

# Repeat mapping


This tutorial will walk through setup, installation and running our reproducible peaks pipeline. This may be run after running the core eCLIP pipeline. For this tutorial, we will be spinning up an AWS ec2 instance with recommended hardware settings.

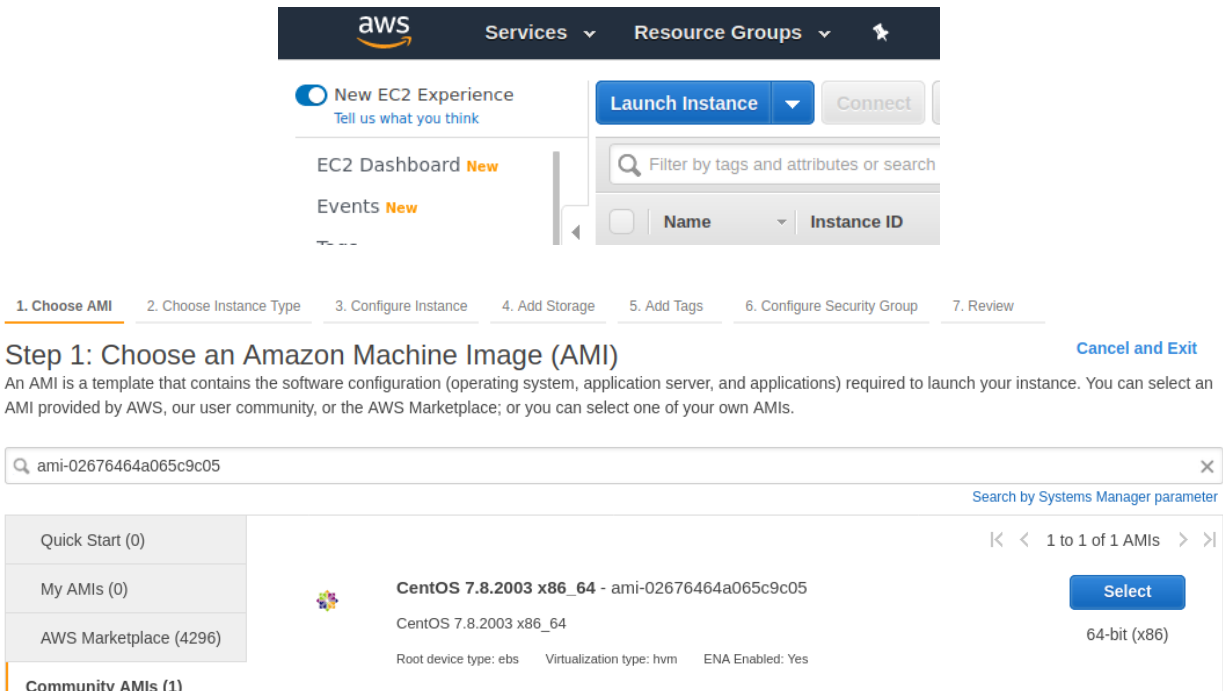
## 1. Setting up AWS

~5 - 20 minutes depending on your level of experience with Amazon web services

### Launch an EC2 instance

The easiest way to launch an instance is through the console. Login and click "Launch Instance" to start the setup wizard. Then, search for the official Centos7 AMI (**ami-02676464a065c9c05**) and configure your instance using the recommended settings. We recommend the **m5.2xlarge** (8 cores, 32G memory) which is plenty for processing a full run. **We will also be attaching 500G** of storage which is more than sufficient for everything we're about to do here. Practically speaking, this number of cores has been set to accommodate bowtie2 during the repeat mapping step, while the memory and storage requirement will scale depending on how many reads are ultimately mapped to repeat elements.

 **Note:** at some point you will need to set up your account credentials and create/download a .pem key. You may read the documentation from AWS [here](#) and will need to refer to this key when you login (see: [Login!](#))



The screenshot shows the AWS Management Console 'Launch Instance' wizard. The 'Choose AMI' step is selected, showing a search bar with the AMI ID 'ami-02676464a065c9c05'. Below the search bar, there are navigation tabs for 'Quick Start (0)', 'My AMIs (0)', 'AWS Marketplace (4296)', and 'Community AMIs (1)'. The 'Community AMIs (1)' tab is active, displaying a search result for 'CentOS 7.8.2003 x86\_64' with the AMI ID 'ami-02676464a065c9c05'. The result includes details like 'Root device type: ebs', 'Virtualization type: hvm', and 'ENA Enabled: Yes'. A 'Select' button is visible next to the result.

We'll be using CentOS 7.8 although other flavors may be used.

## Step 4: Add Storage

Your instance will be launched with the following storage device settings. You can attach additional EBS volumes and instance store volumes to your instance, or edit the settings of the root volume. You can also attach additional EBS volumes after launching an instance, but not instance store volumes. [Learn more](#) about storage options in Amazon EC2.

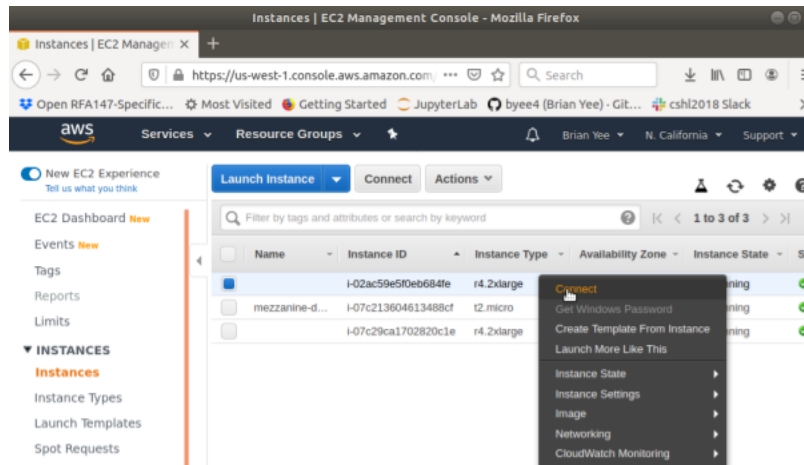
Volume Type	Device	Snapshot	Size (GiB)	Volume Type	IOPS	Throughput (MB/s)	Delete on Termination	Encryption
Root	/dev/sda1	snap-088131aa507fa0d0a	500	General Purpose	1500 / 3000	N/A	<input checked="" type="checkbox"/>	Not Encrypted

[Add New Volume](#)

**(Optional) Setup a "Cost budget" on AWS to remind you of any unexpected charges or caps!**

## Login!

The easiest way to connect is to locate and right click your running instance to connect.



From the console, right-click on your running instance and select 'connect'

## Connect to your instance



### Connection method

- A standalone SSH client
- Session Manager
- EC2 Instance Connect (browser-based SSH connection)

### To access your instance:

1. Open an SSH client. (find out how to [connect using PuTTY](#))
2. Locate your private key file (tmp.pem). The wizard automatically detects the key you used to launch the instance.
3. Your key must not be publicly viewable for SSH to work. Use this command if needed:

```
chmod 400 tmp.pem
```

4. Connect to your instance using its Public DNS:

```
ec2-50-18-8-227.us-west-1.compute.amazonaws.com
```

### Example:

```
ssh -i "tmp.pem" root@ec2-50-18-8-227.us-west-1.compute.amazonaws.com
```

Please note that in most cases the username above will be correct, however please ensure that you read your AMI usage instructions to ensure that the AMI owner has not changed the default AMI username.

If you need any assistance connecting to your instance, please see our [connection documentation](#).

Open your terminal (Mac or Linux) and connect to your ec2 instance using the appropriate DNS and credentials. Windows users may need to download an ssh client (ie. PuTTY, git bash).



**Note:** You may need to change the example command to point to the correct location of your "pem" key, as well as login using "centos" instead of "root." So in this example, we would open our terminal and type:

```
ssh -i "tmp.pem" centos@ec2-50-18-8-227.us-west-1.compute.amazonaws.com
```

## 2. Setting up your initial environment

~0 - 5 minutes

Since we're starting from scratch, we will need to install some basic packages.

Update packages:

```
sudo yum install -y update;
```

Install the screen tool to preserve your session and wget for downloading:

```
sudo yum install -y install screen wget;
```

Install the C compiler (gcc):

```
sudo yum -y group install "Development Tools" # installs gcc
```

Install pip, python, and nodejs which are essential for CWL:

```
wget https://repo.anaconda.com/miniconda/Miniconda3-latest-Linux-x86_64.sh
```

`sh Miniconda3-latest-Linux-x86_64.sh` (and follow onscreen instructions and click "yes" to initialize conda). Use defaults wherever prompted.

Re-run your `~/.bashrc` file and verify successful python/pip installation:

```
source ~/.bashrc;
```

```
which pip;
```

```
which python;
```

### 3. Install Docker and CWL

~0 - 5 minutes

Docker will help us pull pre-built containers with the correct environments and tools for each step of this pipeline. Most of these tools (such as samtools, bedtools, perl) are fairly straightforward on their own to install, but using these containers will simplify the tutorial.

We've elected to use a common workflow language ([CWL](#)) to define each step of our pipeline.

#### Install Docker:

Follow the instructions for installing Docker for [CentOS](#). At the time of writing, the basic commands to install are as follows:

```
sudo yum install -y yum-utils
sudo yum-config-manager \
  --add-repo \
  https://download.docker.com/linux/centos/docker-ce.repo
sudo yum install docker-ce docker-ce-cli containerd.io ### select "y"es to install Docker.
```

```
sudo systemctl start docker
```

**⚠** Make sure you have started Docker with `sudo systemctl start docker` and verify installation with `sudo docker run hello-world`

By default, this will install Docker however we will need to add our user (centos) to the docker group. So we must add ourselves to the group (docker):

```
sudo groupadd docker
sudo usermod -aG docker ${USER}
```

Log out and log back in to ensure you've been added to the `docker` group.

Verify that (centos) belongs to (docker) with `groups`:

```
(base) brian@brian-Inspiron-7577:~$ ssh -i "~/.aws/tmp.pem" centos@ec2-52-53-224-225.us-west-1.compute.amazonaws.com
Last login: Fri Jun 26 18:40:36 2020 from cpe-70-95-26-41.san.res.rr.com
(base) [centos@ip-172-31-29-59 ~]$ groups
centos adm wheel systemd-journal docker
(base) [centos@ip-172-31-29-59 ~]$
```

Make sure your user belongs to the "docker" group.

And you should be able to run docker without sudo. Try it! `docker run hello-world`:

```
(base) [centos@ip-172-31-29-59 cwltool_interm]$ docker run hello-world
Hello from Docker!
This message shows that your installation appears to be working correctly.
```

#### Install CWL:

~0 - 5 minutes

We can use pip to install cwl. Here we are using version (3.0.20200706173533).

```
pip install cwltool==3.0.20200706173533
```

You can verify installation with: `cwltool --help`

## 4. Download data

### Download example eCLIP data into your home directory

~5 - 10 minutes

Download the repeat mapping reference datasets

```
cd ~
wget https://external-collaborator-data.s3-us-west-1.amazonaws.com/reference-data/repeat-family-mapping-grch38.tar.gz
tar -xvf repeat-family-mapping-grch38.tar.gz
```

Download example input data

- This will be similar to direct outputs from our core pipeline and require no pre-processing.

```
cd ~
wget https://external-collaborator-data.s3-us-west-1.amazonaws.com/reference-data/example_data_for_repeat_mapping_hg38.tar.gz
tar -xvf example_data_for_repeat_mapping_hg38.tar.gz
```

### Download the Repeat-family mapping pipeline & add relevant directories to your `$PATH`

```
cd ~
git clone https://github.com/yeolab/repetitive-element-mapping
export PATH=/home/centos/repetitive-element-mapping/cwl:/home/centos/repetitive-element-mapping/wf:$PATH
```

## 5. Run the pipeline

~2-3 hours



**Note:** best to start a [screen session](#) - processing a sample can take up to a full day!

If you've made this far, you should be able to:

1. Verify that Docker and CWL are installed
2. Verify that you have the reference and example data downloaded
3. Verify that the pipeline files can be found in your `$PATH`

At this point, you should be able to navigate to the '`examples/`' directory and run either the `seCLIP_example.yaml` and `peCLIP_example.yaml` manifests (the paths should be set up to point to existing files as outlined from this tutorial, but you may need to double check that the paths do exist).

```
cd /home/centos/repetitive-element-mapping/examples
./repeat_mapping_SE.yaml
```

or

```
./repeat_mapping_PE.yaml
```

### Troubleshooting:

Since our own cluster disallows Docker (common among HPC due to permissions restrictions), it's possible that we may accidentally overwrite the parameter that runs Docker containers. When this happens, you may see your run fail immediately with a message telling you that the pipeline cannot find the tools:

```
ERROR '$TOOL' not found: [Errno 2] No such file or directory:
```

Fortunately, this is an easy fix. Using your favorite text editor, locate the cwl command within `eCLIP_repelement_SE_singleNode` (for single-end eCLIP) or `eCLIP_repelement_PE_singleNode` (for paired-end eCLIP) and remove the `--no-container` parameter:

```
if [ $CWLRunner == cwltool ]
then
  cwltool --debug \
    --outdir ${OUTDIR} \
    --cachedir ${INTERM} \
    --no-container \
    ${WORKFLOW_HOME}/cwl/${WORKFLOW}.cwl \
    ${JOB_FILE} \
    2>${JOB_HOME}/${JOB_NAME}/${PIPELINE}_${JOB_NAME}_LOG.txt | tee ${JOB_HOME}/${JOB_NAME}/${PIPELINE}_${JOB_NAME}_OUTPUT.json
#2>&1 | tee ${JOB_HOME}/${JOB_NAME}/ECLIPREPMAP-JOB-LOG.txt
```

Since this tutorial runs the pipeline with cwltool, you can scroll down to where the command is invoked.

```
if [ $CWLRunner == cwltool ]
then
  cwltool --debug \
    --outdir ${OUTDIR} \
    --cachedir ${INTERM} \
    ${WORKFLOW_HOME}/cwl/${WORKFLOW}.cwl \
    ${JOB_FILE} \
    2>${JOB_HOME}/${JOB_NAME}/${PIPELINE}_${JOB_NAME}_LOG.txt | tee ${JOB_HOME}/${JOB_NAME}/${PIPELINE}_${JOB_NAME}_OUTPUT.json
#2>&1 | tee ${JOB_HOME}/${JOB_NAME}/ECLIPREPMAP-JOB-LOG.txt
```

Removing the `--no-container` parameter will enable pulls from Docker.

## No mandatory DockerRequirement Error:

```
ERROR Workflow error:
No mandatory DockerRequirement, yet path is outside the designated output directory, also know as $(runtime.outdir): MapperEnt(resolved='/h
```

We've noticed that cwltool versions  $\geq 3.0.20200807132242$  is stricter about its Docker requirement (our pipelines suggest "hints" but does not require docker since some clusters (including our own) do not have the capability of running Docker.

### How to fix:

**Option 1.** Install a previous version of cwltool (ie. `pip install cwltool==3.0.20200706173533`)

**Option 2.** Modify the tool definition (`.cwl`) files to make Docker mandatory. Below is an example using the first 20 lines of the `eclip/cwl/makebigwigfiles_PE.cwl` tool. You will simply need to move the Docker definition from `hints` to `requirements`, however be careful to preserve spacing, as these files should be of YAML format (**make sure the ticks line up**):

```
#!/usr/bin/env cwltool

cwlVersion: v1.0

class: CommandLineTool

requirements:
- class: ResourceRequirement
  coresMin: 1
  ramMin: 8000
- class: InitialWorkDirRequirement
  listing:
    - entry: $(inputs.bam)
      writable: true

hints:
- class: DockerRequirement
  dockerPull: brianjee/makebigwigfiles:0.0.3

baseCommand: [makebigwigfiles]
```

## References:

Van Nostrand, Eric L., et al. "Principles of RNA processing from analysis of enhanced CLIP maps for 150 RNA binding proteins." *Genome biology* 21 (2020): 1-26.

Van Nostrand, Eric L., et al. "A large-scale binding and functional map of human RNA-binding proteins." *Nature* 583.7818 (2020): 711-719.

Bao, Weidong, Kenji K. Kojima, and Oleksiy Kohany. "Repbases Update, a database of repetitive elements in eukaryotic genomes." *Mobile Dna* 6.1 (2015): 11.