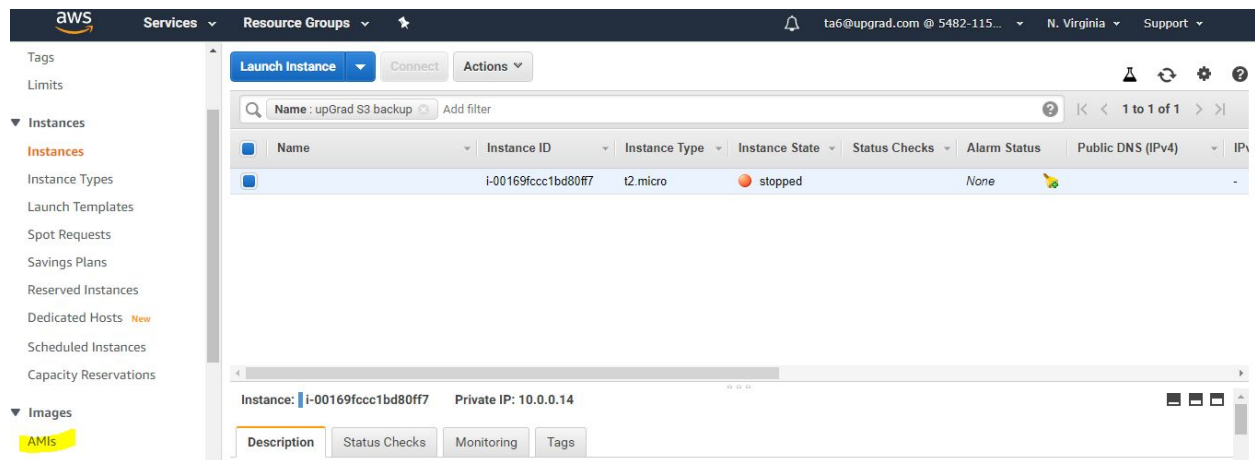# Instance Configuration

Go to the EC2 page on the AWS Management console to launch the instance. As part of this assignment, you are expected to work with the following AMI:
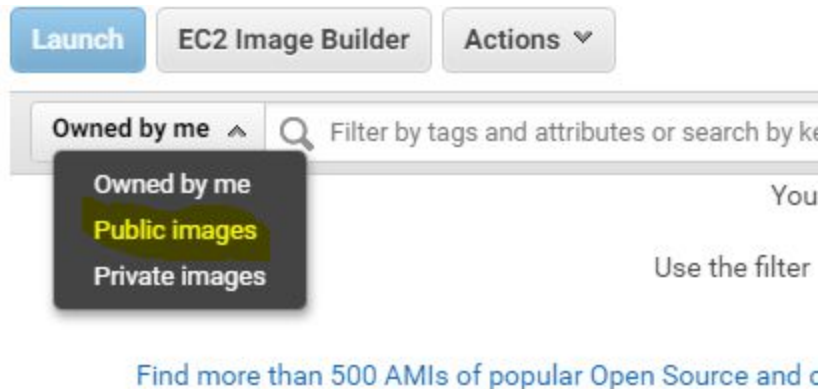
**AMI Name** - spark-jupyter
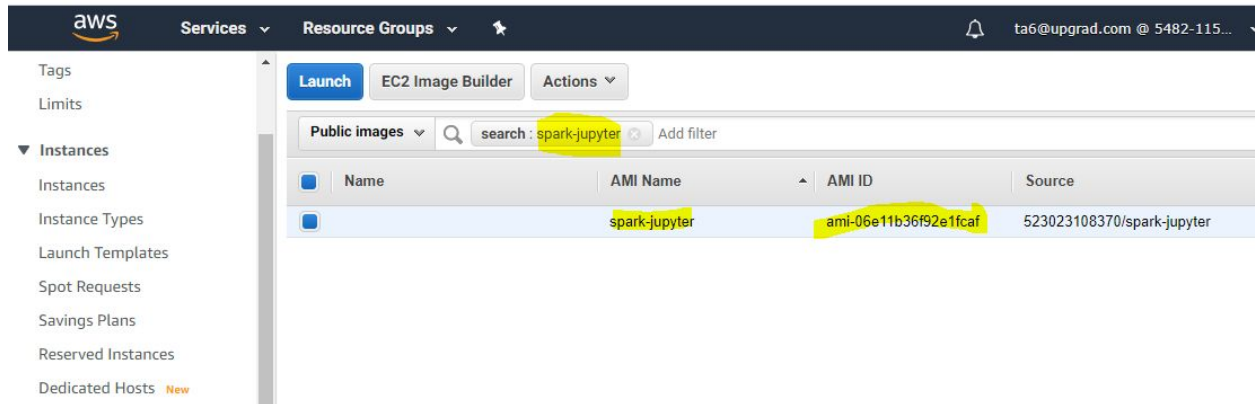**AMI ID** - ami-0437a7d30763297af

To do so, follow the steps given below:

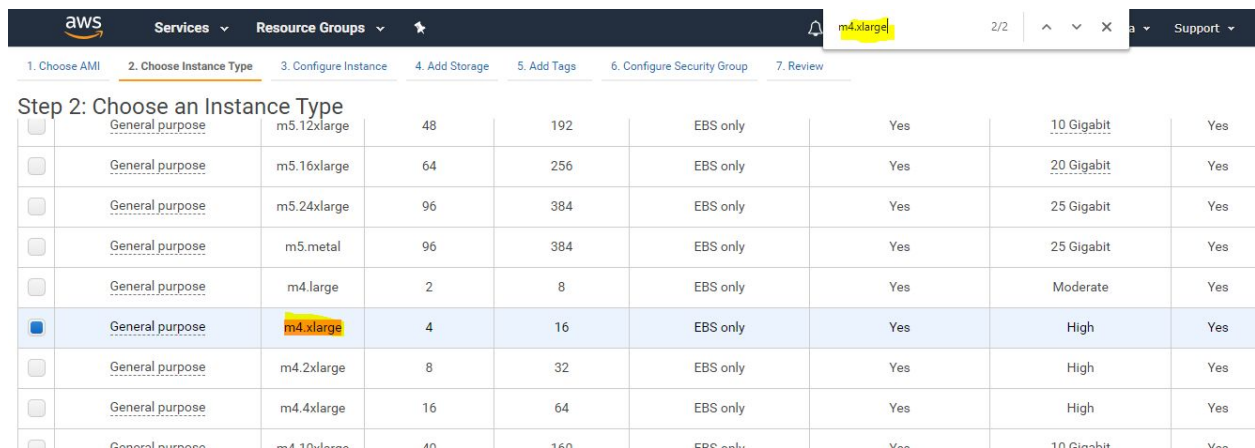1. Go to the EC2 dashboard and open the **AMI** section.



2. Under the AMI section, select the option **Public images** and search for the ami name mentioned above: **spark-jupyter**

3. Select the AMI image and press the **Launch** button. You will move to the next step towards launching the EC2 instance.

4. You are expected to work with the following instance type: **m4.xlarge**. Select the instance-type from the provided list and click on Next.



5. Keep the default following VPC settings under the next section as shown in the image below:

## Step 3: Configure Instance Details

Configure the instance to suit your requirements. You can launch multiple instances from the same AMI, request Spot instances to instance, and more.

| | | |
|---|---|---|
| Number of instances ⓘ | 1 | Launch into Auto Scaling Group ⓘ |
| Purchasing option ⓘ | ☐ Request Spot instances | |
| Network ⓘ | vpc-2bd4f851 (default) | ↻ Create new VPC |
| Subnet ⓘ | No preference (default subnet in any Availability Zone ↕) | Create new sub |
| Auto-assign Public IP ⓘ | Use subnet setting (Enable) ↕ | |
| Placement group ⓘ | ☐ Add instance to placement group | |

Click on Next once done.

6. Next, you are expected to increase the volume size of the instance to **50 GBs**. Keep the Volume Type as a **General Purpose SSD**.

| Volume Type ⓘ | Device ⓘ | Snapshot ⓘ | Size (GiB) ⓘ | Volume Type ⓘ | IOPS ⓘ | Throughput (MB/s) ⓘ | Delete on Termination ⓘ | Encryption ⓘ |
|---|---|---|---|---|---|---|---|---|
| Root | /dev/xvda | snap-0467c9f2f612397c1 | 50 | General Purpose SSD (gp2) ∨ | 150 / 3000 | N/A | ☑ | Not Encrypt ▼ |

**Add New Volume**

Free tier eligible customers can get up to 30 GB of EBS General Purpose (SSD) or Magnetic storage. Learn more about free usage tier eligibility and usage restrictions.

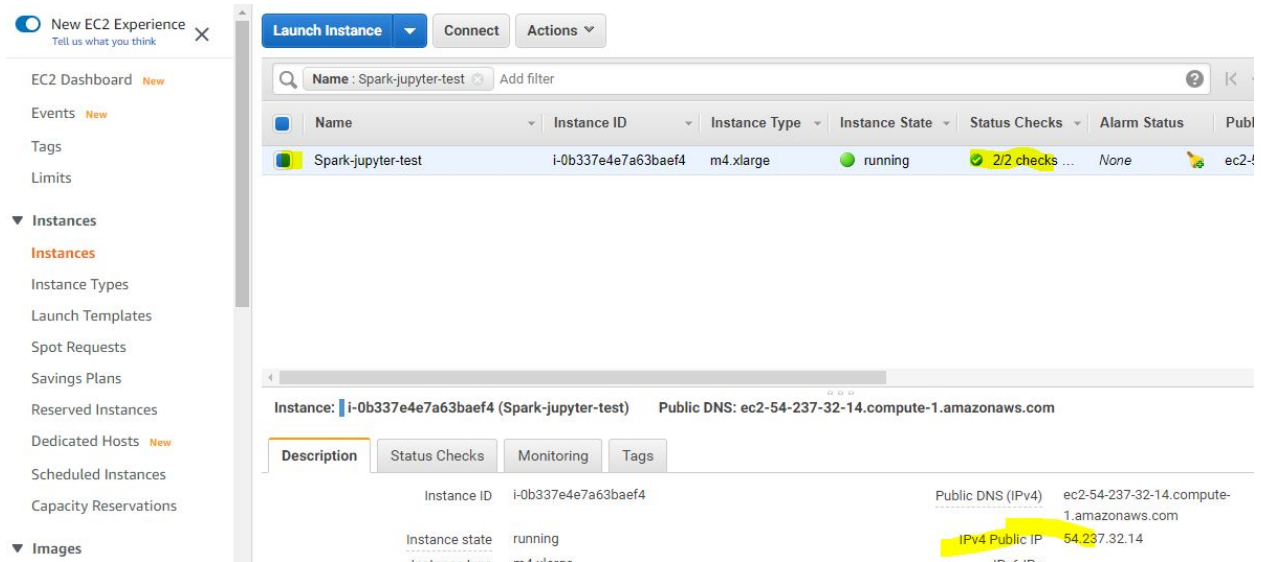Cancel   Previous   **Review and Launch**   Next: Add Tags

7. Under the next section, you can add the tags to your instance. This is not a compulsory step and can be ignored.

8. Under the Configure Security Group section, you must create a new security group that allows the access of the following ports to your IP address: **22** (SSH), **8888** (Custom TCP - Jupyter) and **4040** (Custom TCP - Spark UI).
   **Note**: Remember that this IP must be checked every time you launch the instance.

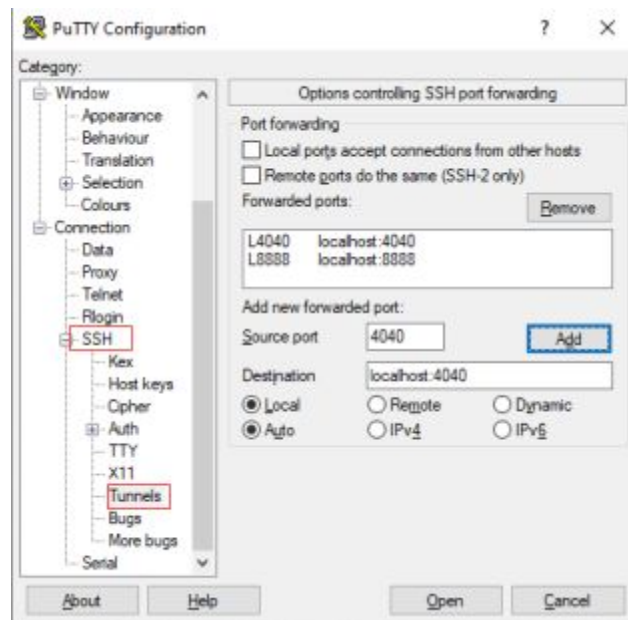   After the group is created, you can click on **Review and Launch**.

9. Next, you will be asked for the key-pair for the instance. You can select an existing key-pair or create a new one if required. Always keep this key-pair safe as it will provide a security layer to the instance.

# Launching the Jupyter Notebook

1. You can use the PuTTYgen software to convert the PEM file into a PPK file if you are using a Windows machine. Ignore this if you are working on a Linux machine or MAC.

2. To connect to the instance:
   a. Windows: Use the PuTTY software to connect to the instance.



You must provide the correct path for the PPK file under the **Auth** section of the **SSH** tab. Moreover, you must also add the required ports under the **Tunnel** section as shown below.
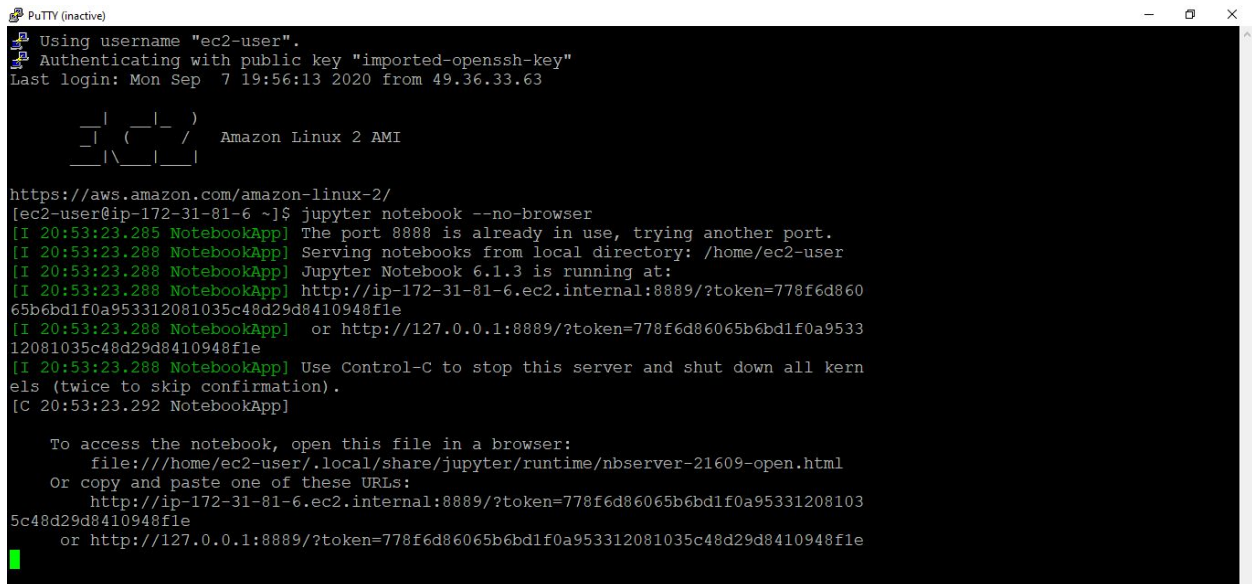
b. Linux/MAC: To connect using your instance's public DNS name, enter the following command:

ssh -i /path/my-key-pair.pem ec2-user@instance-public-dns-name

3. The instance is pre-loaded with all the required tools:
   a. Python
   b. Jupyter
   c. Spark

Hence, you are expected to launch the instance directly using the following command:

**jupyter notebook --no-browser**



You can copy the provided IP address for launching the Jupyter notebook in your local browser. (Here: http://127.0.0.1:8888/?token=778f6d86065b6bd1f0a953312081035c48d29d8410948f1e)

Remember that if you are hosting multiple sessions (8889 in the above image), you must provide access to all the ports under security groups and the tunnel section in PuTTY. You can check the sessions using the following command:

**jupyter notebook list**

4. Once the Jupyter notebook is live, you can use it to write your Python queries directly.



However, to use Spark, you must run the following commands in the notebook to set the environment variables:

```
import os
import sys
os.environ["PYSPARK_PYTHON"]="/usr/bin/python3"
os.environ["PYSPARK_DRIVER_PYTHON"]="/usr/bin/python3"
os.environ["PYSPARK_DRIVER_PYTHON_OPTS"]="notebook --no-browser"
os.environ["JAVA_HOME"] = "/usr/java/jdk1.8.0_161/jre"
os.environ["SPARK_HOME"] = "/home/ec2-user/spark-2.4.4-bin-hadoop2.7"
os.environ["PYLIB"] = os.environ["SPARK_HOME"] + "/python/lib"
sys.path.insert(0, os.environ["PYLIB"] + "/py4j-0.10.7-src.zip")
sys.path.insert(0, os.environ["PYLIB"] + "/pyspark.zip")
```

Now, your instance is ready to perform all the expected tasks.