



the definitive partitionable GPU interface

Alessio Borriero, Gabriele Gaetano Fronzé, Federica Legger and Daniele Monteleone

INFN (Italian National Institution for Nuclear Physics) - Research for Innovation grant - Torino







Theoretical GFLOP/s at base clock







Theoretical GFLOP/s at base clock







Is that too much?

Can we exploit that power?

Is that **efficient**?

From Understanding GPU Architecture: Performance: GPU vs. CPU, Steve Lantz, Cornell Center for Advanced Computing



DIVIDE ET IMPERA!

- Smaller GPUs are easier to use efficiently
- Developing on a small GPU is **accessible**
- Some tasks are by design inefficient on huge GPUs
- Large GPUs can be seen as a **set of smaller ones**









		Nvidia vGPU	Nvidia MIG	AMD MxGPU	PCIe SR-IOV
	Complete and constant API support across possible profiles e.g. all profiles support a complete set of API for compute and graphics				N/A
X	P2P communications between partitions e.g. peer-to-peer communications between multiple virtual partitions for compute tasks	~	0	N/A	N/A
A	Free and easy licensing model e.g. the license is "included" in the hardware vs the licenses are an additional cost/require additional procedures	•		~	
E Contraction of the second	Trivial compatibility matrix e.g. delegated to the OS flavor with no limitations wrt an equivalent physical GPU			~	~
	Certified on any physically compatible host system e.g. the compatibility is based on low-enough standards, therefore any physically and electrically supporting hardware works	—	—		

Both within a vendor's techs or between all vendors!





Same underlying **boilerplate**



Same ideal operations and procedures



Different top level semantics



Proprietary "standards" only partially supported in Linux

Both within a vendor's techs or between all vendors!





The independent inititive to make GPU partitioning:

as easy as using a pair of scissors

vendor independent

mainstream



stands for Open For Better Computing





Single Root I/O Virtualization



Open For Better Computing





- Uniform interface for GPU partitioning
- Expandable toolset for **future new technologies**
- No vendor specificity
- Improved Linux Compatibility





gpu list gpu types gpu partition create gpu partition list gpu partition get



openforbc



gpu list

Lists the available physical GPUs compatible with any partitioning technology





openforbc

o p e n ForBC

gpu types

Lists the available virtual GPU profiles, with plenty of information about memory size (WIP: peak performance).





openforbc

o p e n ForBC

gpu types list -c

Lists the **creatable** virtual GPU profiles, handling the complex compatibility matrices automatically.





openforbc

ре

gpu partition create

Applies one of the available profiles performing the required procedures and ensuring not to disrupt existing instances, providing understandable errors in case of wrong selection



gpu partition destroy



openforbc



gpu partition list

Returns the mode currently set up on the GPU.





openforbc



gpu partition get

Retrieves the information needed to instanciate a new VM or container attached to the virtual GPU instance.





🛑 😑 🔵 🤍 🛑

fish /home/monteleo/openforbc

\$ openforbc gpu list
[nvidia:a100-0] 54c2f5e1-6865-3a7b-93c9-3a6e051ac3f0: NVIDIA A100-PCIE-40GB
\$ openforbc gpu -i nvidia:a100-0 types -c
((2, 0)) + 1400 (0, ((, 0)))

468: GRID A100-4C (4.0GiB) 469: GRID A100-5C (5.0GiB) 470: GRID A100-8C (8.0GiB) 471: GRID A100-10C (10.0GiB) 472: GRID A100-20C (20.0GiB) 473: GRID A100-40C (40.0GiB)

\$ openforbc gpu -i nvidia:a100-0 partition create 471

f74efc9f-d5ea-46db-bf00-ab0a15ecee88

\$ openforbc gpu -i nvidia:a100-0 partition get f74efc9f-d5ea-46db-bf00-ab0a15ecee88
NOTE: please ensure that PCI domain:bus:slot.function is not already used.

<hostdev mode='subsystem' type='mdev' managed='no' model='vfio-pci' display='on'>
 <source>

<address uuid='f74efc9f-d5ea-46db-bf00-ab0a15ecee88'/>

</source>

<address type='pci' domain='0x0000' bus='0x00' slot='0x10' function='0x0'/>
</hostdev>

\$ openforbc gpu -i nvidia:a100-0 partition destroy f74efc9f-d5ea-46db-bf00-ab0a15ecee88





One simple line to install it:

/bin/bash -c "\$(curl -fsSL https://raw.githubusercontent.com/Open-ForBC/OpenForBC/install.sh)"

Every line is open source and hosted on GitHub!



https://github.com/Open-ForBC/OpenForBC/tree/main









Planned in our roadmap We got from AMD an S7150x2 to start thinkering with. Looking forward for v320/v340/v520/v540 **3D tasks:** Top of the line vGPU instances

Up to 128 users per server Down to 1/6 power consuption

Hardware choice: AMD Radeon PRO V340

40500 000000



2D tasks: Up to 32 vGPUs (32 users!)



Is GPU partitioning actually worth it?

We've created a tool to answer this question





BENCHMARKS



Our modular benchmark suite It includes our **custom benchmarks** Is compatible with **Phoronics benchmarks** Is easily **expandable** with additional benchmark definitions





Every line is open source and hosted on GitHub!

Feel free to submit a PR with your own benchmarks!



https://github.com/Open-ForBC/OpenForBC-Benchmark





Let's look at a typical ML task...



In **green** the peak throughput of an A100-40GB PCIe.

In **blue** the average peak throughput of all the creatable* partitions for a given profile.

* All creatable partitions have been allocated and loaded with computation.



...and rescale the peak performance



In green the peak throughput of an A100-40GB PCIe.

In **blue** the peak throughput of each partition topology, computed as the sum of the average throughput of all creatable* partitions given a specific profile.

Maximum speedup of 105%.

* All creatable partitions have been allocated and loaded with computation.



What happens with inference?





More examples:

ForBC





Speedup 305% and 409%

Speedup 461% and 428%



Partners and supporters









WHO ARE WE





INFN's Research for Innovation 2021 grant



Federica Legger Senior researcher – PhD in Physics



Gabriele Gateano Fronzé

Senior researcher – PhD in Physics



Alessio Borriero Physics Master's Student



Daniele Monteleone Computer Engeneering Bachelor's Student

And the whole INFN Torino Computing Group





GPU partitioning shows clear advantages: huge speedups and reduced relative power consuption

Open ForBC introduces a smarter way to use partitionable GPUs on Linux KVM.

Open ForBC handles the complexity due to vendors' choices providing users with a powerful yet simple toolset.

Open ForBC is open source, provides a CLI and a REST API, supports Nvidia GPUs (MIG and vGPU) with AMD support coming next.

Open ForBC Benchmark is an expandable modular benchmark framework.







Thank you for your attention!







Let's look at a typical ML task...



In green the peak performances of an A100-40GB PCIe

In yellow the peak performance of one of N partitions packing 1/N of the whole GPU resources, loaded simultaneously.

In **blue** the peak performance of a single partition packing 1/N of the whole GPU resources, without other running partitions.

BENCHMARKS

...and rescale the peak performance



Same color code as before.

The peak performance of each partition topology has been rescaled by a factor equal to the maximum number of partitons of such kind runnable on a single physical GPU.

Seven 1/7 partitions reach up to 105% speedup with respect to a single A100-40GB PCIe.

BENCHMARKS

Overhead evaluation



Overhead computed as difference between the peak performance of a partitioning topology when all N partitions are loaded simultaneously and when only one is loaded.

This measures the overhead caused by internal scheduling and handling of partitions in the worst case scenario.

Maximum overhead is around 3.7%, with the 4/7 case being caused by the impossibility of allocating more than one 4/7 partition.

