

# Instructions for installing virtual box VM on windows and mac for running Hadoop Lab Software

## Minimum System Requirements:

- 8 GB of RAM
- 20 GB of hard-drive space
- 2 GHz dual core processor
- Windows or Mac

If you do not meet the requirements above or if you use an Apple Silicon device, please use the remote desktop solution; the instruction for that is available on LMS, in the “VM Installation” section.

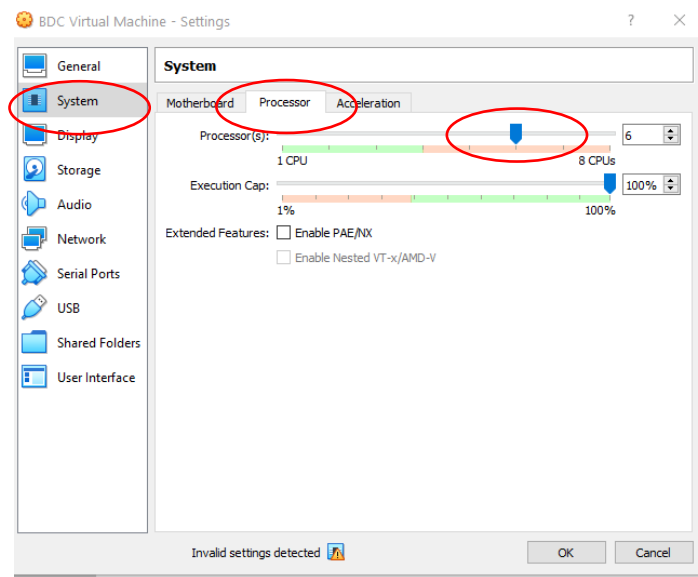
## Installation Guide

Here are the steps to follow to install the required lab software for labs 2 onwards for CSE3BDC.

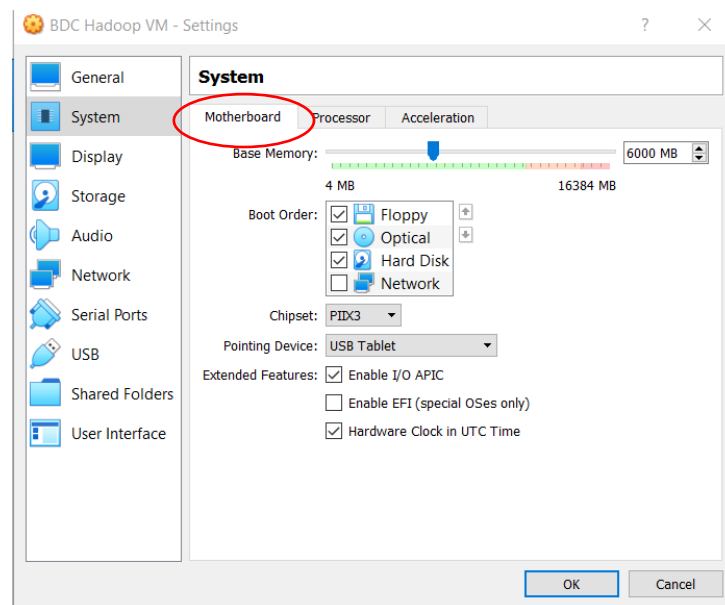
1. Download and unzip the image file from the following link:
  - a. [BDC VM Image](#)**Install VirtualBox.**

Please download the latest version of virtual box from here:

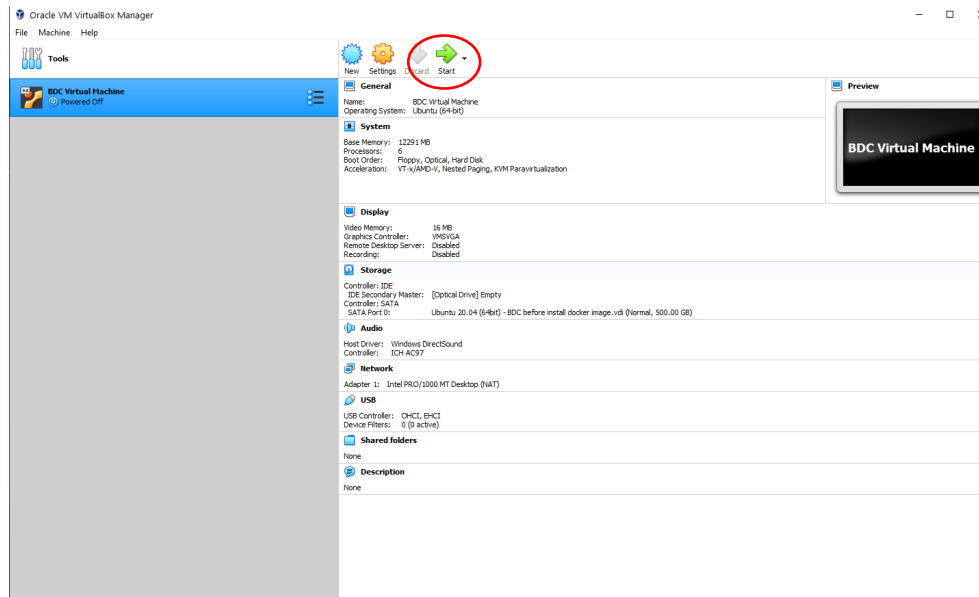
  - o <https://www.virtualbox.org/wiki/Downloads>
2. Go to the unzipped folder *BDC Hadoop VM* folder.
3. Double click on the *BDC Hadoop VM.vbox* file.
4. Go to System settings, and under “Processor” tab, change the number of processors to the following. If you have max of 4 CPUs, then choose 2 processors. If you have a max of 8 CPUs, change to 4. If you have max of 6 CPUs then change to 3. If you find these settings make the rest of your system run really slow then you can reduce the number of CPUs assigned to the VM to be within the green region.



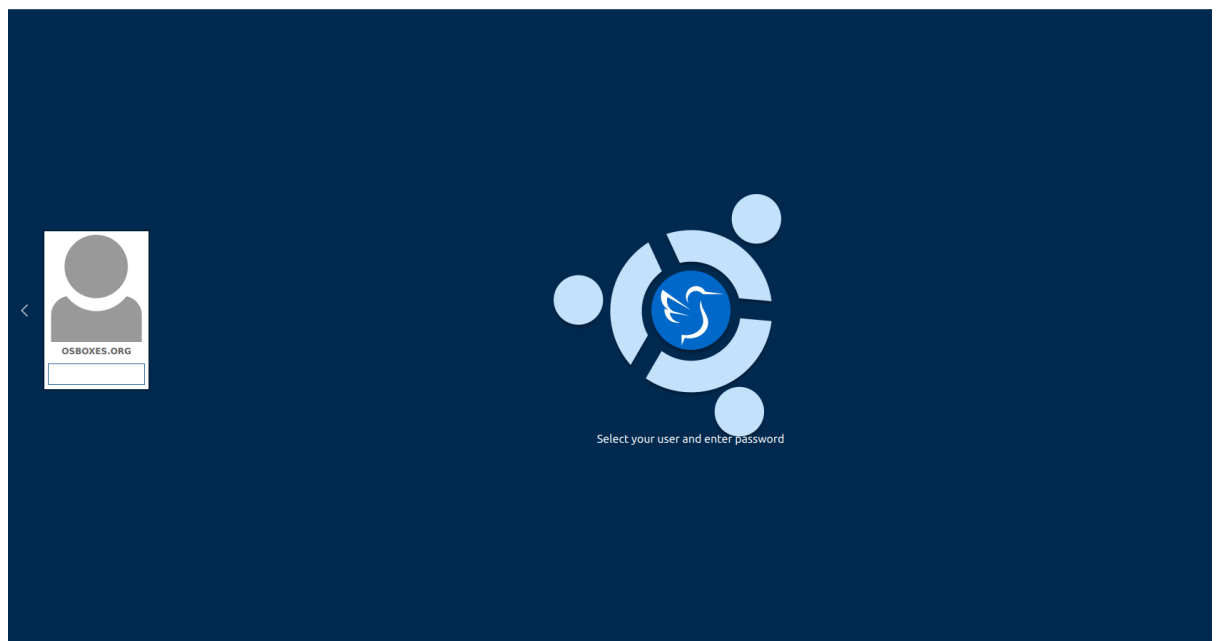
- Next in the motherboard tab select the amount of RAM you want the VM to use. If you have 8 GB of RAM, choose 5 GB for your memory size; if you have 16 GB, choose 8 GB. If you're not sure what your memory size is, just follow the rule of thumb of dragging the slider to be near the left edge of the green part (like below).



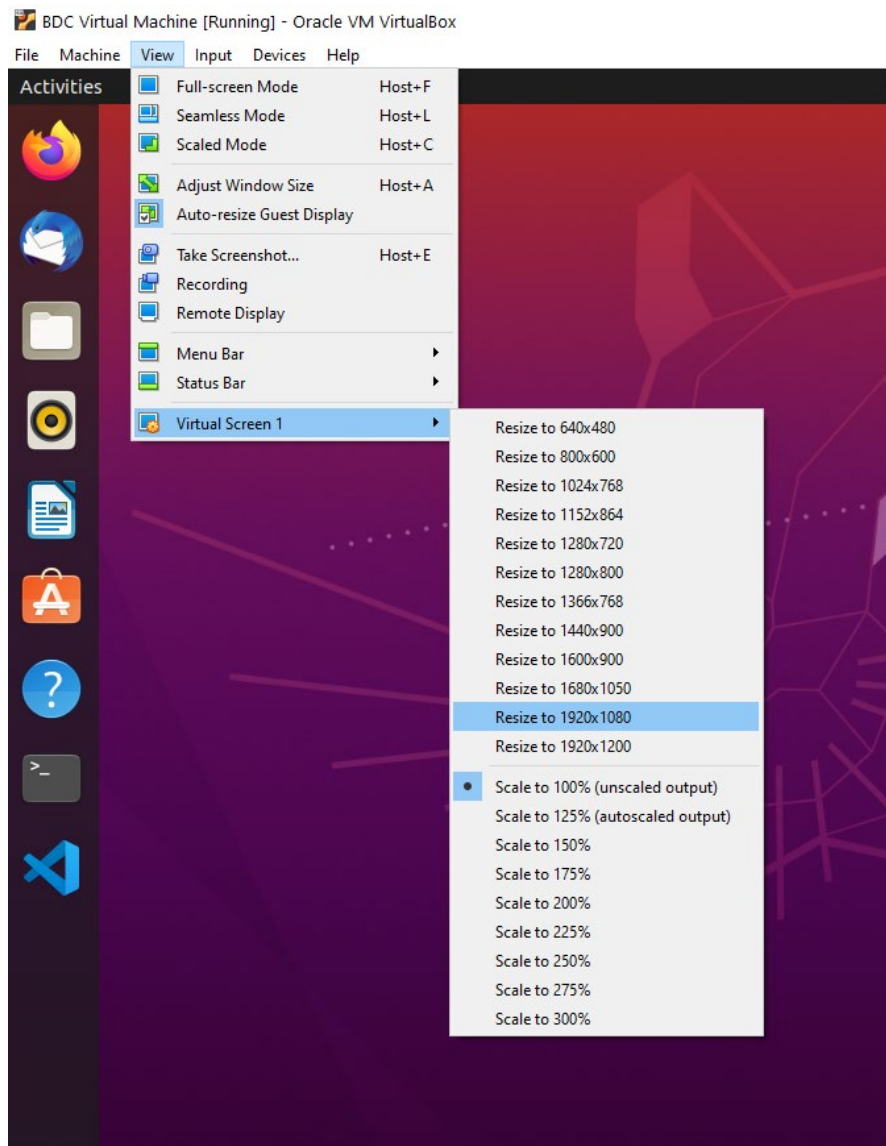
- Next press OK and then start the VM.



- If it asks you which operating system to boot into, please choose Ubuntu. For most people this question will not be asked but if it is asked then choose Ubuntu.
- Login to the osboxes.org account with the password of “osboxes.org”.



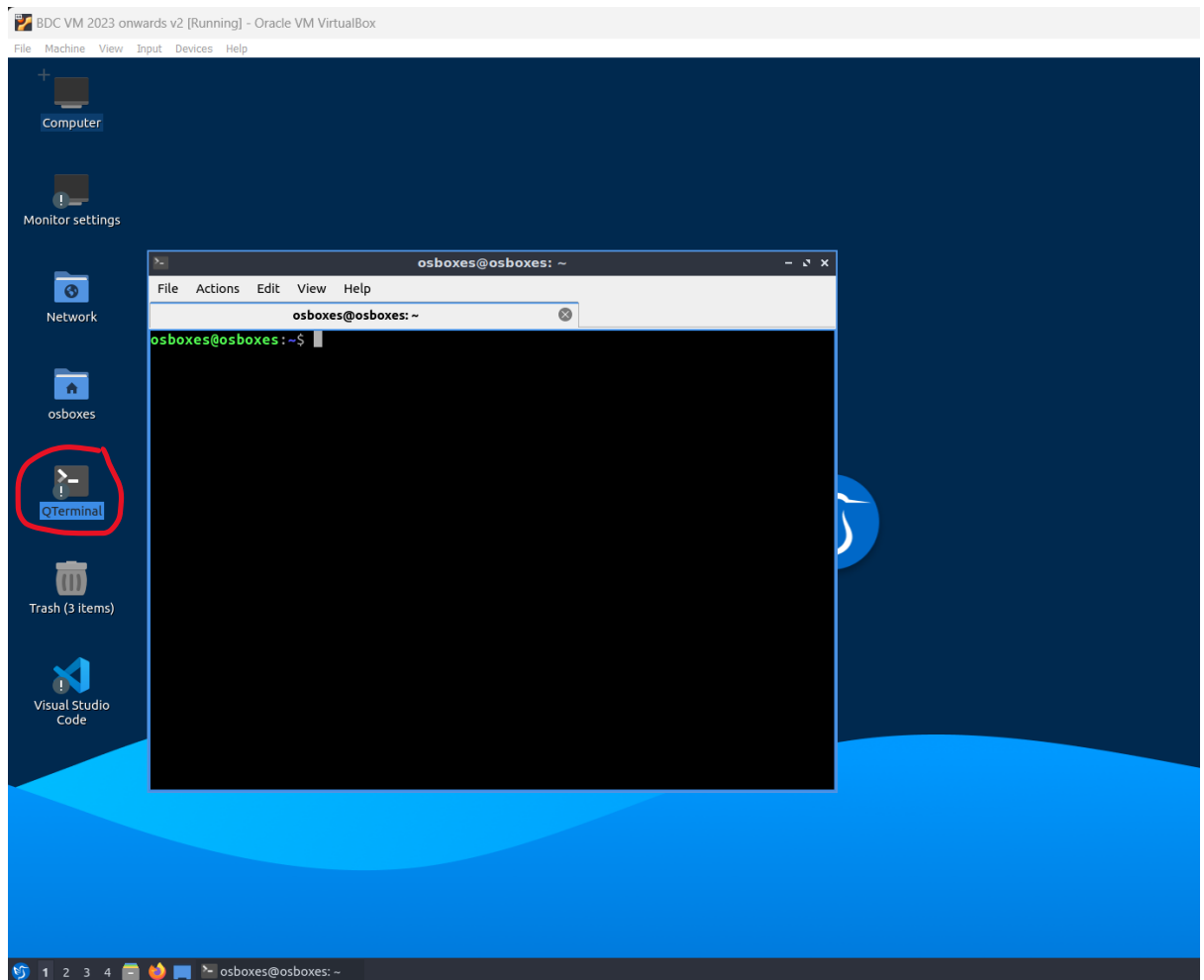
- Sometimes when you resize the VM window **the mouse controls go a bit funny**, if that happens go to *view -> virtual screen 1 -> select* a resize that matches your screen (e.g. 1920 x1080). This will fix any mouse control unaligned issues.



10. Sometimes the virtual machine will ask if you want to update or upgrade the operating. **PLEASE DO NOT UPDATE or UPGRADE** any software packages. In the past, I have updated some software and was then not able to log into the image again!  
**If you really want to update, please first make a backup of the vdi file or a snapshot of the VM and then update. In case you are not able to login again you can at least revert back to your previous version.**

11. Next open a QTerminal and move to the **docker-hive-spark** directory to fire up docker.

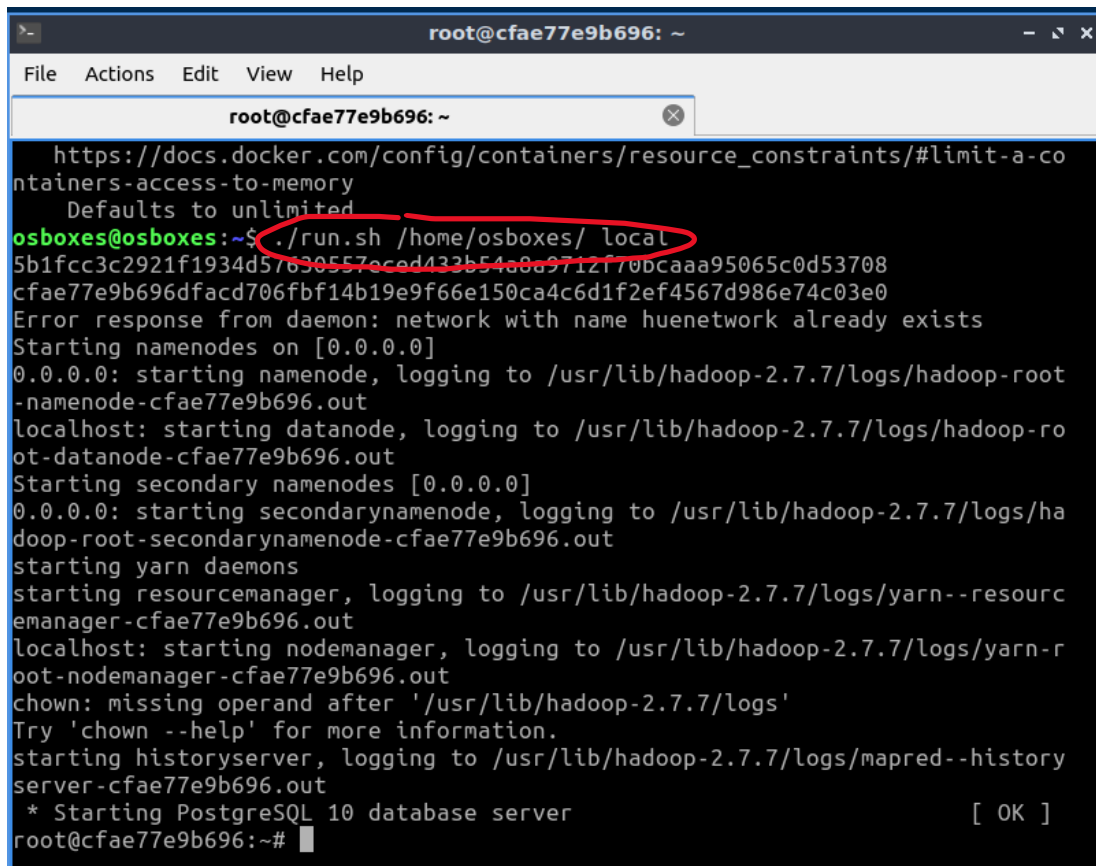
To open a QTerminal, click on the terminal button shown below:



12. Then type the following to see what options you can use when running docker:  
./run.sh

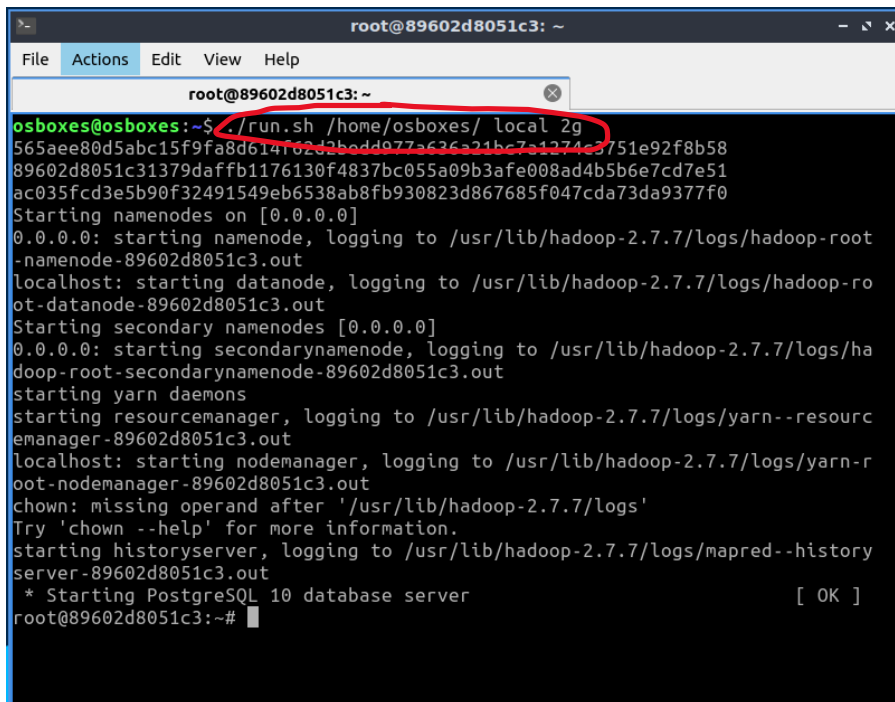
```
osboxes@osboxes: ~  
File Actions Edit View Help  
osboxes@osboxes: ~  
osboxes@osboxes:~$ ./run.sh  
Usage:  
./run.sh <directory-to-map> [MODE] [MEM]  
  
MODE is optional and must be one of "yarnmode"/"y" or "local".  
Defaults to "local"  
MEM is optional and specifies how much RAM is given to the main docker container, provide in the form shown in:  
https://docs.docker.com/config/containers/resource\_constraints/#limit-a-containers-access-to-memory  
Defaults to unlimited  
osboxes@osboxes:~$
```

Use “`./run.sh /home/osboxes/ local`” if you want to run local mode. This will run a lot faster than yarnmode. I suggest you use this mode for everything, except for one of the early exercises in lab 2 where you need to look at the Map Reduce jobs where you should use the `yarnmode` instead.



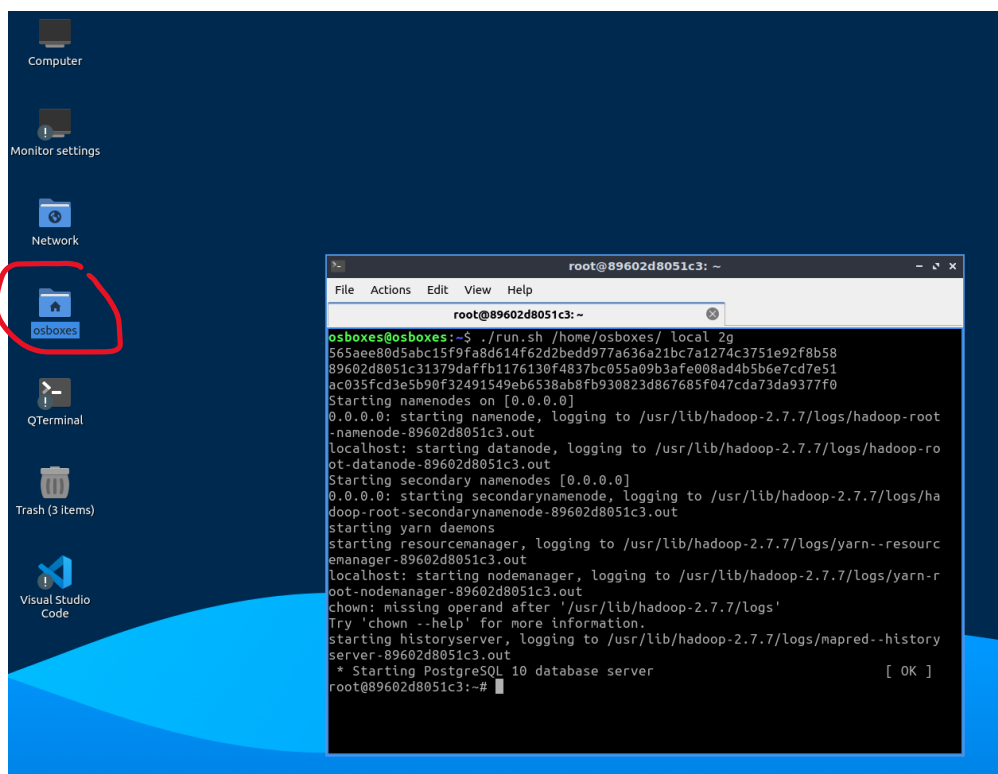
```
root@cfae77e9b696: ~  
File Actions Edit View Help  
root@cfae77e9b696: ~  
https://docs.docker.com/config/containers/resource_constraints/#limit-a-co  
ntainers-access-to-memory  
Defaults to unlimited  
osboxes@osboxes:~$ ./run.sh /home/osboxes/ local  
5b1fcc3c2921f1934d57630557ecad433b54e0a9712f70bcaaa95065c0d53708  
cfae77e9b696dfacd706fbf14b19e9f66e150ca4c6d1f2ef4567d986e74c03e0  
Error response from daemon: network with name huenetwork already exists  
Starting namenodes on [0.0.0.0]  
0.0.0.0: starting namenode, logging to /usr/lib/hadoop-2.7.7/logs/hadoop-root-  
namenode-cfae77e9b696.out  
localhost: starting datanode, logging to /usr/lib/hadoop-2.7.7/logs/hadoop-ro  
ot-datanode-cfae77e9b696.out  
Starting secondary namenodes [0.0.0.0]  
0.0.0.0: starting secondarynamenode, logging to /usr/lib/hadoop-2.7.7/logs/ha  
dooop-root-secondarynamenode-cfae77e9b696.out  
starting yarn daemons  
starting resourcemanager, logging to /usr/lib/hadoop-2.7.7/logs/yarn--resourc  
emanager-cfae77e9b696.out  
localhost: starting nodemanager, logging to /usr/lib/hadoop-2.7.7/logs/yarn-r  
oot-nodemanager-cfae77e9b696.out  
chown: missing operand after '/usr/lib/hadoop-2.7.7/logs'  
Try 'chown --help' for more information.  
starting historyserver, logging to /usr/lib/hadoop-2.7.7/logs/mapred--history  
server-cfae77e9b696.out  
* Starting PostgreSQL 10 database server [ OK ]  
root@cfae77e9b696:~#
```

Finally, the last parameter is if you want to limit the amount of RAM you allocate to the *docker* image. I suggest you select at least 2GB. If you have allocated 4GB of RAM to the VM then I would give between 2GB to 3GB to docker. If you selected 6GB of RAM to the VM then I would give 4GB. Or you can just leave this parameter blank. In the example below I set the VM to use just 2GB of RAM. A tutor was able to do all the labs and the assignment using just 1GB of RAM assigned to docker.

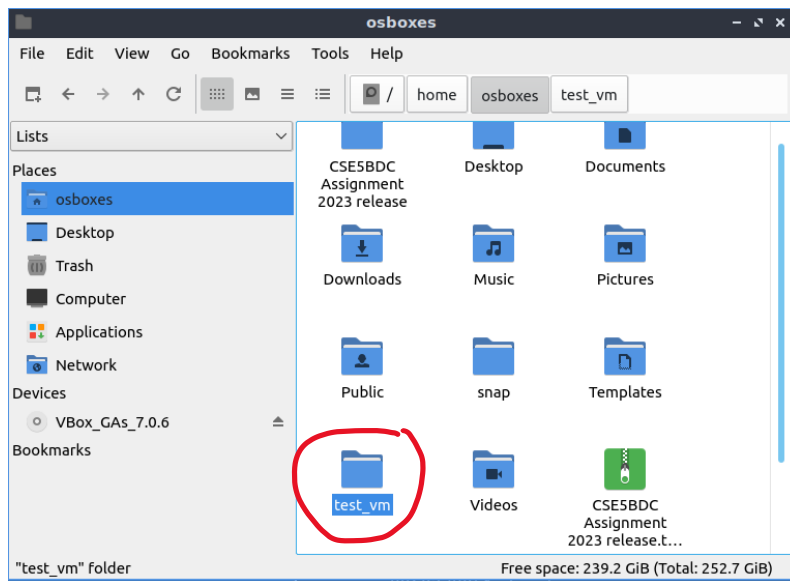


```
root@89602d8051c3: ~  
File Actions Edit View Help  
root@89602d8051c3: ~  
osboxes@osboxes:~$ ./run.sh /home/osboxes/ local 2g  
565aee80d5abc15f9fa8d614f62d2bedd977a636a21bc7a1274c3751e92f8b58  
89602d8051c31379daffb1176130f4837bc055a09b3afe008ad4b5b6e7cd7e51  
ac035fcd3e5b90f32491549eb6538ab8fb930823d867685f047cda73da9377f0  
Starting namenodes on [0.0.0.0]  
0.0.0.0: starting namenode, logging to /usr/lib/hadoop-2.7.7/logs/hadoop-root-  
namenode-89602d8051c3.out  
localhost: starting datanode, logging to /usr/lib/hadoop-2.7.7/logs/hadoop-ro  
ot-datanode-89602d8051c3.out  
Starting secondary namenodes [0.0.0.0]  
0.0.0.0: starting secondarynamenode, logging to /usr/lib/hadoop-2.7.7/logs/ha  
dooop-root-secondarynamenode-89602d8051c3.out  
starting yarn daemons  
starting resource manager, logging to /usr/lib/hadoop-2.7.7/logs/yarn--resourc  
e manager-89602d8051c3.out  
localhost: starting nodemanager, logging to /usr/lib/hadoop-2.7.7/logs/yarn-r  
oot-nodemanager-89602d8051c3.out  
chown: missing operand after '/usr/lib/hadoop-2.7.7/logs'  
Try 'chown --help' for more information.  
starting historyserver, logging to /usr/lib/hadoop-2.7.7/logs/mapred--history  
server-89602d8051c3.out  
* Starting PostgreSQL 10 database server [ OK ]  
root@89602d8051c3:~#
```

13. Now we are in the docker image, we can try out a Hive script to see if everything is working as they should. Let's first open the code in the visual studio code editor to take a look at it. In First click the osboxes folder.



Look for the test\_vm folder then move into the folder and then find file t1-wordcount.hql and then *right click and open with visual studios*.



14. Go back to the terminal installing the docker image. If the install has finished you should be at a **root** prompt. Now change to the **labfiles** directory. This will take you to your home osboxes directory. Now change to **test\_vm** directory. Then type the following command to run hive to test if your machine is fast enough to run the labs: **hive -f t1-wordcount.hql** (command also included in screenshot below)

If everything finishes within 1 or 2 mins then it should be good enough to do the labs. If this take more than 5 mins to finish then your machine maybe too slow.

**Trouble shooting: if the directory **test\_vm** cannot be found. Please double check that the second parameter to **./run.sh /home/osboxes/** is typed in correctly.**

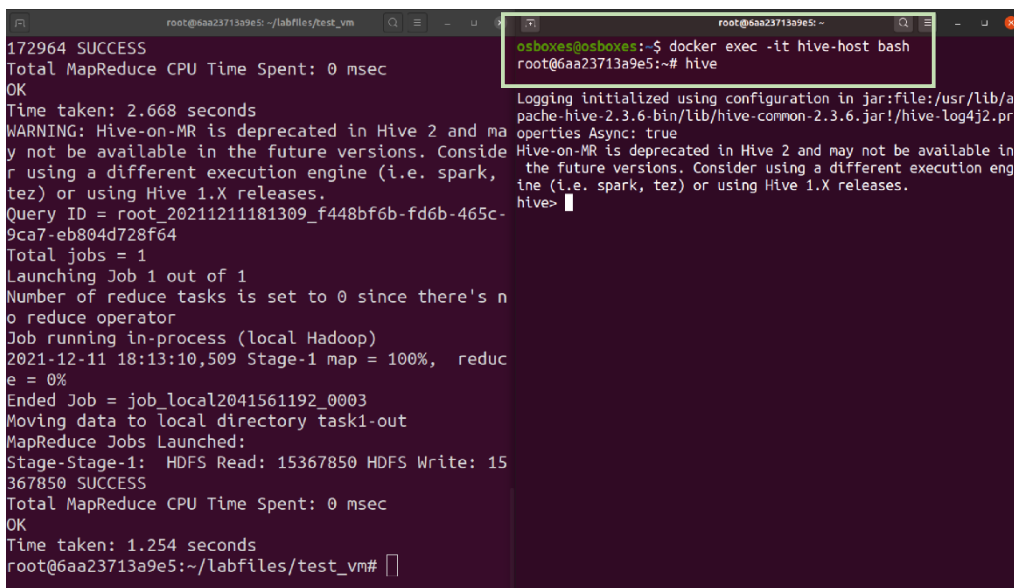
```
root@89602d8051c3:~# ls
labfiles
root@89602d8051c3:~# cd labfiles/
root@89602d8051c3:~/labfiles# ls
'CSE5BDC Assignment 2023 release'  Downloads  Templates  test_vm
'CSE5BDC Assignment 2023 release.tar.gz' Music      Videos
Desktop                          Pictures   run.sh
Documents                        Public    snap
root@89602d8051c3:~/labfiles# cd test_vm
root@89602d8051c3:~/labfiles/test_vm# hive -f t1-wordcount.hql
```

15. Now let's start another terminal window and use that for the hive interpreter. So, we will have two windows open. For the second window we will use the following command to connect into the docker container above.

**docker exec -it hive-host bash**



Then open the hive interpreter using the **hive** command. See the right window below:



The image shows two terminal windows. The left window displays the output of a Hive query, showing success and execution details. The right window shows the Hive interpreter prompt after running the 'hive' command.

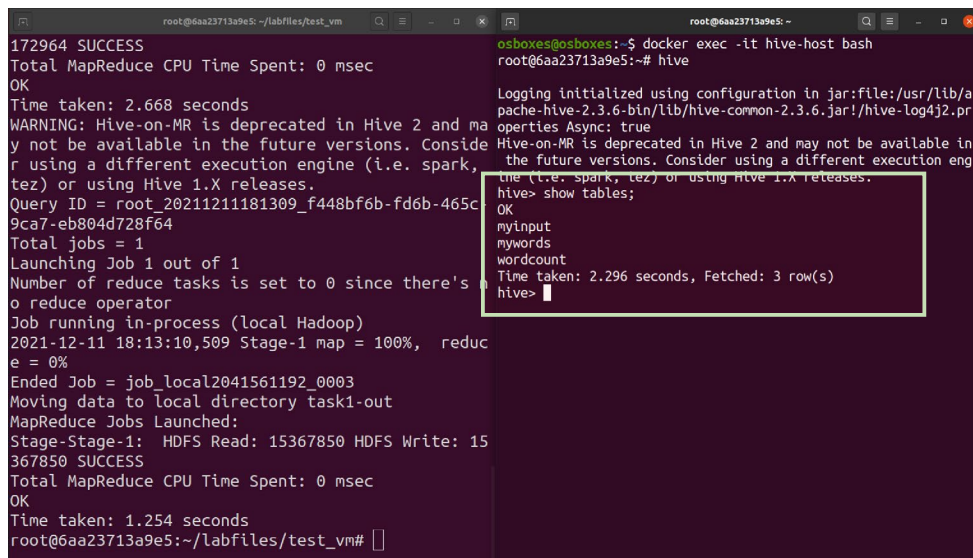
```
root@6aa23713a9e5:~/labfiles/test_vm# 172964 SUCCESS
Total MapReduce CPU Time Spent: 0 msec
OK
Time taken: 2.668 seconds
WARNING: Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a different execution engine (i.e. spark, tez) or using Hive 1.X releases.
Query ID = root_20211211181309_f448bf6b-fd6b-465c-9ca7-eb804d728f64
Total jobs = 1
Launching Job 1 out of 1
Number of reduce tasks is set to 0 since there's no reduce operator
Job running in-process (local Hadoop)
2021-12-11 18:13:10,509 Stage-1 map = 100%, reduce = 0%
Ended Job = job_local2041561192_0003
Moving data to local directory task1-out
MapReduce Jobs Launched:
Stage-Stage-1: HDFS Read: 15367850 HDFS Write: 15367850 SUCCESS
Total MapReduce CPU Time Spent: 0 msec
OK
Time taken: 1.254 seconds
root@6aa23713a9e5:~/labfiles/test_vm#
```

```
osboxes@osboxes:~$ docker exec -it hive-host bash
root@6aa23713a9e5:~# hive
Logging initialized using configuration in jar:file:/usr/lib/apache-hive-2.3.6-bin/lib/hive-common-2.3.6.jar!/hive-log4j2.properties Async: true
Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a different execution engine (i.e. spark, tez) or using Hive 1.X releases.
hive>
```

16. Now on the **right** window you can type the following in the hive interpreter to see what tables were created.

*show tables;*

If everything works correctly you should see three tables (myinput, mywords and wordcount)



The image shows two terminal windows. The left window is the same as before. The right window shows the output of the 'show tables;' command in the Hive interpreter, listing three tables: myinput, mywords, and wordcount.

```
osboxes@osboxes:~$ docker exec -it hive-host bash
root@6aa23713a9e5:~# hive
Logging initialized using configuration in jar:file:/usr/lib/apache-hive-2.3.6-bin/lib/hive-common-2.3.6.jar!/hive-log4j2.properties Async: true
Hive-on-MR is deprecated in Hive 2 and may not be available in the future versions. Consider using a different execution engine (i.e. spark, tez) or using Hive 1.X releases.
hive> show tables;
OK
myinput
mywords
wordcount
Time taken: 2.296 seconds, Fetched: 3 row(s)
hive>
```

Note the first time you use the hive interpreter it will be slow. But if you run it again it should be super fast. Try the above command again in the hive interpreter.

17. Saving via OneDrive. If drag and drop does not work for some reason. You can upload your work to your own OneDrive account for backup. Just need to log into your OneDrive account using the FireFox browser using your La Trobe student account. If you have trouble logging into OneDrive, you can upload your files to your google drive account or anywhere else on the web. Alternatively, if you have a DropBox account, you can use that too.