read distribution output. Type: File.
- **read_quality_out**: RSeQC read quality output. Type: File.

6.6.5 Differential Gene Expression analysis from RNA-Seq data

Our notebooks are designed to execute a Differential Gene Expression analysis using two available tools: [DESeq2](#) and [EdgeR](#). Also, the results for the interception of both tools output is reported with volcano plots, heatmaps and PCA plots.

The workflow use the **sample_sheet.tsv** file to generate an array with all combinations of **conditions**. The code to generate this array is very simple and can be found in the cell number 3 in the **05 - DGA.ipynb** notebook.

```
comparisons = []
for s in itertools.combinations(factors['condition'].unique(), 2):
    comparisons.append(list(s))
```

Let's suppose we have a **factors.txt** file with three conditions: **cond1**, **cond2** and **cond3**. The **comparisons** array will look like:

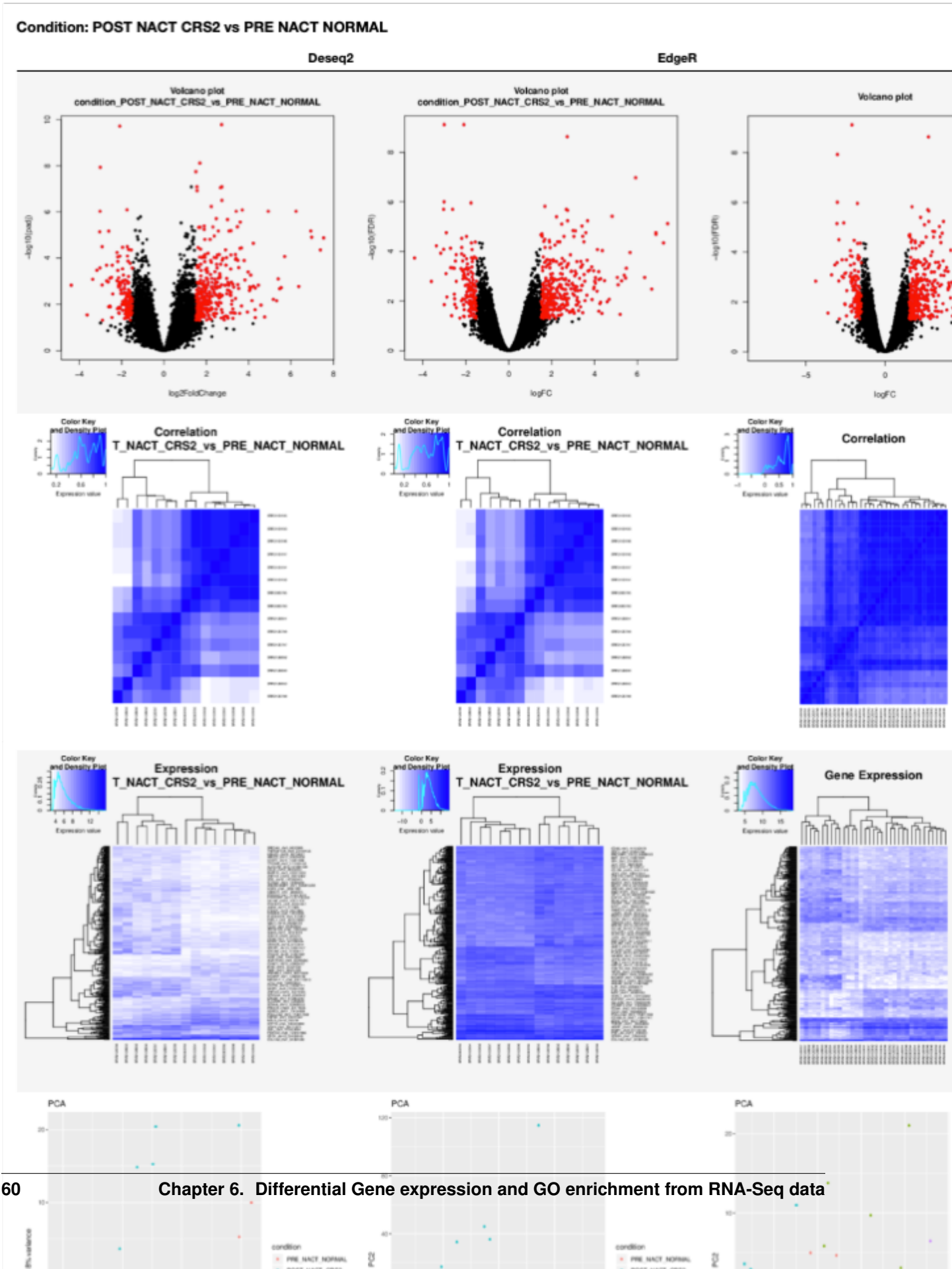
```
comparisons = [  
    ['cond1', 'cond2'],  
    ['cond1', 'cond3'],  
    ['cond2', 'cond3']  
]
```

To avoid this behavior and execute the comparison just in a set of conditions, you should remove the code in the cell number 3 in the **05 - DGA.ipynb** notebook and manually create the array of combinations to be compared as:

```
comparisons = [  
    ['cond1', 'cond3'],  
]
```

The R code used for running DESeq2 is embedded in [deseq2-2conditions.cwl](#) from line 25 to line 169. The R code used for running EdgeR is embedded in [edgeR-2conditions.cwl](#) from line 25 to line 159.

A table with DGA plots is generated for each condition in the **00 - Project Report.ipynb** as shown next.



6.6.6 GO enrichment from RNA-Seq data

The GO enrichment analysis is executed with an *in-house* developed python package named `goenrichment`. This tool uses the `hypergeometric distribution` test to estimate the probability of successes in selecting GO terms from a list of differentially expressed genes. The GO terms are represented as a network using the python library `NetworkX`.

The tool uses a pre-computed database, currently available for `human` and `mouse`, at <https://ftp.ncbi.nlm.nih.gov/pub/goenrichment/>. However, the project web page describe how to create your own database from a set of reference databases.

The workflow uses the `factors.txt` file to generate an array with all combinations of `conditions`. The code to generate this array is very simple and can be found in the cell number 3 in the **06 - GO enrichment.ipynb** notebook.

```
comparisons = []
for s in itertools.combinations(factors['condition'].unique(), 2):
    comparisons.append(list(s))
```

Let's suppose we have a `factors.txt` file with three conditions: `cond1`, `cond2` and `cond3`. The `comparisons` array will look like:

```
comparisons = [
    ['cond1', 'cond2'],
    ['cond1', 'cond3'],
    ['cond2', 'cond3']
]
```

To avoid this behavior and execute the comparison just in a set of conditions, you should remove the code in the cell number 3 in the **06 - GO enrichment.ipynb** notebook and manually create the array of combinations to be compared as:

```
comparisons = [
    ['cond1', 'cond3'],
]
```

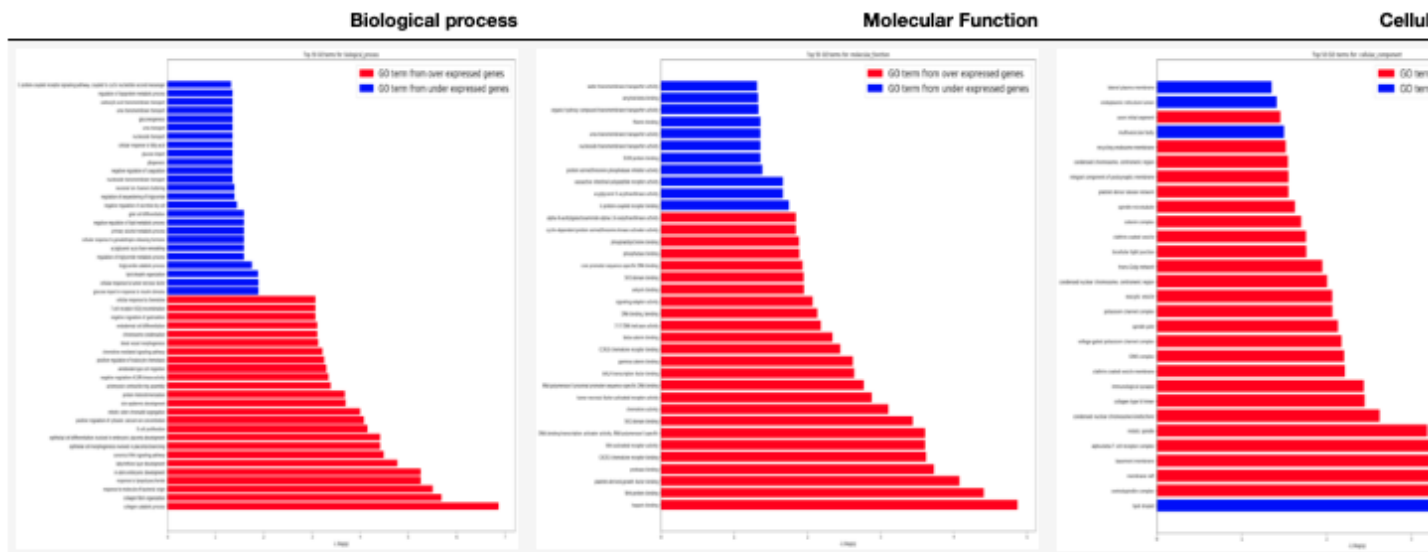
Additionally, the workflow requires three cutoff that are defined in the cell number 5 of the same notebook.

```
min_category_depth=4
min_category_size=3
max_category_size=500
```

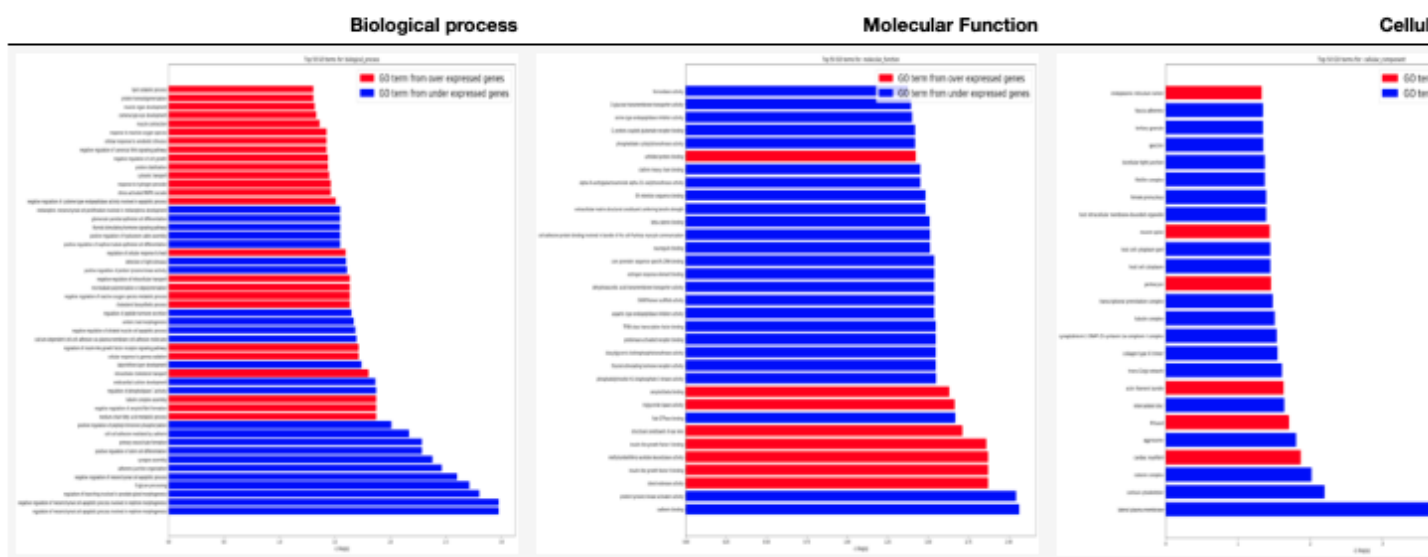
Cutoffs definition

- `min_category_depth`: Min GO term graph depth to include in the report. Default: 4
- `min_category_size`: Min number of gene in a GO term to include in the report. Default: 3
- `max_category_size`: Max number of gene in a GO term to include in the report. Default: 500

A table with GO terms plots is generated for each condition in the **00 - Project Report.ipynb** as shown next. In these plots the red bars are for GO terms selected from the over expressed genes and the blue bars are for GO terms selected from the under expressed genes. It is important to clarify that the two sets of GO terms don't overlap each other.



Condition: POST_NACT_CRIS3 vs POST_NACT_CRIS2



6.7 Demo

PM4NGS includes a demo project that users can use to test the framework. It is pre-configured to use Docker as execution environment.

The RNASeq based demo process samples from the BioProject [PRJNA290924](https://www.ncbi.nlm.nih.gov/bioproject/PRJNA290924).

Use this command to create the project structure in your local computer

```
localhost:~> pm4ngs-rnaseq-demo
```

Once it finish, start the jupyter server and execute the notebooks as it is indicated on them


```
localhost:~> jupyter notebook
[I 14:12:52.956 NotebookApp] Serving notebooks from local directory: /home/veraalva
[I 14:12:52.956 NotebookApp] Jupyter Notebook 6.1.4 is running at:
[I 14:12:52.956 NotebookApp] http://localhost:8888/?
↪token=eae6a8d42ad12d6ace23f5d0923bcec14d0f798127750122
[I 14:12:52.956 NotebookApp] or http://127.0.0.1:8888/?
↪token=eae6a8d42ad12d6ace23f5d0923bcec14d0f798127750122
[I 14:12:52.956 NotebookApp] Use Control-C to stop this server and shut down all
↪kernels (twice to skip confirmatio
n).
[C 14:12:52.959 NotebookApp]

To access the notebook, open this file in a browser:
    file:///home/veraalva/.local/share/jupyter/runtime/nbserver-23251-open.html
Or copy and paste one of these URLs:
    http://localhost:8888/?token=eae6a8d42ad12d6ace23f5d0923bcec14d0f798127750122
    or http://127.0.0.1:8888/?token=eae6a8d42ad12d6ace23f5d0923bcec14d0f798127750122
```

The results of this analysis is <https://ftp.ncbi.nlm.nih.gov/pub/pm4ngs/examples/rnaseq-sra-paired/>

Differential Binding detection from ChIP-Seq data

Differential binding analysis pipeline from ChIP-Seq data.

The workflow is comprised of five steps, as shown in next table. The first step, sample download and quality control, is for downloading samples from the NCBI SRA database, if necessary, or for executing the pre-processing quality control tools on all samples. Sample trimming and alignment are executed next. Post-processing quality control on ChIP-Seq samples is executed using Phantompeakqualtools, as recommended by the ENCODE consortia]. Then, a peak calling step is performed using MACS2, and peak annotation is executed with Homer. Peak reproducibility is explored with IDR. Finally, a differential binding analysis is completed with DiffBind.

Table 1: Differential binding analysis, ChIP-Seq pipeline

Step	Jupyter Notebook	Workflow CWL	Tool	Input	Output	CWL Tool
Sample Download and Quality Control	01 - Pre-processing QC	download_quality_SRA-Tools.cwl	SRA-Tools	SRA accession	Fastq	fastq-dump.cwl
			FastQC	Fastq	FastQC HTML and Zip	fastqc.cwl
Trimming	02 - Samples trimming		Trimmomatic	Fastq	Fastq	trimmomatic-PE.cwl trimmomatic-SE.cwl
Alignments and Quantification	03 - Alignments	chip-seq-alignment.cwl	BWA	Fastq	BAM	bwa-mem.cwl
			Samtools	SAM	BAM Sorted BAM BAM index BAM stats BAM flagstats	samtools-view.cwl samtools-sort.cwl samtools-index.cwl samtools-stats.cwl samtools-flagstat.cwl
Peak Calling and IDR	04 - Peak Calling and IDR	peak-calling-MACS2.cwl	IGVtools	Sorted BAM	TDF	igvtools-count.cw
			MACS2	tagAlign	Peaks (TSV), plots	macs2-callpeak.cwl
			Homer	BAM	Annotation (TSV) FPKM values (TSV)	homer-annotatePeaks.cwl homer-makeTagDirectory.cwl
			IDR	Peaks	Peaks (TSV), plots	idr.cwl
Differential binding Analysis	05 - Differential binding Detection		DiffBind	Read Matrix	Peaks (TSV), plots	diffBind.cwl

7.1 Input requirements

The input requirement for the ChIPSeq pipeline is the *Sample sheet* file.

7.2 Pipeline command line

The ChIPSeq based project can be created using the following command line:

```
localhost:~> pm4ngs-chipseq
usage: Generate a PM4NGS project for ChIPSeq data analysis [-h] [-v]
```

(continues on next page)

(continued from previous page)

```
--sample-sheet
SAMPLE_SHEET
[--config-file CONFIG_FILE]
[--copy-rawdata]
```

Options:

- **sample-sheet:** Sample sheet with the samples metadata
- **config-file:** YAML file with configuration for project creation
- **copy-rawdata:** Copy the raw data defined in the sample sheet to the project structure. (The data can be hosted locally or in an http server)

7.3 Creating the Differential Binding detection from ChIP-Seq data project

The **pm4ngs-chipseq** command line executed with the **--sample-sheet** option will let you type the different variables required for creating and configuring the project. The default value for each variable is shown in the brackets. After all questions are answered, the CWL workflow files will be cloned from the github repo [ncbi/cwl-ngs-workflows-cbb](https://github.com/ncbi/cwl-ngs-workflows-cbb) to the folder **bin/cwl**.

```
localhost:~> pm4ngs-chipseq --sample-sheet my-sample-sheet.tsv
Generating ChIP-Seq data analysis project
author_name [Roberto Vera Alvarez]:
email [veraalva@ncbi.nlm.nih.gov]:
project_name [my_ngs_project]:
dataset_name [my_dataset_name]:
is_data_in_SRA [None]:
Select sequencing_technology:
1 - single-end
2 - paired-end
Choose from 1, 2 [1]: 2
create_demo [y]:
number_spots [1000000]:
organism [human]:
genome_name [hg38]:
genome_dir [hg38]:
aligner_index_dir [hg38/BWA]:
genome_fasta [hg38/genome.fa]:
genome_gtf [hg38/genome.gtf]:
genome_mappable_size [hg38]:
fdr [0.05]:
use_docker [y]:
max_number_threads [16]:
Cloning Git repo: https://github.com/ncbi/cwl-ngs-workflows-cbb to /Users/veraalva/my_
↪ngs_project/bin/cwl
Updating CWLs dockerPull and SoftwareRequirement from: /Users/veraalva/my_ngs_project/
↪requirements/conda-env-dependencies.yaml
bedtools with version 2.29.2 update image to: quay.io/biocontainers/bedtools:2.29.2--
↪hc088bd4_0
/Users/veraalva/my_ngs_project/bin/cwl/tools/bedtools/bedtools-docker.yml with_
↪old image replaced: quay.io/biocontainers/bedtools:2.28.0--hdf88d34_0 (continues on next page)
```

(continued from previous page)

```

bioconductor-diffbind with version 2.16.0 update image to: quay.io/biocontainers/
→bioconductor-diffbind:2.16.0--r40h5f743cb_0
  /Users/veraalva/my_nginx_project/bin/cwl/tools/R/deseq2-pca.cwl with old image_
→replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
  /Users/veraalva/my_nginx_project/bin/cwl/tools/R/macscutoff.cwl with old image_
→replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
  /Users/veraalva/my_nginx_project/bin/cwl/tools/R/dga_heatmaps.cwl with old image_
→replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
  /Users/veraalva/my_nginx_project/bin/cwl/tools/R/diffbind.cwl with old image_
→replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
  /Users/veraalva/my_nginx_project/bin/cwl/tools/R/edgeR-2conditions.cwl with old_
→image replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
  /Users/veraalva/my_nginx_project/bin/cwl/tools/R/volcano_plot.cwl with old image_
→replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
  /Users/veraalva/my_nginx_project/bin/cwl/tools/R/readQC.cwl with old image_
→replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
  /Users/veraalva/my_nginx_project/bin/cwl/tools/R/deseq2-2conditions.cwl with old_
→image replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
bwa with version 0.7.17 update image to: quay.io/biocontainers/bwa:0.7.17--hed695b0_7
  /Users/veraalva/my_nginx_project/bin/cwl/tools/bwa/bwa-docker.yml with old image_
→replaced: quay.io/biocontainers/bwa:0.7.17--h84994c4_5
homer with version 4.11 update image to: quay.io/biocontainers/homer:4.11--
→pl526h9a982cc_2
  /Users/veraalva/my_nginx_project/bin/cwl/tools/homer/homer-docker.yml with old_
→image replaced: quay.io/biocontainers/homer:4.11--pl526h2bce143_2
Copying file /Users/veraalva/Work/Developer/Python/pm4ngs/pm4ngs-chipseq/example/
→pm4ngs_chipseq_demo_sample_data.csv to /Users/veraalva/my_nginx_project/data/my_
→dataset_name/sample_table.csv
6 files loaded
Using table:
  sample_name file condition replicate
0 SRR7549105      ES         1
1 SRR7549106      ES         2
2 SRR7549109      MEF        1
3 SRR7549110      MEF        2
4 SRR7549114      rB         1
5 SRR7549113      rB         2
Done

```

The **pm4ngs-chipseq** command line will create a project structure as:

```

.
├── LICENSE
├── README.md
├── bin
│   └── cwl
├── config
│   └── init.py
├── data
│   └── my_dataset_name
├── doc
├── index.html
├── notebooks
│   ├── 00 - Project Report.ipynb
│   ├── 01 - Pre-processing QC.ipynb
│   ├── 02 - Samples trimming.ipynb
│   └── 03 - Alignments.ipynb

```

(continues on next page)

(continued from previous page)

```

├── 04 - Peak Calling and IDR.ipynb
├── 05 - Differential binding Detection.ipynb
├── requirements
│   └── conda-env-dependencies.yaml
├── results
│   └── my_dataset_name
├── src
└── tmp

12 directories, 11 files

```

Note: ChIP-Seq based project variables

- **author_name:** Default: [Roberto Vera Alvarez]
- **email:** Default: [veraalva@ncbi.nlm.nih.gov]
- **project_name:** Name of the project with no space nor especial characters. This will be used as project folder's name.
Default: [my_nginx_project]
- **dataset_name:** Dataset to process name with no space nor especial characters. This will be used as folder name to group the data. This folder will be created under the **data/{{dataset_name}}** and **results/{{dataset_name}}**.
Default: [my_dataset_name]
- **is_data_in_SRA:** If the data is in the SRA set this to y. A CWL workflow to download the data from the SRA database to the folder **data/{{dataset_name}}** and execute FastQC on it will be included in the **01 - Pre-processing QC.ipynb** notebook.
Set this option to **n**, if the fastq files names and location are included in the sample sheet.
Default: [y]
- **Select sequencing technology:** Select one of the available sequencing technologies in your data.
Values: 1 - single-end, 2 - paired-end
- **create_demo:** If the data is downloaded from the SRA and this option is set to y, only the number of spots specified in the next variable will be downloaded. Useful to test the workflow.
Default: [y]: y
- **number_spots:** Number of sport to download from the SRA database. It is ignored is the **create_demo** is set to **n**.
Default: [1000000]
- **organism:** Organism to process, e.g. human. This is used to link the selected genes to the NCBI gene database.
Default: [human]
- **genome_name:** Genome name , e.g. hg38 or mm10.
Default: [hg38]
- **genome_dir:** Absolute path to the directory with the genome annotation (genome.fa, genes.gtf) to be used by the workflow or the name of the genome.

If the name of the genome is used, PM4NGS will include a cell in the **03 - Alignments.ipynb** notebook to download the genome files. The genome data will be at **data/{{dataset_name}}/{{genome_name}}/**

Default: [hg38]

- **aligner_index_dir:** Absolute path to the directory with the genome indexes for BWA.

If `{{genome_name}}/BWA` is used, PM4NGS will include a cell in the **03 - Alignments.ipynb** notebook to create the genome indexes for BWA.

Default: [hg38/BWA]

- **genome_fasta:** Absolute path to the genome fasta file

If `{{genome_name}}/genome.fa` is used, PM4NGS will use the downloaded fasta file.

Default: [hg38/genome.fa]

- **genome_gtf:** Absolute path to the genome GTF file

If `{{genome_name}}/genes.gtf` is used, PM4NGS will use the downloaded GTF file.

Default: [hg38/gene.gtf]

- **genome_bed:** Absolute path to the genome BED file

If `{{genome_name}}/genes.bed` is used, PM4NGS will use the downloaded BED file.

Default: [hg38/genes.bed]

- **genome_mappable_size:** Genome mappable size used by MACS. For human can be hg38 or in case of other genomes it is a number.

Default: [hg38]

- **fdr:** Adjusted P-Value to be used as cutoff in the DG analysis, e.g. 0.05.

Default: [0.05]

- **use_docker:** Set this to y if you will be using Docker. Otherwise Conda needs to be installed in the computer.

Default: [y]

- **max_number_threads:** Number of threads available in the computer.

Default: [16]

7.4 Jupyter server

PM4NGS uses Jupyter as interface for users. After project creation the jupyter server should be started as shown below. The server will open a browser windows showing the project's structure just created.

```
localhost:~> jupyter notebook
```

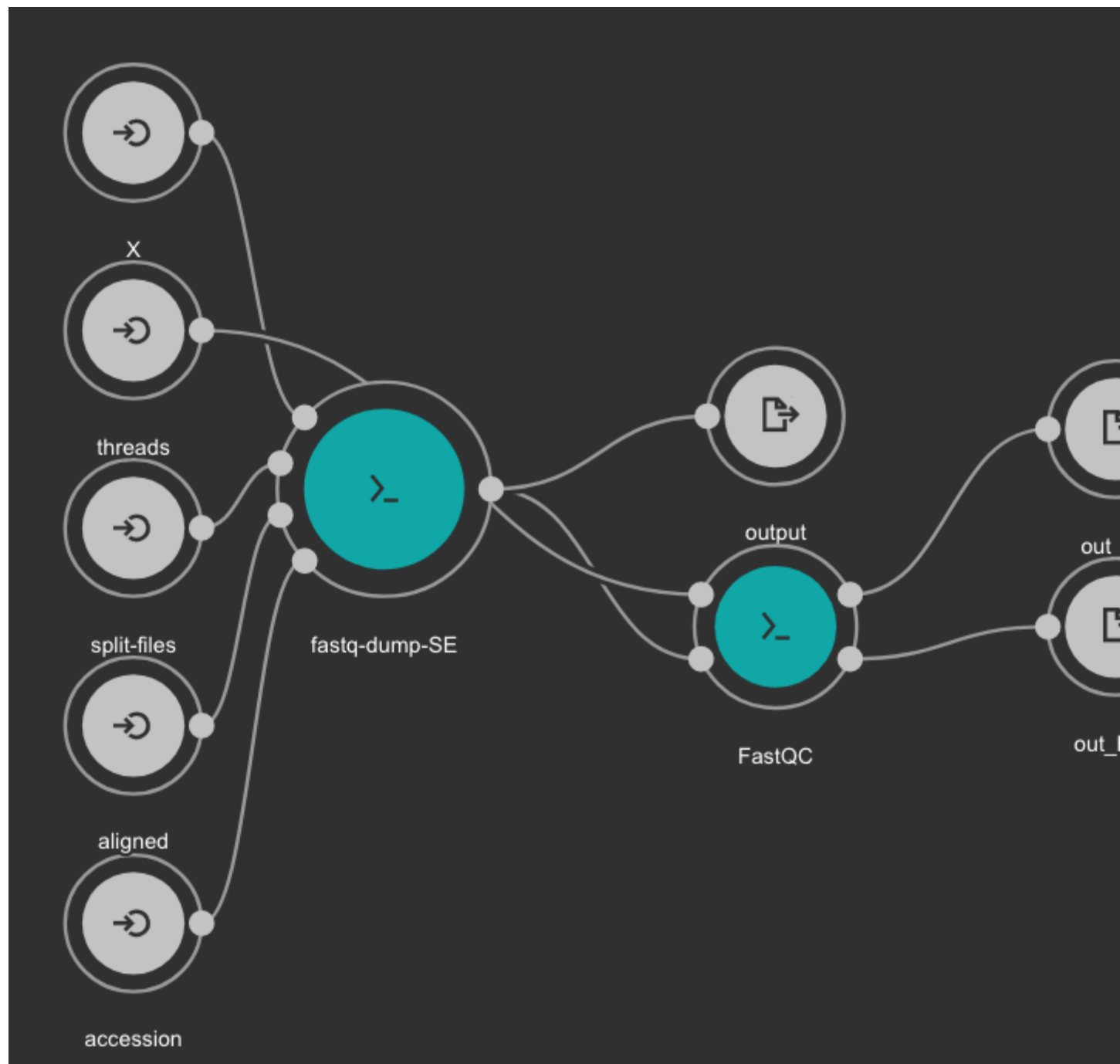
7.5 Data processing

Start executing the notebooks from 01 to 05 waiting for each step completion. The **00 - Project Report.ipynb** notebook can be executed after each notebooks to see the progress in the analysis.

7.6 CWL workflows

7.6.1 SRA download and QC workflow

This CWL workflow is designed to download FASTQ files from the NCBI SRA database using `fastq-dump` and then, execute `fastqc` generating a quality control report of the sample.



Inputs

- **accession**: SRA accession ID. Type: string. Required.
- **aligned**: Used it to download only aligned reads. Type: boolean. Optional.
- **split-files**: Dump each read into separate file. Files will receive suffix corresponding to read number. Type: boolean. Optional.
- **threads**: Number of threads. Type: int. Default: 1. Optional.
- **X**: Maximum spot id. Optional.

Outputs

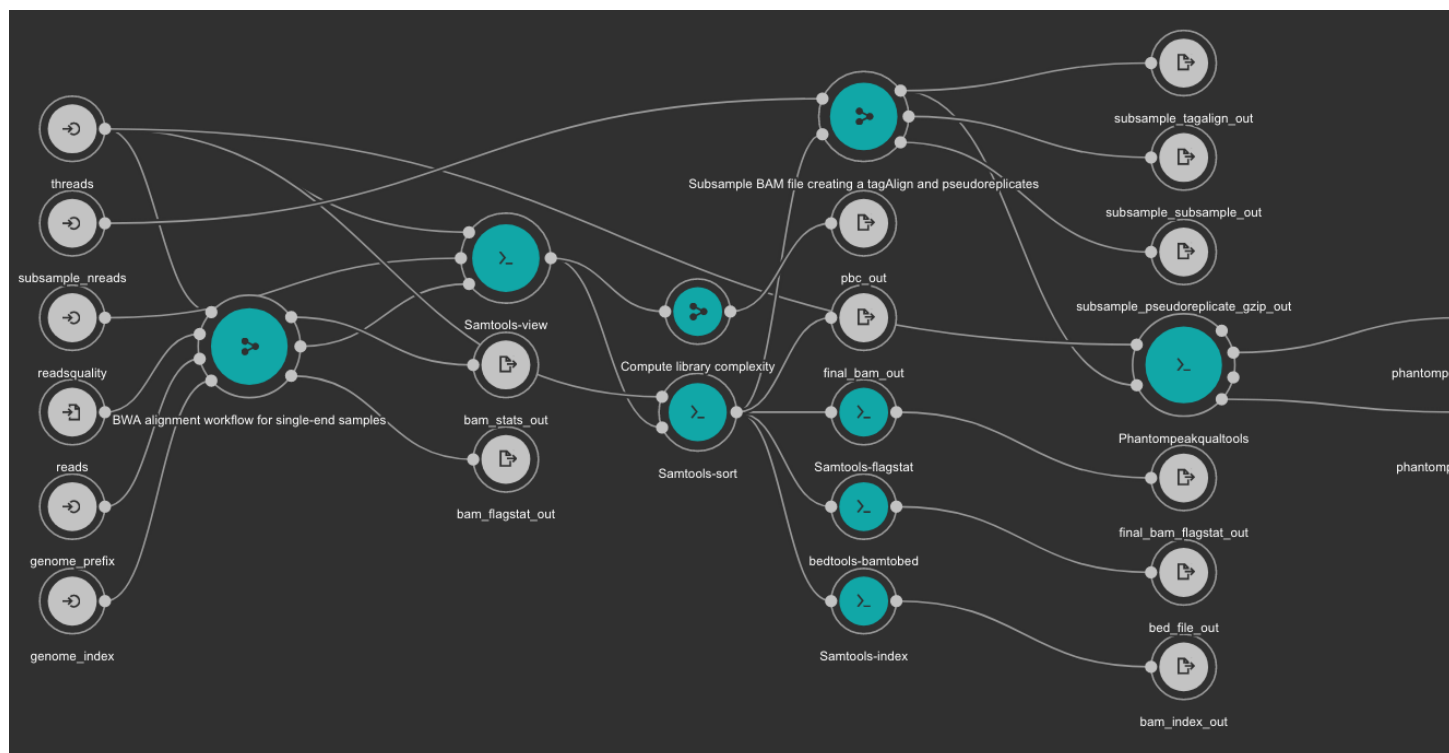
- **output**: Fastq files downloaded. Type: File[]
- **out_zip**: FastQC report ZIP file. Type: File[]
- **out_html**: FastQC report HTML. Type: File[]

7.6.2 Samples Trimming

Our workflows uses [Trimmomatic](#) for read trimming. The Jupyter notebooks uses some basic Trimmomatic options that need to be modified depending on the FastQC quality control report generated for the sample.

7.6.3 BWA based alignment and quality control workflow

This workflow use [BWA](#) as base aligner. It also use [SAMtools](#) and [bedtools](#) for file conversion and statistics report. Finally, [Phantompeakqualtools](#) is used to generate quality control report for the processed samples.



Inputs

- **genome_index**: Aligner indexes directory. Type: Directory. Required. Variable `ALIGNER_INDEX` in the Jupyter Notebooks.
- **genome_prefix**: Prefix of the aligner indexes. Generally, it is the name of the genome FASTA file. It can be used as `os.path.basename(GENOME_FASTA)` in the Jupyter Notebooks. Type: string. Required.
- **readsquality**: Minimum MAPQ value to use reads. We recommend for ChIP_exo data a value of: 30. Type: int. Required.
- **threads**: Number of threads. Type: int. Default: 1. Optional.
- **subsample_nreads**: Number of reads to be used for the subsample. We recommend for ChIP_exo data a value of: 500000. Type: int. Required.
- **reads**: FastQ file to be processed. Type: File. Required.

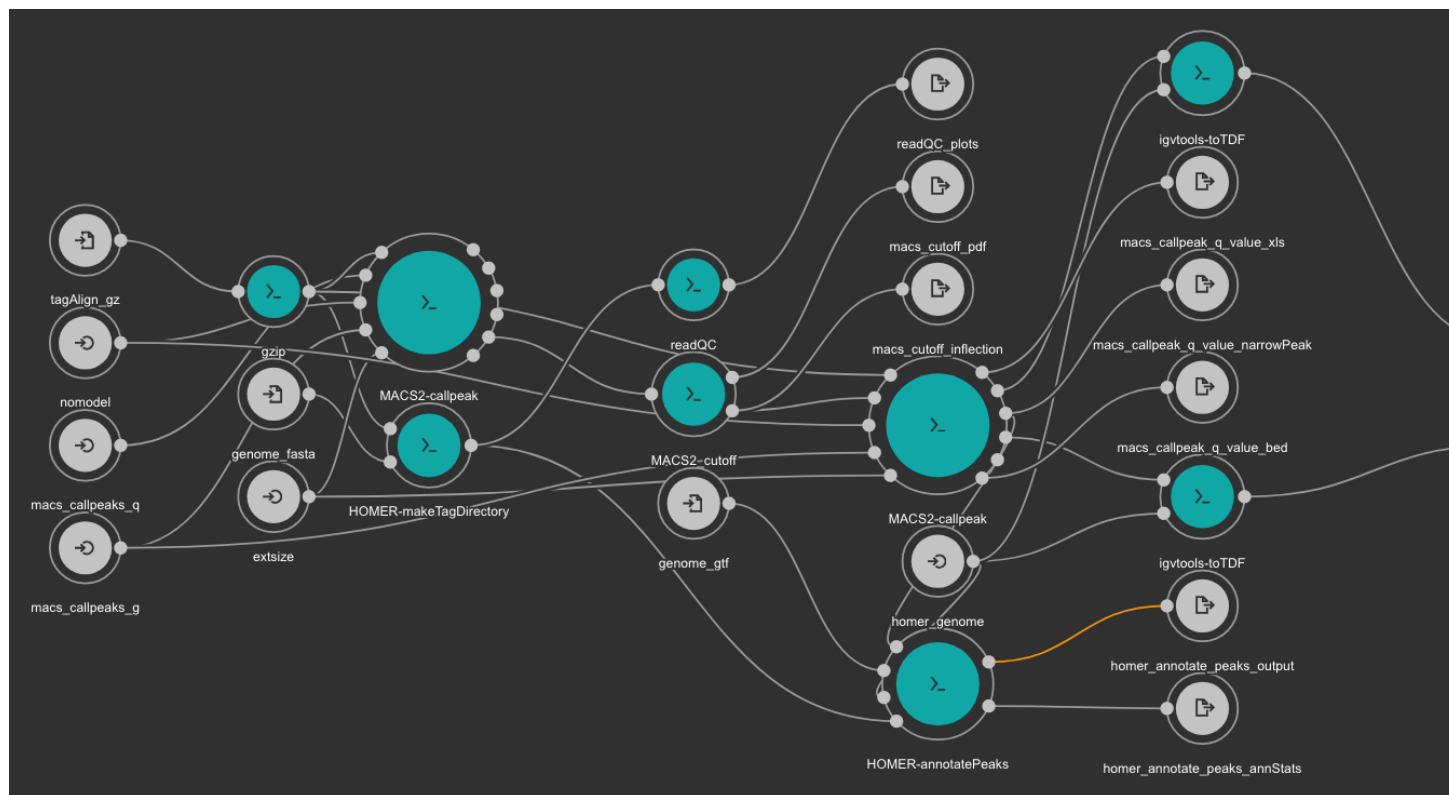
Outputs

- **bam_flagstat_out**: SAMtools flagstats for unfiltered BAM file. Type: File.
- **bam_stats_out**: SAMtools stats for unfiltered BAM file. Type: File.
- **final_bam_flagstat_out**: SAMtools flagstats for filtered BAM file. Type: File.
- **bed_file_out**: Aligned reads in BED format. Type: File.
- **final_bam_out**: Final BAM file filtered and sorted. Type: File.
- **bam_index_out**: BAM index file. Type: File.

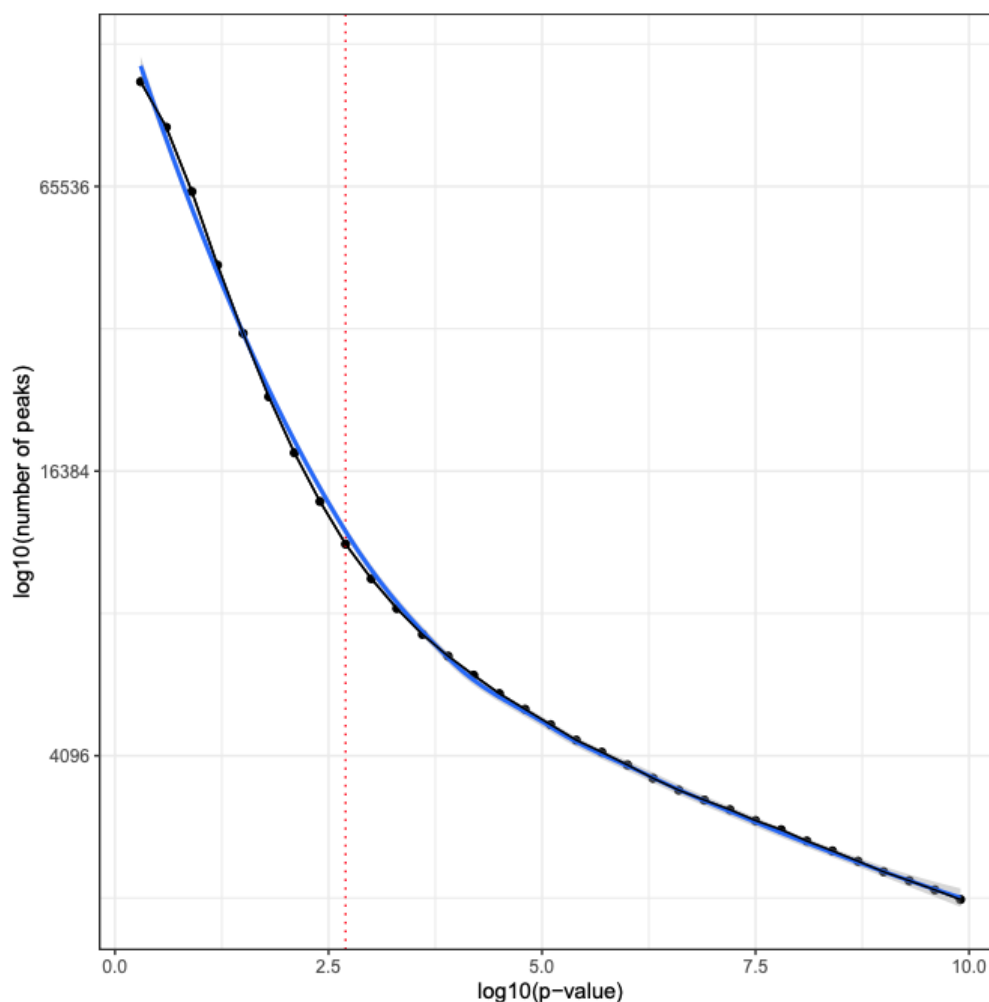
- **pbc_out**: Library complexity report. Type: File.
- **phantompeakqualtools_output_out**: Phantompeakqualtools main output. Type: File.
- **phantompeakqualtools_output_savp**: Phantompeakqualtools SAVP output. Type: File.
- **subsample_pseudoreplicate_gzip_out**: Subsample pseudoreplicates tagAlign gzipped. Type: File[].
- **subsample_tagalign_out**: Subsample tagAlign gzipped. Type: File[].
- **subsample_subsample_out**: Subsample shuffled tagAlign gzipped. Type: File[].

7.6.4 Peak caller workflow using MACS2

This workflow uses [MACS2](#) as peak caller tool. The annotation is created using [Homer](#) and TDF files are created with [igvtools](#).



[MACS2](#) is executed two times. First, the **cutoff-analysis** option is used to execute a cutoff value analysis which is used to estimate a proper value for the p-value used by MACS2 (for more detailed explanation read this [thread](#)).



RSeQC is also executed for quality control.

Inputs

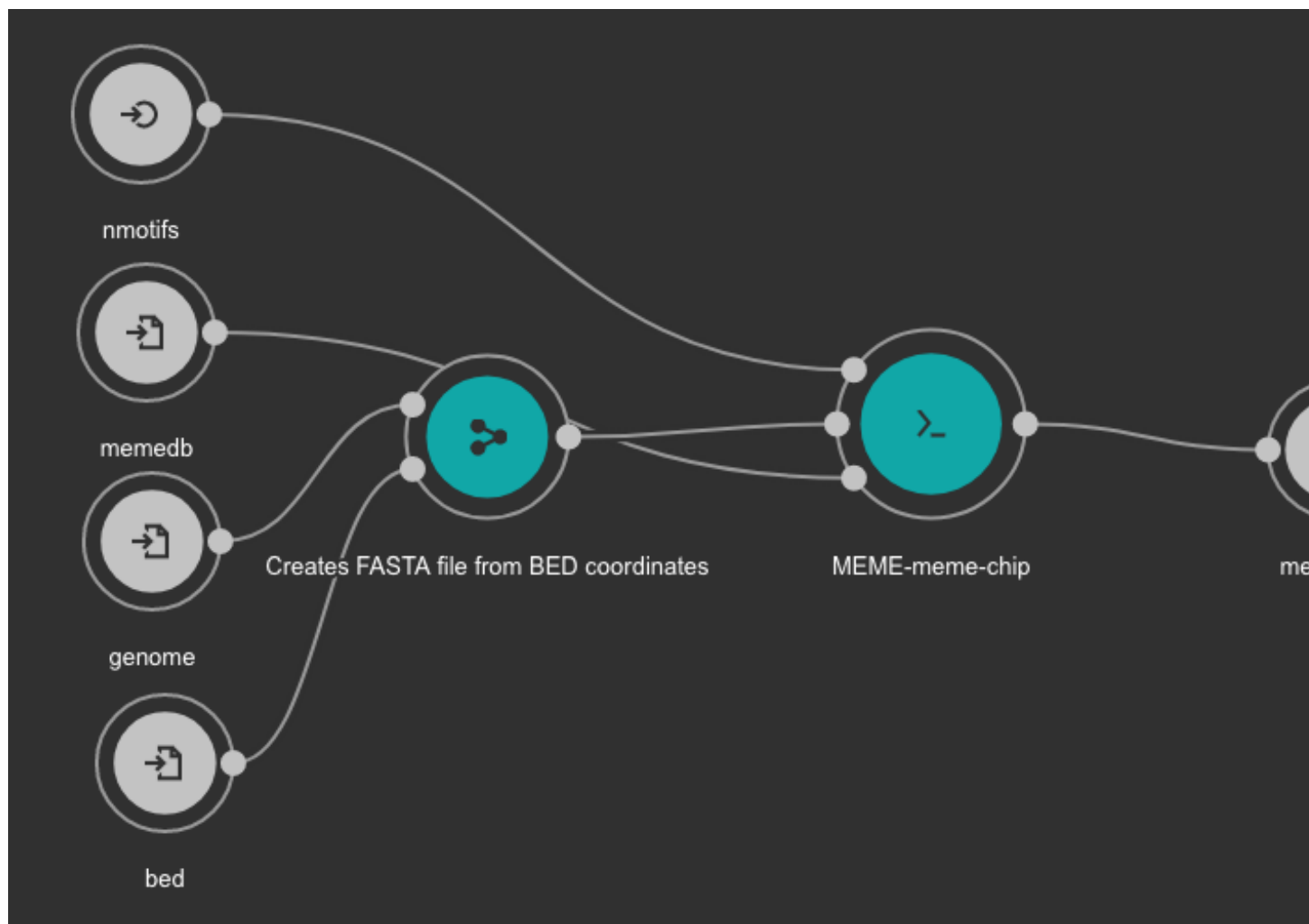
- **homer_genome**: Homer genome name. Type: string. Required.
- **genome_fasta**: Genome FASTA file. Type: File. Required. Variable GENOME_FASTA in the Jupyter Notebooks.
- **genome_gtf**: Genome GTF file. Type: File. Required. Variable GENOME_GTF in the Jupyter Notebooks.
- **tagAlign_gz**: Tag aligned file created with the BWA based alignment and quality control workflow. Type: File. Required.
- **macs_callpeaks_g**: Genome mappable size as defined in MACS2. Type: string. Required. Variable GENOME_MAPPABLE_SIZE in the Jupyter Notebooks.
- **macs_callpeaks_q**: MACS2 **q** option. Starting qvalue (minimum FDR) cutoff to call significant regions. Type: float. Required. Variable fdr in the Jupyter Notebooks.
- **nomodel**: MACS2 nomodel option. MACS will bypass building the shifting model. Type: boolean. Optional.
- **extsize**: MACS2 extsize option. MACS uses this parameter to extend reads in 5'→3' direction to fix-sized fragments. Type: int. Optional.

Outputs

- **readQC_plots**: RSeQC plots. Type: File[]
- **macs_cutoff_pdf**: MACS2 cutoff analysis plot in PDF format. Type: File
- **macs_cutoff_inflection**: MACS2 inflection q value used for the second round. Type: File
- **macs_callpeak_q_value_narrowPeak**: Final MACS2 narrowpeak file. Type: File
- **macs_callpeak_q_value_xls**: Final MACS2 XLS file. Type: File
- **macs_callpeak_q_value_bed**: Final MACS2 BED file. Type: File
- **homer_annotate_peaks_output**: Homer annotated BED file. Type: File
- **homer_annotate_peaks_annStats**: Homer annotation statistics. Type: File
- **lambda_tdf_out**: MACS2 lambda file in TDF format. Type: File
- **pileup_tdf_out**: MACS2 pileup file in TDF format. Type: File

7.6.5 MEME Motif detection workflow

Motif detection is executed using the [MEME](#) suite.



Inputs

- **bed**: BED file with detected peaks. Type: File. Required.
- **memedb**: MEME database for use by Tomtom and CentriMo. Type: File. Required.
- **genome**: Genome FASTA file. Variable GENOME_FASTA in the Jupyter Notebooks. Type: File. Required.
- **nmotifs**: Maximum number of motifs to find. We recommend use 10. Type: int. Required.

Outputs

- **meme_out**: MEME output directory. Type: Directory

7.6.6 MEME databases

MEME workflow depends on the MEME databases. Go to the MEME Suite Download page: <http://meme-suite.org/doc/download.html>

Download the latest version for the Motif Databases and GOMo Databases.

The downloaded files should be uncompressed in a directory **data/meme**. The final directory should be:

```
localhost:~> cd data
localhost:~> mkdir meme
localhost:~> cd meme
localhost:~> wget http://meme-suite.org/meme-software/Databases/motifs/motif_
↳databases.12.18.tgz
localhost:~> wget http://meme-suite.org/meme-software/Databases/gomo/gomo_databases.3.
↳2.tgz
localhost:~> tar xzf motif_databases.12.18.tgz
localhost:~> tar xzf gomo_databases.3.2.tgz
localhost:~> rm gomo_databases.3.2.tgz motif_databases.12.18.tgz
localhost:~> tree -d
.
├── gomo_databases
└── motif_databases
    ├── ARABD
    ├── CIS-BP
    ├── CISBP-RNA
    ├── ECOLI
    ├── EUKARYOTE
    ├── FLY
    ├── HUMAN
    ├── JASPAR
    ├── MALARIA
    ├── MIRBASE
    ├── MOUSE
    ├── PROKARYOTE
    ├── PROTEIN
    ├── RNA
    ├── TFBSShape
    ├── WORM
    └── YEAST

19 directories
```

See also an example in our test project: <https://ftp.ncbi.nlm.nih.gov/pub/pm4ngs/examples/chipexo-single/data/>

7.6.7 Differential binding analysis with Diffbind

Differential binding event is detected with **Diffbind**. This tool will be executed for all comparisons added to the **comparisons** array. See cell number 3 in the notebook **05 - Differential binding analysis.ipynb** (ChIP-Seq workflow).

Inputs

- **bamDir**: Directory with the BAM files. Type: Directory. Required.
- **bedDir**: Directory with BED files created from MACS2 peak calling workflow Type: Directory. Required.

Outputs

- **outpng**: PNG files created from Diffbind. Type File[]

- **outxls**: XLS files created from Diffbind. Type: File[]
- **outbed** BED files created from Diffbind. Type: File[]

7.7 Demo

PM4NGS includes a demo project that users can use to test the framework. It is pre-configured to use Docker as execution environment.

The ChIPSeq based demo process samples from the BioProject [PRJNA481982](#).

Use this command to create the project structure in your local computer

```
localhost:~> pm4ngs-chipseq-demo
```

Once it finish, start the jupyter server and execute the notebooks as it is indicated on them

```
localhost:~> jupyter notebook
[I 14:12:52.956 NotebookApp] Serving notebooks from local directory: /home/veraalva
[I 14:12:52.956 NotebookApp] Jupyter Notebook 6.1.4 is running at:
[I 14:12:52.956 NotebookApp] http://localhost:8888/?
↪token=eae6a8d42ad12d6ace23f5d0923bcec14d0f798127750122
[I 14:12:52.956 NotebookApp] or http://127.0.0.1:8888/?
↪token=eae6a8d42ad12d6ace23f5d0923bcec14d0f798127750122
[I 14:12:52.956 NotebookApp] Use Control-C to stop this server and shut down all
↪kernels (twice to skip confirmatio
n).
[C 14:12:52.959 NotebookApp]

To access the notebook, open this file in a browser:
    file:///home/veraalva/.local/share/jupyter/runtime/nbserver-23251-open.html
Or copy and paste one of these URLs:
    http://localhost:8888/?token=eae6a8d42ad12d6ace23f5d0923bcec14d0f798127750122
    or http://127.0.0.1:8888/?token=eae6a8d42ad12d6ace23f5d0923bcec14d0f798127750122
```

The results of this analysis is <https://ftp.ncbi.nlm.nih.gov/pub/pm4ngs/examples/chipseq-hmgn1/>

Detection of binding events from ChIP-exo data

The workflow is comprised of five steps, as shown in next table. The first step, sample download and quality control, is for downloading samples from the NCBI SRA database, if necessary, or for executing the pre-processing quality control tools on all samples. After that, sample trimming and alignment are executed. Post-processing quality control on the ChIP-exo samples is performed, using Phantompeakqualtools. After that, the peak calling is performed, using MACE. Finally, DNA motif annotation is performed with the MEME Suite.

Table 1: DNA motif identification, ChIP-exo data pipeline

Step	Jupyter Notebook	Workflow CWL	Tool	Input	Output	CWL Tool
Sample Download and Quality Control	01 - Pre-processing QC	download_quality_SRA-Tools.cwl	SRA-Tools	SRA accession	Fastq	fastq-dump.cwl
			FastQC	Fastq	FastQC HTML and Zip	fastqc.cwl
Trimming	02 - Samples trimming		Trimmomatic	Fastq	Fastq	trimmomatic-PE.cwl trimmomatic-SE.cwl
Alignments and Quantification	03 - Alignments	chip-seq-alignment.cwl	BWA	Fastq	BAM	bwa-mem.cwl
			Samtools	SAM	BAM Sorted BAM BAM index BAM stats BAM flagstats	samtools-view.cwl samtools-sort.cwl samtools-index.cwl samtools-stats.cwl samtools-flagstat.cwl
Peak Calling and IDR	04 - Peak Calling	peak-caller-MACE-SE.cwl	IGVtools	Sorted BAM	TDF	igvtools-count.cw
			MACE	tagAlign	Peaks (TSV), plots	mace.cwl
			BAMScale	BAM	TPM values	bamscale-cov.cwl bamscale-scale.cwl
DNA motif identification	05 - MEME Motif		MEME Suite	Peaks	DNA Motif , plots	meme-chip.cwl

8.1 Input requirements

The input requirement for the ChIPexo pipeline is the *Sample sheet* file.

8.2 Pipeline command line

The ChIPexo based project can be created using the following command line:

```
localhost:~> pm4ngs-chipexo
usage: Generate a PM4NGS project for ChIPexo data analysis [-h] [-v]
                                --sample-sheet
                                SAMPLE_SHEET
                                [--config-file CONFIG_FILE]
                                [--copy-rawdata]
```

Options:

- **sample-sheet:** Sample sheet with the samples metadata
- **config-file:** YAML file with configuration for project creation
- **copy-rawdata:** Copy the raw data defined in the sample sheet to the project structure. (The data can be hosted locally or in an http server)

8.3 Creating the Detection of binding events from ChIP-exo data project

The **pm4ngs-chipexo** command line executed with the **--sample-sheet** option will let you type the different variables required for creating and configuring the project. The default value for each variable is shown in the brackets. After all questions are answered, the CWL workflow files will be cloned from the github repo [ncbi/cwl-ngs-workflows-cbb](https://github.com/ncbi/cwl-ngs-workflows-cbb) to the folder **bin/cwl**.

```
localhost:~> pm4ngs-chipexo --sample-sheet my-sample-sheet.tsv
Generating ChIP-exo data analysis project
author_name [Roberto Vera Alvarez]:
email [veraalva@ncbi.nlm.nih.gov]:
project_name [my_ngs_project]:
dataset_name [my_dataset_name]:
is_data_in_SRA [None]:
Select sequencing_technology:
1 - single-end
2 - paired-end
Choose from 1, 2 [1]:
create_demo [y]:
number_spots [5000000]:
organism [human]:
genome_name [hg38]:
genome_dir [hg38]:
aligner_index_dir [hg38/BWA]:
genome_fasta [hg38/genome.fa]:
genome_gtf [hg38/genes.gtf]:
genome_chromsizes [hg38/genome.sizes]:
use_docker [y]:
max_number_threads [16]:
Cloning Git repo: https://github.com/ncbi/cwl-ngs-workflows-cbb to /Users/veraalva/my_
↪ngs_project/bin/cwl
Updating CWLs dockerPull and SoftwareRequirement from: /Users/veraalva/my_ngs_project/
↪requirements/conda-env-dependencies.yaml
bamscale with version 0.0.3 update image to: quay.io/biocontainers/bamscale:0.0.3--
↪ha85820d_0
/Users/veraalva/my_ngs_project/bin/cwl/tools/bamscale/bamscale-docker.yml with
↪old image replaced: quay.io/biocontainers/bamscale:0.0.5--h18f8b1d_1
bedtools with version 2.29.2 update image to: quay.io/biocontainers/bedtools:2.29.2--
↪hc088bd4_0
/Users/veraalva/my_ngs_project/bin/cwl/tools/bedtools/bedtools-docker.yml with
↪old image replaced: quay.io/biocontainers/bedtools:2.28.0--hdf88d34_0
bioconductor-diffbind with version 2.16.0 update image to: quay.io/biocontainers/
↪bioconductor-diffbind:2.16.0--r40h5f743cb_0
/Users/veraalva/my_ngs_project/bin/cwl/tools/R/deseq2-pca.cwl with old image
↪replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
```

(continues on next page)

(continued from previous page)

```

/Users/veraalva/my_nginx_project/bin/cwl/tools/R/macscutoff.cwl with old image_
→replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
/Users/veraalva/my_nginx_project/bin/cwl/tools/R/dga_heatmaps.cwl with old image_
→replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
/Users/veraalva/my_nginx_project/bin/cwl/tools/R/diffbind.cwl with old image_
→replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
/Users/veraalva/my_nginx_project/bin/cwl/tools/R/edgeR-2conditions.cwl with old_
→image replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
/Users/veraalva/my_nginx_project/bin/cwl/tools/R/volcano_plot.cwl with old image_
→replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
/Users/veraalva/my_nginx_project/bin/cwl/tools/R/readQC.cwl with old image_
→replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
/Users/veraalva/my_nginx_project/bin/cwl/tools/R/deseq2-2conditions.cwl with old_
→image replaced: quay.io/biocontainers/bioconductor-diffbind:2.16.0--r40h5f743cb_2
bwa with version 0.7.17 update image to: quay.io/biocontainers/bwa:0.7.17--hed695b0_7
/Users/veraalva/my_nginx_project/bin/cwl/tools/bwa/bwa-docker.yml with old image_
→replaced: quay.io/biocontainers/bwa:0.7.17--h84994c4_5
There is not biocontainer image for gffread version 0.12.1
homer with version 4.11 update image to: quay.io/biocontainers/homer:4.11--
→pl526h9a982cc_2
/Users/veraalva/my_nginx_project/bin/cwl/tools/homer/homer-docker.yml with old_
→image replaced: quay.io/biocontainers/homer:4.11--pl526h2bce143_2
mace with version 1.2 update image to: quay.io/biocontainers/mace:1.2--py27h99da42f_0
/Users/veraalva/my_nginx_project/bin/cwl/tools/mace/mace-docker.yml with old image_
→replaced: quay.io/biocontainers/mace:1.2--py27h99da42f_1
meme with version 5.1.1 update image to: quay.io/biocontainers/meme:5.1.1--
→py37pl526h072abfd_3
/Users/veraalva/my_nginx_project/bin/cwl/tools/meme/meme-docker.yml with old image_
→replaced: quay.io/biocontainers/meme:5.1.1--py27pl526h53063a7_3
Copying file /Users/veraalva/Work/Developer/Python/pm4ngs/pm4ngs-chipexo/example/
→pm4ngs_chipexo_demo_sample_data.csv to /Users/veraalva/my_nginx_project/data/my_
→dataset_name/sample_table.csv
6 files loaded
Using table:
  sample_name file                condition  replicate
0  SRR4011416   Exp_O2_growth_no_rifampicin      1
1  SRR4011417   Exp_O2_growth_no_rifampicin      2
2  SRR4011421   Exp_O2_growth_rifampicin         1
3  SRR4011425   Exp_O2_growth_rifampicin         2
4  SRR4011418   Stat_O2_growth_no_rifampicin      1
5  SRR4011419   Stat_O2_growth_no_rifampicin      2
Done

```

The **pm4ngs-chipexo** command line will create a project structure as:

```

.
├── LICENSE
├── README.md
├── bin
│   └── cwl
├── config
│   └── init.py
├── data
│   └── my_dataset_name
├── doc
├── index.html
└── notebooks

```

(continues on next page)

(continued from previous page)

```

├── 00 - Project Report.ipynb
├── 01 - Pre-processing QC.ipynb
├── 02 - Samples trimming.ipynb
├── 03 - Alignments.ipynb
├── 04 - Peak Calling.ipynb
├── 05 - MEME Motif.ipynb
├── requirements
│   └── conda-env-dependencies.yaml
├── results
│   └── my_dataset_name
├── src
└── tmp

12 directories, 11 files

```

Note: ChIPexo based project variables

- **author_name:** Default: [Roberto Vera Alvarez]
- **email:** Default: [veraalva@ncbi.nlm.nih.gov]
- **project_name:** Name of the project with no space nor especial characters. This will be used as project folder's name.
Default: [my_ngs_project]
- **dataset_name:** Dataset to process name with no space nor especial characters. This will be used as folder name to group the data. This folder will be created under the **data/{{dataset_name}}** and **results/{{dataset_name}}**.
Default: [my_dataset_name]
- **is_data_in_SRA:** If the data is in the SRA set this to y. A CWL workflow to download the data from the SRA database to the folder **data/{{dataset_name}}** and execute FastQC on it will be included in the **01 - Pre-processing QC.ipynb** notebook.
Set this option to **n**, if the fastq files names and location are included in the sample sheet.
Default: [y]
- **Select sequencing technology:** Select one of the available sequencing technologies in your data.
Values: 1 - single-end, 2 - paired-end
- **create_demo:** If the data is downloaded from the SRA and this option is set to y, only the number of spots specified in the next variable will be downloaded. Useful to test the workflow.
Default: [y]: y
- **number_spots:** Number of sport to download from the SRA database. It is ignored is the **create_demo** is set to **n**.
Default: [1000000]
- **organism:** Organism to process, e.g. human. This is used to link the selected genes to the NCBI gene database.
Default: [human]
- **genome_name:** Genome name , e.g. hg38 or mm10.
Default: [hg38]

- **genome_dir:** Absolute path to the directory with the genome annotation (genome.fa, genes.gtf) to be used by the workflow or the name of the genome.

If the name of the genome is used, PM4NGS will include a cell in the **03 - Alignments.ipynb** notebook to download the genome files. The genome data will be at **data/{{dataset_name}}/{{genome_name}}/**

Default: [hg38]
 - **aligner_index_dir:** Absolute path to the directory with the genome indexes for BWA.

If **{{genome_name}}/BWA** is used, PM4NGS will include a cell in the **03 - Alignments.ipynb** notebook to create the genome indexes for BWA.

Default: [hg38/BWA]
 - **genome_fasta:** Absolute path to the genome fasta file

If **{{genome_name}}/genome.fa** is used, PM4NGS will use the downloaded fasta file.

Default: [hg38/genome.fa]
 - **genome_gtf:** Absolute path to the genome GTF file

If **{{genome_name}}/genes.gtf** is used, PM4NGS will use the downloaded GTF file.

Default: [hg38/gene.gtf]
 - **genome_chromsizes:** TSV file with the genome chromosomes name and size.

If **{{genome_name}}/genome.sizes** is used, PM4NGS will use the downloaded file.

Default: [hg38/genome.sizes]
 - **use_docker:** Set this to y if you will be using Docker. Otherwise Conda needs to be installed in the computer.

Default: [y]
 - **max_number_threads:** Number of threads available in the computer.

Default: [16]
-

8.4 Jupyter server

PM4NGS uses Jupyter as interface for users. After project creation the jupyter server should be started as shown below. The server will open a browser windows showing the project's structure just created.

```
localhost:~> jupyter notebook
```

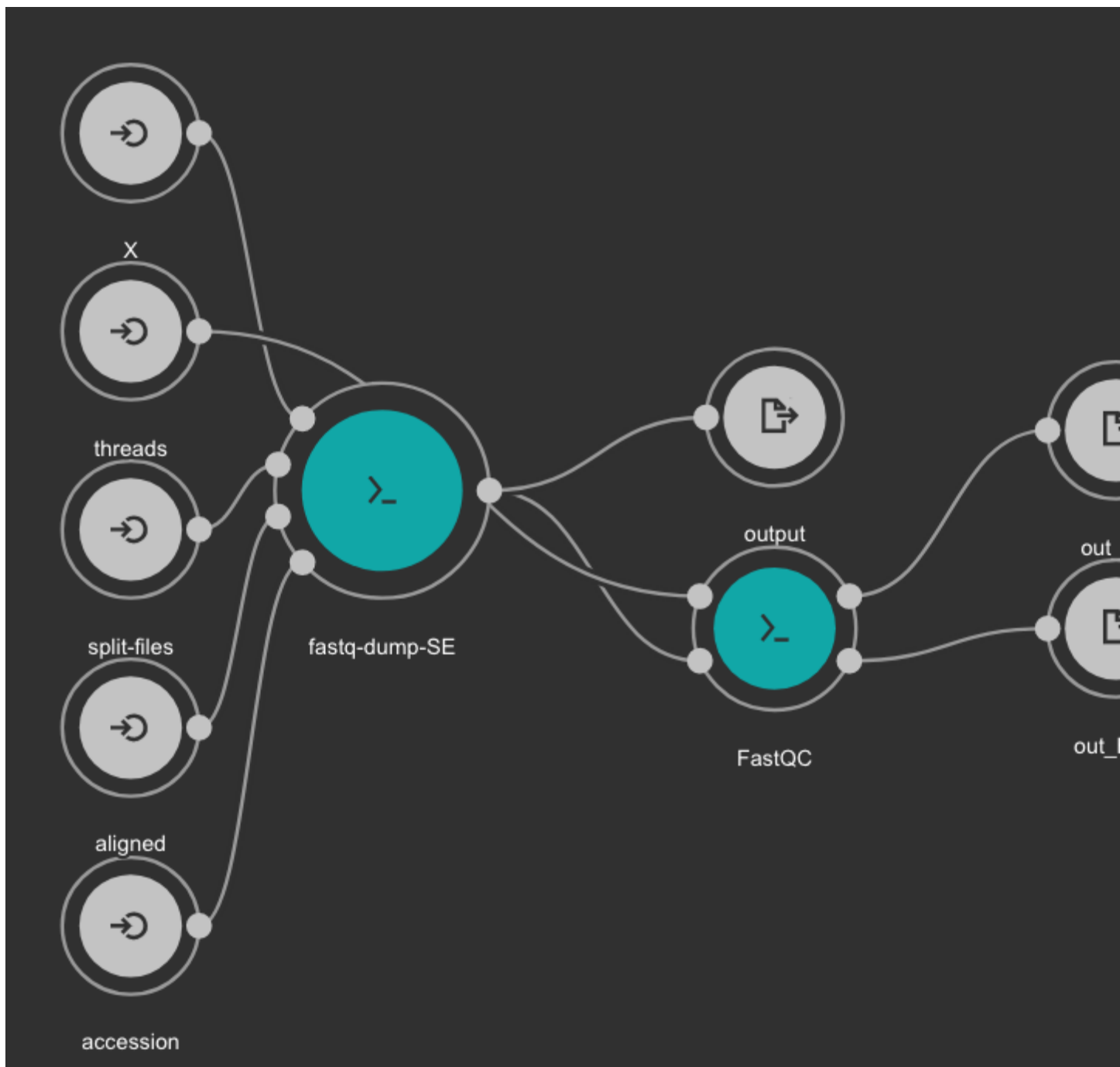
8.5 Data processing

Start executing the notebooks from 01 to 05 waiting for each step completion. The **00 - Project Report.ipynb** notebook can be executed after each notebooks to see the progress in the analysis.

8.6 CWL workflows

8.6.1 SRA download and QC workflow

This CWL workflow is designed to download FASTQ files from the NCBI SRA database using `fastq-dump` and then, execute `fastqc` generating a quality control report of the sample.



Inputs

- **accession**: SRA accession ID. Type: string. Required.
- **aligned**: Used it to download only aligned reads. Type: boolean. Optional.
- **split-files**: Dump each read into separate file. Files will receive suffix corresponding to read number. Type: boolean. Optional.
- **threads**: Number of threads. Type: int. Default: 1. Optional.
- **X**: Maximum spot id. Optional.

Outputs

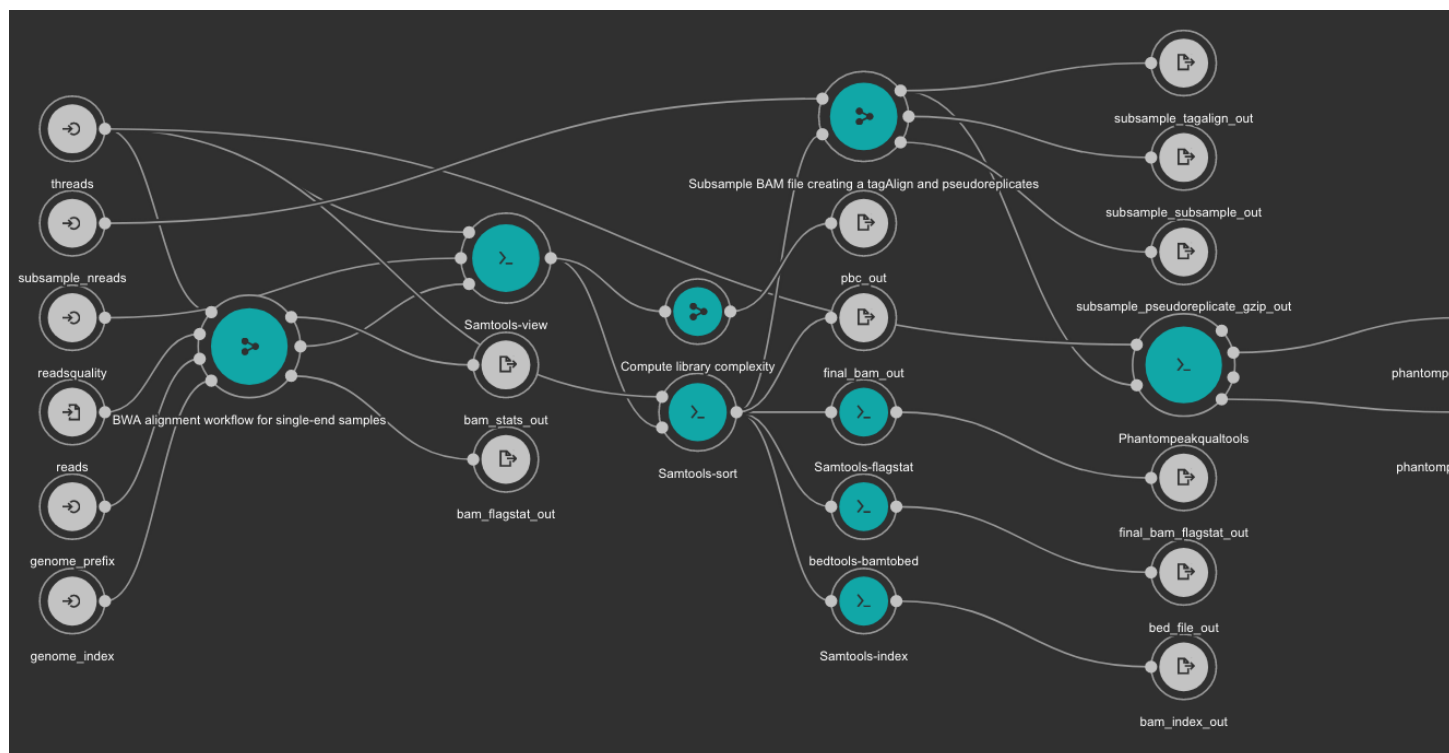
- **output**: Fastq files downloaded. Type: File[]
- **out_zip**: FastQC report ZIP file. Type: File[]
- **out_html**: FastQC report HTML. Type: File[]

8.6.2 Samples Trimming

Our workflows uses [Trimmomatic](#) for read trimming. The Jupyter notebooks uses some basic Trimmomatic options that need to be modified depending on the FastQC quality control report generated for the sample.

8.6.3 BWA based alignment and quality control workflow

This workflow use [BWA](#) as base aligner. It also use [SAMtools](#) and [bedtools](#) for file conversion and statistics report. Finally, [Phantompeakqualtools](#) is used to generate quality control report for the processed samples.



Inputs

- **genome_index**: Aligner indexes directory. Type: Directory. Required. Variable `ALIGNER_INDEX` in the Jupyter Notebooks.
- **genome_prefix**: Prefix of the aligner indexes. Generally, it is the name of the genome FASTA file. It can be used as `os.path.basename(GENOME_FASTA)` in the Jupyter Notebooks. Type: string. Required.
- **readsquality**: Minimum MAPQ value to use reads. We recommend for ChIP_exo data a value of: 30. Type: int. Required.
- **threads**: Number of threads. Type: int. Default: 1. Optional.
- **subsample_nreads**: Number of reads to be used for the subsample. We recommend for ChIP_exo data a value of: 500000. Type: int. Required.
- **reads**: FastQ file to be processed. Type: File. Required.

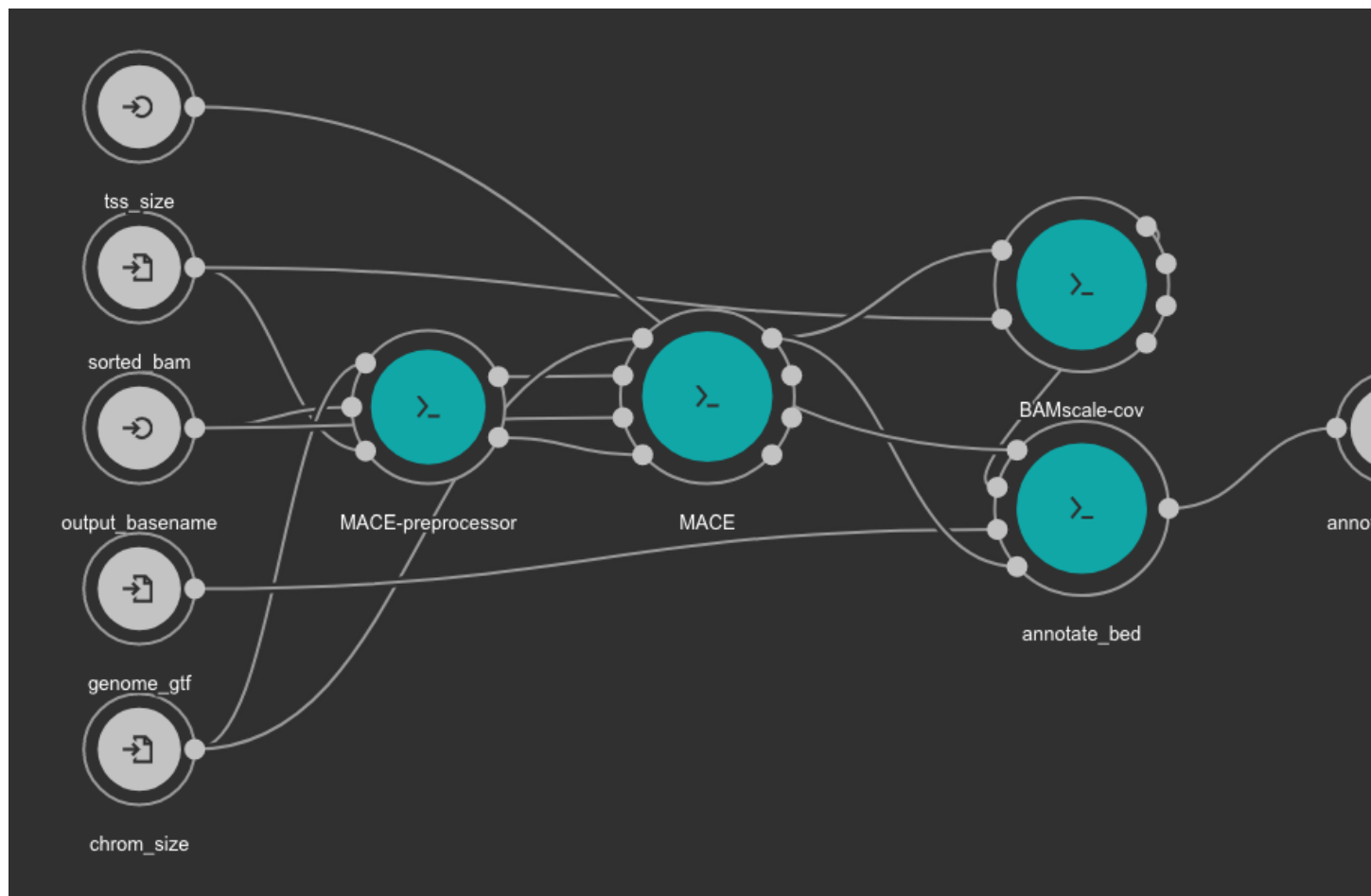
Outputs

- **bam_flagstat_out**: SAMtools flagstats for unfiltered BAM file. Type: File.
- **bam_stats_out**: SAMtools stats for unfiltered BAM file. Type: File.
- **final_bam_flagstat_out**: SAMtools flagstats for filtered BAM file. Type: File.
- **bed_file_out**: Aligned reads in BED format. Type: File.
- **final_bam_out**: Final BAM file filtered and sorted. Type: File.
- **bam_index_out**: BAM index file. Type: File.

- **pbc_out**: Library complexity report. Type: File.
- **phantompeakqualtools_output_out**: Phantompeakqualtools main output. Type: File.
- **phantompeakqualtools_output_savp**: Phantompeakqualtools SAVP output. Type: File.
- **subsample_pseudoreplicate_gzip_out**: Subsample pseudoreplicates tagAlign gzipped. Type: File[].
- **subsample_tagalign_out**: Subsample tagAlign gzipped. Type: File[].
- **subsample_subsample_out**: Subsample shuffled tagAlign gzipped. Type: File[].

8.6.4 Peak caller workflow using MACE

This workflow uses [MACE](#) as peak caller tool. The annotation is created from the GTF file using a *in-house* python script available [here](#). [BAMscale](#) is used for the quantification of the resulting peaks.



Inputs

- **chrom_size**: Chromosome size file. Tab or space separated text file with 2 columns: first column is chromosome name, second column is size of the chromosome. Type: File. Required. Variable GENOME_CHROMSIZES in the Jupyter Notebooks.

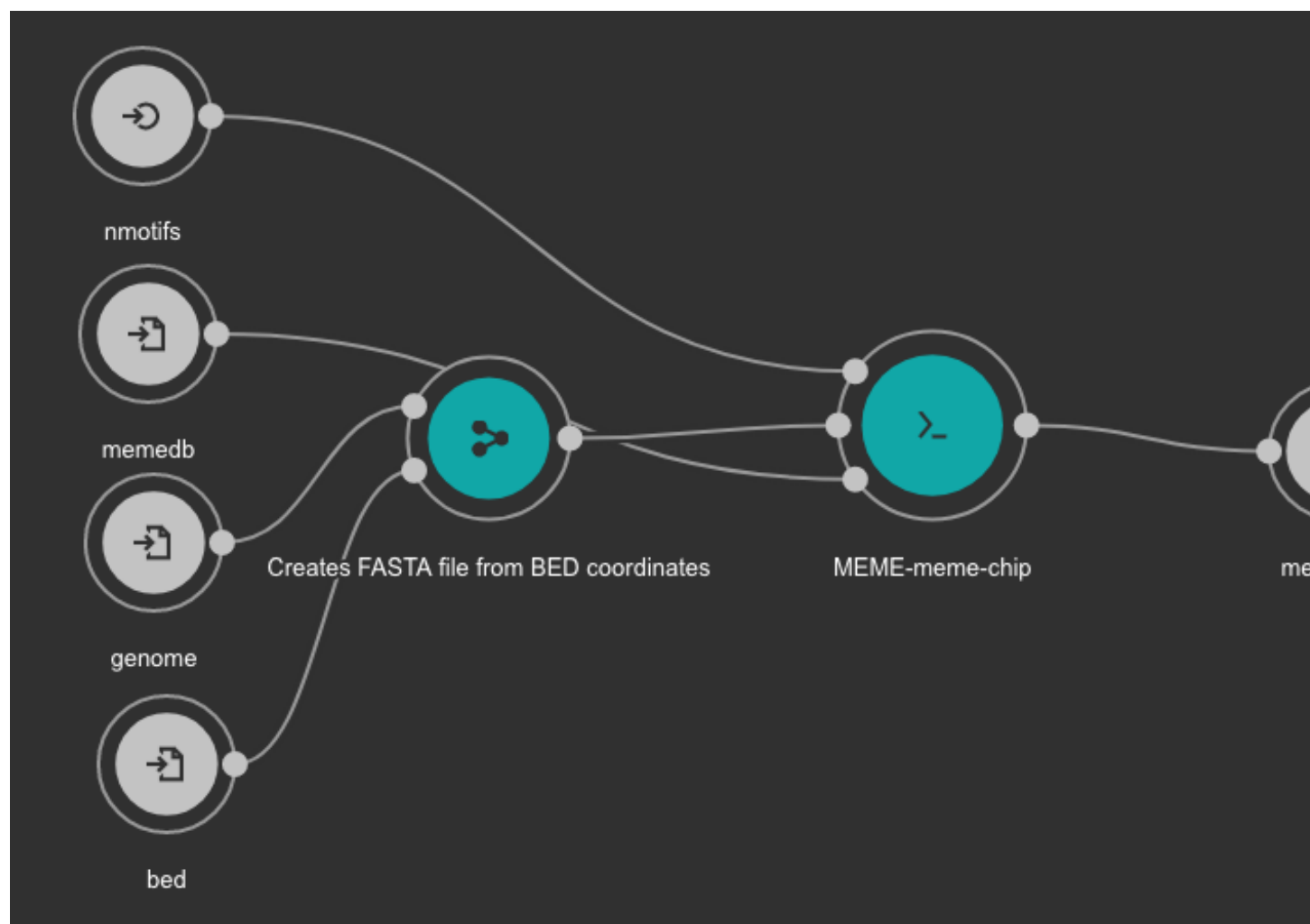
- **output_basename**: Prefix for the output file. Type: string. Required.
- **genome_gtf**: Genome GTF file. Variable `GENOME_GTF` in the Jupyter Notebooks. Type: File. Required.
- **tss_size**: Number of bp to use for TSS definition. We recommend use 1000. Type: int. Required.

Outputs

- **annotated_bed**: Annotated detected peaks in BED format. Type: File

8.6.5 MEME Motif detection workflow

Motif detection is executed using the [MEME](#) suite.



Inputs

- **bed**: BED file with detected peaks. Type: File. Required.

- **memedb**: MEME database for use by Tomtom and CentriMo. Type: File. Required.
- **genome**: Genome FASTA file. Variable GENOME_FASTA in the Jupyter Notebooks. Type: File. Required.
- **nmotifs**: Maximum number of motifs to find. We recommend use 10. Type: int. Required.

Outputs

- **meme_out**: MEME output directory. Type: Directory

8.6.6 MEME databases

MEME workflow depends on the MEME databases. Go to the MEME Suite Download page: <http://meme-suite.org/doc/download.html>

Download the latest version for the Motif Databases and GOMo Databases.

The downloaded files should be uncompressed in a directory **data/meme**. The final directory should be:

```
localhost:~> cd data
localhost:~> mkdir meme
localhost:~> cd meme
localhost:~> wget http://meme-suite.org/meme-software/Databases/motifs/motif_
↳databases.12.18.tgz
localhost:~> wget http://meme-suite.org/meme-software/Databases/gomo/gomo_databases.3.
↳2.tgz
localhost:~> tar xzf motif_databases.12.18.tgz
localhost:~> tar xzf gomo_databases.3.2.tgz
localhost:~> rm gomo_databases.3.2.tgz motif_databases.12.18.tgz
localhost:~> tree -d
.
├── gomo_databases
└── motif_databases
    ├── ARABD
    ├── CIS-BP
    ├── CISBP-RNA
    ├── ECOLI
    ├── EUKARYOTE
    ├── FLY
    ├── HUMAN
    ├── JASPAR
    ├── MALARIA
    ├── MIRBASE
    ├── MOUSE
    ├── PROKARYOTE
    ├── PROTEIN
    ├── RNA
    ├── TFBSShape
    ├── WORM
    └── YEAST

19 directories
```

See also an example in our test project: <https://ftp.ncbi.nlm.nih.gov/pub/pm4ngs/examples/chipexo-single/data/>

8.7 Demo

PM4NGS includes a demo project that users can use to test the framework. It is pre-configured to use Docker as execution environment.

The ChIPexo based demo process samples from the BioProject [PRJNA338159](#).

Use this command to create the project structure in your local computer

```
localhost:~> pm4ngs-chipexo-demo
```

Once it finish, start the jupyter server and execute the notebooks as it is indicated on them

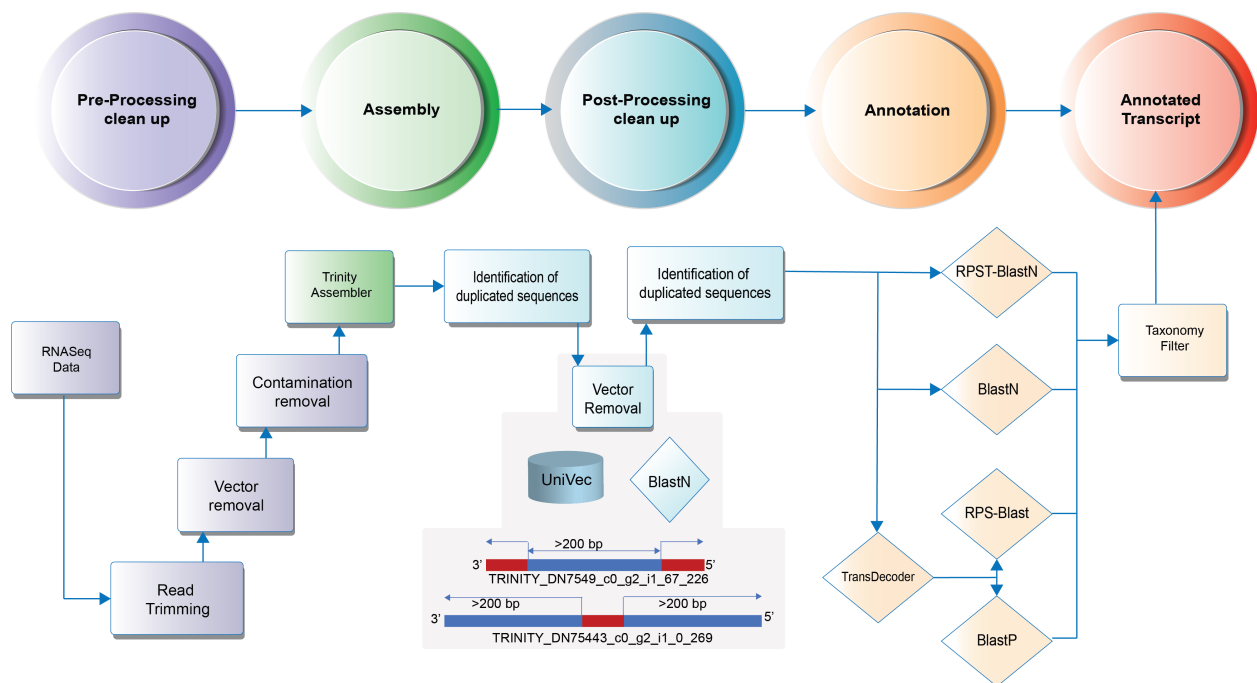
```
localhost:~> jupyter notebook
[I 14:12:52.956 NotebookApp] Serving notebooks from local directory: /home/veraalva
[I 14:12:52.956 NotebookApp] Jupyter Notebook 6.1.4 is running at:
[I 14:12:52.956 NotebookApp] http://localhost:8888/?
↪token=eae6a8d42ad12d6ace23f5d0923bcec14d0f798127750122
[I 14:12:52.956 NotebookApp] or http://127.0.0.1:8888/?
↪token=eae6a8d42ad12d6ace23f5d0923bcec14d0f798127750122
[I 14:12:52.956 NotebookApp] Use Control-C to stop this server and shut down all
↪kernels (twice to skip confirmatio
n).
[C 14:12:52.959 NotebookApp]

To access the notebook, open this file in a browser:
    file:///home/veraalva/.local/share/jupyter/runtime/nbserver-23251-open.html
Or copy and paste one of these URLs:
    http://localhost:8888/?token=eae6a8d42ad12d6ace23f5d0923bcec14d0f798127750122
    or http://127.0.0.1:8888/?token=eae6a8d42ad12d6ace23f5d0923bcec14d0f798127750122
```

The results of this analysis is <https://ftp.ncbi.nlm.nih.gov/pub/pm4ngs/examples/chipexo-single/>

Transcriptome Annotation pipeline for non-model organisms

9.1 Annotation Workflow



Attention: This annotation pipeline uses Google Cloud platform for all computing tasks. Users should have installed and configured [Cloud SDK](#).

9.2 GCP configuration

Cloud SDK should be installed and configured. The pipelines are based on the [Cloud Life Sciences \(beta\) API](#).

To install the beta commands use this command line:

```
localhost:~> gcloud components install beta
```

kubectl command is also required. This is used by Elastic-blast to configure the Kubernetes cluster on GCP.

It can be installed:

```
localhost:~> gcloud components install kubectl
```

9.3 Input requirements

The input requirement for this pipeline is the *Sample sheet* file.

9.4 Pipeline command line

The annotation project can be created using the following command line:

```
localhost:~> pm4ngs-transcriptome-annotation
usage: Generate a PM4NGS project for Transcriptome-Annotation data analysis [-h] [-v]
                                     --sample-sheet SAMPLE_SHEET [--config-file CONFIG_
                                     ↪FILE] [--copy-rawdata]

Generate a PM4NGS project for Transcriptome-Annotation data analysis:
error: the following arguments are required: --sample-sheet
```

Options:

- **sample-sheet:** Sample sheet with the samples SRA run IDs in the first column
- **config-file:** YAML file with configuration for project creation
- **copy-rawdata:** Copy the raw data defined in the sample sheet to the project structure. (The data can be hosted locally or in an http server)

9.5 Creating the annotation project

The **pm4ngs-transcriptome-annotation** command line executed with the **--sample-sheet** option will let you type the different variables required for creating and configuring the project. The default value for each variable is shown in the brackets.

```
localhost:~> pm4ngs-transcriptome-annotation --sample-sheet my-sample-sheet.tsv
Generating Transcriptome-Annotation data analysis project
author_name [Roberto Vera Alvarez]:
email [veraalva@ncbi.nlm.nih.gov]:
```

(continues on next page)

(continued from previous page)

```

project_name [my_ngs_project]: nopal-annotation
dataset_name [my_dataset_name]: PRJNA320545
Select sequencing_technology:
1 - single-end
2 - paired-end
Choose from 1, 2 [1]: 1
Copying file /home/veraalva/my-sample-sheet.tsv to /home/veraalva/nopal-annotation/
↳ data/PRJNA320545/sample_table.csv
Done

```

The **pm4ngs-rnaseq** command line will create a project structure as:

```

.
├── LICENSE
├── README.md
├── bin
│   └── gcp
│       ├── pipeline-blastn.json
│       ├── pipeline-contamination-cleanup.json
│       ├── pipeline-download-sra.json
│       ├── pipeline-read-assignment.json
│       ├── pipeline-split-fasta.json
│       ├── pipeline-transcriptome-annotation-rpsblast.json
│       ├── pipeline-transcriptome-annotation-rpstablastn.json
│       ├── pipeline-transcriptome-annotation.json
│       ├── pipeline-transcriptome-cleanup.json
│       ├── pipeline-transcriptome-fastq-cleanup.json
│       ├── pipeline-trimming-fastq-pe.json
│       ├── pipeline-trimming-fastq-se.json
│       └── pipeline-trinity.json
├── config
│   └── init.py
├── data
│   └── PRJNA320545
│       └── sample_table.csv
├── doc
├── notebooks
│   ├── 01 - Download and pre-processing quality control.ipynb
│   ├── 02 - Sample Trimming.ipynb
│   ├── 03 - Vector removal.ipynb
│   ├── 04 - Detecting Contamination.ipynb
│   ├── 05 - Trinity assembly.ipynb
│   ├── 06 - Vector Detection and data Partitioning.ipynb
│   ├── 07 - Transcriptome annotation.ipynb
│   ├── 08 - Transcript Annotation - Blast.ipynb
│   ├── 09 - Transcript Annotation - CDD.ipynb
│   ├── 10 - Transcript Submission to TSA.ipynb
│   ├── 11 - Alignment of raw read to the transcriptome.ipynb
│   └── 12 - Quantifying transcripts.ipynb
├── requirements
│   └── conda-env-dependencies.yaml
├── results
│   └── PRJNA320545
├── src
└── tmp

12 directories, 29 files

```

Note: RNASeq based project variables

- **author_name:** Default: [Roberto Vera Alvarez]
 - **email:** Default: [veraalva@ncbi.nlm.nih.gov]
 - **project_name:** Name of the project with no space nor especial characters. This will be used as project folder's name.
Default: [my_nginx_project]
 - **dataset_name:** Dataset to process name with no space nor especial characters. This will be used as folder name to group the data. This folder will be created under the **data/{{dataset_name}}** and **results/{{dataset_name}}**.
Default: [my_dataset_name]
 - **Select sequencing_technology:** Select one of the available sequencing technologies in your data.
Values: 1 - single-end, 2 - paired-end
-

9.6 Jupyter server

PM4NGS uses Jupyter as interface for users. After project creation the jupyter server should be started as shown below. The server will open a browser windows showing the project's structure just created.

```
localhost:~> jupyter notebook
```

9.7 Data processing

Start executing the notebooks from 01 to 12 waiting for each step completion.

9.8 Demo

PM4NGS includes a demo project that users can use to test the framework. It is pre-configured to use Docker as execution environment.

The annotated based demo process samples from the BioProject [PRJNA320545](#).

Use this command to create the project structure in your local computer

```
localhost:~> pm4ngs-transcriptome-annotation-demo
```

Once it finish, start the jupyter server and execute the notebooks as it is indicated on them

```
localhost:~> jupyter notebook
[I 14:12:52.956 NotebookApp] Serving notebooks from local directory: /home/veraalva
[I 14:12:52.956 NotebookApp] Jupyter Notebook 6.1.4 is running at:
[I 14:12:52.956 NotebookApp] http://localhost:8888/?
↪token=eae6a8d42ad12d6ace23f5d0923bcec14d0f798127750122
[I 14:12:52.956 NotebookApp] or http://127.0.0.1:8888/?
↪token=eae6a8d42ad12d6ace23f5d0923bcec14d0f798127750122
```

(continues on next page)

(continued from previous page)

```
[I 14:12:52.956 NotebookApp] Use Control-C to stop this server and shut down all
↪ kernels (twice to skip confirmatio
n).
[C 14:12:52.959 NotebookApp]

To access the notebook, open this file in a browser:
    file:///home/veraalva/.local/share/jupyter/runtime/nbserver-23251-open.html
Or copy and paste one of these URLs:
    http://localhost:8888/?token=eae6a8d42ad12d6ace23f5d0923bcec14d0f798127750122
    or http://127.0.0.1:8888/?token=eae6a8d42ad12d6ace23f5d0923bcec14d0f798127750122
```


CHAPTER 10

Create a new PM4NGS pipeline

The recommended way to create a new pipeline is modifying one of the current available pipelines. Before reading this instructions, the user should understand the basics concepts of the [Cookiecutter](#).

10.1 PM4NGS based pipeline folder structure

The PM4NGS based pipeline folder structure is:

```
.
├── LICENSE
├── README.md
├── cookiecutter.json
├── example
│   ├── pm4ngs_chipseq_demo_config.yaml
│   └── pm4ngs_chipseq_demo_sample_data.csv
├── hooks
│   ├── post_gen_project.py
│   └── pre_gen_project.py
├── {{cookiecutter.project_name}}
│   ├── LICENSE
│   ├── README.md
│   ├── bin
│   ├── config
│   │   └── init.py
│   ├── data
│   │   └── {{cookiecutter.dataset_name}}
│   ├── doc
│   ├── index.html
│   └── notebooks
│       ├── 00 - Project Report.ipynb
│       ├── 01 - Pre-processing QC.ipynb
│       ├── 02 - Samples trimming.ipynb
│       └── 03 - Alignments.ipynb
```

(continues on next page)

(continued from previous page)

```

├── 04 - Peak Calling and IDR.ipynb
├── 05 - Differential binding Detection.ipynb
├── requirements
│   └── conda-env-dependencies.yaml
├── results
│   └── {{cookiecutter.dataset_name}}
├── src
└── tmp

14 directories, 18 files

```

Folders:

- **cookiecutter.json**: the default cookiecutter JSON file where all input variables are defined
- **example**: this folder should include the sample sheet and yaml config file for the pipeline demo
- **hooks**: Cookiecutter pre and post execution scripts. PM4NGS uses these pre and post scripts to execute tasks required to configure the pipeline like cloning the CWL repo, copying the raw data and the sample sheet.

In the **post_gen_project.py** is where the CWL repository is defined. Users should modify this variable with the localization of their own CWL workflows.

```
CWL_WORKFLOW_REPO = 'https://github.com/ncbi/cwl-ngs-workflows-cbb'
```

This repo will be clone to a folder named **cwl** in the project **bin** directory.

- **{{cookiecutter.project_name}}**: this is a folder with the cookiecutter variable for project name.

This folder store all pipeline's files.

- **config/init.py**: This is the configuration file for the workflow that should be loaded at the beginning of all notebooks. Any function for creating tables, figures and plots not included in PM4NGS should be included in the **init.py** file. Users can also do PR on the PM4NGS repo to the package: `pm4ngs.jupyterngsplugin` to include their functions in the PM4NGS package but this is not required.
- **notebooks**: In this folder all jupyter notebook are stored.

The specific requirement for the notebooks are:

1. Run the **init.py** at the first cell of the notebook

```
%run ../config/init.py
```

2. Use relative paths all the time.
3. Global variables should be defined in the **init.py** file.
4. HTML code can be used to display information on the notebook:

To show a link to the sample sheet this can be used on the notebook cell:

```
filename = os.path.relpath(os.path.join(DATA, DATASET, 'sample_sheet.tsv'))
html_link = '<a href="/" target="_blank">sample_sheet.tsv</a>'.
format(filename.replace(' ', '%20'))
display(Markdown(html))
```

5. PM4NGS include functions to create links from an image or PDF file: Note that this functions require Poppler (<https://poppler.freedesktop.org/>) installed.


```
from pm4ngs.jupyterngsplugin.markdown.utils import get_link_image

width = 450
height = 450
filename = os.path.relpath(os.path.join(RESULTS, DATASET, 'dga', 'condition_
↳ POST_NACT_CRIS2_vs_PRE_NACT_NORMAL_deseq2_pca.pdf'))
html_link = get_link_image(filename, width, height, ' --- ')
display(Markdown(html))
```

- **requirements/conda-env-dependencies.yaml**: Define the conda packages and versions to run the pipeline.

10.2 CWL tools and workflows specifications

The workflow repository should include two main directories: **tools** and **workflows**. The first directory, **tools**, includes all computational tools used by the workflows in CWL format. The second folder, **workflows**, includes all workflows.

Each CWL tool should include two YAML files with suffixes **bioconda.yaml** and **docker.yaml** that are imported in the **hints** block.

```
hints:
  - $import: bwa-docker.yaml
  - $import: bwa-bioconda.yaml
```

See example [bwa-mem.cwl](#)

As it is indicated by its names, the **bioconda.yaml** file stores the software requirements for executing the CWL tool using Conda. The files specify the package name and version. The CWL runner will create a Conda environment and install the package, if it doesn't exist, at runtime.

```
class: SoftwareRequirement
packages:
  - package: 'bwa'
    version:
      - '0.7.17'
  specs:
    - https://anaconda.org/bioconda/bwa
```

See example [bwa-bioconda.yaml](#)

The **docker.yaml** file defines the Biocontainers docker image to be used. This image will be pulled by the CWL runner at runtime.

```
class: DockerRequirement
dockerPull: quay.io/biocontainers/bwa:0.7.17--h84994c4_5
```

See example [bwa-docker.yaml](#)

PM4NGS uses the [Biocontainers Registry](#) through its python interface named [bioconda2biocontainer](#) to keep CWL docker images defined in the **docker.yaml** file updated to its latest tag. The Bioconda package name and version defined in the **bioconda.yaml** file is passed as an argument to the [update_cwl_docker_from_tool_name](#) method in [bioconda2biocontainer](#) returning the latest docker image available for the tool. PM4NGS after cloning the CWL repository, reads the Bioconda package names and version from the **bioconda.yaml** files and updates all defined docker images to its latest tags modifying all **docker.yaml** files.

11.1 STAR Genomic Indexes

The **genome.fa** and **genes.gtf** files should be copied to the genome directory.

```
localhost:~> mkdir genome
localhost:~> cd genome
localhost:~> mkdir STAR
localhost:~> cd STAR
localhost:~> cwl-runner https://github.com/ncbi/cwl-ngs-workflows-cbb/blob/master/
↳tools/star/star-index.cwl --runThreadN 16 --genomeDir . --genomeFastaFiles ../
↳genome.fa --sjdbGTFfile ../genes.gtf
localhost:~> cd ..
localhost:~> tree
.
├── genes.gtf
├── genome.fa
└── STAR
    ├── chrLength.txt
    ├── chrNameLength.txt
    ├── chrName.txt
    ├── chrStart.txt
    ├── exonGeTrInfo.tab
    ├── exonInfo.tab
    ├── geneInfo.tab
    ├── Genome
    ├── genomeParameters.txt
    ├── Log.out
    ├── SA
    ├── SAindex
    ├── sjdbInfo.txt
    ├── sjdbList.fromGTF.out.tab
    ├── sjdbList.out.tab
    └── transcriptInfo.tab
```

(continues on next page)

(continued from previous page)

```
1 directory, 18 files
```

Here all files inside the directory **STAR** are created by the workflow.

11.2 BWA Genomic Indexes

The **genome.fa** file should be copied to the genome directory.

```
localhost:~> mkdir genome
localhost:~> cd genome
localhost:~> mkdir BWA
localhost:~> cd BWA
localhost:~> cwl-runner --no-container https://github.com/ncbi/cwl-ngs-workflows-cbb/
↳blob/master/tools/tools/BWA/bwa-index.cwl --sequences ../genome.fa
localhost:~> cd ..
localhost:~> tree
.
├── BWA
│   ├── genome.fa
│   ├── genome.fa.amb
│   ├── genome.fa.ann
│   ├── genome.fa.bwt
│   ├── genome.fa.pac
│   └── genome.fa.sa
└── genome.fa

1 directory, 7 files
```

11.3 Create BED from GTF

For generating a BED file from a GTF.

The **genes.gtf** file should be copied to the genome directory.

```
localhost:~> mkdir genome
localhost:~> cd genome
localhost:~> cwl-runner --no-container https://github.com/ncbi/cwl-ngs-workflows-cbb/
↳blob/master/workflows/UCSC/gtftobed.cwl --gtf genes.gtf
localhost:~> tree
.
├── genes.bed
├── genes.genePred
└── genome.gtf

0 directory, 3 files
```

Here the files **genes.bed** and **genes.genePred** are created from the workflow.

CHAPTER 12

PM4NGS available genomes

A set of genomes are available for download with the required files already configured. If the genome name is used, PM4NGS will include in the alignment notebook a cell to download, uncompress the genome file and create the aligner indexes.

Note: We recommend to move out of the first project the generated genome folder `data/{{genome_name}}/` to a different folder so it can be reutilized for other projects. This will avoid the extra time for the index creation and will save space.

CHAPTER 13

Introduction

PM4NGS was designed to generate a standard organizational structure for Next Generation Sequencing (NGS) data analysis. It includes a directory structure for the project, several Jupyter notebooks for data management and CWL workflows for pipeline execution.

Our work was inspired by a manuscript by Prof. William Noble in 2009: [A Quick Guide to Organizing Computational Biology Projects](#). We recommend reading this paper for a better understanding of the guiding principles of our project.

The project is composed of three main parts.

1. A project organizational structure which define standard files and directories for the project
2. Jupyter Notebooks as user interfaces for data management and visualization
3. CWL workflows that execute the data analysis

PM4NGS source code includes the templates used by **cookiecutter** to generate the project organizational structure and the Jupyter notebooks. The CWL workflows are defined in a separate GitHub project named: [cwl-ngs-workflows-cbb](#).

CHAPTER 14

Features

- NGS data integration, management and analysis uses Jupyter notebooks, CWL workflows and cookiecutter project templates
- Easy installation and use with a minimum command line interaction
- Data analysis CWL workflows executed from the Jupyter notebook with automatic failing detection and can be validated with published data
- CWL workflows and Jupyter Notebooks are ready for cloud computing
- Project reports and dynamic content creation after data processing using CWL workflows are included
- Optional use of Docker/Biocontainers or Conda/Bioconda for Bioinformatics tool installations and managements are also included

1. Vera Alvarez R, Pongor LS, Mariño-Ramírez L and Landsman D. PM4NGS, a project management framework for next-generation sequencing data analysis, GigaScience, Volume 10, Issue 1, January 2021, giaa141, <https://doi.org/10.1093/gigascience/giaa141>
2. Vera Alvarez R, Mariño-Ramírez L and Landsman D. Transcriptome annotation in the cloud: complexity, best practices, and cost, GigaScience, Volume 10, Issue 2, February 2021, giaa163, <https://doi.org/10.1093/gigascience/giaa163>
3. Vera Alvarez R, Pongor LS, Mariño-Ramírez L and Landsman D. Containerized open-source framework for NGS data analysis and management [version 1; not peer reviewed]. F1000Research 2019, 8(ISCBCOMM J):1229 (poster) (doi: 10.7490/f1000research.1117155.1)

CHAPTER 16

Help and Support

For query/questions regarding PM4NGS, please write veraalva@ncbi.nlm.nih.gov

For feature requests or bug reports, please open an issue on our [GitHub Repository](#).

Public Domain notice

National Center for Biotechnology Information.

This software is a "United States Government Work" under the terms of the United States Copyright Act. It was written as part of the authors' official duties as United States Government employees and thus cannot be copyrighted. This software is freely available to the public for use. The National Library of Medicine and the U.S. Government have not placed any restriction on its use or reproduction.

Although all reasonable efforts have been taken to ensure the accuracy and reliability of the software and data, the NLM and the U.S. Government do not and cannot warrant the performance or results that may be obtained by using this software or data. The NLM and the U.S. Government disclaim all warranties, express or implied, including warranties of performance, merchantability or fitness for any particular purpose.

Please cite NCBI in any work or product based on this material.