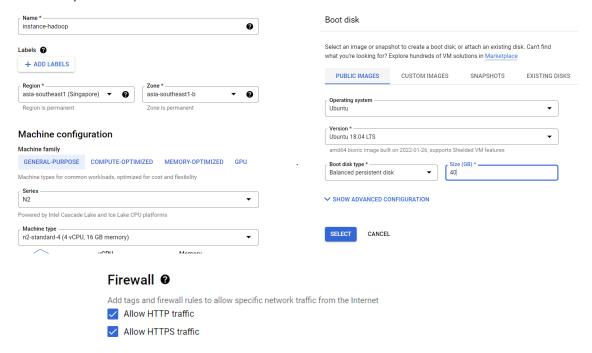
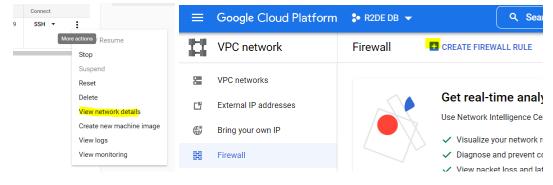
Special Live by Shane: Hadoop Docker Installation Guide

การ Set-up Google Compute Engine เพื่อทำโปรเจ็ค Hadoop

1) Provision GCP compute engine for Docker (configuration ที่นอกเหนือจากนี้ ให้ปล่อยเป็นค่า default)

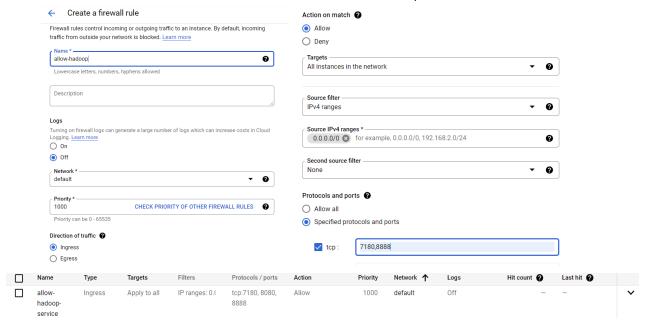


2) หลังจาก provision VM เรียบร้อยแล้ว ให้ทำการ set-up Network Firewall ของ VM ตัวนั้น



3) Set-up Firewall Network ตาม configuration ข้างต้น (หมายเหตุ การ set-up Firewall ข้างต้นเป็นการเปิด Firewall for all inbound เพื่อลดความ ซับซ้อน ของ configuration สำหรับผู้เริ่มต้นใช้งาน Networking ควรใช้ในกรณีศึกษาเท่านั้น ไม่ ควรใช้ในงาน Production หรืองานที่มีความสำคัญทางข้อมูล อาจทำให้ผิด policy ด้าน security ได้)

(จุดประสงค์หลักเพื่อเป็นการเปิด tcp port : 7180 สำหรับ Cloudera Manager และ tcp port : 8888 สำหรับ Cloudera HUE ให้เข้าถึงผ่านหน้า browser ได้)



4) หลังจาก set-up เสร็จเรียบร้อย กด ปุ่ม SSH ที่หน้า Compute Engine ของตัวที่ create ขึ้นมา เพื่อเข้าสู่หน้า SSH terminal

การเริ่มต้นใช้งาน Cloudera Hadoop via Docker container

- 1) ดิดตั้ง docker บน Ubuntu 18 ตาม **Step 1-4 จ่าก guide ของ Digital Ocean**https://www.digitalocean.com/community/tutorials/how-to-install-and-use-docker-on-ubuntu-18-04
- 2) หลังจากติดตั้งเสร็จเรียบร้อย ให้ ใช้คำสั่ง docker pull mikelemikelo/cloudera-spark:latest เพื่อ pull docker image มาใช้ที่ VM ของคุณ
- 3) หลังจาก pull image มาสำเร็จเรียบร้อยแล้ว สามารถเช็คได้ด้วยการใช้คำสั่ง docker images
- 4) ใช้คำสั่ง docker run --hostname=quickstart.cloudera --privileged=true -it -p 8888:8888 -p 8080:8080 -p 7180:7180 -p 88:88/udp -p 88:88 mikelemikelo/cloudera-spark:latest /usr/bin/docker-quickstart-light

```
nuttatun c@hadoop-sandbox:~$ docker run --hostname=quickstart.cloudera --privileged=true -i
t -p 8888:8888 -p 8080:8080 -p 7180:7180 -p 88:88/udp -p 88:88 mikelemikelo/cloudera-spark:
latest /usr/bin/docker-quickstart-light
```

5) รอจนกว่า container จะรัน Hadoop fundamental commands จนเสร็จสิ้น

```
-Doozie.https.port=11443 -Doozie.https.keystore.pass=password
 -Xmx1024m -Doozie.https.port=11443 -Doozie.https.keystore.pass=password -Xmx1024m -Dderby.
stream.error.file=/var/log/oozie/derby.log
Adding to CATALINA OPTS:
                             -Doozie.home.dir=/usr/lib/oozie -Doozie.config.dir=/etc/oozie/
conf -Doozie.log.dir=/var/log/oozie -Doozie.data.dir=/var/lib/oozie -Doozie.instance.id=qui
ckstart.cloudera -Doozie.config.file=oozie-site.xml -Doozie.log4j.file=oozie-log4j.properti
es -Doozie.log4j.reload=10 -Doozie.http.hostname=quickstart.cloudera -Doozie.admin.port=110
01 -Doozie.http.port=11000 -Doozie.https.port=11443 -Doozie.base.url=http://quickstart.clou
dera:11000/oozie -Doozie.https.keystore.file=/var/lib/oozie/.keystore -Doozie.https.keystor
e.pass=password -Djava.library.path=:/usr/lib/hadoop/lib/native:/usr/lib/hadoop/lib/native
Using CATALINA_BASE: /var/lib/oozie/tomcat-deployment
Using CATALINA_HOME: /usr/lib/bigtop-tomcat
Using CATALINA_TMPDIR: /var/lib/oozie
Using JRE HOME:
                      /usr/lib/jvm/jre-1.8.0-openjdk.x86_64
Using CLASSPATH:
                      /usr/lib/bigtop-tomcat/bin/bootstrap.jar
Using CATALINA_PID:
                       /var/run/oozie/oozie.pid
[root@quickstart cloudera]#
[root@quickstart cloudera] # sudo /home/cloudera/cloudera-manager --express && service ntpd
[QuickStart] Shutting down CDH services via init scripts...
kafka-server: unrecognized service
```

6) ใช้คำสั่ง sudo /home/cloudera/cloudera-manager --express && service ntpd start เพื่อเริ่ม ตันการใช้งาน Cloudera Manager

```
[QuickStart] Configuring deployment...

Submitted jobs: 14
[QuickStart] Deploying client configuration...

Submitted jobs: 16
[QuickStart] Starting Cloudera Management Service...

Submitted jobs: 24
[QuickStart] Enabling Cloudera Manager daemons on boot...

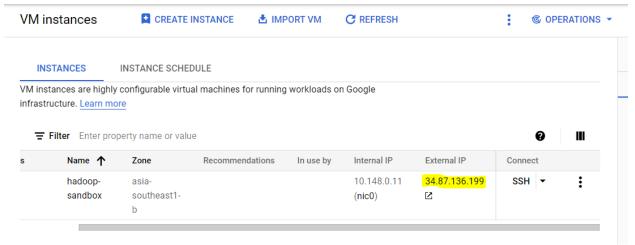
Success! You can now log into Cloudera Manager from the QuickStart VM's browser:

http://quickstart.cloudera:7180

Username: cloudera
Password: cloudera

Starting ntpd:
[OK]
[root@quickstart cloudera]#
```

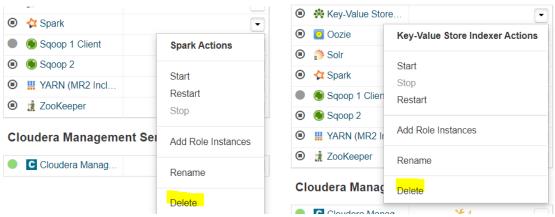
 ใช้ Public / External IP ของ VM ในการเข้า Cloudera Manager ด้วย port 7180 ยกตัวอย่างเช่น http://34.87.136.199:7180/
 -> user:cloudera, pw:cloudera



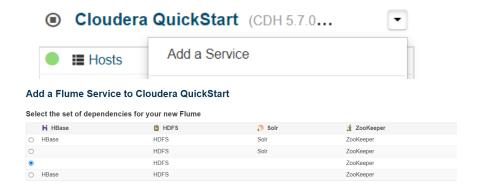
8) ก่อนจะเริ่มตัน start cluster ให้ไปที่เชอร์วิส Hive ทั้งด้าน ซ้ายมือ set ค่า ของ configuration ข้าง ตัน เป็น none



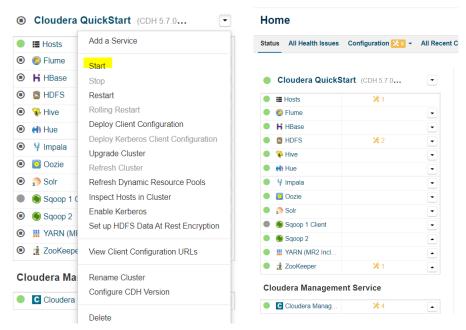
9) ลบ Service ที่ไม่ได้ใช้งานและไม่จำเป็นออก เช่น Spark และ KV Indexer (Spark ของ CDH นี้เป็น version 1.x.x ที่ไม่สามารถใช้งาน spark structured streaming ได้ ให้ ลบ service ตรงนี้ แล้วไปใช้ Spark Local ที่ติดตั้งไว้แล้วใน docker container แทน)



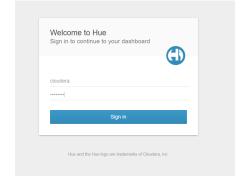
10) ทำการ Add service "Flume"



11) Start cluster และ รอจนกว่า service จะเข้า status all-up-and-running



12) ทำการเข้าสู่ระบบด้วย Cloudera HUE โดยใช้ Public / External IP ของ VM ด้วย port 8888 ยก ตัวอย่างเช่น <u>http://34.87.136.199:8888/</u> -> user:cloudera, pw:cloudera



การ Clone Git Repository เพื่อใช้ Materials จากตัวอย่างใน Special Live

- 1) กด Ctrl + P + Q เพื่อออกจาก docker container มาที่ VM terminal หลัก
- ใช้คำสั่ง git clone https://github.com/nuttatunc/hadoop-r2de-data.git เพื่อ clone repository จาก special live
- 3) ใช้คำสั่ง docker container ls เพื่อ ดู container_id ของ container ที่รัน Hadoop Service อยู่
- 4) ใช้คำสั่ง docker cp hadoop-r2de-data/data <container_id>:/ เพื่อนำ ไฟล์ที่ได้จาก git repo ก็ อปปี้เข้าไปใน container
- 5) ใช้คำสั่ง docker exec -it <container_id> /bin/bash เพื่อกลับไปเข้าใช้งาน container ตามปกติ โดยจะมี folder ที่ชื่อว่า data ที่มี code และ ส่วนประกอบอื่นๆเหมือนภายใน Special Live

```
nuttatun c@hadoop-sandbox:~$ docker container ls
 CONTAINER ID
                                                                                              CREATED
                                                                                                                        STATUS
  PORTS
                                             NAMES
 bdf9a0c19f38 r2de-hadoop:latest "/usr/bin/docker-qui..." 5 seconds ago Up 4 seconds 0.0.0.0:88->88/tcp, :::88->88/tcp, 0.0.0.0:7180->7180/tcp, :::7180->7180/tcp, 0.0.0.0:8880
bdf9a0c19f38 r2de-hadoop:latest
 ->8080/tcp, :::8080->8080/tcp, 0.0.0.0:8888->8888/tcp, 0.0.0.0:88->88/udp, :::8888->8888/tc
 p, :::88->88/udp, 4040/tcp mystifying_faraday
  nuttatun c@hadoop-sandbox:~$ docker cp data/data bdf9a0c19f38:/
  nuttatun c@hadoop-sandbox:~$ docker exec -it bdf9a0c19f38 /bin/bash
 [root@quickstart cloudera]# cd ../..
[root@quickstart /]# ls -lrt
total 80
drwxr-xr-x 2 root root 4096 Sep 23 2011 srv
drwxr-xr-x 2 root root 4096 Sep 23 2011 mnt
drwxr-xr-x 2 root root 4096 Sep 23 2011 media
drwxr-xr-x 2 root root 4096 Mar 4 2015 lost+found
drwxr-xr-x 3 root root 4096 Mar 4 2015 selinux
drwxrwxr-x 1 root root 4096 Apr 5 2016 home
drwxr-xr-x 2 root root 4096 Apr 6 2016 packer-files
dr-xr-xr-x 9 root root 4096 Apr 6 2016 lib
dr-xr-xr-x 2 root root 4096 Apr 6 2016 sbin
drwxr-xr-x 3 root root 4096 Apr 6 2016 boot
drwxrwxr-x 1 root root 4096 Apr 6 2016 usr
drwxr-xr-x 1 root root 4096 Apr 6 2016 var
dr-xr-xr-x 2 root root 4096 Apr 6 2016 bin
dr-xr-xr-x 1 root root 4096 Feb 3 2021 lib64
drwxr-xr-x 1 root root 4096 Feb 3 2021 opt
dr-xr-x--- 1 root root 4096 Feb 3 2021 root
drwxrwxr-x 7 1002 1003 4096 Jan 28 16:29 data
drwxr-xr-x 1 root root 4096 Jan 29 11:38 etc
dr-xr-xr-x 13 root root 0 Jan 29 11:38 sys
dr-xr-xr-x 220 root root 0 Jan 29 11:38 proc
drwxr-xr-x 13 root root 3720 Jan 29 11:38 dev
drwxrwxrwt 1 root root 4096 Jan 29 11:38
[root@quickstart /]# [
```