

# CSCI4180 Tutorial-2

# Hadoop Setup on OpenStack

# Windows Azure Guide

ZHANG, Mi

[mzhang@cse.cuhk.edu.hk](mailto:mzhang@cse.cuhk.edu.hk)

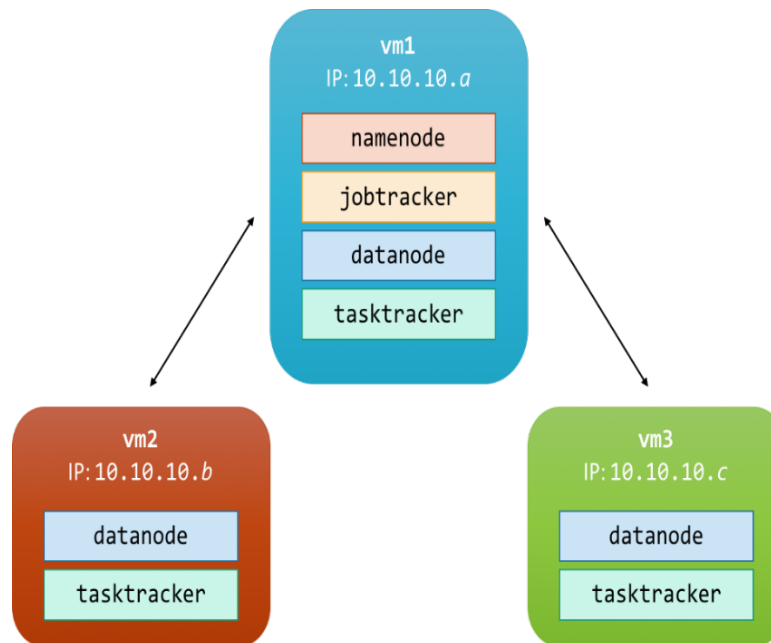
Sep. 24, 2015

# Outline

- Hadoop setup on OpenStack
  - Set up Hadoop cluster
  - Manage Hadoop cluster
  - WordCount Example
- Windows Azure guide
  - Access Azure
  - Create VMs
  - Install Hadoop

# Set up Hadoop Cluster

- We've created three VM instances of our own.
  - Architecture



# Set up Hadoop Cluster

- We'll set up small-scale Hadoop cluster using these VM instances.
- What you've done in tutorial-1:
  - Setting up HTTP proxy.
  - Installing Java.
  - Configuring **/etc/hosts**.

# Set up Hadoop Cluster

- Switch to normal user “hadoop”
  - `su - hadoop`
- If you do not have user “hadoop”
  - `adduser hadoop`
  - enter your password when necessary...
  - `su - hadoop`

# Set up Hadoop Cluster

- Download Hadoop on EACH node
  - wget  
<http://archive.apache.org/dist/hadoop/core/hadoop-0.20.203.0/hadoop-0.20.203.0rc1.tar.gz>
- Place Hadoop in home directory on EACH node
  - tar xzf hadoop-0.\*.\*.tar.gz
  - mv hadoop-0.\*.\* hadoop

# Set up Hadoop Cluster

- Set environment variable on EACH node.
  - I recommend you put them in `~/.bashrc`
    - `export HADOOP_HOME=~/.hadoop`
    - `export PATH=$PATH:$HADOOP_HOME/bin`
- Set hadoop environment on EACH node.
  - Append the following lines to `~/hadoop/conf/hadoop-env.sh`
    - `export JAVA_HOME=/usr/lib/jvm/java-7-oracle`
    - `#depends on where you put the jvm`
    - `export HADOOP_OPTS=-Djava.net.preferIPv4Stack=true`
- Set path for HDFS storage on EACH node.
  - `#under HOME directory`
  - `mkdir hadoop/tmp`

# Set up Hadoop Cluster

- Configure SSH on EACH node
  - `ssh-keygen -t rsa -P ""`
  - `cat $HOME/.ssh/id_rsa.pub >> $HOME/.ssh/authorized_keys`
- Configure SSH on namenode only
  - `ssh-copy-id -i $HOME/.ssh/id_rsa.pub hadoop@vm1`
  - `ssh-copy-id -i $HOME/.ssh/id_rsa.pub hadoop@vm2`
  - `ssh-copy-id -i $HOME/.ssh/id_rsa.pub hadoop@vm3`
- Check SSH configuration
  - whether namenode can ssh all the datanodes without typing password. E.g.,
    - `ssh vm2`



# Set up Hadoop Cluster

- Set hadoop core on EACH node
  - Add property in `~/hadoop/conf/core-site.xml`

```
<property>
  <name>hadoop.tmp.dir</name>
  <value>/home/hadoop/hadoop/tmp</value>
</property>
<property>
  <name>fs.default.name</name>
  <value>hdfs://vm1:54310</value>
</property>
```

# Set up Hadoop Cluster

- Set hadoop mapreduce on EACH node
  - Add property in `~/hadoop/conf/mapred-site.xml`

```
<property>  
  <name>mapred.job.tracker</name>  
  <value>vm1:54311</value>  
</property>
```

# Set up Hadoop Cluster

- Set hadoop HDFS on EACH node
  - Add property in `~/hadoop/conf/hdfs-site.xml`

```
<property>  
  <name>dfs.replication</name>  
  <value>3</value>  
</property>
```

# Set up Hadoop Cluster

- Set hadoop master on namenode
  - Add hostname which is supposed to run *NameNode* and *JobTracker* in `~/hadoop/conf/masters`
    - vm1
- Set hadoop slaves on namenode
  - Add hostname which is supposed to run *DataNode* and *TaskTracker* in `~/hadoop/conf/slaves`
    - vm1
    - vm2
    - vm3

# Set up Hadoop Cluster

- Format namenode on namenode
  - `hadoop namenode -format`
- Start hadoop on namenode
  - `start-dfs.sh`
  - `start-mapred.sh`
  - # you can type “jps” to see whether the startup is successful.
- Stop hadoop on namenode
  - `stop-mapred.sh`
  - `stop-dfs.sh`

# Set up Hadoop Cluster

- Some operations related to HDFS
  - From Local to HDFS
    - `hadoop dfs -copyFromLocal <local dir/file> <hdfs URI>`  
( for user home URI: /home/hadoop )
  - From HDFS to Local
    - `hadoop dfs -copyToLocal <hdfs URI> <local dir/file>`
  - List files in HDFS
    - `hadoop dfs -ls <hdfs URI>`
  - Cat files in HDFS
    - `hadoop dfs -cat <hdfs URI>`

# Manage Hadoop Cluster

- Add one more instance into cluster
  - Stop Hadoop services on namenode
  - For the new instance, repeat steps from slide 4 to slide 11
  - Add IP of new instance in `~/hadoop/conf/slaves` on namenode
  - Format namenode and start Hadoop

# Manage Hadoop Cluster

- Remove one instance from cluster
  - Stop Hadoop services on namenode
  - Remove IP of the instance from `~/hadoop/conf/slaves`
  - Format namenode and start Hadoop



# WordCount Example

- Download the java source code from course website, say, WordCount.java, to your namenode, home directory
  - Compile and run the program
    - `mkdir wordcount`
    - `javac -classpath $HADOOP_HOME/hadoop-core-0.20.203.0.jar WordCount.java -d wordcount`
    - `jar -cvf wordcount.jar -C wordcount/ .`
    - `hadoop jar wordcount.jar org.myorg.WordCount */HDFS URI/to/input/file* */HDFS URI/to/output/directory*`
- Note that the part-r-00000 is the actual output.

# Windows Azure platform

- Windows Azure guide
  - Access Azure
  - Create VMs
  - Install Hadoop

# Overview

- Get the 14-character code before you start.
- Redeem your Windows Azure at  
<https://www.microsoftazurepass.com/azureu>

# Redeem the pass

Microsoft Azure

All Microsoft Sites

Search Microsoft Azure

## Try Microsoft Azure

The Microsoft Educator Grant Program provides access to Azure for use in the classroom by university students and their professor.



### ELIGIBILITY REQUIREMENTS:

- Requester must be a faculty member of an accredited university
- Passes must be used for a specific class
  - a. Multiple classes should have multiple requests
  - b. Multiple sections of same class may use one request
- Passes must be used by students
- Faculty must provide all information requested
  - a. If you are unable to provide any of the requests, please attach a document explaining reason for omission, and another method to verify.

### APPLICATION IS FOR EDUCATORS/FACULTY ONLY.

If you do not meet the eligibility requirements, the free trial program is also available for 1-month at no cost. ([Free Trial Link](#))

All fields are required.

Full Name:

University Faculty Email:

Country:

University (full name):

Course Name:

Course Description:  
( Please specifically indicate how you intend to leverage Microsoft Azure; 4000 character limit )

Course Start Date:

# of Student Passes Requested:

Course URL:

Faculty URL:

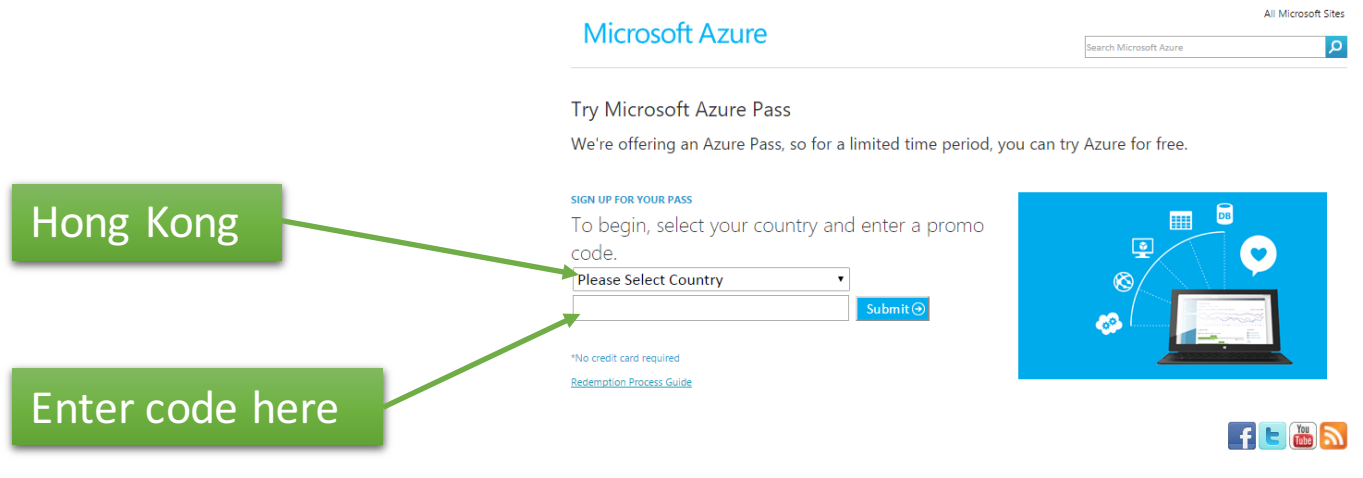
Course Syllabus:  No file chosen

University ID:  No file chosen

Click here

Have a code already?  
Redeem it here.

# Redeem the pass



The screenshot shows the Microsoft Azure Pass redemption page. At the top, the 'Microsoft Azure' logo is on the left, and 'All Microsoft Sites' with a search bar is on the right. The main heading is 'Try Microsoft Azure Pass', followed by the text 'We're offering an Azure Pass, so for a limited time period, you can try Azure for free.' Below this, a section titled 'SIGN UP FOR YOUR PASS' instructs users to 'To begin, select your country and enter a promo code.' There is a dropdown menu labeled 'Please Select Country' and a 'Submit' button. A green box labeled 'Hong Kong' has an arrow pointing to the country dropdown. Another green box labeled 'Enter code here' has an arrow pointing to the promo code input field. Below the input fields, it says '\*No credit card required' and provides a link to the 'Redemption Process Guide'. To the right of the form is a blue graphic showing a laptop with various cloud service icons (mail, calendar, database, etc.) connected to it. At the bottom right are social media icons for Facebook, Twitter, YouTube, and RSS.

- Followed with some register information.
- To redeem the pass, you also need a windows live account.

# Redeem the pass

Microsoft Azure

All Microsoft Sites

Search Microsoft Azure



Login your windows account

Please sign in using your Organizational account or Microsoft account, fill out the remaining information, then click submit.

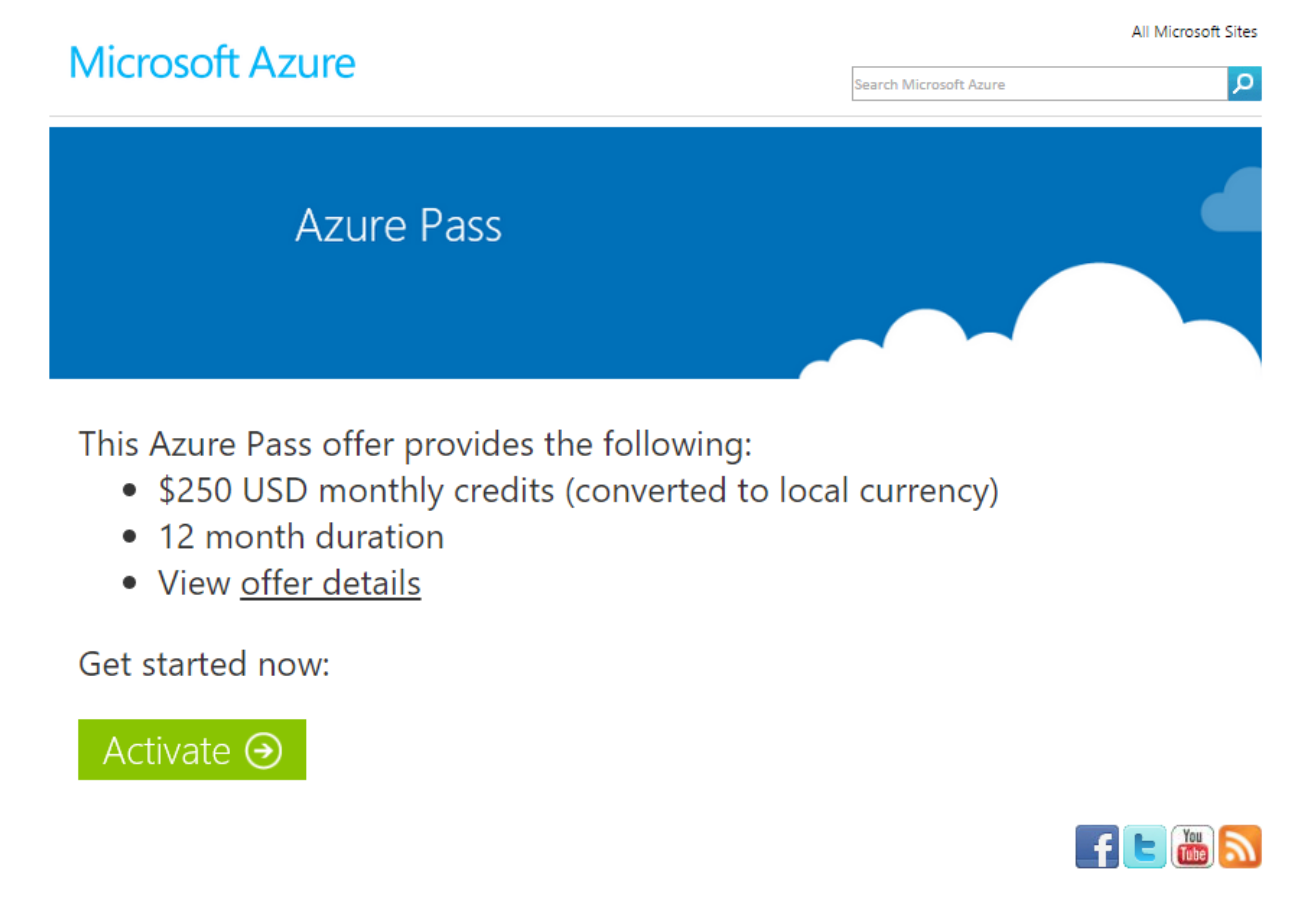
Sign in

Create a Microsoft Account: <https://signup.live.com>

Create an Organizational Account: <https://account.windowsazure.com/organization>



# Redeem the pass



The screenshot shows the Microsoft Azure website with a search bar at the top right. The main heading is "Azure Pass" on a blue background with white clouds. Below this, a list of benefits is provided, followed by a "Get started now:" section with a green "Activate" button. Social media icons for Facebook, Twitter, YouTube, and RSS are at the bottom right.

Microsoft Azure

All Microsoft Sites

Search Microsoft Azure

## Azure Pass

This Azure Pass offer provides the following:

- \$250 USD monthly credits (converted to local currency)
- 12 month duration
- View [offer details](#)

Get started now:

Activate ➔

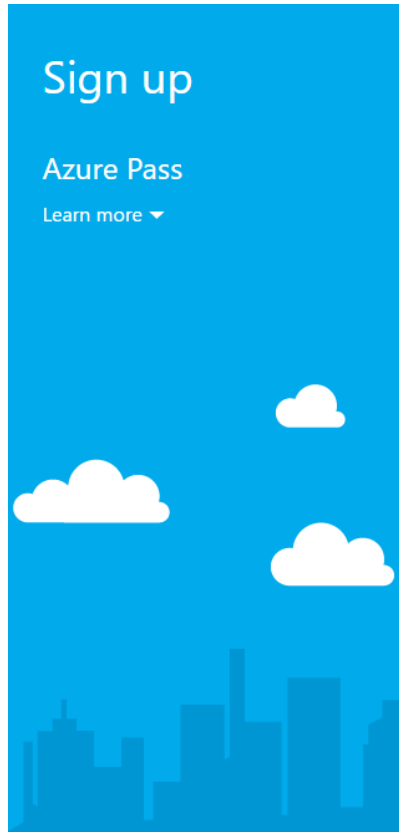
Facebook Twitter YouTube RSS

# Redeem the pass

## Sign up

### Azure Pass

[Learn more](#) ▼



## Microsoft Azure

1155074923@link.cuhk.edu.hk ▼

- ### About you

FIRST NAME	LAST NAME	COUNTRY/REGION ?
<input type="text" value="Mi"/>	<input type="text" value="ZHANG"/>	<input td="" type="text" value="Hong Kong SAR" ▼<=""/>
CONTACT EMAIL ?	COMPANY/SCHOOL	
<input type="text" value="mzhang@cse.cuhk.edu.hk"/>	<input type="text" value="- Optional -"/>	
- ### Contact phone number ?
- ### Agreement

☒ I agree to the [subscription agreement](#), [offer details](#), and [privacy statement](#).

☒ Microsoft may use my email and phone to provide special Microsoft Azure offers.

[Sign up](#) ➔

phone number

English

© 2015 Microsoft [Privacy & Cookies](#) [Trademarks](#) [Legal](#) [Support](#) [Give Us Feedback](#)

Microsoft



# Redeem the pass

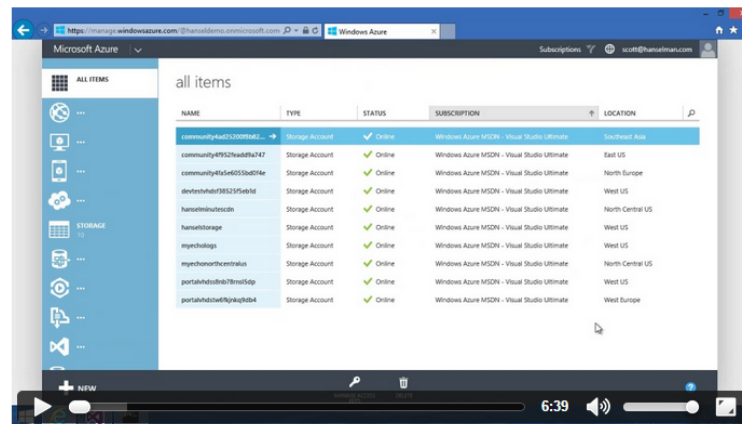
Welcome to Microsoft Azure!

Your subscription - Azure Pass

Your subscription is ready for you!

Start managing my service >

Take a tour of the management experience while you wait.



## Tutorials

Get started with...

[Web Apps](#)

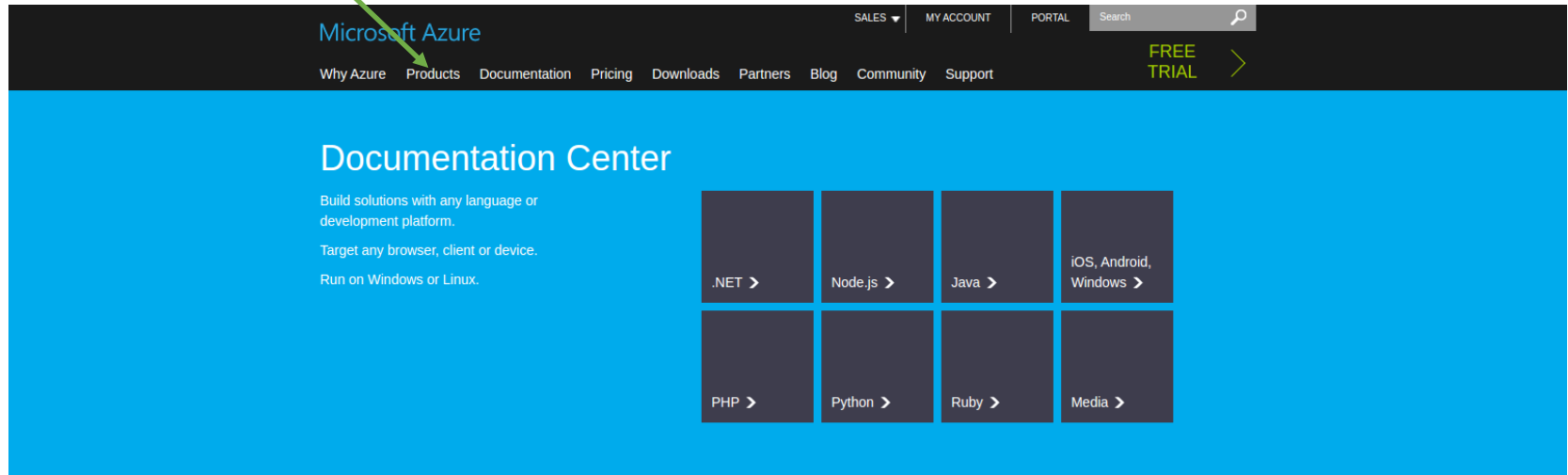
[Virtual Machines](#)

[SQL Database](#)

[Storage](#)

# Create VMs

Click Products



Microsoft Azure

SALES ▼ | MY ACCOUNT | PORTAL | Search

Why Azure | **Products** | Documentation | Pricing | Downloads | Partners | Blog | Community | Support

**FREE TRIAL** >

## Documentation Center

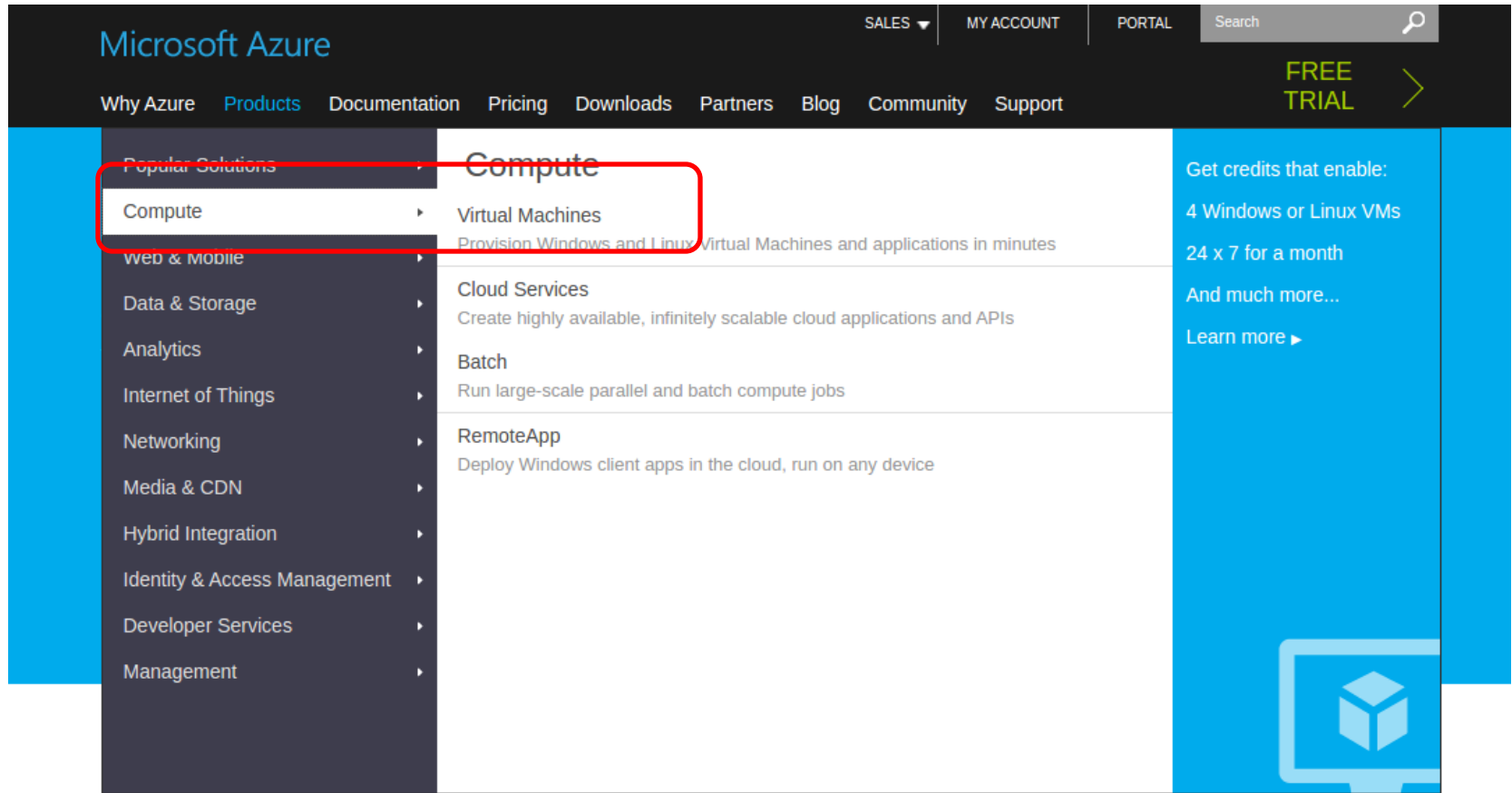
Build solutions with any language or development platform.  
Target any browser, client or device.  
Run on Windows or Linux.

<a href="#">.NET &gt;</a>	<a href="#">Node.js &gt;</a>	<a href="#">Java &gt;</a>	<a href="#">iOS, Android, Windows &gt;</a>
<a href="#">PHP &gt;</a>	<a href="#">Python &gt;</a>	<a href="#">Ruby &gt;</a>	<a href="#">Media &gt;</a>

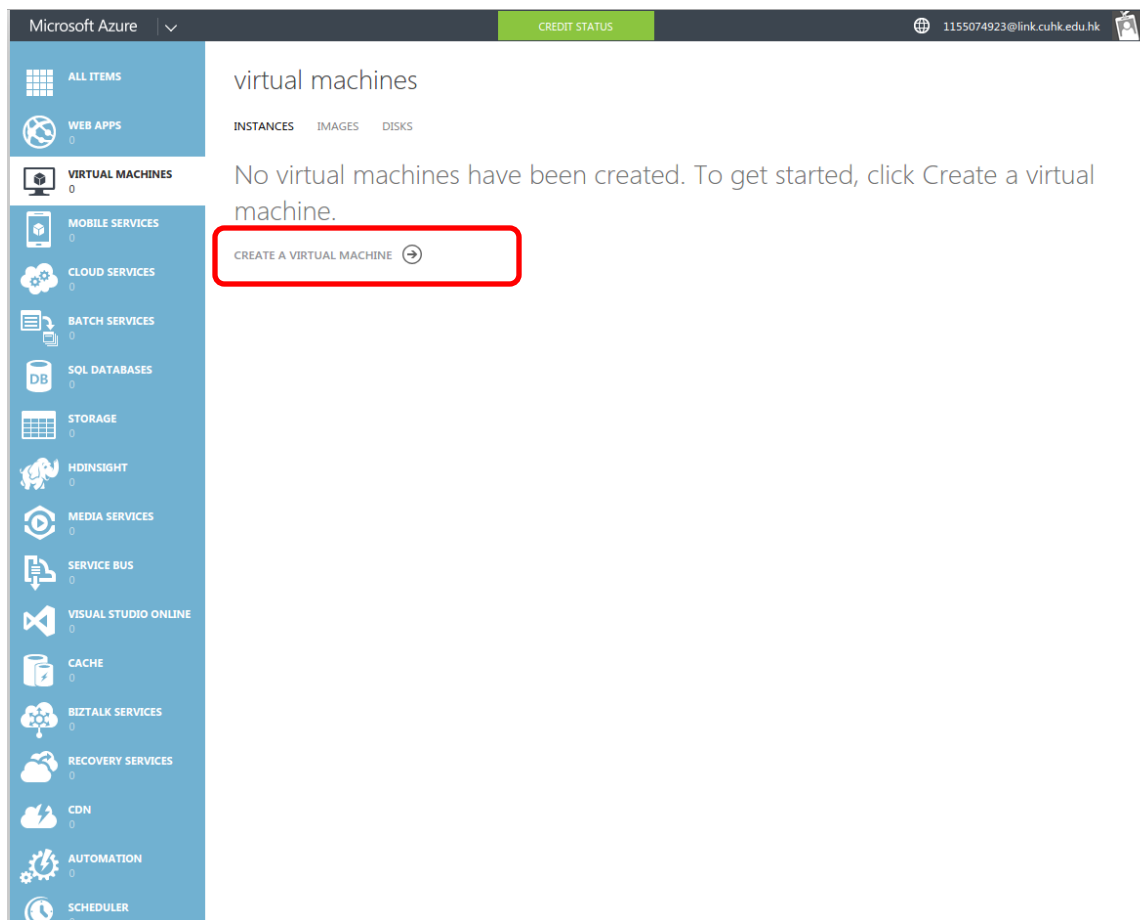
## Documentation by service

Compute	Web & Mobile	Data & Storage	Analytics	Networking
<a href="#">Virtual Machines</a>	<a href="#">App Service</a>	<a href="#">DocumentDB</a>	<a href="#">HDInsight</a>	<a href="#">Virtual Network</a>
<a href="#">Cloud Services</a>	<a href="#">Websites</a>	<a href="#">SQL Database</a>	<a href="#">Machine Learning</a>	<a href="#">ExpressRoute</a>
<a href="#">Batch</a>	<a href="#">Mobile Services</a>	<a href="#">Redis Cache</a>	<a href="#">Stream Analytics</a>	<a href="#">Application Gateway</a>
<a href="#">RemoteApp</a>	<a href="#">API Management</a>	<a href="#">Storage</a>	<a href="#">Data Factory</a>	<a href="#">Traffic Manager</a>
<a href="#">Service Fabric</a>	<a href="#">Push Notifications</a>	<a href="#">StorSimple</a>	<a href="#">Event Hubs</a>	<a href="#">DNS</a>
	<a href="#">Mobile Engagement</a>	<a href="#">Azure Search</a>	<a href="#">Data Catalog</a>	<a href="#">Load Balancer</a>
		<a href="#">SQL Data Warehouse</a>		<a href="#">VPN Gateway</a>

# Create VMs



# Create VMs



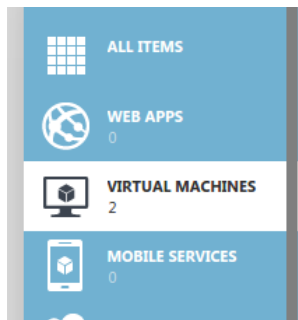
# Create VMs

The screenshot shows the Azure portal's 'NEW' page. On the left, the 'COMPUTE' category is selected, and 'VIRTUAL MACHINE' is chosen from the list. The 'QUICK CREATE' section on the right contains the following fields:

- DNS NAME:** csci4180vm1
- IMAGE:** Ubuntu Server 12.04 L' (highlighted with a red box)
- SIZE:** A1 (1 core, 1.75 GB me) (highlighted with a red box)
- USER NAME:** azureuser
- NEW PASSWORD:** [Redacted]
- CONFIRM:** [Redacted]
- REGION/AFFINITY GROUP:** Southeast Asia

At the bottom right, the 'CREATE A VIRTUAL MACHINE' button is highlighted with a red box. Below the form, a disclaimer states: 'By clicking Create, I agree this software is from Canonical and Canonical's legal terms apply to it. Microsoft doesn't provide rights for third party software.'

# Create VMs

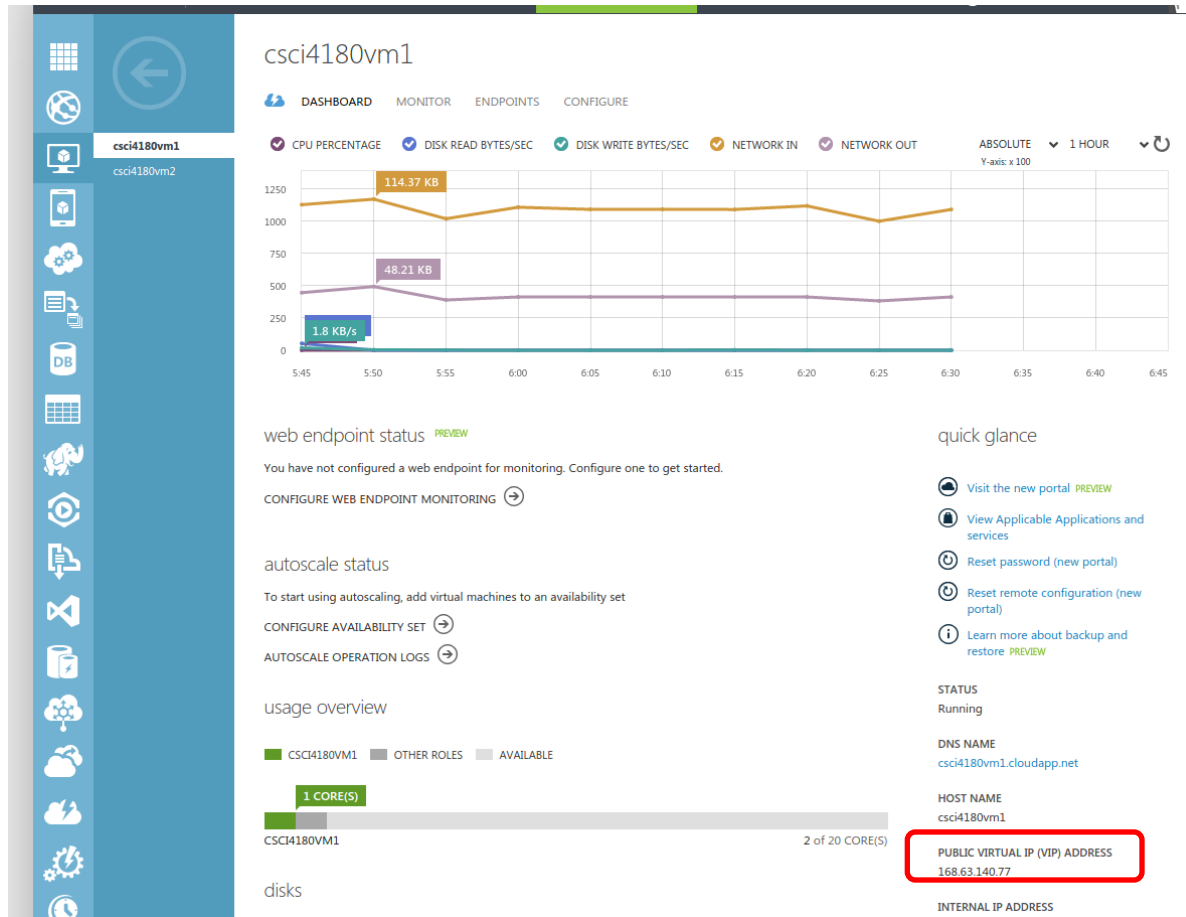


## virtual machines

INSTANCES IMAGES DISKS

NAME		STATUS	SUBSCRIPTION	LOCATION	DNS NAME	
csci4180vm1	→	✓ Running	Azure Pass	Southeast Asia	csci4180vm1.cloudapp.net	
csci4180vm2		✓ Running	Azure Pass	Southeast Asia	csci4180vm2.cloudapp.net	

# Create VMs



You can ssh  
to your VM  
using this IP.

# Create VMs

- In your terminal, **ssh azureuser@\*your vm IP\***

```
azureuser@csci4180vm1.cloudapp.net's password:
Welcome to Ubuntu 12.04.5 LTS (GNU/Linux 3.13.0-63-generic x86_64)

* Documentation:  https://help.ubuntu.com/

System information as of Wed Sep 16 11:27:08 UTC 2015

System load:  0.04               Processes:            217
Usage of /:   3.8% of 28.80GB     Users logged in:     0
Memory usage: 7%                 IP address for eth0: 10.62.162.156
Swap usage:   0%

Graph this data and manage this system at:
https://landscape.canonical.com/

Get cloud support with Ubuntu Advantage Cloud Guest:
http://www.ubuntu.com/business/services/cloud

0 packages can be updated.
0 updates are security updates.

New release '14.04.3 LTS' available.
Run 'do-release-upgrade' to upgrade to it.
```



# Create VMs

- If you SSH the VM via cse wired network, you may need to configure your ssh setting.
  - Append the following lines to `~/.ssh/config`
    - Host \*your hostname\*
    - User azureuser
    - HostName \*your hostname\*
    - ProxyCommand nc -x socks.cse.cuhk.edu.hk:1080 %h %p
- Then you can login to the vm using ssh
  - ssh \*your hostname\*

# Install Hadoop

- Install Java on EACH VM:
  - `sudo apt-get update`
  - `sudo apt-get upgrade`
  - `sudo add-apt-repository ppa:webupd8team/java`
  - `sudo apt-get update`
  - `sudo apt-get install oracle-java7-installer`
- You could follow the instruction in Tutorial 1.
  - <http://mtyiou.github.io/csci4180-fall15/>

# Install Hadoop

- Repeat the process of installing hadoop on OpenStack from slide 5 to slide 11.
  - Slide 8: replace vm1, vm2 with their respective public IP.
  - Slide 9, 10: change "vm1" to "127.0.0.1" when editing .xml files.
- Set hadoop masters on namenode
  - Edit `~/hadoop/conf/masters`
    - 127.0.0.1
- Set hadoop slaves on namenode
  - Edit `~/hadoop/conf/slaves`
    - 127.0.0.1
    - \*Another vm public IP \*
- After starting HDFS and MapReduce, you can run the WordCount example.

**Thank you!**