

Docker task in HPC Pack

We introduced docker task in HPC Pack 2016 Update1.

To use this feature, set the environment variable `CCP_DOCKER_IMAGE` of a task so that it could be run in a docker container on Linux compute node.

Environment variable `CCP_DOCKER_IMAGE` indicates the docker image to use by this task, the format is like running a docker command in Linux:

`CCP_DOCKER_IMAGE=[Docker Registry]<Repository>[:Tag]`

Besides, there are several environment variables could be used to enhance this feature.

- `CCP_DOCKER_NVIDIA` to indicate if using command 'nvidia-docker', instead of using 'docker', to start docker container.
 - `CCP_DOCKER_VOLUMES` to set the directories to be mounted from host to docker container as volumes.
 - `CCP_DOCKER_DEBUG` to indicate if leaving the container alive for debugging after the command in it finishes, the container needs to be removed manually later.
- ❖ To run docker task, docker should be installed on Linux compute nodes first.
 - ❖ Currently, the docker image used in HPC Pack should have `/bin/bash`, or the docker task would fail. We will try to remove this restriction later.
 - ❖ A docker task can be allocated with multiple nodes to run MPI application, nevertheless, one node should not be allocated to multiple docker tasks.
 - ❖ To run MPI application in Linux nodes and docker containers in HPC Pack, a SSH key pair should be prepared and configured to Linux nodes for setting the mutual trust between them. The key pair could be generated by `ssh-keygen` on a Linux node or by `puttygen.exe`(PuTTY Key Generator) on a windows node.
 - ❖ To run MPI application in docker containers on Linux nodes, the docker image should have `sudo`, `ssh` service and `MPI` installed.

Run MPI in docker container step by step

1. Deploy cluster with ARM template

Use the template “[NOT Active Directory Domain integrated - Single head node cluster for Linux workloads](#)”
Use WindowsServer2012R2 as Head Node OS and use Ubuntu_16.04 as Compute Node Image.

BASICS

* Subscription ✓

* Resource group Create new Use existing
 ✓

* Location ✓

SETTINGS

* Cluster Name ✓

Head Node OS ✓

Head Node Disk Type ✓

Head Node VM Size

Compute Node Image ✓

Compute Node Name Prefix

Compute Node Number ✓

Compute Node Disk Type ✓

Compute Node VM Size

Admin Username ✓

* Admin Password ✓

Hpc Pack Version

* Vault Name ✓

* Vault Resource Group ✓

* Certificate Url ✓

* Cert Thumbprint ✓

2. Install docker on Linux compute nodes by clusrun

Connect to the cluster, open Command Prompt and type below command:

```
clusrun /nodegroup:linuxnodes apt update^; apt -y install docker.io
```

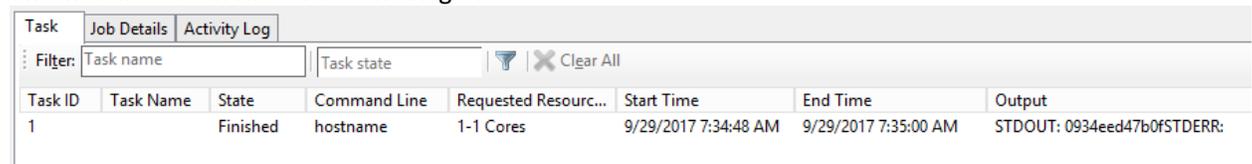
3. Run command in container as docker task

- Submit a job containing 1 docker task:

Take the Linux compute nodes online and type below command in Command Prompt:

```
job submit /env:ccp_docker_image=docker.io/library/ubuntu:16.04 hostname
```

Check task result in HPC Cluster Manager:



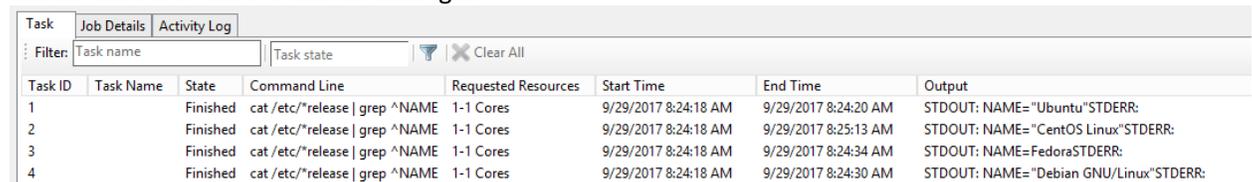
Task ID	Task Name	State	Command Line	Requested Resourc...	Start Time	End Time	Output
1		Finished	hostname	1-1 Cores	9/29/2017 7:34:48 AM	9/29/2017 7:35:00 AM	STDOUT: 0934eed47b0fSTDERR:

- Submit a job containing multiple docker tasks:

Type below command in Command Prompt:

```
job new
job add !! /env:ccp_docker_image=ubuntu cat /etc/*release ^| grep ^^NAME
job add !! /env:ccp_docker_image=centos cat /etc/*release ^| grep ^^NAME
job add !! /env:ccp_docker_image=debian cat /etc/*release ^| grep ^^NAME
job add !! /env:ccp_docker_image=fedora cat /etc/*release ^| grep ^^NAME
job submit /id:!!
```

Check task result in HPC Cluster Manager:



Task ID	Task Name	State	Command Line	Requested Resources	Start Time	End Time	Output
1		Finished	cat /etc/*release grep ^NAME	1-1 Cores	9/29/2017 8:24:18 AM	9/29/2017 8:24:20 AM	STDOUT: NAME="Ubuntu"STDERR:
2		Finished	cat /etc/*release grep ^NAME	1-1 Cores	9/29/2017 8:24:18 AM	9/29/2017 8:25:13 AM	STDOUT: NAME="CentOS Linux"STDERR:
3		Finished	cat /etc/*release grep ^NAME	1-1 Cores	9/29/2017 8:24:18 AM	9/29/2017 8:24:34 AM	STDOUT: NAME="Fedora"STDERR:
4		Finished	cat /etc/*release grep ^NAME	1-1 Cores	9/29/2017 8:24:18 AM	9/29/2017 8:24:30 AM	STDOUT: NAME="Debian GNU/Linux"STDERR:

- ❖ Tasks would inherit the environment variables of their job if they don't have the same ones, thus the docker image can also be assigned in job environment variables.

Type below command in Command Prompt:

```
job new /jobenv:ccp_docker_image=ubuntu
job add !! hostname^; cat /etc/*release ^| grep ^^NAME
job add !! hostname^; cat /etc/*release ^| grep ^^NAME
job add !! hostname^; cat /etc/*release ^| grep ^^NAME
job add !! /env:ccp_docker_image=centos hostname^; cat /etc/*release ^| grep ^^NAME
job submit /id:!!
```

Check task result in HPC Cluster Manager:

Task ID	Task Name	State	Command Line	Requested Resources	Start Time	End Time	Output
1		Finished	hostname; cat /etc/*release ...	1-1 Cores	9/30/2017 2:53:56 AM	9/30/2017 2:53:58 AM	STDOUT: 33ad8282744eNAME="Ubuntu"STDERR:
2		Finished	hostname; cat /etc/*release ...	1-1 Cores	9/30/2017 2:53:56 AM	9/30/2017 2:53:58 AM	STDOUT: 3f87543d0a5fNAME="Ubuntu"STDERR:
3		Finished	hostname; cat /etc/*release ...	1-1 Cores	9/30/2017 2:53:56 AM	9/30/2017 2:53:58 AM	STDOUT: fbe26a50c8edNAME="Ubuntu"STDERR:
4		Finished	hostname; cat /etc/*release ...	1-1 Cores	9/30/2017 2:53:56 AM	9/30/2017 2:54:11 AM	STDOUT: 4b989e0cdabdNAME="CentOS Linux"STDERR:

4. Run MPI task in container as docker task

- Build customized docker image with mpich installed

❖ Perform this step in any Linux node with docker installed.

Start a container with docker image ubuntu:

```
docker run -it ubuntu
```

Install sudo, ssh, vim and mpich with apt-get:

```
apt update; apt -y install sudo ssh vim mpich
```

Write a simple mpi program:

```
mkdir /mpisample
chmod o+w /mpisample
cd /mpisample
vim helloMpi.c
```

Edit helloMpi.c with below content:

```
#include<mpi.h>
#include<stdio.h>

int main(int argc, char** argv)
{
    int rank, size, processor_name_length;
    char processor_name[1000];
    MPI_Init(NULL, NULL);
    MPI_Comm_rank (MPI_COMM_WORLD, &rank);
    MPI_Comm_size (MPI_COMM_WORLD, &size);
    MPI_Get_processor_name(processor_name, &processor_name_length);
    printf("Hello from %s, rank %d out of %d processors.\n", processor_name, rank, size);
    MPI_Finalize();
}
```

Compile helloMpi.c with mpicc and create a shell script run.sh:

```
mpicc helloMpi.c -o helloMpi
vim run.sh
```

Edit run.sh with below content:

```
#!/bin/bash
echo $CCP_NODES | tr " " "\n" | sed "1d;n;d" | cat > host_file
```

```
num=$1
[ -z "$num" ] && num=4
mpirun -n $num -f host_file ./helloMpi
```

Set the execution permission of run.sh and exit the docker container:

```
chmod +x run.sh
exit
```

- ❖ A docker hub account is needed to perform this operation.

Commit and push the docker image to docker hub:

```
docker commit $(docker ps -qa -n 1) <docker hub account>/mpich
docker login -u <docker hub account> -p <password>
docker push <docker hub account>/mpich
```

- ❖ This is a way to build docker image manually, alternative is using Dockerfile.

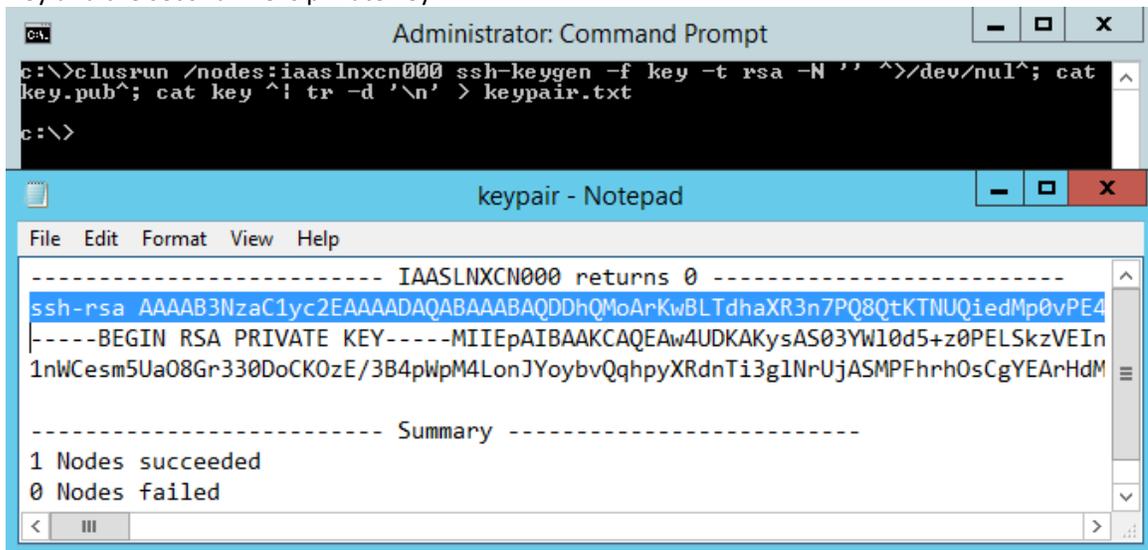
- [Run MPI task as root](#)

Generate key pair by command ssh-keygen on Linux node with clusrun.

Open a Command Prompt and run below command:

```
clusrun /nodes:iaaslnxcn000 ssh-keygen -f key -t rsa -N '' ^>/dev/nul^; cat key.pub^; cat key ^| tr -d '\n' > keypair.txt
```

A key pair will be generated and stored in the file keypair.txt, the first line of clusrun result is public key and the second line is private key.



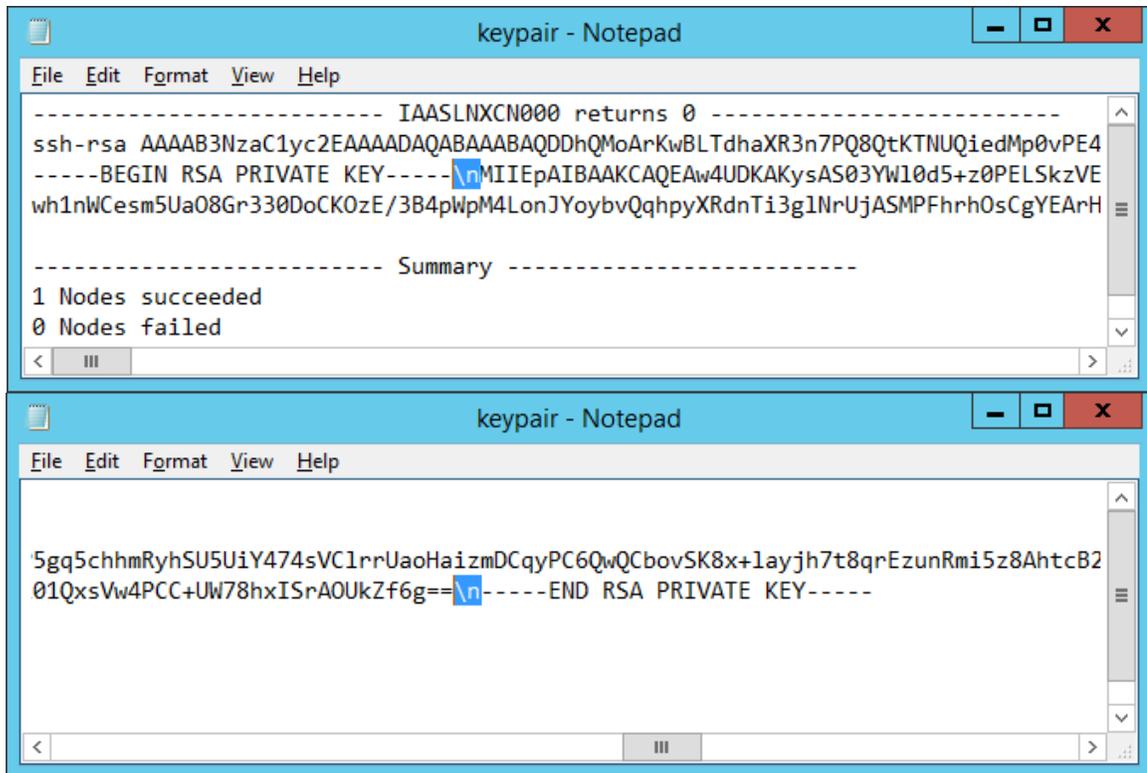
The screenshot shows two windows. The top window is a Command Prompt titled "Administrator: Command Prompt" with the following command and output:

```
c:\>clusrun /nodes:iaaslnxcn000 ssh-keygen -f key -t rsa -N '' ^>/dev/nul^; cat key.pub^; cat key ^| tr -d '\n' > keypair.txt
c:\>
```

The bottom window is a Notepad window titled "keypair - Notepad" showing the output of the command:

```
----- IAASLNXCN000 returns 0 -----
ssh-rsa AAAAB3NzaC1yc2EAAAADAQABAAQDDhQMoArKwBLTdhaxR3n7PQ8QtKTNUQiedMp0vPE4
-----BEGIN RSA PRIVATE KEY-----MIIEpAIBAAKCAQEaw4UDKAKysAS03YWl0d5+z0PELSkzVEIn
1nWCesm5Ua08Gr330DoCK0zE/3B4pWpM4LonJYoybvQqhpYXRdnTi3g1NrUjASMPFhrh0sCgYEAhHdM
-----
Summary
1 Nodes succeeded
0 Nodes failed
```

Insert 2 newline characters '\n' in public key, one after "-----BEGIN RSA PRIVATE KEY-----" and the other before "-----END RSA PRIVATE KEY-----".



Set mutual trust between Linux compute nodes for root.

Creating ssh key file id_rsa with clusrun in Command Prompt:

```

clusrun /nodegroup:linuxnodes mkdir -p /root/.ssh
clusrun /nodegroup:linuxnodes echo '<public key>' ^> /root/.ssh/authorized_keys
clusrun /nodegroup:linuxnodes echo -e '<private key>' ^> /root/.ssh/id_rsa
clusrun /nodegroup:linuxnodes chmod 600 /root/.ssh/id_rsa

```

- ❖ The mutual trust could also be set by other ways like login Linux compute nodes with putty and configuring manually. For normal user, a credential file containing the key pair could be passed with command hpccred so that the mutual trust could be set automatically by HPC Pack.

Submit a job including a docker task to run the MPI application in the docker image we built in above step:

```

job submit /env:ccp_docker_image=<docker hub account>/mpich /numnodes:4
/workdir:/mpisample ./run.sh 16

```

```

c:\>clusrun /nodes:iaaslnxcn000 ssh-keygen -f key -t rsa -N '' ^>/dev/nul^; cat
key.pub^; cat key ^| tr -d '\n' > keypair.txt

c:\>clusrun /nodegroup:linuxnodes mkdir -p /root/.ssh
----- IAASLNXCN002 returns 0 -----
----- IAASLNXCN003 returns 0 -----
----- IAASLNXCN000 returns 0 -----
----- IAASLNXCN001 returns 0 -----

----- Summary -----
4 Nodes succeeded
0 Nodes failed

c:\>clusrun /nodegroup:linuxnodes echo 'ssh-rsa AAAAB3NzaC1yc2EAAAADAQABAAQDD
hQMoArKwBLTidhaXR3n7PQ8QtKTNUQiedMpvPE4DkPsmMMFW/REZgo1yyh73BUIq3AULdmDEUu0gScAG
/ppmfKhjT1tLKrPbuwax93+4kZ/oFwIsKfeNducDduFq1LLbc8dpdXgo1df0gLPiTU+Hx06TBMzwQ9fH
QeAtyrmFeklfOcUZXFB3n56NHgwb0a8WpSmzwdRH iWzFR9Ah6moQh4PpbGe l0+6tSWqkXYS p5f lSUJX
dQReso+ehWEUwLfM7UZlpcEMlJso/jKH4hyUltUvpcyeghMhBF64KbMB4XAFZCv9C2iMp321blYI7iQSU
JpqD6yH/eNa9cCWx2eS3 root@IaaSLnxCN000' ^> /root/.ssh/authorized_keys
----- IAASLNXCN000 returns 0 -----
----- IAASLNXCN001 returns 0 -----
----- IAASLNXCN003 returns 0 -----
----- IAASLNXCN002 returns 0 -----

----- Summary -----
4 Nodes succeeded
0 Nodes failed

c:\>clusrun /nodegroup:linuxnodes echo -e '-----BEGIN RSA PRIVATE KEY-----\nMIIE
pAIBAAKCAQEAw4UDKAKysAS03YWl0d5+z0PELSkzUEInTKdLzX0A5D0pjDBUv0RGYKncsoe9wUcKtwF
S3ZgxFbtIEAnBv6aZnYoY09bSyyqz278GsfD/uJGf6BcCLCn3jXbnA3bhatZS23PHaXU6qNXX9IC6Yk1P
h8dOkwTM8EPXx0HgLcQ5hXpJXznFGUxW95+eJR4MLgdGvFj7Js8HUR41sXUfQIEpqqEIE D6WxnpdPurU1
qpF2EgeX5U1SU3UEXrKPnm1hFzi3z01WZaXBDJSbKP4yh+IcLJbVaXmnoITIQReuCMzAeFwBWqr/QtoJ
Kd9tWymI04kE1Caag+sh/3jWvXA119nktwIDAQABAoIBADEFjgdqRawUP+YDfMQCg8a/3sRtEipD3m2sn
qFeRyW2eBc6q3DyKl6+Ynyl laEje+nU5HN+NXd6alJgSMP/RB08Qr2UeMwYf y9pD1UQvBtUFRfIXrkUZ
ow1Du+WgACFQwJxDS e/YRb6fGaps0+uu0Hkofzf18efEtItv8/6uP5gq5c hhmRyhSU5Uiy474sUC1rrU
aoHaizmDCgyPC6QwQCbovSK8x+lajjh7t8qrEzunRmi5z8AhtcB20hPcA0bs8h+sW7ksZ1us7vzfzRAsI
c8b0Uhp0q4TixjUvz0gyywlbedW+K5F/14ajTs/t0ajmQdErYYqzQyM4fS2xAA34SKfKcGyEA9bXA+K5C
5EbAzeDe0Wr1JFPcYjPg1385xWNIp4pfGG0ohZiQLt2Ymdr1SEn2LatpUidsM7rsOPLoWTEMTLo1UR7C
alsEsUlg7gfuFM1Ldc5p0Qig0uttYqTPJJPfNKb6mmNTEBvhlK28UKt0Y0171topmcAArBhyjGdGXHJ+
XGUCyYEAy7UrX7H7f6jGckGCjze+LlKdvoq/2tQ0KxmwFo1R86Y5UydmKoRmRsm6IdIFpDeFGTD7vXki
aRDn+yWeThUza7Tb5eDhWYT00szwh1nWCesm5Ua08Gr330DoCK0zE/3B4pWpM4LonJYoybvQghpyXRdn
Tig1NrUjASMPFhrhOsCgYEAhHdMIQGoQDUNhjuOEw53s100tv07MHHm15S/KmagDMLieKBafPu9id4
gpovquTDNk7rmw/cWUjf1CPWX1Dgtb1if+atI4J8Plf17/MutF0q7Hak9CF+IRxPrXvokzqUcQK/HfJ/
2UdtHC4kvsMXOdaeFqZ3k1M1G00TJrbSzsUCgYEAAGTzSMCqrwTWx48I1nlg+Zty93UGhc4ce4H8bxt
ne+INCQiHexDdpyRj93zU1T0HD0miQEawrQn0AMe60IrdhpuVnNUGAZ27HubLEv2iKP6nT1hQ023E1R
sphTN+d60/Fuqxv4tzPh1fz+ASkYKhqLLz51Tf07eGEpbqOgaS8CgYACm6+4AoAiyRbYcOsduXSWnPCo
DSFmcUzZGQc6nkh/33L4PG8XmnD0m0DSMC49mKUENIGCB0NuUcoLYIvTNj4PERhJWUCeQY30AdivzCzn
0AIFJ8/t//KxtSverHj35po2GJwTWkAb69HHNz01QxsUw4PCC+UW78hxISrAOUkZf6g==\n-----END R
SA PRIVATE KEY-----' ^> /root/.ssh/id_rsa
----- IAASLNXCN000 returns 0 -----
----- IAASLNXCN002 returns 0 -----
----- IAASLNXCN003 returns 0 -----
----- IAASLNXCN001 returns 0 -----

----- Summary -----
4 Nodes succeeded
0 Nodes failed

c:\>clusrun /nodegroup:linuxnodes chmod 600 /root/.ssh/id_rsa
----- IAASLNXCN003 returns 0 -----
----- IAASLNXCN000 returns 0 -----
----- IAASLNXCN001 returns 0 -----
----- IAASLNXCN002 returns 0 -----

----- Summary -----
4 Nodes succeeded
0 Nodes failed

c:\>job submit /env:ccp_docker_image=zclok/mpich /numnodes:4 /workdir:/mpisample
./run.sh 16
Job has been submitted. ID: 37.

c:\>task view 37.1

```

```
c:\>job submit /env:ccp_docker_image=zclock/mpich /numnodes:4 /workdir:/mpisample
./run.sh 16
Job has been submitted. ID: 37.
```

```
c:\>task view 37.1
```

```
Task Id           : 37.1
State             : Finished
Task Name        :
Command Line     : ./run.sh 16
Resource Request : 4-4 nodes
Task Type        : Basic Task
Submit Time      : 9/30/2017 6:13:44 AM
Start Time       : 9/30/2017 6:13:44 AM
End Time         : 9/30/2017 6:15:24 AM
Elapsed Time     : 00:00:01:39
Total Kernel Time : 00:00:00:00:180 d:h:m:s:ms
Total User Time  : 00:00:00:00:430 d:h:m:s:ms
Working Set      : 21 MB
Processes        :
Required Nodes   :
Allocated Nodes  :
IAASLNXCN000, IAASLNXCN001, IAASLNXCN002, IAASLNXCN003
Pending Reason   :
Exit Code        : 0
Error Message    :
```

```
Output
STDOUT: Hello from IaaSLnxCN001, rank 5 out of 16 processors.
Hello from IaaSLnxCN001, rank 1 out of 16 processors.
Hello from IaaSLnxCN001, rank 9 out of 16 processors.
Hello from IaaSLnxCN001, rank 13 out of 16 processors.
Hello from IaaSLnxCN002, rank 10 out of 16 processors.
Hello from IaaSLnxCN002, rank 14 out of 16 processors.
Hello from IaaSLnxCN000, rank 4 out of 16 processors.
Hello from IaaSLnxCN000, rank 8 out of 16 processors.
Hello from IaaSLnxCN000, rank 12 out of 16 processors.
Hello from IaaSLnxCN003, rank 7 out of 16 processors.
Hello from IaaSLnxCN003, rank 11 out of 16 processors.
Hello from IaaSLnxCN003, rank 15 out of 16 processors.
Hello from IaaSLnxCN003, rank 3 out of 16 processors.
Hello from IaaSLnxCN002, rank 6 out of 16 processors.
Hello from IaaSLnxCN002, rank 2 out of 16 processors.
Hello from IaaSLnxCN000, rank 0 out of 16 processors.
```

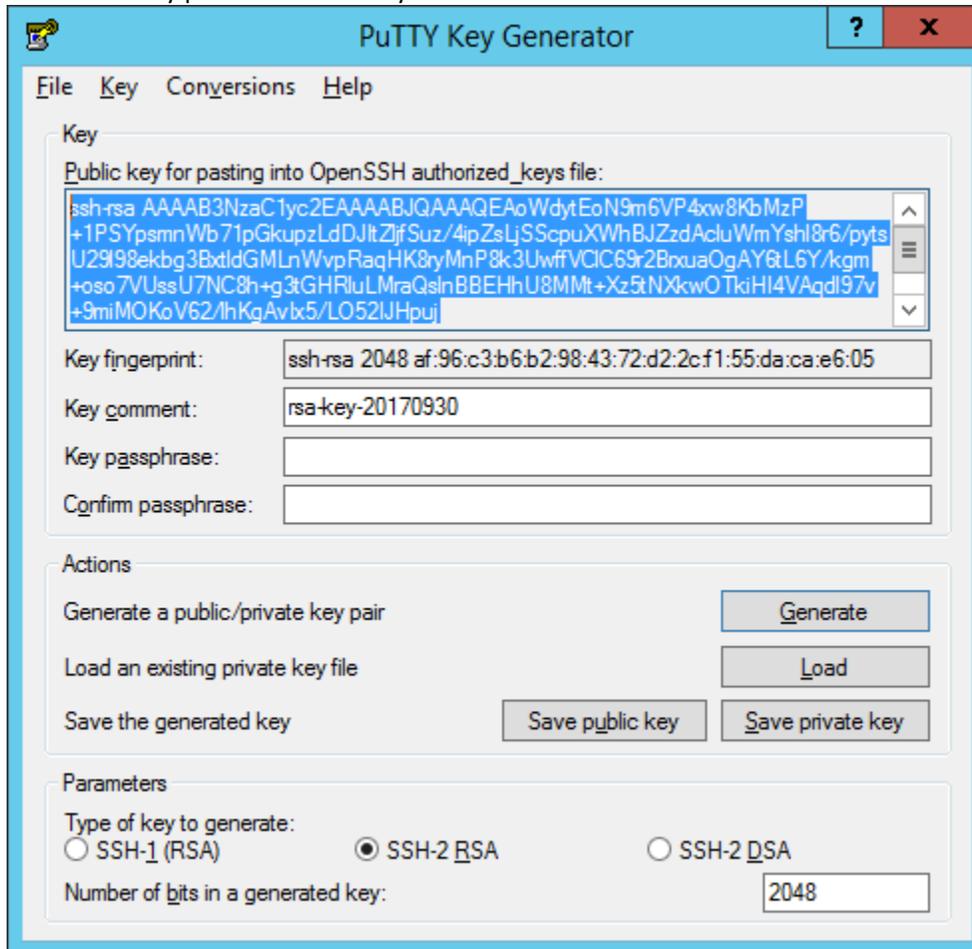
```
STDERR:
```

```
c:\>
```

- Run MPI task as a normal user

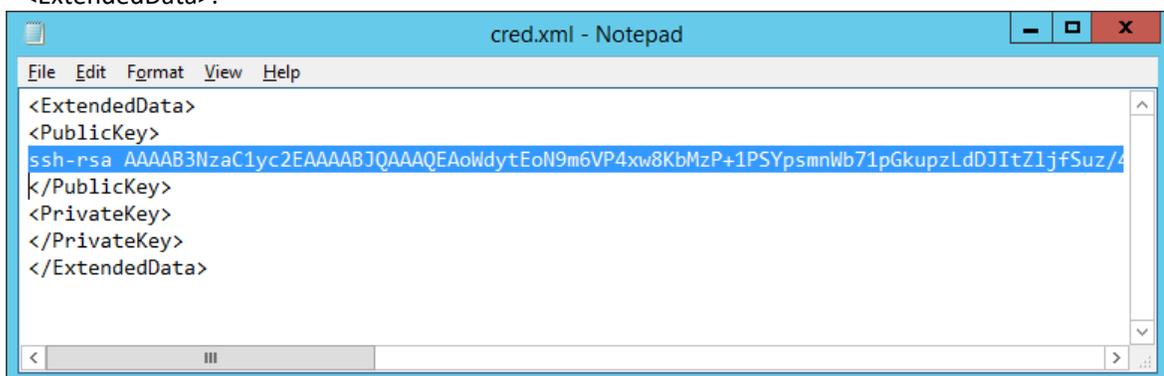
create a normal user, like “user1”.

Generate a key pair with PuTTY Key Generator.

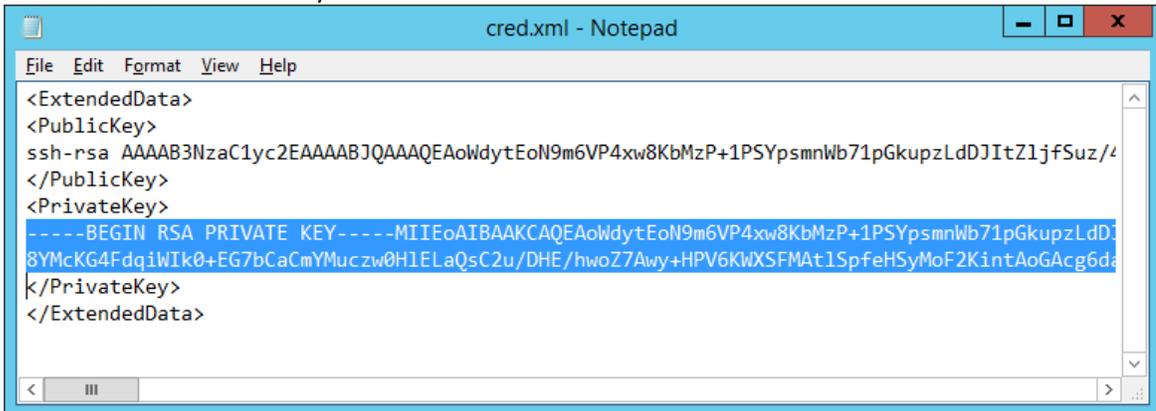


Create a xml file named “cred.xml” in C:\.

Copy public key from the window to file “cred.xml” as the content of xml element <PublicKey> in <ExtendedData>.



Export OpenSSH Key in PuTTY Key Generator as private key, and copy it to "cred.xml" as the content of xml element <PrivateKey> in <ExtendedData>.



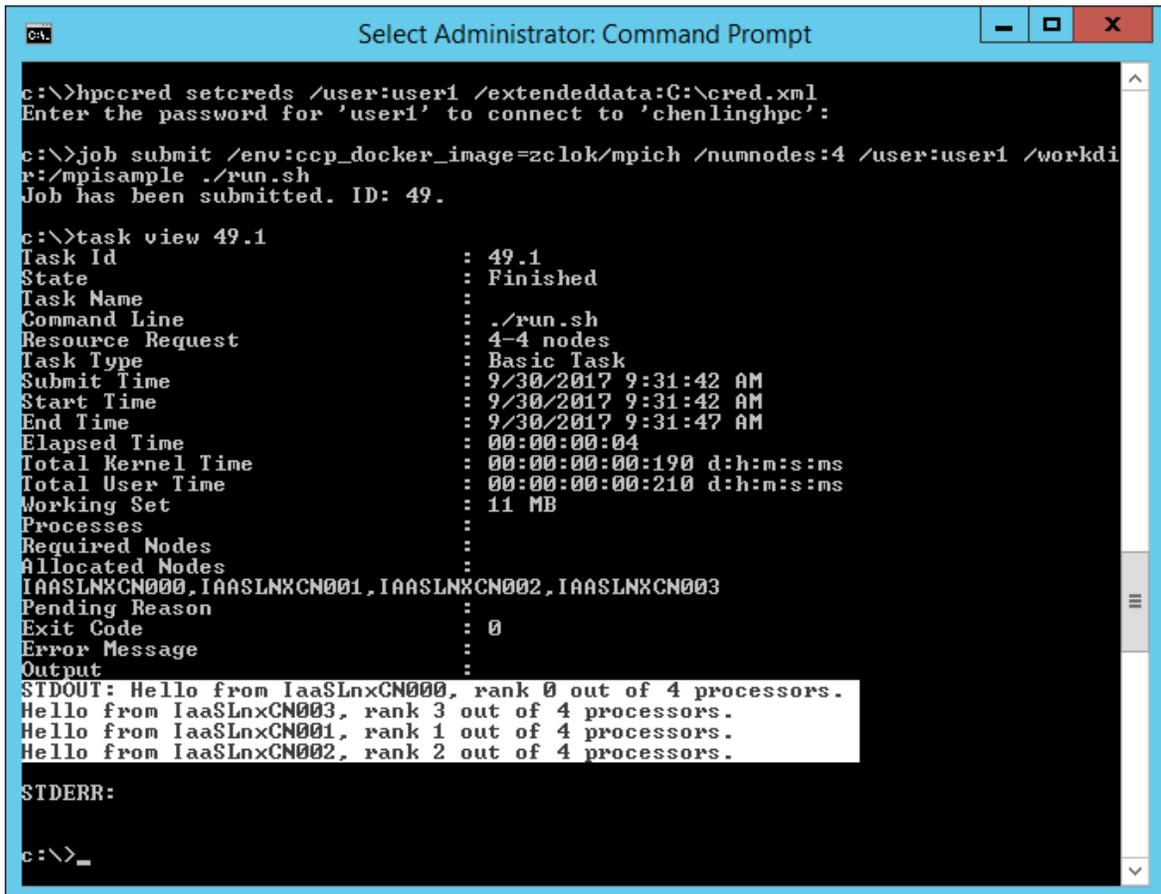
```
cred.xml - Notepad
File Edit Format View Help
<ExtendedData>
<PublicKey>
ssh-rsa AAAAB3NzaC1yc2EAAAABJQAAAQEAAoWdytEoN9m6VP4xw8KbMzP+1PSYpsmnWb71pGkupzLdDJItZ1jfSuz/4
</PublicKey>
<PrivateKey>
-----BEGIN RSA PRIVATE KEY-----MIIEOAIBAACAQEAAoWdytEoN9m6VP4xw8KbMzP+1PSYpsmnWb71pGkupzLdDJ
8YMcK64FdqiWIk0+EG7bCaCmYMuczw0H1ELaQsC2u/DHE/hwoZ7Awy+HPV6KWXSFMat1Spf0HSyMoF2KintAoGAcg6d
</PrivateKey>
</ExtendedData>
```

Set credential for normal user 'user1' in Command Prompt:

```
hpccred setcreds /user:user1 /extendeddata:C:\cred.xml
```

Submit job to run MPI application in docker task as user1:

```
job submit /env:ccp_docker_image=<docker hub account>/mpich /numnodes:4 /user:user1
/workdir:/mpisample ./run.sh
```



```
Select Administrator: Command Prompt
c:\>hpccred setcreds /user:user1 /extendeddata:C:\cred.xml
Enter the password for 'user1' to connect to 'chenlinghpc':

c:\>job submit /env:ccp_docker_image=zclock/mpich /numnodes:4 /user:user1 /workdi
r:/mpisample ./run.sh
Job has been submitted. ID: 49.

c:\>task view 49.1
Task Id                : 49.1
State                  : Finished
Task Name              :
Command Line           : ./run.sh
Resource Request       : 4-4 nodes
Task Type              : Basic Task
Submit Time            : 9/30/2017 9:31:42 AM
Start Time             : 9/30/2017 9:31:42 AM
End Time               : 9/30/2017 9:31:47 AM
Elapsed Time           : 00:00:00:04
Total Kernel Time     : 00:00:00:00:190 d:h:m:s:ms
Total User Time        : 00:00:00:00:210 d:h:m:s:ms
Working Set            : 11 MB
Processes              :
Required Nodes         :
Allocated Nodes       :
IAASLNXCN000,IAASLNXCN001,IAASLNXCN002,IAASLNXCN003
Pending Reason         :
Exit Code              : 0
Error Message          :
Output
STDOUT: Hello from IaaSLnxCN000, rank 0 out of 4 processors.
Hello from IaaSLnxCN003, rank 3 out of 4 processors.
Hello from IaaSLnxCN001, rank 1 out of 4 processors.
Hello from IaaSLnxCN002, rank 2 out of 4 processors.

STDERR:

c:\>_
```