# Carnegie Mellon University

# Introduction to AWS/GCP

16-824 Visual Learning and Recognition (Fall 2022)

Presented by Vanshaj (vanshajc@)

# Course Logistics

- You will receive $50 of AWS credit per homework

  - We will send this credit out soon via email

- $50 GCP credit sent for HW1

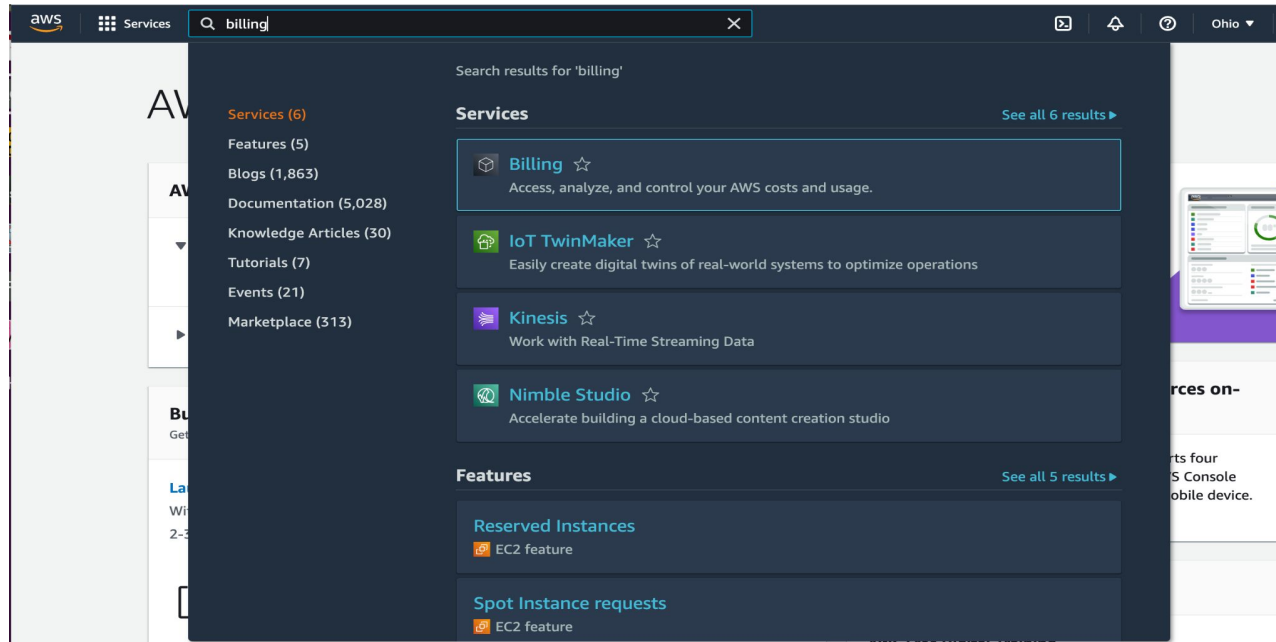  - We will send $100 more GCP credits later

# Other Resources

- Personal Laptop / Lab GPU Server

  - We cannot provide installation support

- Google Colab

  - Free GPU access (time limited)

  - Very easy to set up

  - Might need pro ($10/mo) for smoother experience

**Carnegie Mellon University**

# AWS Account Setup

- **Use andrew email address**

- https://aws.amazon.com/console/

- Redeem credits via billing console

- Apply the credits to your AWS account
  - AWS console - Billing - Credits - Redeem Credit - Enter Promo Code

- Apply the credits to your AWS account
  - AWS console - Billing - Credits - Redeem Credit - Enter Promo Code

# Amazon Elastic Compute Cloud (EC2)

- Visual learning is computation-intensive, requiring special hardware: GPU

- EC2 provides <u>resizable compute capacity</u> in the cloud

- Easy deployment of hardware, system softwares and data storage

  according to your demand and budget

**Carnegie
Mellon
University**

# EC2 Instances

- An instance is a virtual machine with specific hardware configuration

- Life cycle of an instance:
  - Launch
  - **Stop - when you are not working on EC2, only charges storage**
  - (Re)start
  - Terminate - be careful, the instance and root storage will be deleted.

# EC2 Instances Types

- Various instance types available - different in RAMs, CPUs, GPUs

- t2.micro: baseline CPU performance

  - 750 hours of free usage per month

  - (Optional) write code; test part of code

- p2.xlarge: what we use for assignments and projects

  - 4-core CPU, Nvidia Tesla K80 (11,441MB memory)

  - **costs $0.9 per Hour**

- g4dn instances: alternative for initial testing with GPU

# EC2 Dashboard

# EC2 Regions

# EC2 Limits



- Increase limit for P instances and G instances
- You won't be able to launch any instances with GPU until you request this service limit increase
- **Do this ASAP - it takes a while for AWS to approve**

# EC2 Limits



- Increase limit for P instances and G instances
- You won't be able to launch any instances with GPU until you request this service limit increase
- **Do this ASAP - it takes a while for AWS to approve**

**EC2 Dashboard** New

Events

Tags

Reports

Limits

▼ **INSTANCES**

Instances

Instance Types

Launch Templates New

Spot Requests

Savings Plans

Reserved Instances

Dedicated Hosts

Capacity Reservations

▼ **IMAGES**

AMIs

Bundle Tasks

▼ **ELASTIC BLOCK STORE**

Volumes

Snapshots

Lifecycle Manager

▼ **NETWORK & SECURITY**

Security Groups

Elastic IPs New

We're redesigning the EC2 console to make it easier to use and improve performance. We'll release new screens periodically. We encourage you to try them and let us know where console and the new console, use the New EC2 Experience toggle.

EC2

## Resources

You are using the following Amazon EC2 resources in the US East (Ohio) Region:

| | | | | | |
|---|---|---|---|---|---|
| Running instances | 1 | Elastic IPs | 0 | Dedicated Hosts | 0 |
| Snapshots | 0 | Volumes | 1 | Load balancers | 0 |
| Key pairs | 1 | Security groups | 7 | Placement groups | 0 |

ⓘ Easily size, configure, and deploy Microsoft SQL Server Always On availability groups on AWS using the AWS Launch Wizard for SQL Server. Learn more ✕

## Launch instance

To get started, launch an Amazon EC2 instance, which is a virtual server in the cloud.

**Launch instance** ▼

Note: Your instances will launch in the US East (Ohio) Region

## Scheduled events

## Service health

**Service Health Dashboard** ↗

Region
US East (Ohio)

Status
⊘ This service is operating normally

## Availability Zone status

Zone
us-east-2a (use2-az1)

Status
⊘ Availability Zone is operating normally

## Step 2: Choose an Instance Type

Amazon EC2 provides a wide selection of instance types optimized to fit different use cases. Instances are virtual servers that can run applications. They have varying combinations of CPU, memory, storage, and networking capacity, and give you the flexibility to choose the appropriate mix of resources for your applications. Learn more about instance types and how they can meet your computing needs.

Filter by: **p2** ▾   **Current generation** ▾   **Show/Hide Columns**

**Currently selected:** p2.xlarge (- ECUs, 4 vCPUs, 2.7 GHz, -, 61 GiB memory, EBS only)

| | Family ▾ | Type ▾ | vCPUs ⓘ ▾ | Memory (GiB) ▾ | Instance Storage (GB) ⓘ ▾ | EBS-Optimized Available ⓘ ▾ | Network Performance ⓘ ▾ | IPv6 Support ⓘ ▾ |
|---|---|---|---|---|---|---|---|---|
| ☑ | p2 | p2.xlarge | 4 | 61 | EBS only | Yes | High | Yes |
| ☐ | p2 | p2.8xlarge | 32 | 488 | EBS only | Yes | 10 Gigabit | Yes |
| ☐ | p2 | p2.16xlarge | 64 | 732 | EBS only | Yes | 25 Gigabit | Yes |

Cancel   **Previous**   **Review and Launch**   **Next: Configure Instance Details**

Choose t2.micro for CPU,
or if P instance type limits aren't set

**Carnegie
Mellon
University**

# Step 3: Configure Instance Details

Configure the instance to suit your requirements. You can launch multiple instances from the same AMI, request Spot instances to take advantage of the lower pricing, assign an access management role to the instance, and more.

| | | |
|---|---|---|
| **Number of instances** ⓘ | `1` | Launch into Auto Scaling Group ⓘ |
| **Purchasing option** ⓘ | ☐ Request Spot instances | |
| **Network** ⓘ | vpc-a688b4ce (default) ⬍ | C  Create new VPC |
| **Subnet** ⓘ | No preference (default subnet in any Availability Zone ⬍ | Create new subnet |
| **Auto-assign Public IP** ⓘ | Use subnet setting (Enable) ⬍ | |
| **Hostname type** ⓘ | Use subnet setting (IP name) ⬍ | |
| **DNS Hostname** ⓘ | ☑ Enable IP name IPv4 (A record) DNS requests | |
| | ☑ Enable resource-based IPv4 (A record) DNS requests | |
| | ☐ Enable resource-based IPv6 (AAAA record) DNS requests | |
| **Placement group** ⓘ | ☐ Add instance to placement group | |
| **Capacity Reservation** ⓘ | Open ⬍ | |
| **Domain join directory** ⓘ | No directory ⬍ | C  Create new directory |

Cancel   **Previous**   **Review and Launch**   **Next: Add Storage**

Just click Add Storage
Do NOT change any settings

**Carnegie Mellon University**

# Step 4: Add Storage

Your instance will be launched with the following storage device settings. You can attach additional EBS volumes and instance store volumes to your instance, or edit the settings of the root volume. You can also attach additional EBS volumes after launching an instance, but not instance store volumes. Learn more about storage options in Amazon EC2.

| Volume Type ⓘ | Device ⓘ | Snapshot ⓘ | Size (GiB) ⓘ | Volume Type ⓘ | IOPS ⓘ | Throughput (MB/s) ⓘ | Delete on Termination ⓘ | Encryption ⓘ |
|---|---|---|---|---|---|---|---|---|
| Root | /dev/sda1 | snap-0afe27d2ce58f1639 | 130 | General Purpose SSD (gp2) | 390 / 3000 | N/A | ☑ | Not Encrypte ▼ |

**Add New Volume**

Free tier eligible customers can get up to 30 GB of EBS General Purpose (SSD) or Magnetic storage. Learn more about free usage tier eligibility and usage restrictions.

▼ Shared file systems ⓘ

You currently don't have any file systems on this instance. Select "Add file system" button below to add a file system.

**Add file system**

Cancel    Previous    **Review and Launch**    Next: Add Tags

Carnegie
Mellon
University

# Step 7: Review Instance Launch

▼ AMI Details

**Deep Learning AMI (Ubuntu 18.04) Version 54.0 - ami-0476bba883df7cca6**

MXNet-1.8.0 & 1.7.0, TensorFlow-2.4.3, 2.3.4 & 1.15.5, PyTorch-1.7.1 & 1.8.1, Neuron, & others. NVIDIA CUDA, cuDNN, NCCL, Intel MKL-DNN, Docker, NVIDIA-Docker & EFA support. For fully managed experience, check: https://aws.amazon.com/sagemaker

Root Device Type: ebs    Virtualization type: hvm

▼ Instance Type                                                                        Edit instance type

| Instance Type | ECUs | vCPUs | Memory (GiB) | Instance Storage (GB) | EBS-Optimized Available | Network Performance |
|---|---|---|---|---|---|---|
| p2.xlarge | - | 4 | 61 | EBS only | Yes | High |

▼ Security Groups                                                                      Edit security groups

**Security group name**    launch-wizard-20
**Description**            launch-wizard-20 created 2022-02-08T19:31:08.619-05:00

| Type ⓘ | Protocol ⓘ | Port Range ⓘ | Source ⓘ | Description ⓘ |
|---|---|---|---|---|
| SSH | TCP | 22 | 0.0.0.0/0 | |

▶ Instance Details                                                                     Edit instance details

▶ Storage                                                                              Edit storage

▶ Tags                                                                                 Edit tags

Just click Launch
Do NOT change any settings

Cancel    Previous    **Launch**

**Carnegie Mellon University**

**1. Name your key pair**
**2. Download key pair**
**3. Launch**

Comments:
1. Make sure that you store key pair file safely and privately. You won't be able to log into an instance if its key pair is lost.
2. You may choose to reuse an existing key pair when you launch a second instance. Make sure you can still find the key pair.

Carnegie
Mellon
University

# Check your instance status

# Instructions for login

# Login via SSH

```
chmod 0400 <path_to_key_file>
```

```
ssh -i <path_to_key_file> <username>@<public_ip_from_EC2>
```

**Carnegie
Mellon
University**

```
=======================================================================
       __|  __|_  )
       _|  (     /   Deep Learning Base AMI (Ubuntu 16.04) Version 21.0
      ___|\___|___|
=======================================================================

Welcome to Ubuntu 16.04.6 LTS (GNU/Linux 4.4.0-1098-aws x86_64v)

Nvidia driver version: 418.87.01
CUDA versions available: cuda-10.0 cuda-10.1 cuda-8.0 cuda-9.0 cuda-9.2
Default CUDA version is 10.0
Libraries: cuDNN, NCCL, Intel MKL-DNN

AWS Deep Learning AMI Homepage: https://aws.amazon.com/machine-learning/amis/
Developer Guide and Release Notes: https://docs.aws.amazon.com/dlami/latest/devguide/what-is-dlami.html
Support: https://forums.aws.amazon.com/forum.jspa?forumID=263
For a fully managed experience, check out Amazon SageMaker at https://aws.amazon.com/sagemaker
When using INF1 type instances, please update regularly using the instructions at: https://github.com/aws/aws-neuron-sdk/tree/master/release-notes
=======================================================================

 * Documentation:  https://help.ubuntu.com
 * Management:      https://landscape.canonical.com
 * Support:         https://ubuntu.com/advantage

 * Multipass 1.0 is out! Get Ubuntu VMs on demand on your Linux, Windows or
   Mac. Supports cloud-init for fast, local, cloud devops simulation.

     https://multipass.run/

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

59 packages can be updated.
37 updates are security updates.



The programs included with the Ubuntu system are free software;
the exact distribution terms for each program are described in the
individual files in /usr/share/doc/*/copyright.

Ubuntu comes with ABSOLUTELY NO WARRANTY, to the extent permitted by
applicable law.

ubuntu@ip-172-31-25-48:~$ ls
Nvidia_Cloud_EULA.pdf  README  src  tools
```

# Configure Anaconda Environment

```
$ wget https://repo.continuum.io/archive/Anaconda3-2019.10-Linux-x86_64.sh
$ bash Anaconda3-2019.10-Linux-x86_64.sh
```

Agree to terms and conditions; yes to everything (including conda init); use default installation path.

```
$ source ~/.bashrc
$ conda create -n pytorch_p36 python=3.6
```
Create an environment called pytorch_p36, with the Python 3.6 interpreter.

```
$ conda activate pytorch_p36
$ conda install pytorch=1.3.0 torchvision cudatoolkit=10.1 -c pytorch
```
Now you should be able to use PyTorch on your instance.

Carnegie
Mellon
University

# Amazon Machine Images (AMI)

- A template of pre-installed system softwares and libraries

  - For deep learning: Linux OS, Nvidia GPU driver, CUDA, cuDNN, etc.

- When an instance is launched, a root storage is created with a copy of AMI.

- **Deep Learning AMI - we prefer <u>base AMI</u> w/ Ubuntu 18**

**Deep Learning AMI (Ubuntu 18.04) Version 54.0** - ami-0476bba883df7cca6

MXNet-1.8.0 & 1.7.0, TensorFlow-2.4.3, 2.3.4 & 1.15.5, PyTorch-1.7.1 & 1.8.1, Neuron, & others. NVIDIA CUDA, cuDNN, NCCL, Intel MKL-DNN, Docker, NVIDIA-Docker & EFA support. For fully managed experience, check: https://aws.amazon.com/sagemaker

Root device type: ebs     Virtualization type: hvm     ENA Enabled: Yes

Select

64-bit (x86)

Carnegie Mellon University

# Storage - Amazon Elastic Block Store (EBS)

- General Purpose SSD (gp2)
  - **$0.1 per GB-month**
  - 30 GB-month per month of free usage for the first 12 months
- Incurs charges as long as the EBS is **active**
  - "Active" means that it is allocated to your account, until it is deleted.
  - Even if it stores no data
  - Even if it is not attached to your instance

# Root Storage vs. Additional Storage

- **Root storage**
- Automatically created when an instance is launched
- Stores OS, libraries
- Automatically deleted when an instance is terminated
- Instance cannot boot when root storage is detached
- 50GB for the deep learning AMI

- **Additional storage**
- Need to create manually
- **Stores your code, model checkpoints & dataset**
- Need to delete manually
- Can be detached from or attached to different instances
- 20 GB for the 1st assignment
- Must be in the same availability zone as the instance

**Carnegie Mellon University**

# Python/Pytorch Development on AWS

- Deep Learning Base AMI provides almost everything except Python

- Anaconda enables easy Python library management

- Create an environment with Python 3.6

- Alternative: virtualenv & pip

  - we don't provide instruction

- Use IDE like VSCode to write code w/ a nice UI (some conda integration)

**Carnegie Mellon University**

# Create additional storage (Optional)



Important: check the availability zone of your instance; your additional storage must be in the same zone.

# Create Volume

| | |
|---|---|
| **Volume Type** | General Purpose SSD (gp2) ▼ ❶ |
| **Size (GiB)** | 20   (Min: 1 GiB, Max: 16384 GiB) ❶ |
| **IOPS** | 100 / 3000   (Baseline of 3 IOPS per GiB with a minimum of 100 IOPS, burstable to 3000 IOPS) ❶ |
| **Availability Zone\*** | us-east-2b ▼ ❶ |
| **Throughput (MB/s)** | Not applicable ❶ |
| **Snapshot ID** | Select a snapshot ▼ ↻ ❶ |
| **Encryption** | ☐ Encrypt this volume |

# Initialize & mount additional storage

- `lsblk`
  - show the attached devices (find the name of the attached EBS volume)
- `sudo file -s /dev/<device_name>`
  - show whether the device has a filesystem
- `sudo mkfs -t ext4 /dev/<device_name>`
  - e.g. `sudo mkfs -t ext4 /dev/xvdf`
  - create the ext4 file system on the device (do this only if the previous command returned "data")
- `sudo mkdir <mount_point>`
  - e.g. `sudo mkdir /mnt/data`
  - make a folder for mounting the volume
- `sudo mount <device_name> <mount_point>`
  - e.g. `sudo mount /dev/xvdf /mnt/data`
  - mount the volume
- `sudo chmod 777 -R <mount_point>`
  - add permissions to access the folder

```
(base) ubuntu@ip-172-31-25-48:~$ lsblk
NAME      MAJ:MIN RM   SIZE RO TYPE MOUNTPOINT
xvda      202:0    0    50G  0 disk
└─xvda1   202:1    0    50G  0 part /
 xvdf     202:80   0    20G  0 disk
```

/dev/sdf in previous slide renamed to /dev/xvdf

**Carnegie Mellon University**

# Test additional storage

- `cd <mount point>`
  - go to the mount point of the additional storage, e.g., /mnt/data

- `wget` http://host.robots.ox.ac.uk/pascal/VOC/voc2007/VOCtest_06-Nov-2007.tar
  - download the dataset required for the 1st assignment
  - also store your code and model checkpoints on additional EBS storage

- `df -h`
  - show the usage on the disks
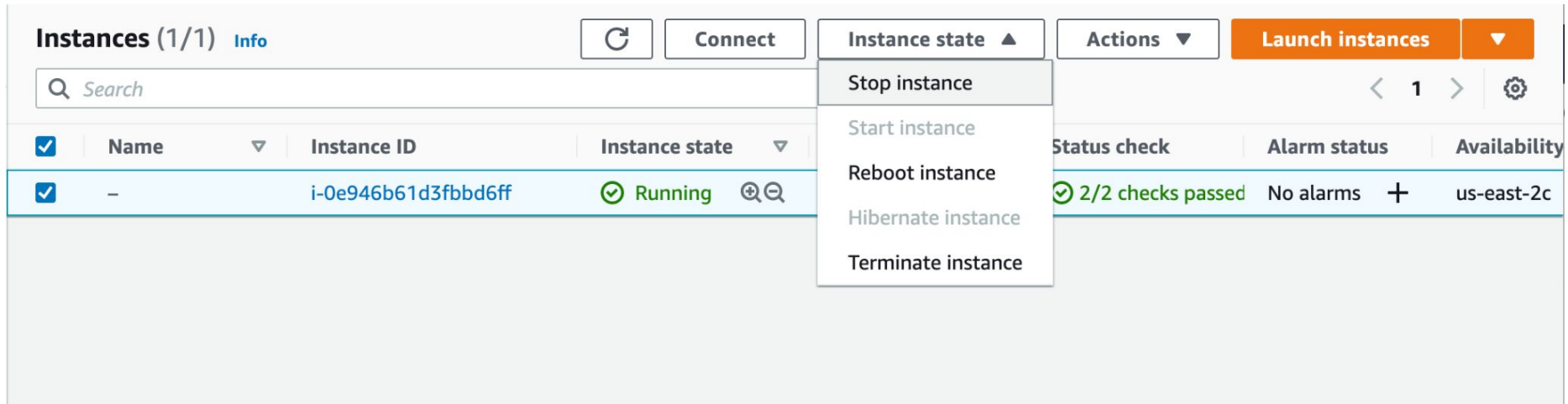  - monitor root storage and additional storage

```
(base) ubuntu@ip-172-31-25-48:/mnt/data$ df -h
Filesystem      Size  Used Avail Use% Mounted on
udev             30G     0   30G   0% /dev
tmpfs           6.0G  8.8M  6.0G   1% /run
/dev/xvda1       49G   36G   13G  74% /
tmpfs            30G     0   30G   0% /dev/shm
tmpfs           5.0M     0  5.0M   0% /run/lock
tmpfs            30G     0   30G   0% /sys/fs/cgroup
/dev/loop0       88M   88M     0 100% /snap/core/5742
/dev/loop1       17M   17M     0 100% /snap/amazon-ssm-agent/784
/dev/loop2       18M   18M     0 100% /snap/amazon-ssm-agent/1480
/dev/loop3       90M   90M     0 100% /snap/core/8268
tmpfs           6.0G     0  6.0G   0% /run/user/1000
/dev/xvdf        20G  475M   19G   3% /mnt/data
```

root storage

additional storage

# Final Word: Stop your instance! For now..

# Google Cloud Platform

Same Idea

- structure
- billing strategy
- increase quotas
- create instance (pay attention to many details).

**Carnegie Mellon University**

# GCP Account Setup

- **Use personal email address**

  - Do not use andrew account (Will likely face organization error)

- https://cloud.google.com/

- Redeem credits at: https://console.cloud.google.com/education

# Google Cloud Platform

- Redeem coupon to your personal account

  -

# Google Cloud Platform (Basic steps)

- Upgrade to paid account (not free account)
- Enable "Compute Engine API"
- Create a project
- **Increase your quota for "gpus_all_region"**
- Create "VM Instance"
- Select the zone you are residing in

# Increase Quotas (otherwise you cannot use all types of GPUs)

# Increase Quotas (otherwise you cannot use all types of GPUs)

Quotas for project "16726"    ✏ EDIT QUOTAS

| Near the limit | Low usage | All quotas |
|---|---|---|
| 0 | 5,375 | 5,568 |
| View quotas | View quotas | |

≡ Filter   Metric : compute.googleapis.com/gpus_all_regions ⊗   Enter property name or value   ✕  ❓  ▥

| ☑ | Service | Quota | Dimensions (e.g. location) | Limit | Current usage percentage ↓ | 7 day peak usage percentage |
|---|---|---|---|---|---|---|
| ☑ | Compute Engine API | GPUs (all regions) | | 0 | 0% | 0% |

# You must have a paid account before increasing Quotas (otherwise you cannot use all types of GPUs)



Filter   Metric : compute.googleapis.com/gpus_all_regions ⊗   Enter property name or value

| | Service | Quota | Dimensions (e.g. location) | Limit | Current usage percentage |
|---|---|---|---|---|---|
| ☐ | Compute | GPUs (all | | 0 | |

Please be aware that free trial accounts for Google Cloud Platform have limited quota during their trial period. In order to increase your quota, please upgrade to a paid account by clicking "Upgrade my account" from the top of any page once logged in to Google Cloud Console.

**Carnegie Mellon University**

Now you can proceed to:
- Compute Engine -> VM instances -> Create Instance
- Only Tesla P4 is available (as of 2022.02)
- You can use your preferred GPUs (we won't be making recommendations).

Carnegie Mellon University

- Set SSH Keys: https://cloud.google.com/compute/docs/instances/adding-removing-ssh-keys#create sshkeys
- add disk: youtube
- Pricing: https://cloud.google.com/compute/gpus-pricing

# Final words

- **Stop an instance if you are not working on it!**

- $150 AWS and GCP for all 3 assignments and course project

- Use AWS/GCP wisely & responsibly

- Enjoy the deep learning alchemy!

Carnegie
Mellon
University