## **Virtual Machine is 19 GB. Don't use cell phone connection to download. Use wifi/adsl/fiber connection.**

## GİZLİLİK VE GÜVENLİK POLİTİKASI

Download Linki:

https://drive.google.com/file/d/1qWPMeRpsOAKZsIQDK4xh8DuZCj5WExTJ/view?usp=sharing

User: train
Password: Ankara06
root password: Ankara06

# Hardware Requirements

Your pc/laptop should have at least **8 GB RAM** (greater is better), **4 cpu cores** and **60 GB** free disk space. **SSD disk highly recommended**. If you use hdd give vm higher memory  e.g. 10 GB.

# Software Requirements

Windows 7-10 or MacOS operating system. (not M1 cpu)
VirtualBox or VmWare Player or Workstation.

Optionally recommended installations: MobaXTerm or gitbash to connect VM (Windows users only).

-------------------------------------------------------------------------------------------------------------------

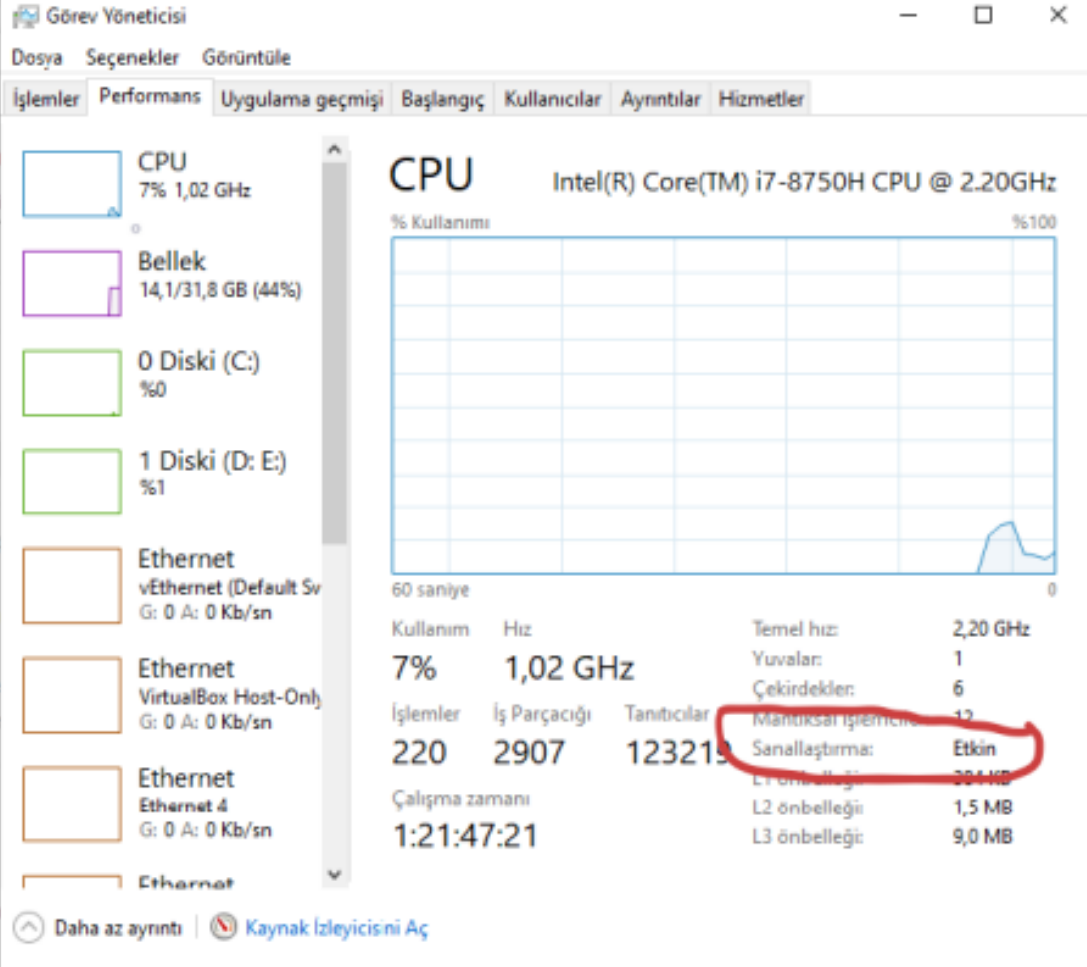# - Possible MacOS Virtualbox Installation Error and Solution

If you get above error, Solution#1 can help you in this post:
https://medium.com/@DMeechan/fixing-the-installation-failed-virtualbox-error-on-mac-high sierra-7c421362b5b5

-----------------------------------------------------------------------------------------------------------------------
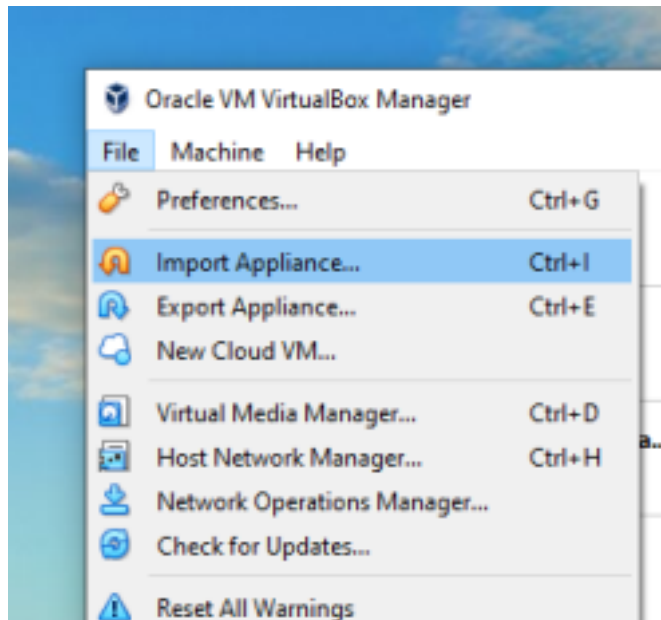
# CPU virtualization must be enabled

CPU virtualization is enabled (Generally it is enabled, if not, you must activate from BIOS). Below you see how to check for Windows computers.
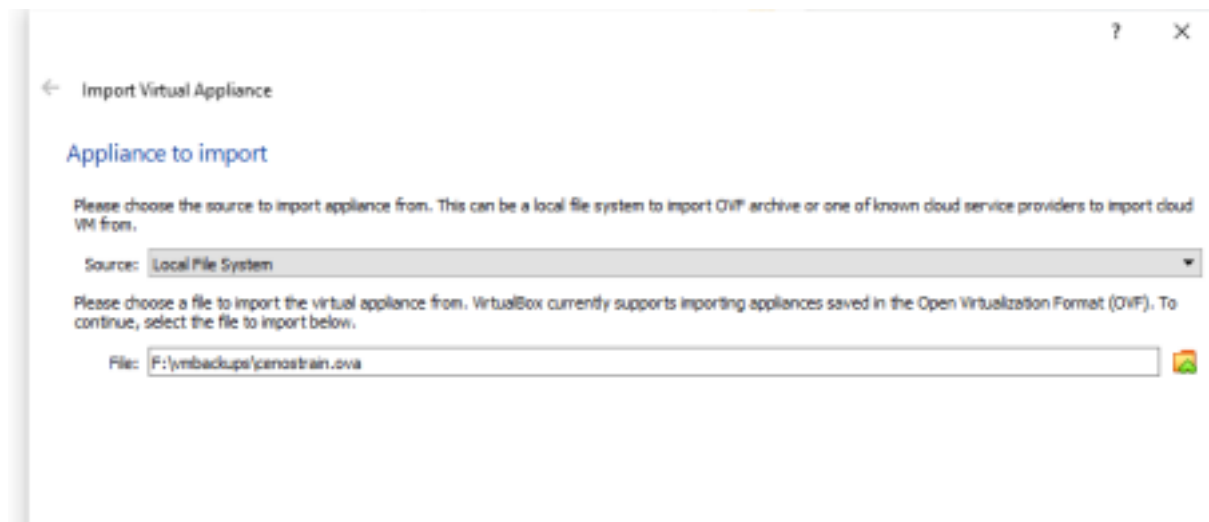
# Importing vm into VirtualBox

## 1. Open VirtualBox and Import Appliance



-------------------------------------------------------------------------------------------------------------------

## 2. Locate your VM



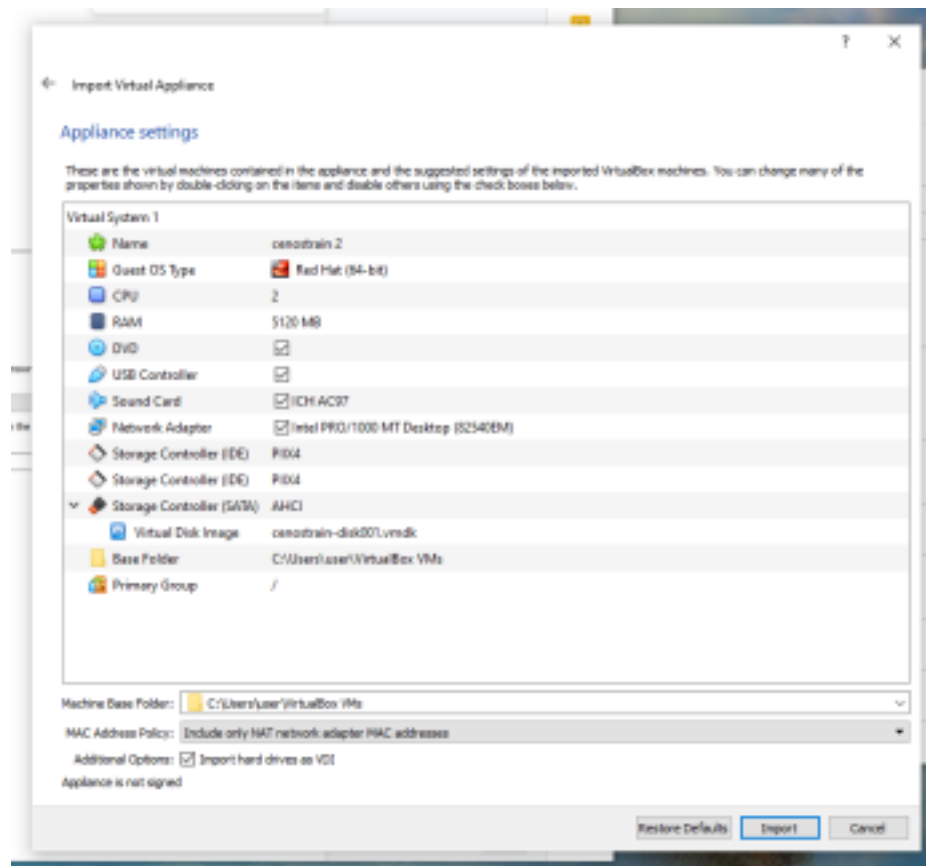-------------------------------------------------------------------------------------------------------------------
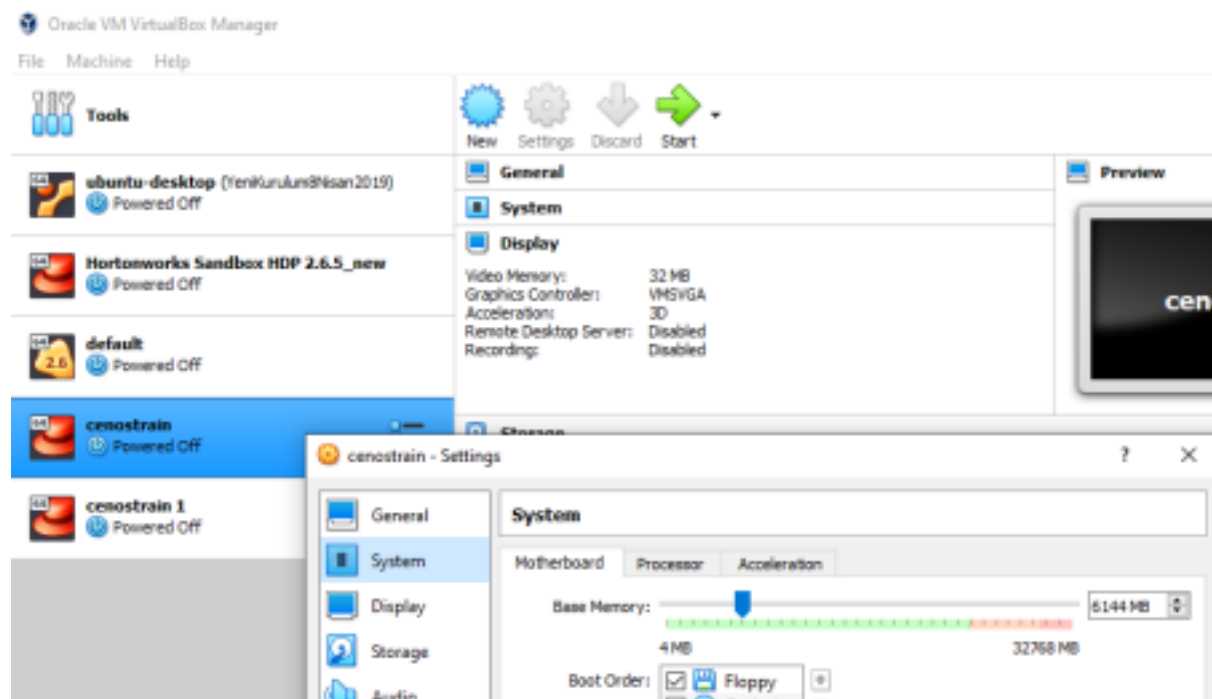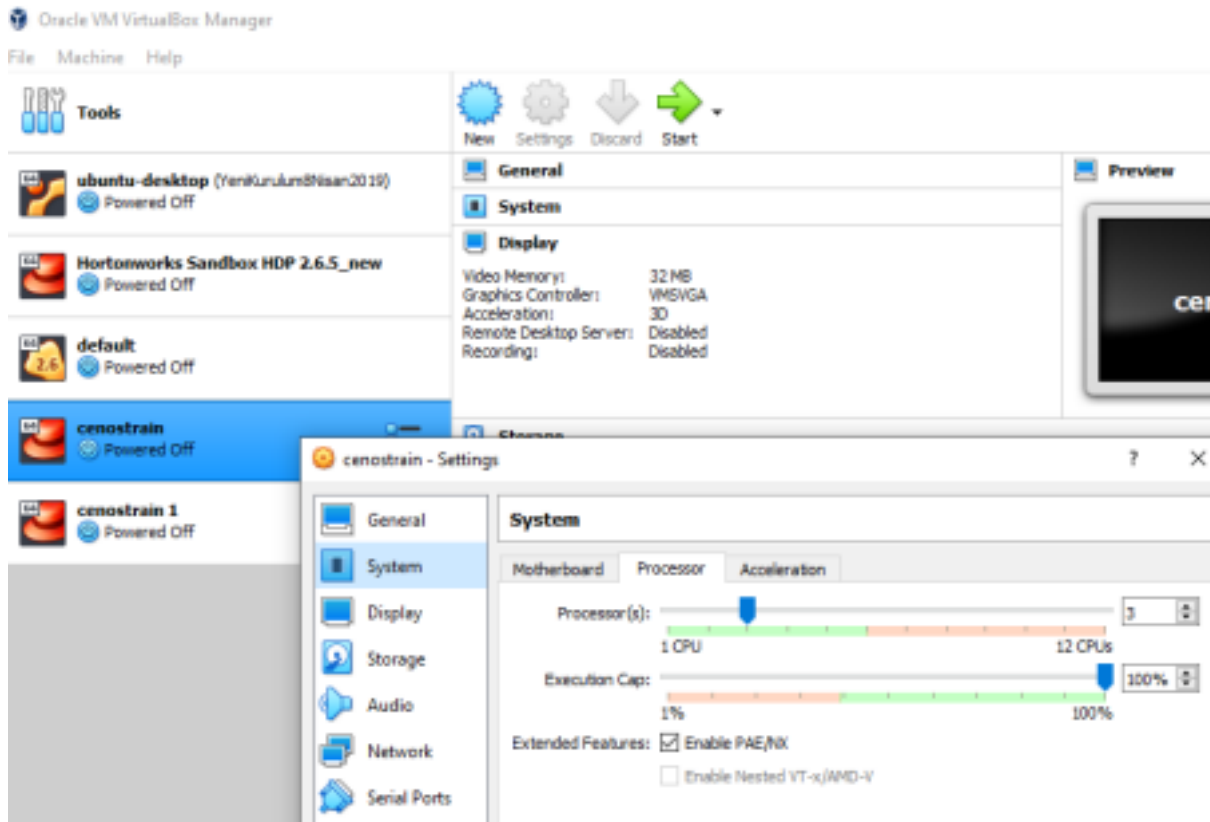-

## 3. Import VM

If necessary change Machine Base Folder (must be at least 60 Gb free disk space) click Import

Give vm higher memory and cpu core. The higher the better. But don't forget that your host machine also needs some resources, find the balance.



-------------------------------------------------------------------------------------------------------------------

-
CPU

# Test Services

**Caution!! Some commands exceed the second line. Copy both lines at the same time.**

## 1. Test Kafka

**Start Zookeeper and Kafka**
```
[train@localhost ~]$ sudo systemctl start zookeeper
[train@localhost ~]$ sudo systemctl status zookeeper
        Must be running

[train@localhost ~]$ sudo systemctl start kafka
[train@localhost ~]$ sudo systemctl status kafka
        Must be running
```

**Kafka topic create**
```
[train@localhost ~]$ kafka-topics.sh --bootstrap-server localhost:9092 --create
-- topic test5 --partitions 3 --replication-factor 1
```

**Expected output**
```
Created topic test5.
```

**Kafka topic list**
```
[train@localhost ~]$ kafka-topics.sh --bootstrap-server localhost:9092 --list
```

**Expected output**
```
__consumer_offsets
test
test1
test2
```

```
test5
```

**Stop Kafka**

```
[train@localhost ~]$ sudo systemctl stop kafka
[train@localhost ~]$ sudo systemctl stop zookeeper
```

First start zookeeper, First stop kafka.

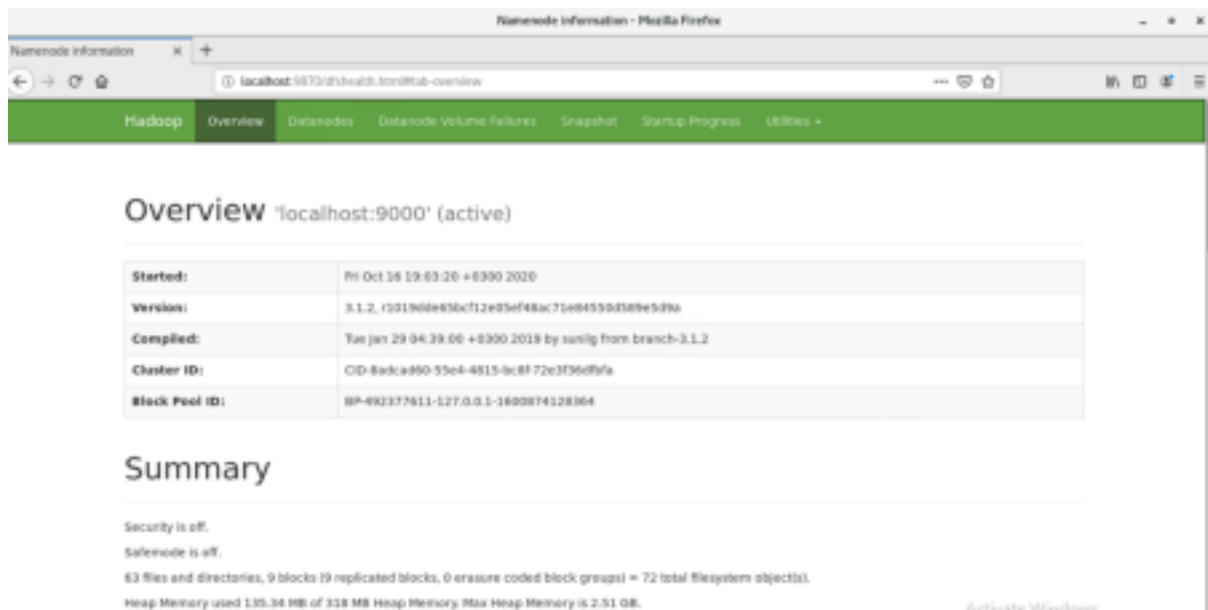# 2. Test Hadoop, YARN and Hive

```
[train@localhost ~]$ start-all.sh
```
Wait till warnings end, it will take some time.

**Hadoop Test**
```
[train@localhost ~]$ jps
9552 SecondaryNameNode
9825 ResourceManager
10338 RunJar
9142 NameNode
10393 RunJar
10825 Jps
9276 DataNode
9950 NodeManager
```

**Namenode Web UI Test**
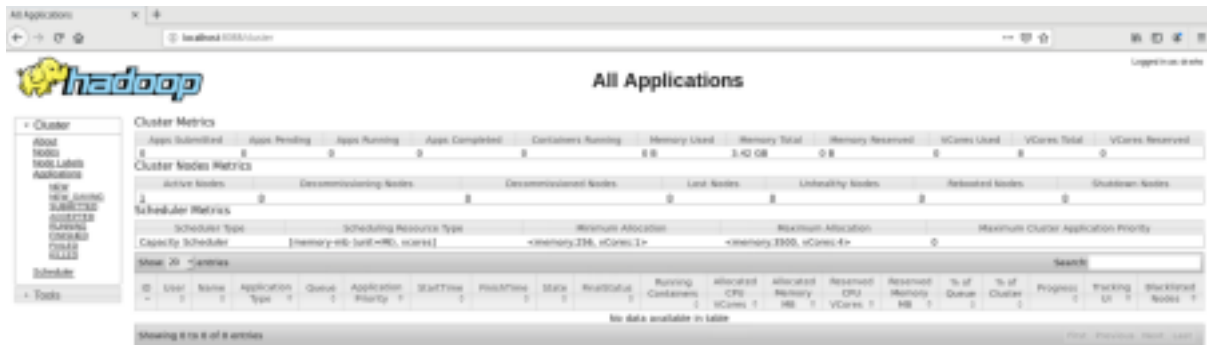Open browser (Applications -> Firefox) and enter `http://localhost:9870/` see namenode
ui



**See: Live Nodes 1 (Decommissioned: 0, In Maintenance: 0)**

**YARN Test**
On browser `http://localhost:8088/cluster` you should see the resource manager page

## HDFS Write Test

```
[train@localhost ~]$ hdfs dfs -put ~/datasets/Advertising.csv
hdfs://localhost:9000/tmp

[train@localhost ~]$ hdfs dfs -ls hdfs://localhost:9000/tmp
-rw-r--r-- 1 train supergroup 4556 2021-05-09 11:54
hdfs://localhost:9000/tmp/Advertising.csv
. . .
. . .
```

. . . represents other files

# 3. Hive and beeline test

After starting Hadoop **wait a while (at least 30 secs to up Hive services)** then test the beeline.

**Is Hive2server running?**
```
pgrep -f org.apache.hive.service.server.HiveServer2
```
**Expected output**
```
<pid>
```

**Is Hive Metastore running?**
```
pgrep -f org.apache.hadoop.hive.metastore.HiveMetaStore
```
**Expected output**
```
<pid>
```

**Beeline Test**
**Connect hive through beeline**
```
[train@localhost ~]$ beeline -n train -u jdbc:hive2://127.0.0.1:10000
```

**Stop beeline warnings**
```
0: jdbc:hive2://127.0.0.1:10000> set
```

```
hive.server2.logging.operation.level=NONE;
```
**Create database table insert and**

**select**

```
0: jdbc:hive2://127.0.0.1:10000> create database if not exists bookstore;
```

```
0: jdbc:hive2://127.0.0.1:10000> use bookstore;
```

```
0: jdbc:hive2://127.0.0.1:10000> create table if not exists bookstore.books(id
int,book name string,isbn bigint,book id bigint,price float,price currency
string,rating count int,author id bigint,publisher id bigint);
```

```
0: jdbc:hive2://127.0.0.1:10000> insert into bookstore.books values(13,"Madam
Bovary (Ciltli)",6050948752,489127179,25.115735,"TRY",5,4098249,46868),
(22,"Mai ve Siyah (Eleştirel
Basım)",9750523533,492625951,25.349610000000002,"TRY",17,3066057,63217)
, (27,"Nutuk",9759914288,9927355,11.48147,"TRY",23,9705003,46868),
 (34,"Devlet",9754734263,395307782,27.9994,"TRY",0,8978000,20709);
```

```
0: jdbc:hive2://127.0.0.1:10000> select id, book_name from
bookstore.books; OK
+-----+------------------------------+
| id  | book_name |
+-----+------------------------------+
| 13  | Madam Bovary (Ciltli) |
| 22  | Mai ve Siyah (Eleştirel Basım) |
| 27  | Nutuk |
| 34  | Devlet |
+-----+------------------------------+
4 rows selected (0.369 seconds)
```

**Exit beeline**
```
!q
```

# 4. Postgresql test

**Command**
```
[train@localhost ~]$ systemctl status postgresql-10
```
**Expected output**
```
● postgresql-10.service - PostgreSQL 10 database server
 Loaded: loaded (/usr/lib/systemd/system/postgresql-10.service; enabled; vendor
preset: disabled)
 Active: active (running) since Wed 2020-09-23 16:49:49 +03; 1h 38min ago
Docs: https://www.postgresql.org/docs/10/static/
 Process: 1134 ExecStartPre=/usr/pgsql-10/bin/postgresql-10-check-db-dir ${PGDATA}
(code=exited, status=0/SUCCESS)
 Main PID: 1156 (postmaster)
```

# 5. Sqoop test

```
[train@localhost ~]$ sqoop version | grep Sqoop
```
**Output**
```
2020-09-23 19:54:35,388 INFO sqoop.Sqoop: Running Sqoop version:
1.4.6 Sqoop 1.4.6
```

```
[train@localhost ~]$ sqoop job --list | grep -A10 'Available
jobs' Output
```
```
2021-05-14 10:24:46,960 INFO sqoop.Sqoop: Running Sqoop version:
1.4.6 Available jobs:
```

# 6. Pyspark test

**Local mode**
```
[train@localhost ~]$ pyspark --master local
```

## Output

```
Python 3.6.8 (default, Nov 16 2020, 16:55:22)
[GCC 4.8.5 20150623 (Red Hat 4.8.5-44)] on linux
Type "help", "copyright", "credits" or "license" for more information.
2021-08-14 23:28:17,984 WARN util.Utils: Your hostname, localhost.localdomain
resolves to a loopback address: 127.0.0.1; using 10.0.2.15 instead (on interface
enp0s3)
2021-08-14 23:28:17,988 WARN util.Utils: Set SPARK_LOCAL_IP if you need to bind to
another address
2021-08-14 23:28:19,249 WARN util.NativeCodeLoader: Unable to load
native-hadoop library for your platform... using builtin-java classes where
applicable Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use
setLogLevel(newLevel).
Welcome to

      ____              __
     / __/__  ___ _____/ /__
    _\ \/ _ \/ _ `/ __/  '_/
   /__ / .__/\_,_/_/ /_/\_\   version 3.1.1
      /_/

Using Python version 3.6.8 (default, Nov 16 2020 16:55:22)
Spark context Web UI available at http://10.0.2.15:4040
Spark context available as 'sc' (master = local, app id =
local-1628972904777). SparkSession available as 'spark'.
```

From http://localhost:4040 check spark web ui.

## Local mode exit
```
>>> exit()
```

## YARN mode
```
[train@localhost ~]$ pyspark --master yarn
```

## Output

```
Python 3.6.8 (default, Nov 16 2020, 16:55:22)
[GCC 4.8.5 20150623 (Red Hat 4.8.5-44)] on linux
Type "help", "copyright", "credits" or "license" for more information.
2021-08-14 23:29:58,266 WARN util.Utils: Your hostname, localhost.localdomain
resolves to a loopback address: 127.0.0.1; using 10.0.2.15 instead (on interface
enp0s3)
2021-08-14 23:29:58,269 WARN util.Utils: Set SPARK_LOCAL_IP if you need to bind to
another address
2021-08-14 23:29:59,451 WARN util.NativeCodeLoader: Unable to load native-hadoop
library for your platform... using builtin-java classes where applicable Setting
default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use
setLogLevel(newLevel).
2021-08-14 23:30:06,471 WARN yarn.Client: Neither spark.yarn.jars nor
spark.yarn.archive is set, falling back to uploading libraries under SPARK_HOME.
Welcome to

      ____              __
     / __/__  ___ _____/ /__
    _\ \/ _ \/ _ `/ __/  '_/
   /__ / .__/\_,_/_/ /_/\_\   version 3.1.1
      /_/

Using Python version 3.6.8 (default, Nov 16 2020 16:55:22)
Spark context Web UI available at http://10.0.2.15:4040
Spark context available as 'sc' (master = yarn, app id =
```

```
application_1628972365632_0002).
SparkSession available as 'spark'.
From http://localhost:8088/cluster/apps/RUNNING check spark have resources
```

| | ID | User | Name | Application Type | Queue | Application Priority | StartTime | FinishTime | State | FinalStatus | Running Containers | A V |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Show 20 entries | | | | | | | | | | | | |
| | application_1600874193447_0002 | train | PySparkShell | SPARK | default | 0 | Wed Sep 23 19:58:26 +0300 2020 | N/A | RUNNING | UNDEFINED | 2 | 2 |

## Write with Pyspark to hive

```
>>>df=spark.read.option("header",True).option("inferSchema",True).csv("/tmp/Advert
i sing.csv")

>>> df.write.mode("overwrite").saveAsTable("advertising")

>>> spark.sql("select * from advertising").show(3)
+---+-----+-----+---------+-----+
| ID| TV|Radio|Newspaper|Sales|
+---+-----+-----+---------+-----+
|  1|230.1| 37.8|     69.2| 22.1|
|  2| 44.5| 39.3|     45.1| 10.4|
|  3| 17.2| 45.9|     69.3|  9.3|
+---+-----+-----+---------+-----+
only showing top 3 rows
```
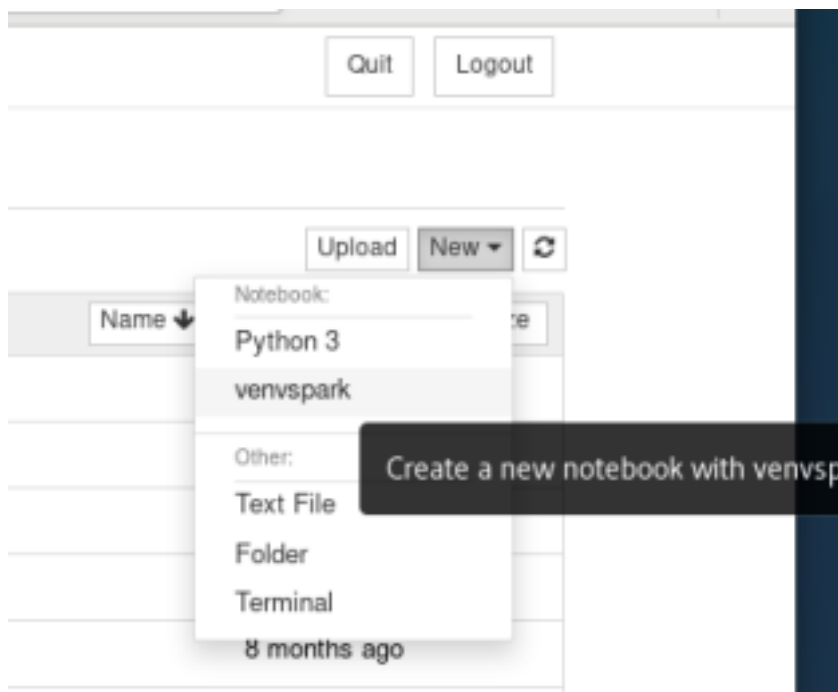
## YARN mode exit

```
>>> exit()
```

## Spark python virtual environment

```
[train@localhost ~]$ source venvspark/bin/activate
```

## Open Jupyter notebook

```
(venvspark) [train@localhost ~]$ jupyter notebook > notebook.log 2>&1
```

& Jupyter will automatically open up in the browser.

Create **venvspark** notebook

**Close Jupyter notebook**
Save your notebook and close the browser.
On terminal stop jupyter notebook
```
 (venvspark) [train@localhost ~]$ jupyter notebook stop

Shutting down server on port 8888 ...
 [1]+ Done jupyter notebook > notebook.log 2>&1
```

**Close spark virtual environment**
```
 (venvspark) [train@localhost ~]$ deactivate
```

# 7. Stop Hadoop, YARN and Hive

```
stop-all.sh
```

# 8. Airflow test

```
[train@localhost ~]$ sudo systemctl start airflow
```

```
[train@localhost ~]$ sudo systemctl start airflow-scheduler
```

Wait 30 secs for the Airflow web server spin up. Open browser
http://127.0.0.1:1502 and see airflow web ui. **username: admin, password: admin**

**Stop Airflow**

```
[train@localhost ~]$ sudo systemctl stop airflow-scheduler
[train@localhost ~]$ sudo systemctl stop airflow
```

# 9. Docker test

```
[train@localhost ~]$ sudo systemctl start docker
```

```
[train@localhost ~]$ sudo systemctl status docker
```

# 10. Docker-compose test

```
[train@localhost ~]$ docker-compose version

docker-compose version 1.29.2, build 5becea4c
docker-py version: 5.0.0
CPython version: 3.7.10
OpenSSL version: OpenSSL 1.1.0l 10 Sep 2019
```

# 11. Minikube (Kubernetes) test

```
[train@localhost ~]$ minikube delete
```

```
[train@localhost ~]$ minikube start
```

```
[train@localhost ~]$ kubectl get all
```

```
NAME TYPE CLUSTER-IP EXTERNAL-IP PORT(S) AGE service/kubernetes ClusterIP
10.96.0.1 <none> 443/TCP 10m
```

**Stop minikube**
```
[train@localhost ~]$ minikube stop
```
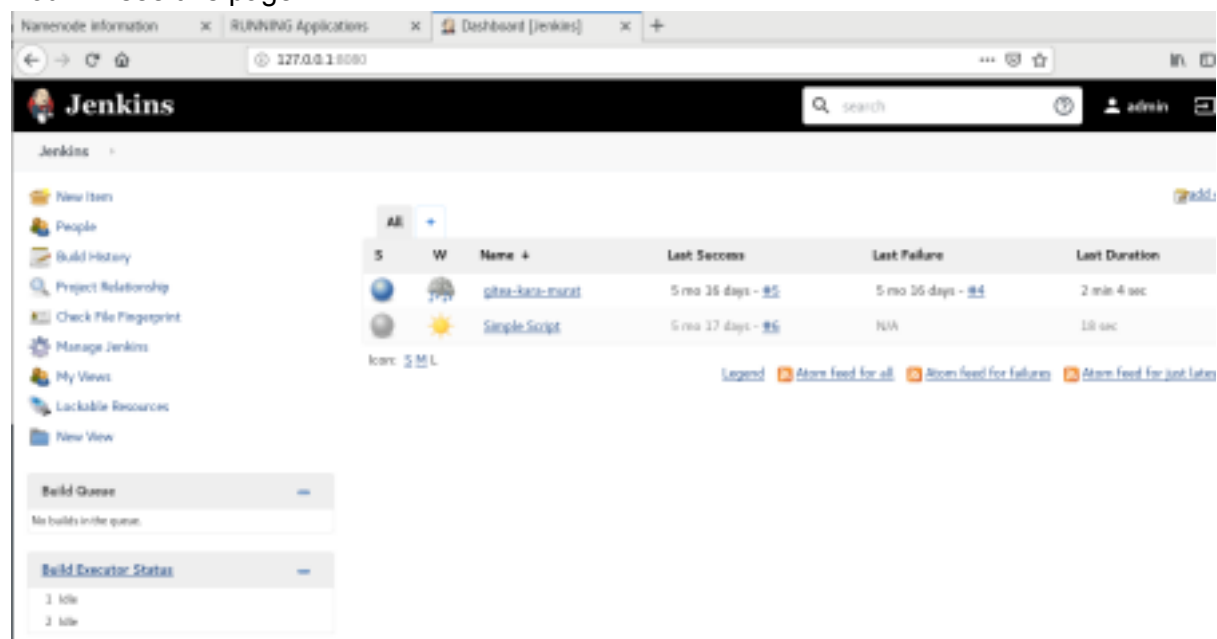
## 12. Jenkins Test

**Start Jenkins**
```
[train@localhost ~]$ sudo systemctl start jenkins
```

Open Jenkins UI
On browser open `localhost:8080`

**username: admin, password: Ankara06**

You will see this page



**Stop Jenkins**
```
[train@localhost ~]$ sudo systemctl stop jenkins
```

## 13. Gitea test

**Start Gitea**
```
[train@localhost ~]$ sudo systemctl start gitea
```
Open Gitea UI
On browser open `localhost:3000`
Sign in with **username: jenkins, password Ankara_06**

**Stop Gitea**
```
[train@localhost ~]$ sudo systemctl stop gitea
```

## 14. Close virtual machine

Check if any docker container is running. If there are any, stop them.

stop kafka

`[train@localhost ~]$` `sudo systemctl stop kafka`

stop zookeeper

`[train@localhost ~]$` `sudo systemctl stop zookeeper`

stop Hadoop services

`[train@localhost ~]$` `stop-all.sh`

Stop any other apps, services e.g. **docker, airflow, gitea** etc. like `sudo systemctl stop <xxxxxxx>`

Shutdown the machine

`[train@localhost ~]$` `sudo shutdown now`

# Congrats VM Test has finished !!!

See Troubleshooting in the following pages.

# Common Errors and Fixes

## 1. YARN Spark Accepted State Disk usage

If your YARN apllication stuck at ACCEPTED state your disk usage probably exceeded %90. To clean it run this command and

`rm -rf /tmp/hadoop-train/nm-local-dir/*`

check

`df -h /`

```
Filesystem Size Used Avail Use% Mounted on
/dev/mapper/centos-root 38G 30G 8.0G 79% /
```

Your application will switch to RUNNING state after a while.

## 2. Windows Users Additional Setup

Install MobaXTerm to connect virtual machine and file transfer between host and guest machine.

## 3. Windows Users ssh connection to Vm

Using gitbash or putty

`$ ssh train@localhost`

If you get an error like below

```
@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@
@ WARNING: REMOTE HOST IDENTIFICATION HAS CHANGED! @
@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@@
IT IS POSSIBLE THAT SOMEONE IS DOING SOMETHING NASTY!
```

```
Someone could be eavesdropping on you right now (man-in-the-middle attack)!
It is also possible that a host key has just been changed.
The fingerprint for the ECDSA key sent by the remote host is
SHA256:hqpjsGjn9WsPsE+KNr87ozyUqhfyNBj7YBT0DrC+PNQ.
Please contact your system administrator.
Add correct host key in /c/Users/user/.ssh/known_hosts to get rid of this message.
Offending ECDSA key in /c/Users/user/.ssh/known_hosts:49
ECDSA host key for localhost has changed and you have requested strict checking.
Host key verification failed.
```
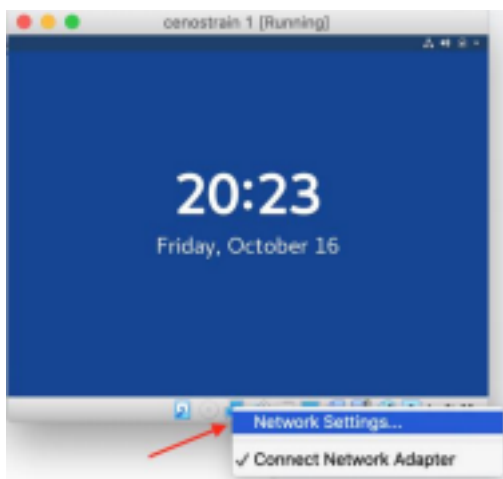
Run
```
$ ssh-keygen.exe -R localhost
```
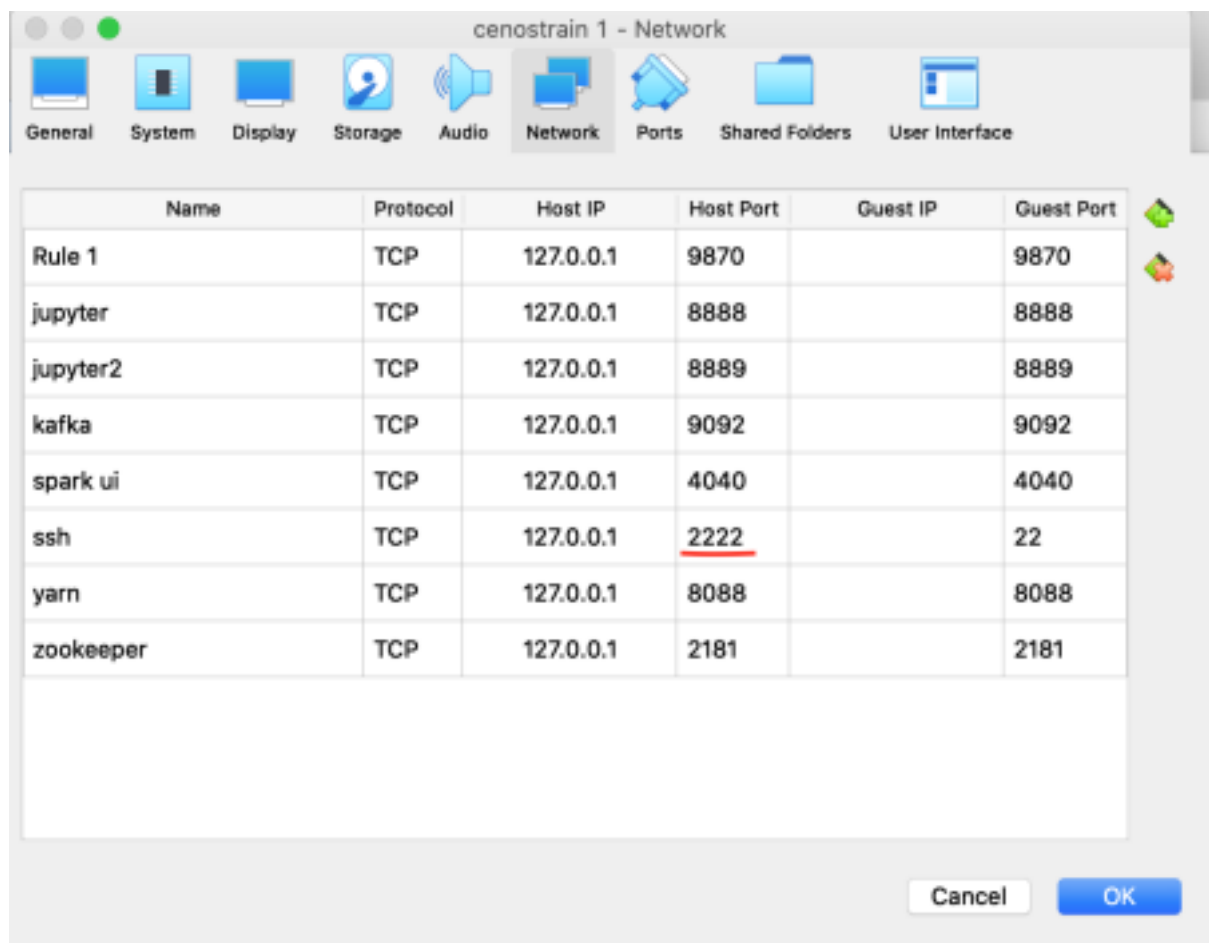Then retry to connect.

# 4. MacOS VM connection

MacOS users can use MacOS terminal to virtual machine.
Port forwarding over 22 is not stable on mac however it works on windows. So we need to change 22 port to 2222 to be able to make ssh connection
**ssh connection settings for mac**
**1-** Open advanced network setting and click port forwarding



change 22 -> 2222

click OK
2- open mac terminal and enter ssh connection command below
`ssh train@127.0.0.1 -p 2222`



write yes



enter password for train user Ankara06

## 5. Jenkins password issue

Learn Jenkins admin password

```
[train@localhost]$ sudo cat
/var/lib/jenkins/secrets/initialAdminPassword
260cd473916f4c01a0e6969857c55128
```



## 6. Additional Settings

1- Disable Screen Lock.

2- Language - Keyboard settings:
Default Keyboard is Turkish for this Virtual Machine. You can add additional keyboard layouts



7. windows vmware setup error

to solve this problem go to control panel and click turn windows features on or off



check windows hypervisor platform and virtual machine platform

Hyper-V Disable

If Hyper-V is enabled, the virtualbox will work slowly. To close Hyper-V completely follow the instructions on this link: https://forums.virtualbox.org/viewtopic.php?f=25&t=99390