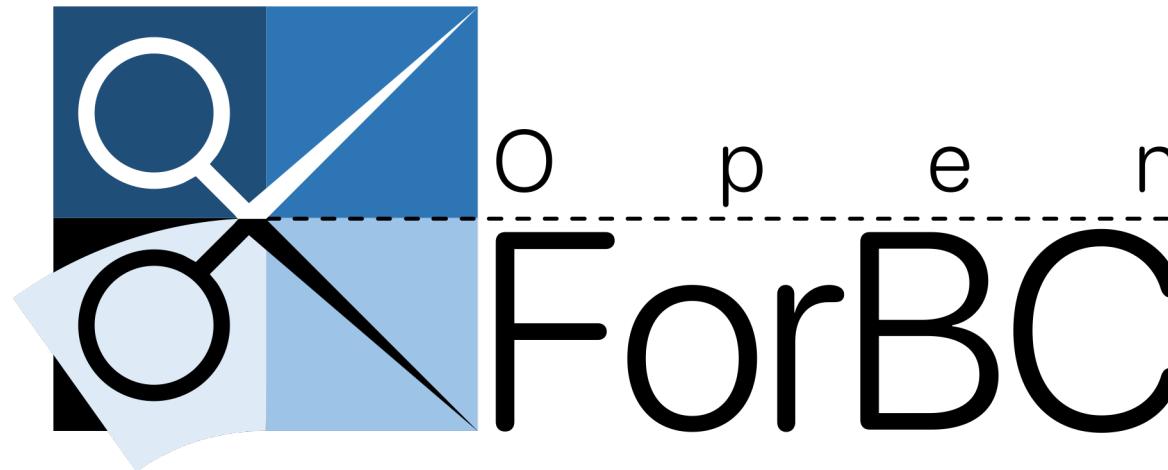


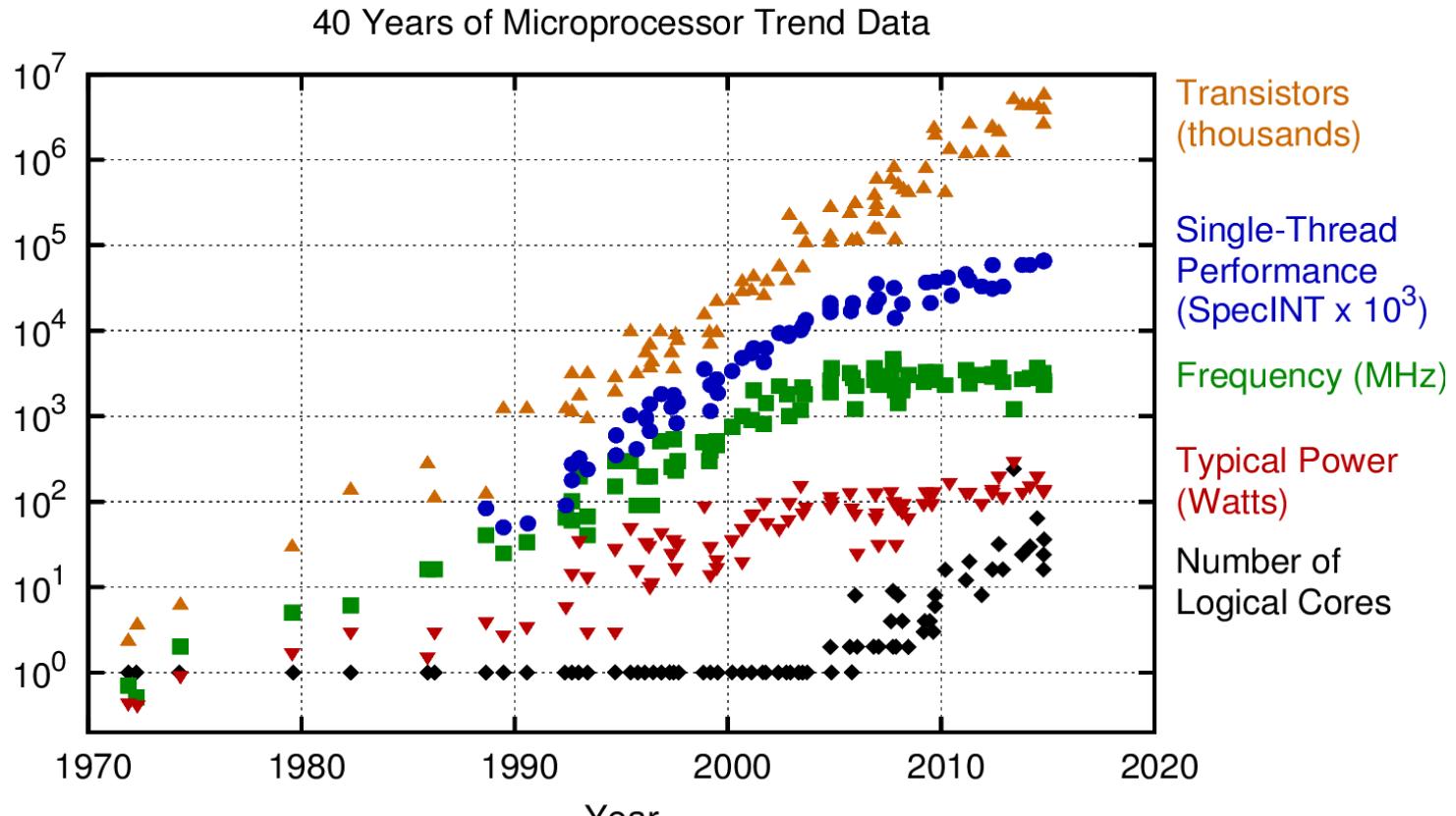


Istituto Nazionale di Fisica Nucleare



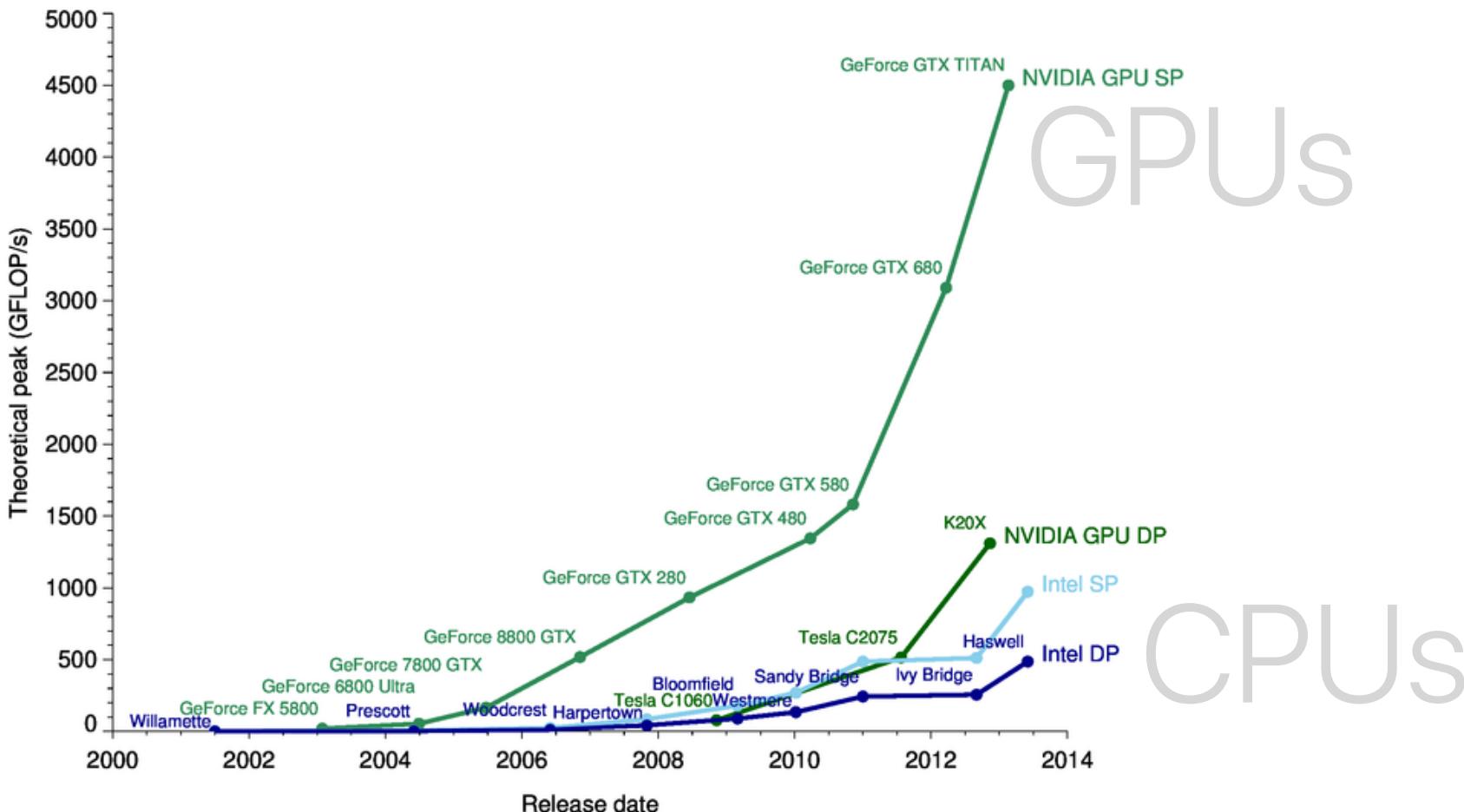
the definitive partitionable GPU interface

RATIONALE

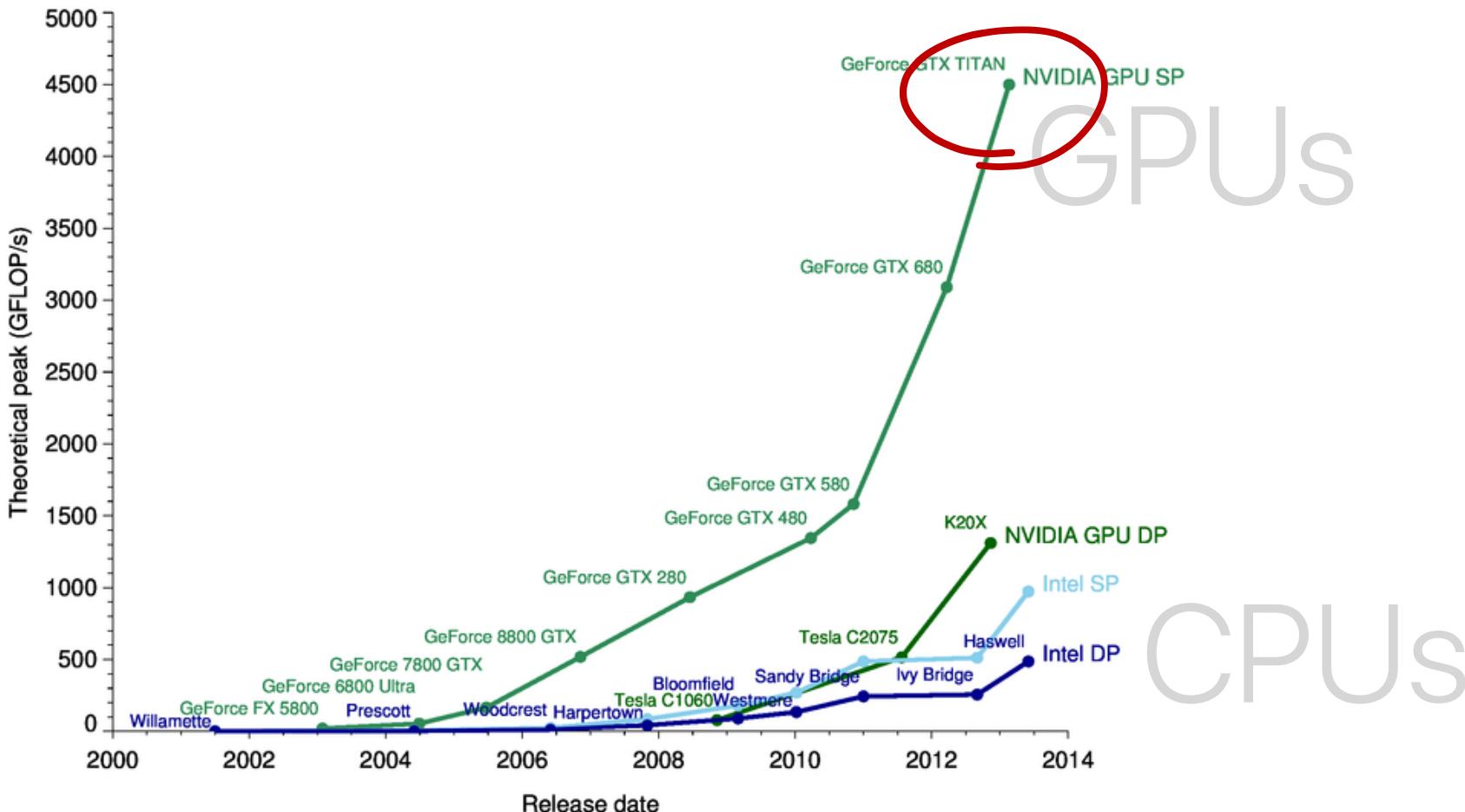


Original data up to the year 2010 collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond, and C. Batten
New plot and data collected for 2010-2015 by K. Rupp

RATIONALE



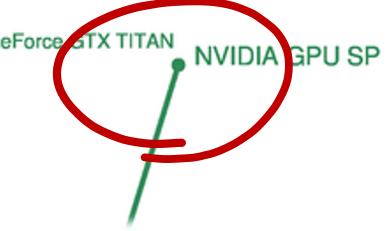
RATIONALE



RATIONALE

Is that too much?

Can we **exploit** that power?



Is that **efficient**?

RATIONALE

Is that too much?

Can we exploit that power?

Is that efficient?



STATE OF ART



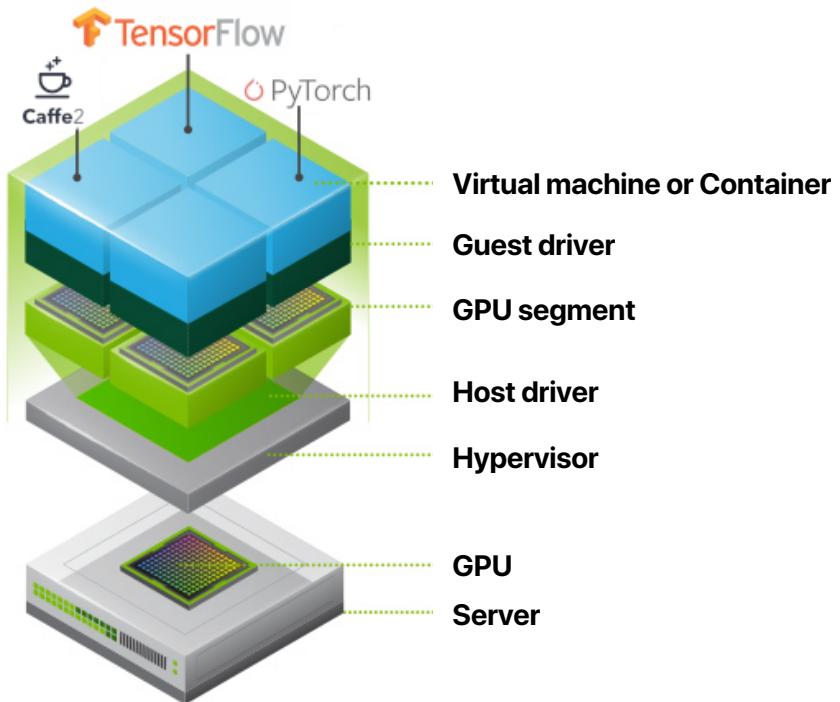
USING VIRTUAL MACHINES
DIVIDE ET IMPERA!

- 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**

STATE OF ART



USING VIRTUAL MACHINES
DIVIDE ET IMPERA!

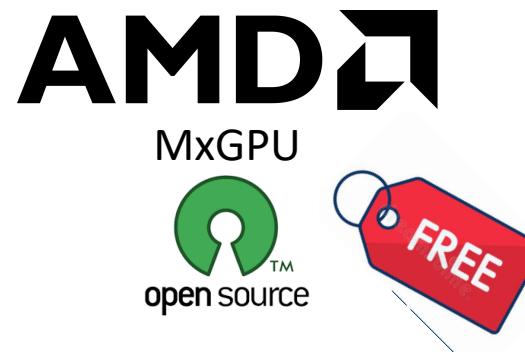
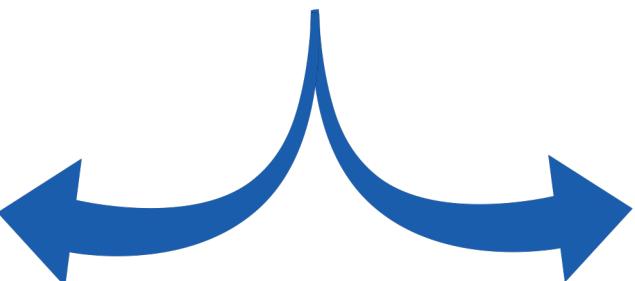
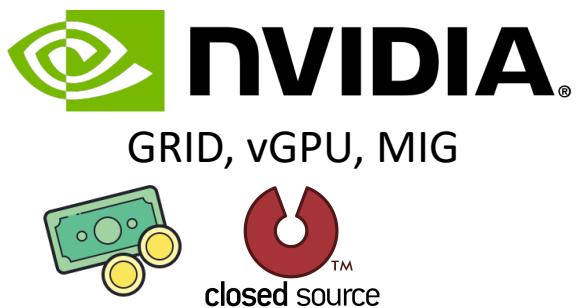


AUTEM...

STATE OF ART



Single Root I/O Virtualization



Everyone is adopting its
own “standard” derived
from the only real one

STATE OF ART



Single Root I/O Virtualization



Same underlying standard



Same technical operations and procedures



Different top level semantics

STATE OF ART



Single Root I/O Virtualization

Proprietary “standards” only partially implemented in Linux



OPEN ForBC



Single Root I/O Virtualization

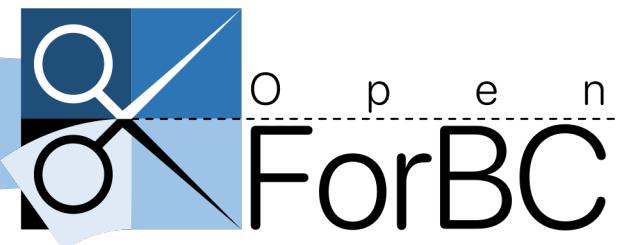


GRID, vGPU, MIG



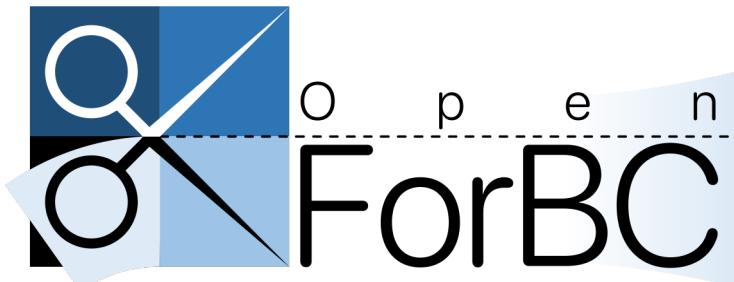
MxGPU

Open For Better Computing



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

OPENForBC



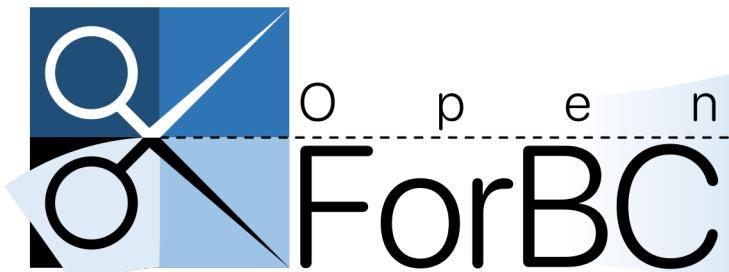
`get_available_modes`

`set_mode`

`get_current_mode`

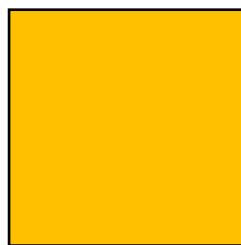
`get_instance`

OPENForBC

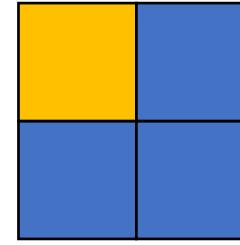


get_available_modes

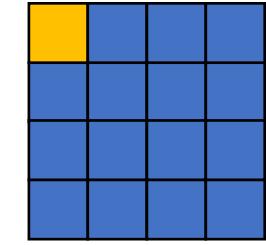
Lists the available virtual GPU profiles, with plenty of information about peak performances and memory size.



1 vGPU

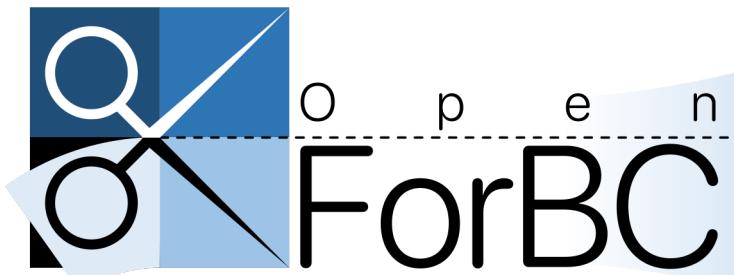


4 vGPUs



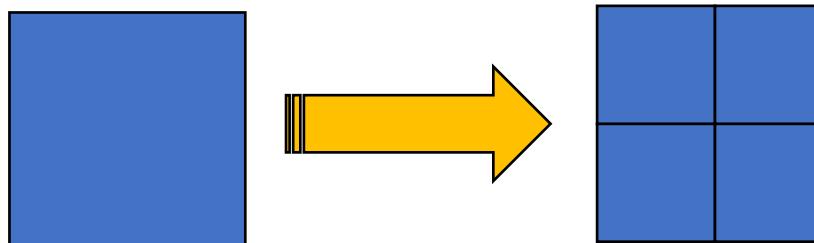
16 vGPUs

OPENForBC

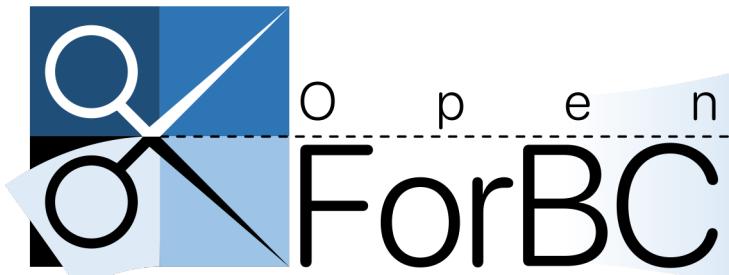


set_mode

Applies one of the available profiles performing the required procedures and ensuring not to disrupt existing instances

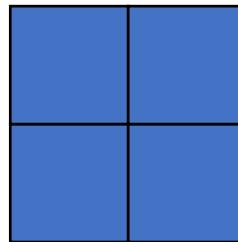


OPENForBC



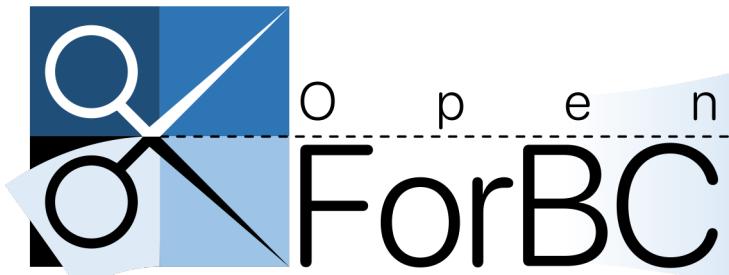
`get_current_mode`

Returns the mode currently set up on the GPU.



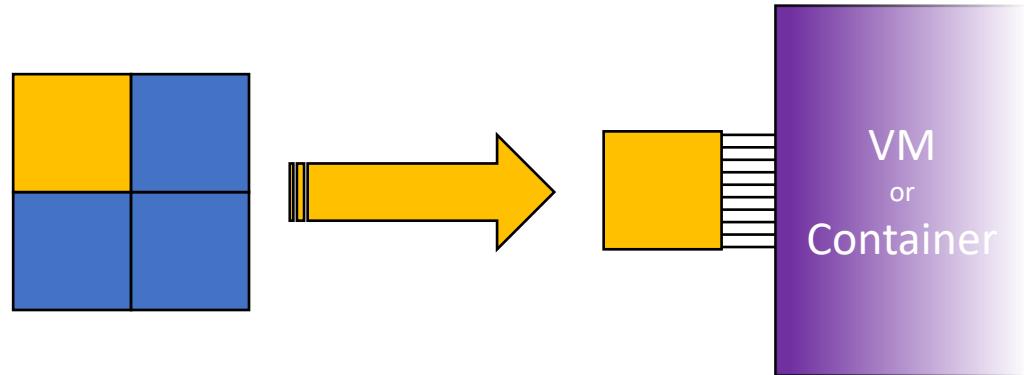
4 vGPUs

OPENForBC



get_instance

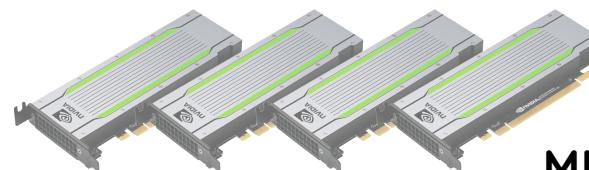
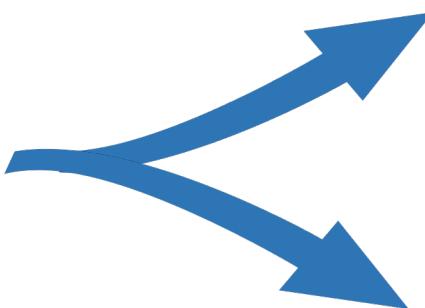
Retrieves the information needed to instantiate a new VM or container attached to the vGPU instance.



OPENForBC

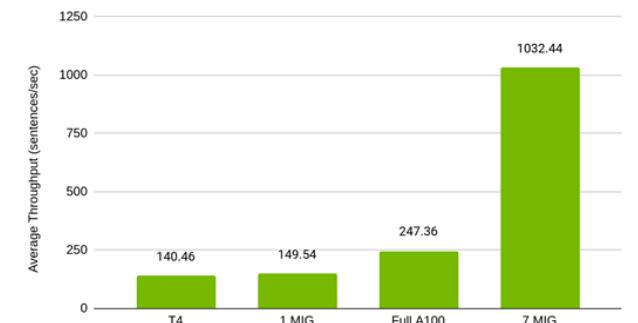


Hardware choice:
Nvidia V100 32GB GPU



ML training:
huge GPU with lots of memory

Benchmark: BERT large TensorFlow Inference (SQuAD, BS=1)

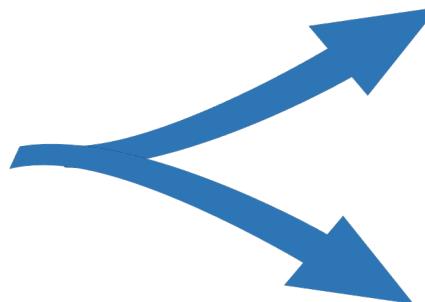


ML inference:
smaller GPUs for higher throughput

OPENForBC



Hardware choice:
AMD Radeon PRO V340



3D tasks:
Top of the line vGPU instances

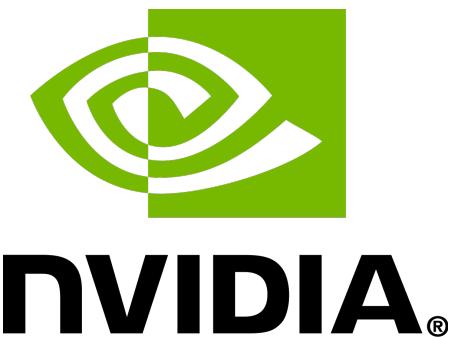


Up to 128 users per server
Down to 1/6 power consumption

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

OPENForBC

Partner tecnici e *sample provider*:



OPENForBC

Federica Legger: 0.2 FTE

Gabriele Gaetano Fronzé: 0.2 FTE

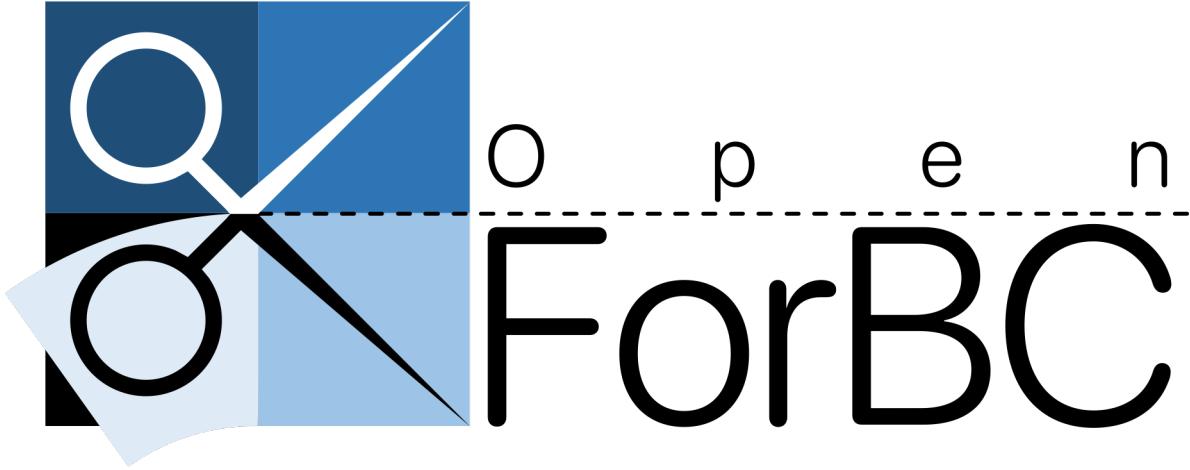
AdR 1 anno: 1.0 FTE

OPEN ForBC

Open ForBC introduces a smarter way to use partitionable GPUs.

Optimizing your hardware deployment for your tasks is few clicks away.

Never depend on manufacturer choices: with Open ForBC learn one and use forever.



Thank you for your attention!

