

GPU as a first class feature

A proposal and an appeal for help

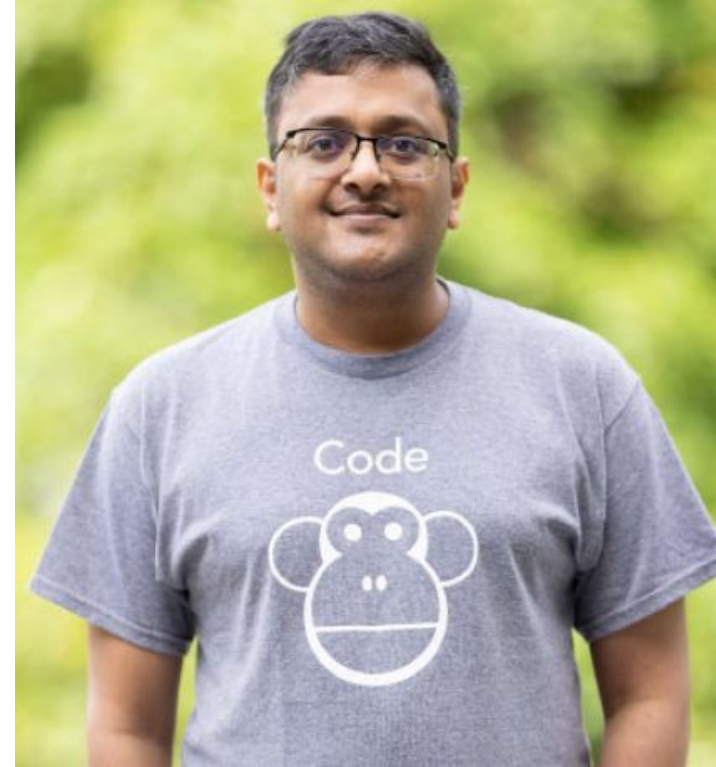
Vishesh Jindal
Software Engineer @ ShapeBlue



November 20 - 22, 2024 | Madrid, Spain

About me

- Contributing to ACS since 2023
- Committer @ Apache CloudStack since 2024
- Software Engineer @ ShapeBlue
- Worked as a Backend & Infrastructure engineer in the past
- Experience managing and administering kubernetes, AWS, PostgreSQL in production
- Passionate about cloud infrastructure, security, and AI

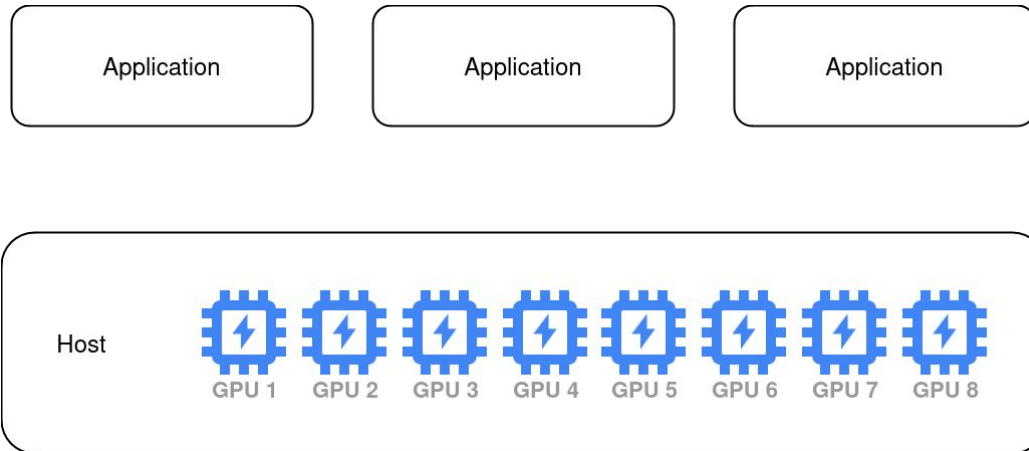


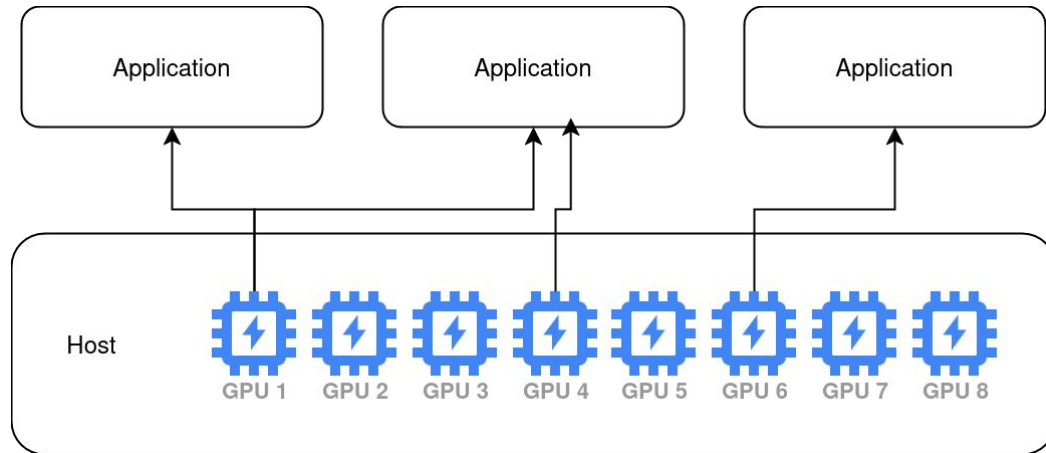
What's in it for you?

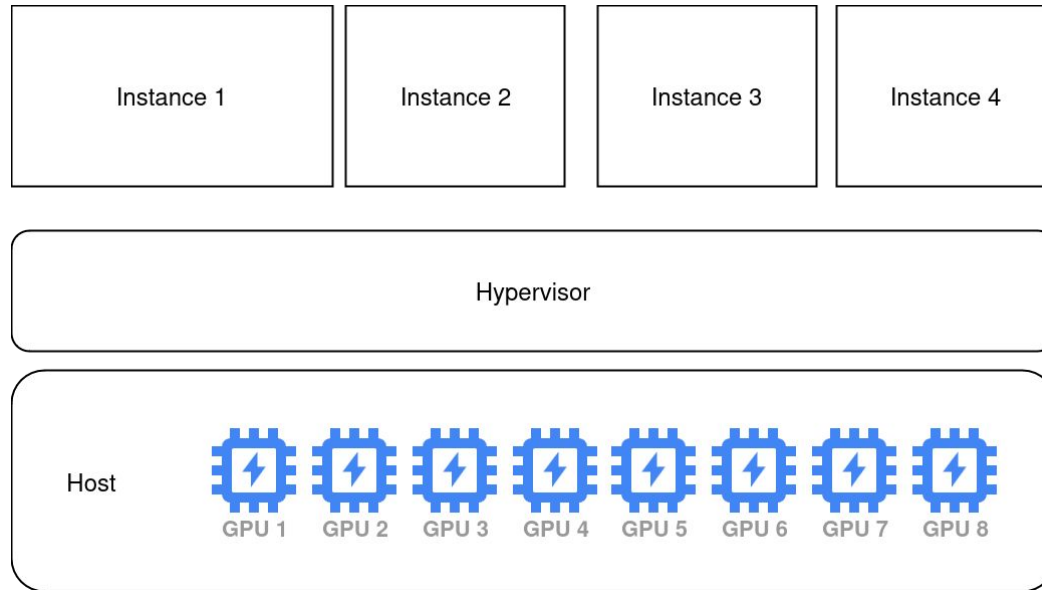


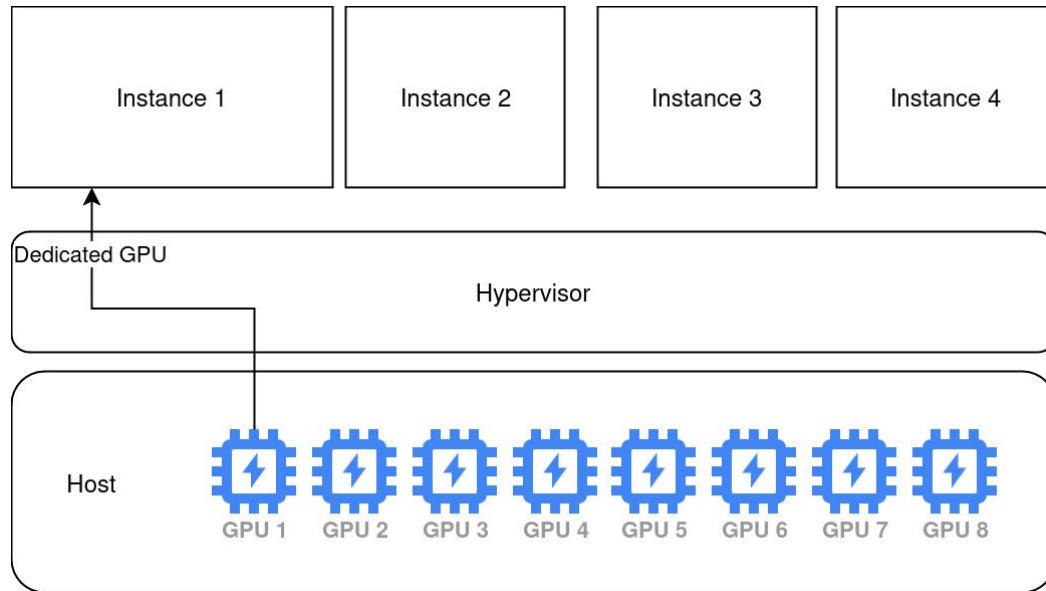
How are GPUs used in the cloud?

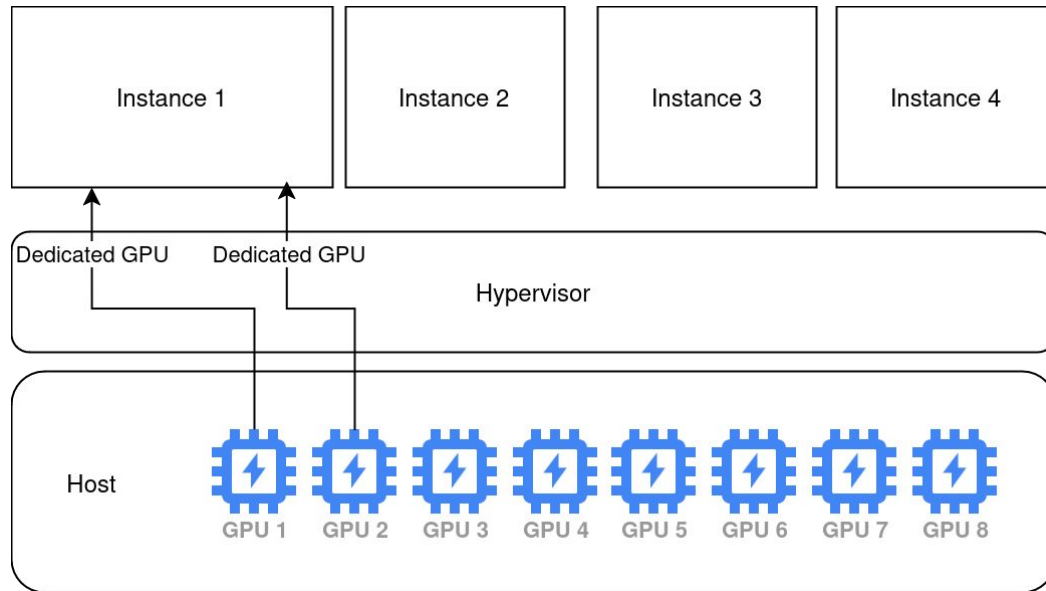


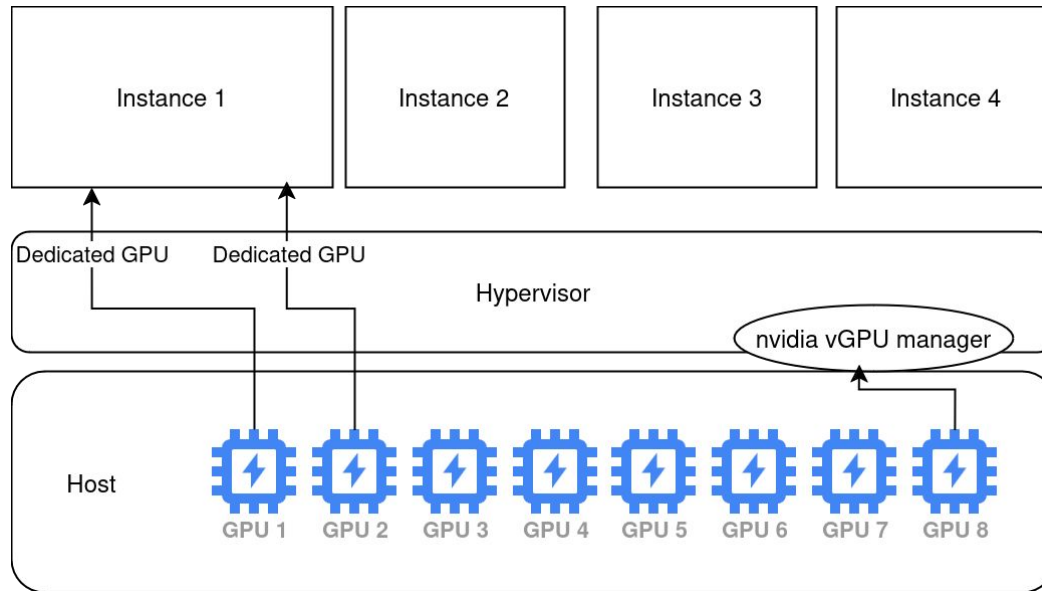


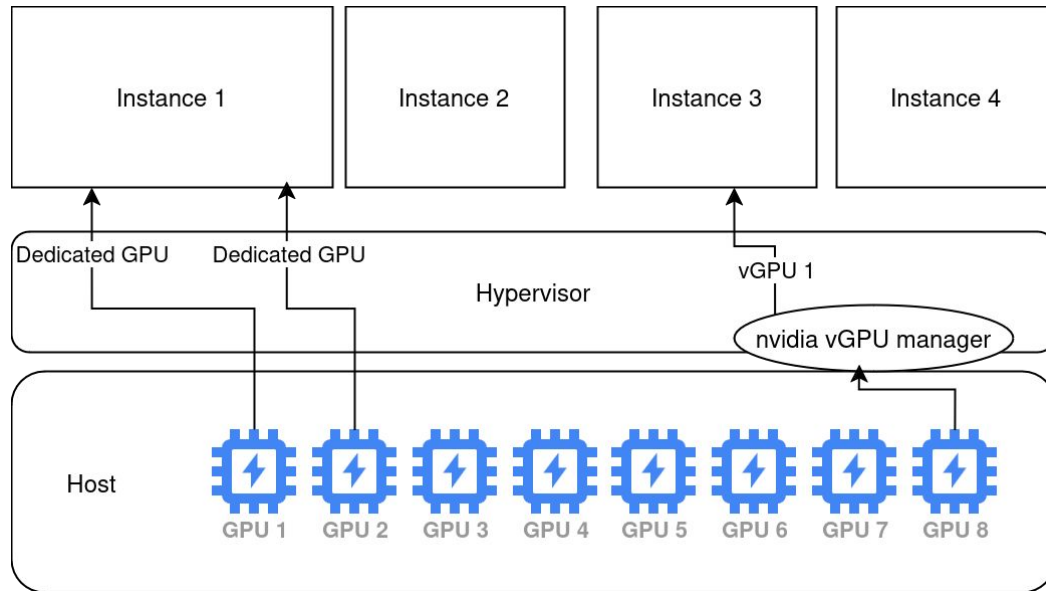


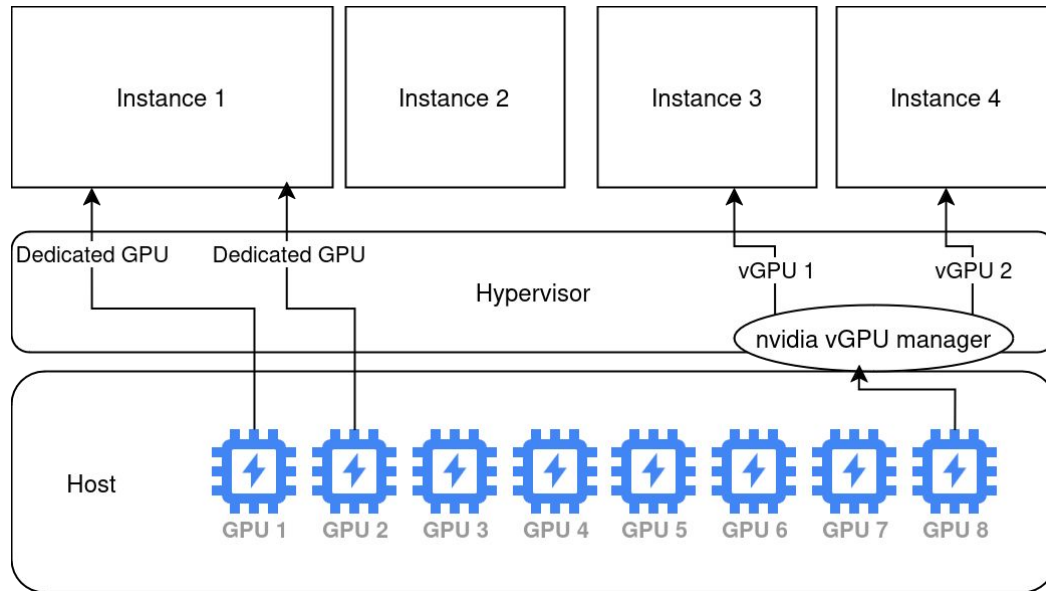


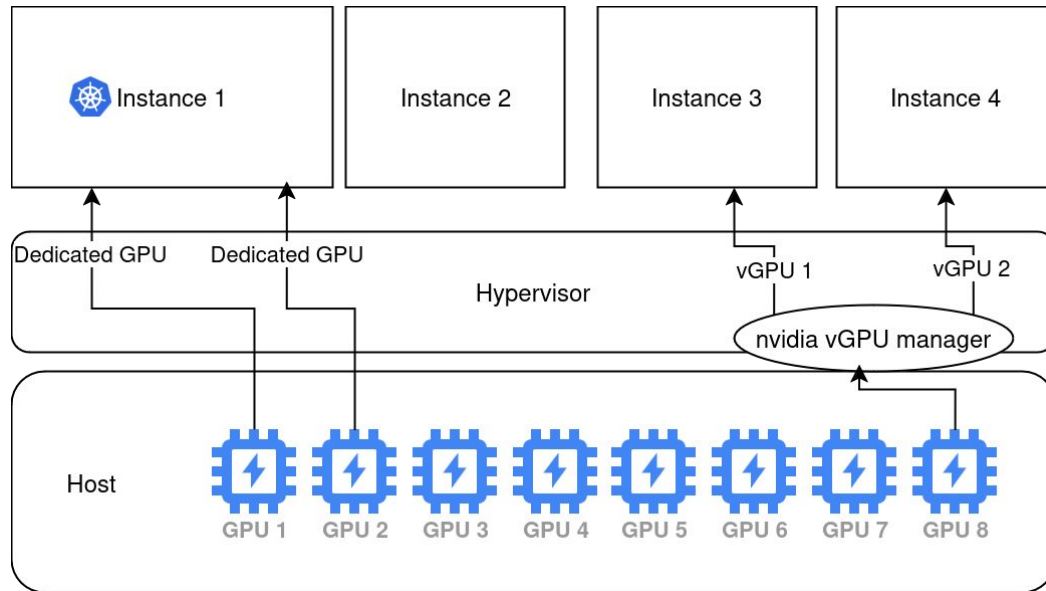


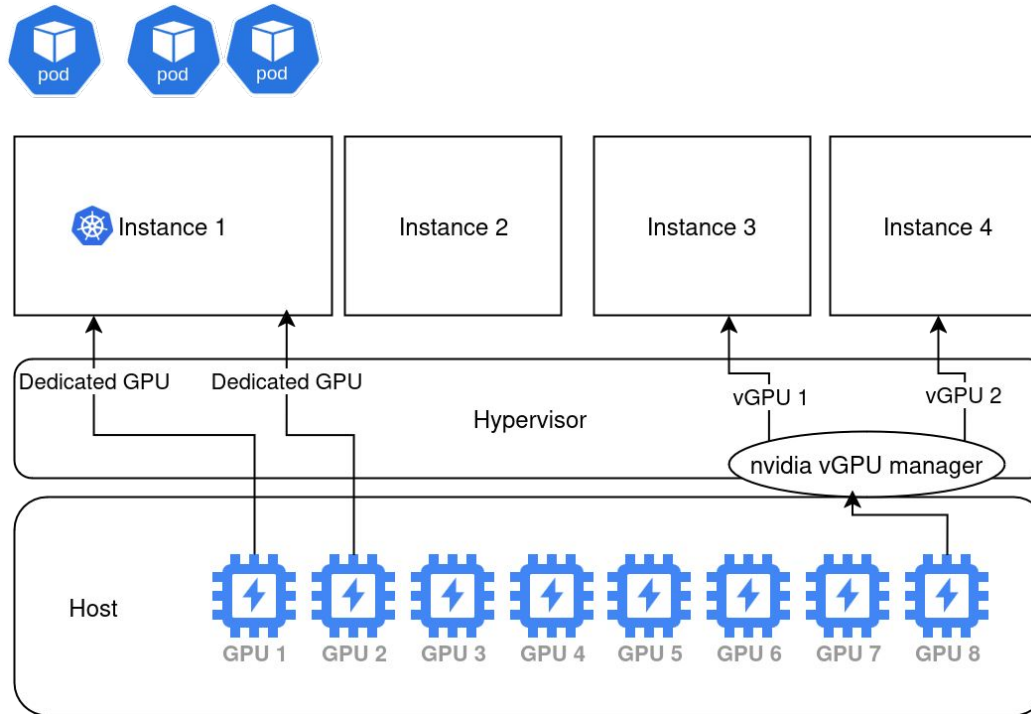


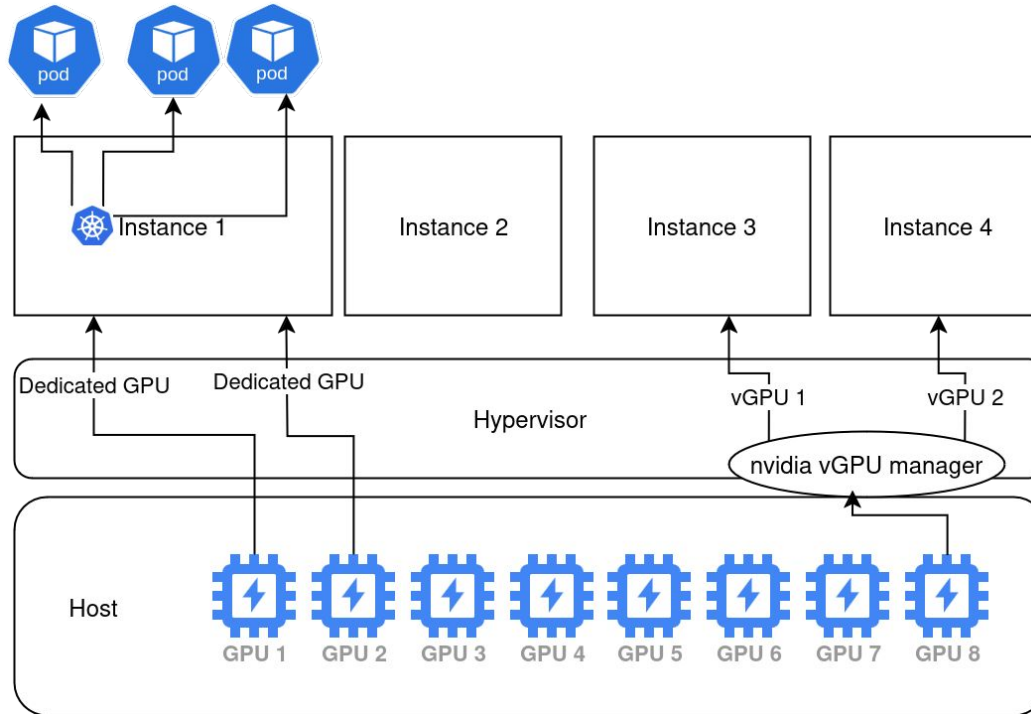












Current State of GPU support



Add compute offering ? ×

* Name ?

Name

Description ?

Description

Compute offering type

Fixed offering Custom constrained **Custom unconstrained**

Host tags ? Network rate (Mb/s) ?

Offer HA ? Dynamic scaling enabled ?

CPU cap ? Volatile ?

Deployment planner ?

GPU

vGPU type

Public



Workarounds with KVM



Configuration

```
[root@helium ~]# lspci -nn|grep NVIDIA
61:00.0 3D controller [0302]: NVIDIA Corporation GV100GL [Tesla V100 SXM2 32GB
62:00.0 3D controller [0302]: NVIDIA Corporation GV100GL [Tesla V100 SXM2 32GB
89:00.0 3D controller [0302]: NVIDIA Corporation GV100GL [Tesla V100 SXM2 32GB
8a:00.0 3D controller [0302]: NVIDIA Corporation GV100GL [Tesla V100 SXM2 32GB
[root@helium ~]#
```

```
[root@helium ~]# cat /etc/modprobe.d/vfio.conf
options vfio-pci ids=10de:1db5,10de:1db5,10de:1db5,10de:1db5
```

Source: <https://lab.piszki.pl/cloudstack-kvm-and-running-vm-with-vgpu/>



Additional libvirt XML for device passthrough

```

<devices>
  <hostdev mode='subsystem' type='pci' managed='yes'>
    <driver name='vfio' />
    <source>
      <address domain='0x0000' bus='0x61' slot='0x00' function='0x0' />
    </source>
    <alias name='nvidia0' />
    <address type='pci' domain='0x0000' bus='0x00' slot='0x00' function='0x0' />
  </hostdev>
  <hostdev mode='subsystem' type='pci' managed='yes'>
    <driver name='vfio' />
    <source>
      <address domain='0x0000' bus='0x62' slot='0x00' function='0x0' />
    </source>
    <alias name='nvidia1' />
    <address type='pci' domain='0x0000' bus='0x00' slot='0x00' function='0x1' />
  </hostdev>
</devices>
<devices>
  <hostdev mode='subsystem' type='pci' managed='yes'>
    <driver name='vfio' />
    <source>
      <address domain='0x0000' bus='0x89' slot='0x00' function='0x0' />
    </source>
    <alias name='nvidia2' />
    <address type='pci' domain='0x0000' bus='0x00' slot='0x00' function='0x2' />
  </hostdev>
  <hostdev mode='subsystem' type='pci' managed='yes'>
    <driver name='vfio' />
    <source>
      <address domain='0x0000' bus='0x8a' slot='0x00' function='0x0' />
    </source>
    <alias name='nvidia3' />
    <address type='pci' domain='0x0000' bus='0x00' slot='0x00' function='0x3' />
  </hostdev>
</devices>

```

Source:
<https://lab.piszki.pl/cloudstack-kvm-and-running-vm-with-vgpu/>



Using libvirt hooks

```

7 class GpuDeviceAdder {
8     String gpuXml = """
9     <devices>
10        <hostdev mode='subsystem' type='pci' managed='yes'>
11            <driver name='vfio'/>
12            <source>
13                <address domain='0x0000' bus='0x61' slot='0x00' function='0x0'/>
14            </source>
15            <alias name='nvidia0'/>
16            <address type='pci' domain='0x0000' bus='0x00' slot='0x00' function='0x0'/>
17        </hostdev>
18    </devices>
19    """
20
21    String transform(Object logger, String xml) {
22        def vmDef = new XmlParser().parseText(xml)
23
24        // Parse GPU XML
25        def gpuDevices = new XmlParser().parseText(gpuXml)
26
27        // Append GPU devices to the VM definition
28        gpuDevices.hostdev.each { gpuDevice ->
29            vmDef.devices[0].append(gpuDevice)
30        }
31
32        // Return updated XML definition
33        return XmlUtil.serialize(vmDef)
34    }
35

```

Source:
<https://gist.github.com/rajujith/f3b3854ed77f2ca68dc4fb5e3ee260c4>



Is that good enough?



What's missing?



User

"I need to deploy instances with multiple GPUs to efficiently run demanding AI and ML workloads."

Cloud Service Provider (CSP)

"I want to offer instances with multiple GPUs to my users, complete with usage tracking and customizable limits, across different hypervisors."



Capabilities

- ❑ Discovery and Inventory Management
- ❑ Groupings and Offerings
- ❑ Allocation and Assignment
- ❑ Usage and Limits

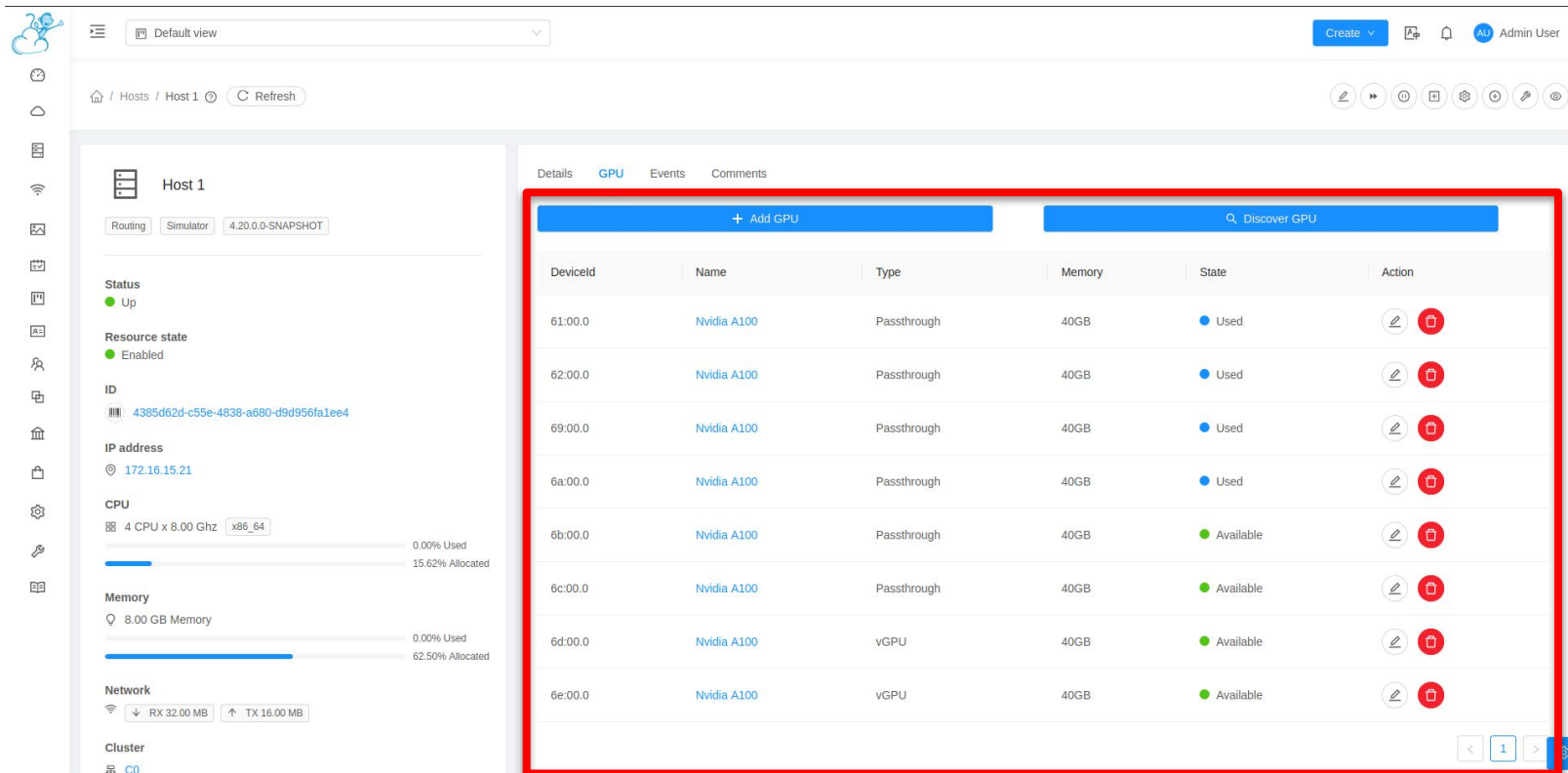


Capabilities

- **Discovery and Inventory Management**
 - ❑ Groupings and Offerings
 - ❑ Allocation and Assignment
 - ❑ Usage and Limits



GPUs on a Host



The screenshot shows the CloudStack management console for a host named 'Host 1'. The interface is divided into two main sections. On the left, the host's general information is displayed, including its status (Up), resource state (Enabled), ID (4385d62d-c55e-4838-a680-d9d956fa1ee4), IP address (172.16.15.21), and CPU/Memory usage. On the right, the 'GPU' tab is active, showing a table of GPUs attached to the host. The table has columns for DeviceId, Name, Type, Memory, State, and Action. A red box highlights the GPU table area. At the top of the GPU table, there are buttons for '+ Add GPU' and 'Discover GPU'. The table lists 10 GPUs, with 4 in a 'Used' state and 6 in an 'Available' state. The 'Action' column for each GPU contains edit and delete icons.

Host 1

Routing Simulator 4.20.0.0-SNAPSHOT

Status
● Up

Resource state
● Enabled

ID
4385d62d-c55e-4838-a680-d9d956fa1ee4

IP address
172.16.15.21

CPU
4 CPU x 8.00 Ghz x86_64
0.00% Used
15.62% Allocated

















Memory
8.00 GB Memory
0.00% Used
62.50% Allocated

Network
RX 32.00 MB TX 16.00 MB

Cluster
CO

Details GPU Events Comments

+ Add GPU Discover GPU

DeviceId	Name	Type	Memory	State	Action
61:00:0	Nvidia A100	Passthrough	40GB	Used	 
62:00:0	Nvidia A100	Passthrough	40GB	Used	 
69:00:0	Nvidia A100	Passthrough	40GB	Used	 
6a:00:0	Nvidia A100	Passthrough	40GB	Used	 
6b:00:0	Nvidia A100	Passthrough	40GB	Available	 
6c:00:0	Nvidia A100	Passthrough	40GB	Available	 
6d:00:0	Nvidia A100	vGPU	40GB	Available	 
6e:00:0	Nvidia A100	vGPU	40GB	Available	 

GPUs on a Host



Default view Create

Home / Hosts Refresh All Metrics Add host + Search

Name	State	Resource state	IP Address	Hypervisor	Instances	Power state	Version	CPU cores	CPU total	CPU used	CPU allocated	Memory total	Memory used	Memory allocated	GPU devices	GPU devices used	GPU memory total	GPU memory used	Network read
Host 1	Up	Enabled	172.16.15.21	Simulator	3 / 3	Disabled	4.21.0.0	4	32.00 Ghz	0.00 Ghz	5.00 Ghz	8.00 GB	0.00 GB	5.00 GB	8	4	320 GB	160 GB	0.03 G
Host 2	Up	Enabled	172.16.15.8	Simulator	0 / 0	Disabled	4.21.0.0	4	32.00 Ghz	0.00 Ghz	0.50 Ghz	8.00 GB	0.00 GB	0.50 GB	8	0	320 GB	0 GB	0.03 G
Host 3	Up	Enabled	172.16.15.23	Simulator	1 / 1	Disabled	4.21.0.0	4	32.00 Ghz	0.00 Ghz	3.00 Ghz	8.00 GB	0.00 GB	3.00 GB	0	0	0 GB	0 GB	0.03 G
Host 4	Up	Enabled	172.16.15.3	Simulator	0 / 0	Disabled	4.21.0.0	4	32.00 Ghz	0.00 Ghz	0.50 Ghz	8.00 GB	0.00 GB	1.00 GB	0	0	0 GB	0 GB	0.03 G

Showing 1-4 of 4 items 1 / 20 / page



Capabilities

- ✓ Discovery and Inventory Management
- **Groupings and Offerings**
- ❑ Allocation and Assignment
- ❑ Usage and Limits



GPU Offerings

Home / GPU Offerings Mine

Search

Name	Description	State	Total memory	GPU Model	Count	Domain	Zone	Order
8x NVIDIA Tesla V100 32GB	NVIDIA Tesla V100	Active	256 GB	NVIDIA Tesla V100 32GB	8			<input type="button" value="edit"/> <input type="button" value="delete"/> <input type="button" value="up"/> <input type="button" value="down"/>
4x NVIDIA Tesla V100 32GB	NVIDIA Tesla V100	Active	128 GB	NVIDIA Tesla V100 32GB	4			<input type="button" value="edit"/> <input type="button" value="delete"/> <input type="button" value="up"/> <input type="button" value="down"/>
2x NVIDIA Tesla V100 32GB	NVIDIA Tesla V100	Active	64 GB	NVIDIA Tesla V100 32GB	2			<input type="button" value="edit"/> <input type="button" value="delete"/> <input type="button" value="up"/> <input type="button" value="down"/>
1x NVIDIA Tesla V100 32GB	NVIDIA Tesla V100	Active	32 GB	NVIDIA Tesla V100 32GB	1			<input type="button" value="edit"/> <input type="button" value="delete"/> <input type="button" value="up"/> <input type="button" value="down"/>

Showing 1-1 of 1 items < 1 > 20 / page

Licensed under the Apache License, Version 2.0.



Add GPU Offering

Add GPU offering ? ×

* Name i

8x NVIDIA A100 80GB

Description i

8x NVIDIA A100 80GB

* GPU Model

NVIDIA A100 80GB ▾

vGPU Disabled

* Number of GPUs i

8

Cancel OK



Capabilities

- ✓ Discovery and Inventory Management
- ✓ Groupings and Offerings
- Allocation and Assignment
- ☐ Usage and Limits



Deploying an Instance

Override root disk offering

Total 4 items < 1 > 10 / page v

5 Data disk

Search

Disk offering	Disk size (in GB)	Min IOPS / Max IOPS
<input checked="" type="radio"/> No thanks	-	-
<input type="radio"/> Small	5 GB	-
<input type="radio"/> Medium	20 GB	-
<input type="radio"/> Large	100 GB	-
<input type="radio"/> Custom	Custom disk size	-

6 GPU Offering

Enable GPU

Name	Total Memory	GPU Model
<input type="radio"/> 8x NVIDIA Tesla V100 32GB	256 GB	NVIDIA Tesla V100 32GB
<input type="radio"/> 4x NVIDIA Tesla V100 32GB	128 GB	NVIDIA Tesla V100 32GB
<input type="radio"/> 2x NVIDIA Tesla V100 32GB	64 GB	NVIDIA Tesla V100 32GB
<input type="radio"/> 1x NVIDIA Tesla V100 32GB	32 GB	NVIDIA Tesla V100 32GB

7 Networks

Simulator

OS type

- CentOS 5.6 (64-bit)

CPU

- 1 CPU x 1.00 Ghz

Memory

- 1024 MB memory

Disk size (in GB)

- 5 GB (Root)

Networks

- Test isolated (Default)

Template

- CentOS 5.6 (64-bit) no GUI (Simulator)

Compute offering

- Test

Zone

- Sandbox-simulator



Service Offering with GPU

Add compute offering ⓘ

* Name ⓘ

Description ⓘ

Compute offering type

Fixed offering Custom constrained Custom unconstrained

Host tags ⓘ

Network rate (Mb/s) ⓘ

Offer HA ⓘ

CPU cap ⓘ

Deployment planner ⓘ

Dynamic scaling enabled ⓘ

Volatile ⓘ

Attach GPU Offering ⓘ

Name	Total Memory	GPU Model
<input type="radio"/> 8x NVIDIA Tesla V100 32GB	256 GB	NVIDIA Tesla V100 32GB
<input type="radio"/> 4x NVIDIA Tesla V100 32GB	128 GB	NVIDIA Tesla V100 32GB
<input type="radio"/> 2x NVIDIA Tesla V100 32GB	64 GB	NVIDIA Tesla V100 32GB
<input type="radio"/> 1x NVIDIA Tesla V100 32GB	32 GB	NVIDIA Tesla V100 32GB



GPU details for an Instance

Instances / QA-b739940f-46fb-4aff-8871-7f4ae8451a8e Refresh

QA-b739940f-46fb-4aff-8871-7f4ae8451a8e

1-2-51-QA Simulator

Status
Running

ID
b739940f-46fb-4aff-8871-7f4ae8451a8e

OS type
CentOS 5.6 (64-bit)

IP address
172.31.0.171

CPU
1 CPU x 1.00 Ghz

Memory
1024 MB memory

Network
1 NIC(s)
eth0 172.31.0.171 Default
Test

Template
CentOS 5.6 (64-bit) no GUI (Simulator)

Compute offering
Test

Host

Details	Name	GPU Offering	Type	Memory
Metrics	NVIDIA Tesla V100	4x NVIDIA Tesla V100 32GB	Passthrough	32GB
GPU devices	NVIDIA Tesla V100	4x NVIDIA Tesla V100 32GB	Passthrough	32GB
Volumes	NVIDIA Tesla V100	4x NVIDIA Tesla V100 32GB	Passthrough	32GB
NICs	NVIDIA Tesla V100	4x NVIDIA Tesla V100 32GB	Passthrough	32GB
Instance Snapshots	< 1 >			
Backups				
Security groups				
Schedules				
Settings				
Events				
Comments				



Capabilities

- ✓ Discovery and Inventory Management
- ✓ Groupings and Offerings
- ✓ Allocation and Assignment
- Usage and Limits



GPU Usage

Home / Usage Refresh

Server: usage-server/10.0.35.234 | Last heartbeat: 28 Oct 2024 11:10:35 (a few seconds ago) | Last successful job: 27 Oct 2024 23:59:59 (11 hours ago)

Domain: Account: **GPU Usage** 41a0-a4fa-f7b80993a960 2024-10-27 → 2024-10-27 [Show usage records](#) [Download CSV](#) [Clear](#)

Fetch usage records for child domains

Account	Domain	Usage type	Resource ID	Start date and time	End date and time	Raw usage (in hours)	Description	View
admin	ROOT	GPU Usage	6e1fdf4f-31f1-41a0-a4fa-f7b80993a960	27 Oct 2024 00:00:00	27 Oct 2024 23:59:59	24	GPU usage for 8x Nvidia Tesla V100 32GB (6e1fdf4f-31f1-41a0-a4fa-f7b80993a960) with 8 Nvidia Tesla V100 32GB card(s).	<input type="text"/>

Showing 1-1 of 1 items | / page



GPU Limits

Domain Details **Limits** Configure limits Settings Events Comments

Instance limits (Unlimited Available)

Used / Limit : 3 / Unlimited

CPU limits (Unlimited Available)

Used / Limit : 6 / Unlimited

Memory limits (MiB) (Unlimited Available)

Used / Limit : 6656 / Unlimited

GPU limits (Unlimited Available)

Used / Limit : 6 / Unlimited

GPU Memory limits (MiB) (Unlimited Available)

Used / Limit : 66560 / Unlimited

Primary storage limits (GiB) (Unlimited Available)

Used / Limit : 108 / Unlimited



Capabilities

- ✓ Discovery and Inventory Management
- ✓ Groupings and Offerings
- ✓ Allocation and Assignment
- ✓ Usage and Limits



This is just a proposal



Why hasn't the community built this yet?



Where do we go from here?



Questions?



Thank you!

#CSCollab24
@CloudStack

