



# ELKH Cloud vGPU challenges and implementation steps

**Ádám Pintér**

Wigner Adatközpont - IT leader



**ELKH** | Eötvös Loránd  
Research Network

# About ELKH Cloud project



- Wigner RCP site & SZTAKI site
- OpenStack cloud solution for
  - research institutes
  - universities
  - industrial research applications
- Supported by Eötvös Loránd Research Network Secretariat
- Officially from 15th of February, 2022



# Available new capacity Wigner Datacenter



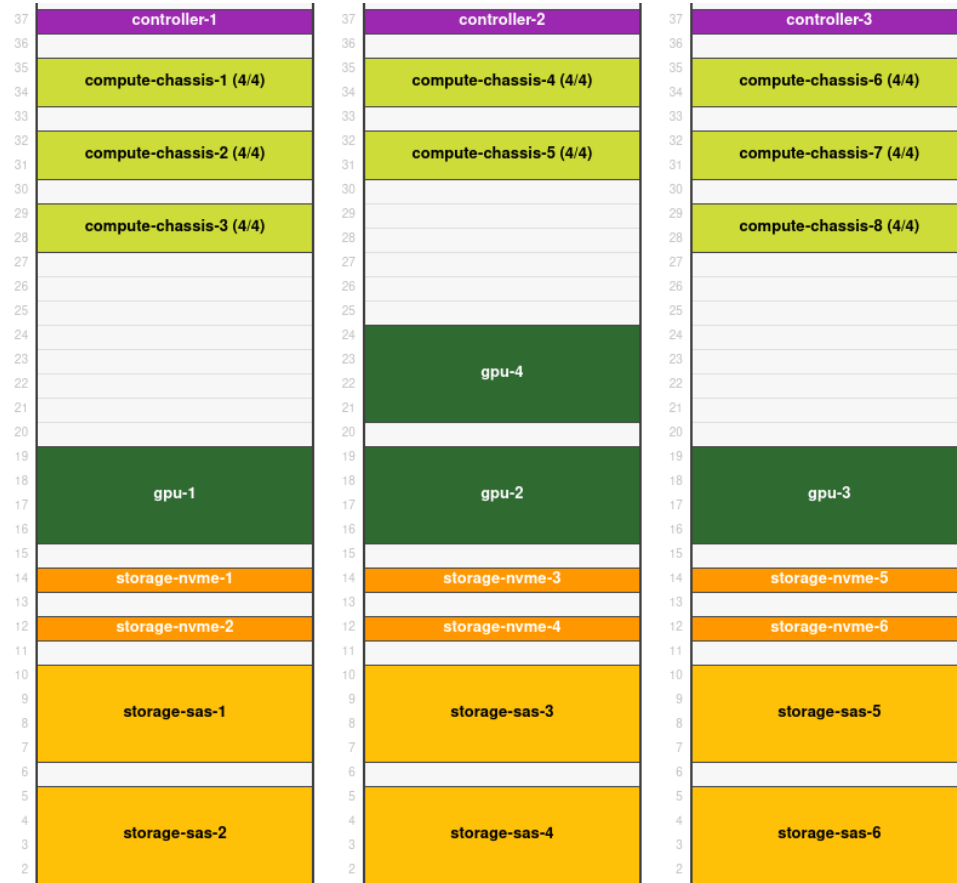
capacity	SZTAKI	Wigner Datacenter	ELKH Cloud full
compute server	10	40	50
GPU server	0	4	4
HDD storage server	6	6	12
SSD storage server	6	6	12
controller server	3	3	6
vCPU	1040	5888	6928
memory (RAM, TB)	7,68	24,576	32,256
HDD (replicated, TB)	384	864	1248
SSD (replicated, TB)	184,32	153,6	337,92

# Available new capacity Wigner Datacenter



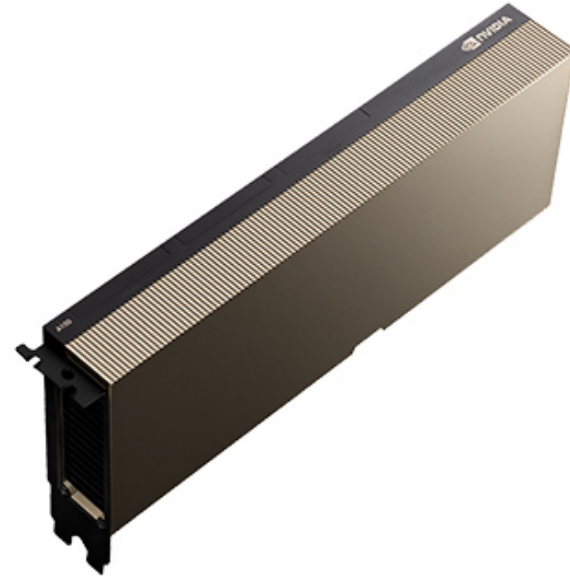
capacity	SZTAKI	Wigner Datacenter	ELKH Cloud full
internal network (Gbit/s)	100	100	100
GPU cards	40	32	72
GPU memory (RAM, GB)	1280	1280	2560
GPU double TFLOPS	312	310,4	622,4
GPU single TFLOPS	628	624	1252
GPU FP 16 tensor TFLOPS	5000	9984	14984

# Rack concept Wigner Datacenter



2022: extra  
8 compute nodes

- Nvidia A100 GPU
  - 6 912 CUDA cores
  - PCIe 4
  - 40 GB GPU RAM
  - 312 FP16 tensor TFLOPS
  - 19.5 single TFLOPS
  - 9.7 double TFLOPS



# Specialities

## Wigner Datacenter



- AMD CPU type (EPYC 7002)
- SSD storage type available in OpenStack
- vGPU (partitioning physical GPU cards)
- RAID 1 system disk
- General compute hypervisor (128 vCPU + 1/2 TB RAM)
- GPU hypervisor (192 vCPU + 1 TB RAM)
- Storage solution: Ceph Pacific version
  - HDD pool
  - SSD pool
  - rack based replication

# Specialities – tape library

## Wigner Datacenter



- Automated hierarchical tape library system
  - With robotic arm
    - Modular
      - 4 tape drives
      - 200 raw tapes
    - For cloud projects - unique requests
    - Long term storage
    - Store data on different hardware





# OpenStack flavor list



- Provide good frame for research applications
- Similar concept

géptípus	vCPU	memória (GB)	vGPU memória (GB) SZTAKI / Wigner	felhasználás
m2.tiny	1	1		általános célú
m2.small	1	2		általános célú
m2.medium	2	4		általános célú
m2.large	4	8		általános célú
m2.xlarge	8	16		általános célú
m2.2xlarge	16	32		általános célú
m2.4xlarge	32	64		általános célú
r2.medium	2	8		memória-optimalizált
r2.large	4	16		memória-optimalizált
r2.xlarge	8	32		memória-optimalizált
r2.2xlarge	16	64		memória-optimalizált
g2.medium	2	8	4 / 5	GPU-gyorsított
g2.large	4	16	8 / 10	GPU-gyorsított
g2.xlarge	8	32	16 / 20	GPU-gyorsított
g2.2xlarge	16	64	32 / 40	GPU-gyorsított

- SZTAKI and Wigner RCP similar concept
- OpenStack quota system
- Basic user support (ticketing, helpdesk mail)
- Elevated level user support – IT engineer consulting
- Elevated level user support – complex problems (Data Science for example)
- Elevated level user support – daily operation (ticketing, helpdesk mail, phone, meetings)

- OpenStack API for better automation
- VPN service
  - two factor authentication
  - software/hardware token
- Supporting vGPU concept in the cloud

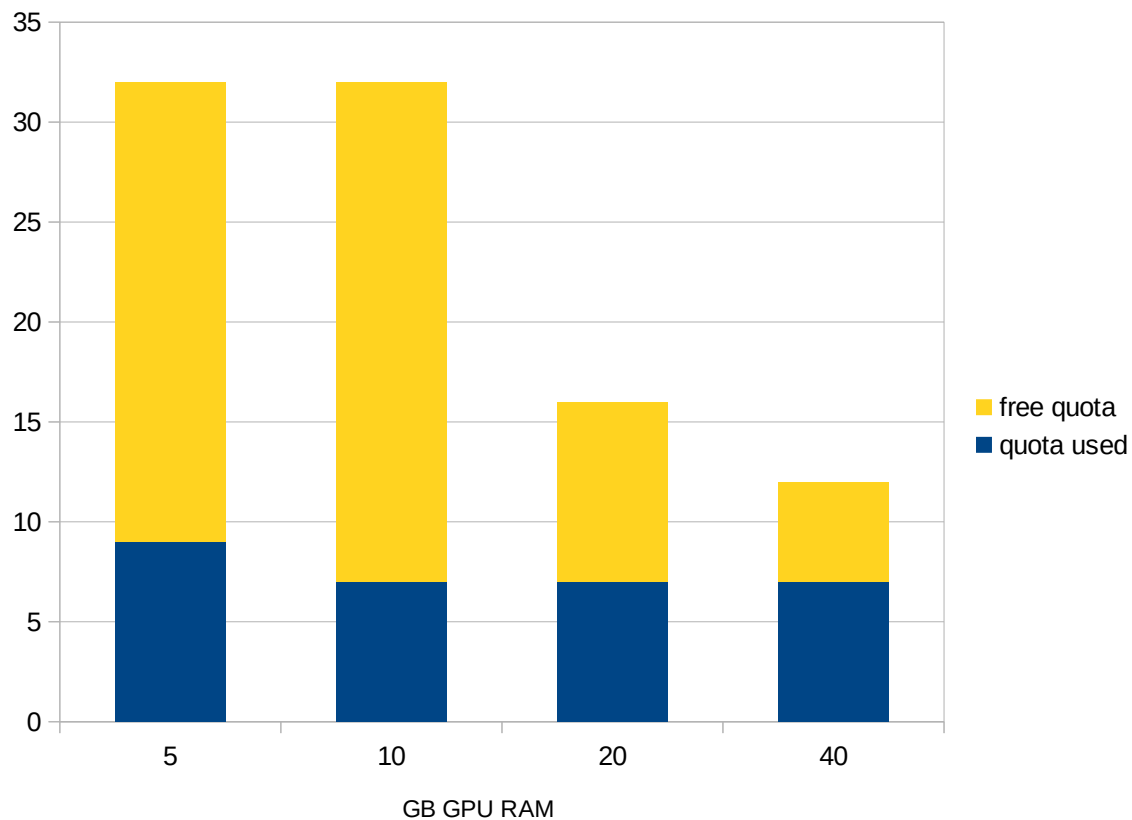
# vGPU - Concept

## Wigner Datacenter



- 4 GPU hosts each with 8 A100 cards
- GPU flavor already very popular
- Only PCIe passthrough at beginning – not flexible
- vGPU concept
- NVIDIA Virtual Compute Server (vCS)

# vGPU - Resources Wigner Datacenter



# vGPU - Hypervisor configuration

## Wigner Datacenter



```
[devices]
enabled_mdev_types = nvidia-469,nvidia-471,nvidia-472,nvidia-473

# A100-5C
[mdev_nvidia-469]
device_addresses = 0000:e1:00.4,0000:e1:00.5,0000:e1:00.6,0000:e1:00.7,0000:e1:01.0,0000:e1:01.1,0000:e1:01.2,0000:e1:01.3

# A100-10C
[mdev_nvidia-471]
device_addresses = 0000:81:00.4,0000:81:00.5,0000:81:00.6,0000:81:00.7,0000:c1:00.4,0000:c1:00.5,0000:c1:00.6,0000:c1:00.7

# A100-20C
[mdev_nvidia-472]
device_addresses = 0000:41:00.4,0000:41:00.5,0000:61:00.4,0000:61:00.5

# A100-40C
[mdev_nvidia-473]
device_addresses = 0000:01:00.4,0000:25:00.4,0000:a1:00.4
```

- Need to be very careful
- Hypervisor => Ubuntu 20.04
- A100 card – linux driver from Nvidia
- nvidia-smi => 470.82.01
- vGPU => 13.1
- CUDA => 11.4
- from VM side
  - 00:05.0 3D controller: NVIDIA Corporation Device 20f1 (rev a1)
  - device ID 20f1 => NVIDIA Corporation GA100 [GRID A100 PCIe 40GB]

# vGPU – Installation steps

## Wigner Datacenter



- Nvidia portal => GPU driver for host and virtual machine
- Licence pool, tokens per project
- Each hypervisor => /usr/lib/nvidia/sriov-manage -e ALL



# vGPU – Installation steps

## Wigner Datacenter



- `openstack flavor create --vcpus 2 --ram 8192 --disk 80 g2.medium`
- `openstack flavor set g2.medium --property "resources:VGPU=1"`
- `openstack trait create CUSTOM_A100_5G`
- `openstack resource provider trait set --trait CUSTOM_A100_5G [UUID]`
- `openstack flavor set --property trait:CUSTOM_A100_5G=required g2.medium`
  
- Deploy configuration
- One vGPU type per card ~ iteration
  
- Driver installation script & guide available
- CUDA installation guide

# vGPU – Running VM from HV Wigner Datacenter



##### NEXT: 0 #####

MDEV\_UUID=2ac2ca31-f78c-4796-b277-5610945f10f7

PCI id 00:01:00.4

vgpu\_type\_id=473

##### NEXT: 1 #####

MDEV\_UUID=2e0d0673-da9f-4c05-bc80-8bd5c5c2ee91

PCI id 00:e1:01.2

vgpu\_type\_id=469

# vGPU – Running VM from HV Wigner Datacenter



```
root@gpu-3:~# mdevctl list
```

```
45187cbe-0b0a-451a-891d-a61e7bf66ada 0000:c1:02.3 nvidia-471  
b694383f-b92c-4ede-bc21-9dd851c7571a 0000:81:00.4 nvidia-471  
5d794584-3db3-421b-b251-2990f482f50f 0000:81:01.5 nvidia-471  
036e3ac0-daba-4982-b709-addf5b46fb19 0000:c1:01.5 nvidia-471  
156653da-670b-462c-8200-f06d383eba98 0000:81:00.6 nvidia-471
```

- CUDA examples test (samples/7\_CUDALibraries/cannyEdgeDetectorNPP)
- Full allocation test
- Host is down? Restore from logs/DB
- Evacuation => feels as a restart
- No licence server / invalid token => 24h CUDA functionalities are stopping

- More application examples as new cloud projects coming
- Better evacuation process
- Slurm Workload Manager + GPU jobs
- Multiple vGPU resource for virtual machine
- Horovod (distributed deep learning training framework)

Thank you for you attention

Ádám Pintér  
pinter.adam@wigner.hu



[www.elkh.org](http://www.elkh.org)