



Contribution ID: 175

Type: Oral

## Pinpoint resource allocation for GPU batch applications

*Thursday 14 March 2024 15:10 (20 minutes)*

With the increasing usage of Machine Learning (ML) in High Energy Physics (HEP), the breadth of new analyses with a large spread in compute resource requirements, especially when it comes to GPU resources. For institutes, like the Karlsruhe Institute of Technology (KIT), that provide GPU compute resources to HEP via their batch systems or the Grid, a high throughput, as well as energy efficient usage of their systems is of the essence. With low intensity GPU analyses specifically, inefficiencies are created by the standard scheduling, as resources are over-assigned to such workflows. An approach that is flexible enough to cover the entire spectrum, from multi-process per GPU, to multi-GPU per process, is necessary. As a follow-up to the techniques presented at the 2022 ACAT, this time we study Nvidia's multi-process service (MPS), its ability to securely distribute device memory and its interplay with the KIT HTCondor batch system. A number of ML applications were benchmarked using this less demanding and more flexible approach to illustrate the performance implications regarding throughput and energy efficiency.

### Significance

Batch systems are crucial for the efficient and high-throughput computing that is required in modern high energy physics. Often, these batch systems are limited by their coarse granularity. Especially for GPU resources, the safe sharing of high performance datacenter GPUs is necessary to avoid gross over-allocation of costly hardware, while still allowing for workflows that require multiple GPUs at once. Nvidia's multi-process service (MPS) enables this kind of flexibility, and we therefore consider it a valuable tool for our goal of high throughput and high energy efficiency.

### References

<https://indico.cern.ch/event/1106990/contributions/4991345/>

### Experiment context, if any

CMS

**Primary author:** VOIGTLAENDER, Tim (KIT - Karlsruhe Institute of Technology (DE))

**Co-authors:** QUAST, Gunter (KIT - Karlsruhe Institute of Technology (DE)); GIFFELS, Manuel (KIT - Karlsruhe Institute of Technology (DE)); SCHNEPF, Matthias Jochen; WOLF, Roger (KIT - Karlsruhe Institute of Technology (DE))

**Presenter:** VOIGTLAENDER, Tim (KIT - Karlsruhe Institute of Technology (DE))

**Session Classification:** Track 1: Computing Technology for Physics Research

**Track Classification:** Track 1: Computing Technology for Physics Research