Services for Machine Learning applications (part 2 of 3)

Diana Gaponcic, IT-CD-PI

GPU clusters at CERN

How to create a GPU cluster

$\bullet\bullet\bullet$

\$ openstack coe cluster create digaponc-gpu-004 --merge-labels --labels nvidia_gpu_enabled=true

How to create a GPU cluster

1. By default 2 nodes are deployed: the master and the default worker node

How to create a GPU cluster

- 1. By default 2 nodes are deployed: the master and the default worker node
- **2. No GPU yet**
	- a. the cluster is configured to manage GPUs, but we don't get a GPU by default

GPU flavors

Consult https://clouddocs.web.cern.ch/gpu_overview.html for an up-to-date list of GPU flavors

Add a GPU node

$\bullet\bullet\bullet$

\$ openstack coe nodegroup create digaponc-gpu-004 gpu-t4 --flavor g2.5xlarge --node-count 1

 \sim \sim

NVIDIA GPU operator

$\bullet\bullet\bullet$

\$ kubectl get pod -n kube-system | grep nvidia nvidia-container-toolkit-daemonset-8hfwn nvidia-cuda-validator-dlpmt nvidia-dcgm-exporter-lm4kn nvidia-device-plugin-daemonset-9w9xk nvidia-driver-daemonset-sqs5c nvidia-operator-validator-7scl5

nvidia-driver-daemonset

Loads the drivers on the node

nvidia-container-toolkit-ctr

The toolkit includes a container runtime library and utilities to automatically configure containers to leverage NVIDIA GPUs.

nvidia-dcgm-exporter + nvidia-operator-validator

NVIDIA Data Center GPU Manager (DCGM) is a suite of tools for managing and monitoring NVIDIA datacenter GPUs. It exposes GPU metrics exporter for Prometheus leveraging NVIDIA DCGM.

nvidia-device-plugin-daemonset

Allows to automatically:

- 1. Expose the number of GPUs on each nodes of your cluster
- 2. Keep track of the health of your GPUs
- 3. Run GPU enabled containers in your Kubernetes cluster.

This is what allows NVIDIA GPUs to be requested by a container using the **nvidia.com/gpu** resource type.

nvidia-cuda-validator

Validates that the stack installation worked

Node feature discovery

$\bullet\bullet\bullet$

\$ kubectl get pod -n kube-system | grep node-feature-discovery cern-magnum-node-feature-discovery-gc-7985cbd94b-q499t cern-magnum-node-feature-discovery-master-7bbccf9b68-fjpp8 cern-magnum-node-feature-discovery-worker-5qjzq cern-magnum-node-feature-discovery-worker-qhbrc

 $1/1$

 $1/1$

 $1/1$

 $1/1$

Node feature discovery

$\bullet\bullet\bullet$

\$ kubectl get pod -n kub cern-magnum-node-feature cern-magnum-node-feature cern-magnum-node-feature cern-magnum-node-feature

local.feature: elements: nvidia.com/cuda.driver-version.full: 550.54.15 nvidia.com/cuda.driver-version.major: "550" nvidia.com/cuda.driver-version.minor: "54" nvidia.com/cuda.driver-version.revision: "15" nvidia.com/cuda.driver.major: "550" nvidia.com/cuda.driver.minor: "54" nvidia.com/cuda.driver.rev: "15" nvidia.com/cuda.runtime-version.full: "12.4" nvidia.com/cuda.runtime-version.major: "12" nvidia.com/cuda.runtime-version.minor: "4" nvidia.com/cuda.runtime.major: "12" nvidia.com/cuda.runtime.minor: "4" nvidia.com/gfd.timestamp: "1728992460" nvidia.com/gpu.compute.major: "7" nvidia.com/gpu.compute.minor: "5" nvidia.com/qpu.count: "1" nvidia.com/qpu.family: turing nvidia.com/qpu.machine: OpenStack-Compute nvidia.com/qpu.memory: "15360" nvidia.com/qpu.mode: compute nvidia.com/gpu.product: Tesla-T4 nvidia.com/gpu.replicas: "1" nvidia.com/gpu.sharing-strategy: none nvidia.com/mig.capable: "false" nvidia.com/mig.strategy: mixed nvidia.com/mps.capable: "false" nvidia.com/vqpu.present: "false"

 $^{\circ}$

 N/A

Tainting

Taint Nodes

With kubernetes templates 1.24+, the gpu-operator helm chart does not taint GPU nodes which will allow all workloads to run in this nodes. We suggest to taint the nodes explicitly by adding the following taint to the GPU nodegroups:

node-role.kubernetes.io/gpu=true:NoSchedule

Disclaimer: We will have automatic tainting in the next release

Let's run some workloads

Example Use Cases (very different GPU consumption behaviour)

Badly coded simulation job:

- Low average GPU usage (CPU dependant workload)
- Needs 10 GiB VRAM $(8 + 2)$ dynamic)
- Long running process

Never know what to expect from a notebook user:

- Potential memory leaks
- Poorly considered batch size
- GPU memory locked by an idle notebook

An inference service which is occasionally triggered by outside events:

- Spiky and unpredictable execution
- Mostly sits idle
- Saturates the GPU cores
- Max 10 GiB VRAM ($2 + 8$ dynamic)

* All use cases were run on a CERN Kubernetes cluster with 1 NVIDIA A100 40GB GPU

Onboard **Only** Use Case 1 = Dedicated GPU

Badly coded simulation job:

- Low average GPU usage (CPU dependant workload)
- Needs 10 GiB VRAM ($8 + 2$ dynamic)
- Long running process

- **GPU underutilized**
- Steady memory utilization \sim 20%

Dedicated GPU drawbacks

- Dedicated GPUs => small/limited GPU offering
- Some use cases cannot fully utilize a GPU => idle time

Dedicated GPU drawbacks

- Dedicated GPUs => small/limited GPU offering
- Some use cases cannot fully utilize a GPU => idle time

How to improve?

GPU Sharing

1. Time-slicing

Time-slicing

- The scheduler gives an equal share of time to all GPU processes and alternates them in a round-robin fashion.
- The memory is shared between the processes
- The compute resources are assigned to one process at a time

GPU

kubectl label node <node-name> nvidia.com/device-plugin.config=slice-4

- GPU underutilized
- Steady memory utilization \sim 20%

- **GPU underutilized**
- Steady memory utilization \sim 20%

- Improved GPU utilization
- **● Better memory consumption (~ 50 %)**

Time-Slicing

Advantages Disadvantages

Works on a wide range of NVIDIA architectures

An easy way to set up GPU concurrency

No process/memory isolation

No ability to set priorities

An unlimited number of partitions Inappropriate for latency-sensitive applications (ex: desktop rendering for CAD workloads)

GPU Sharing

2. Multi Instance GPU

Multi Instance GPU

Multi Instance GPU (MIG) can partition the GPU into up to seven instances, each fully isolated with its own high-bandwidth memory, cache, and compute cores.

[MIG Profiles on A100](https://docs.nvidia.com/datacenter/tesla/mig-user-guide/)

NVIDIA MIG provides multiple strategies for allowing users to reference the graphic card resources:

- **mixed**: Different resource types are enumerated for every MIG device available. Ex: nvidia.com/mig-3g.20gb
- **single**: MIG devices are enumerated as nvidia.com/gpu, and map to the MIG devices available on that node, instead of the full GPUs.
- **none**: No distinction between GPUs with MIG or without. The available devices are listed as nvidia.com/gpu.

kubectl label nodes <node-name> nvidia.com/mig.config=3g.20gb-2x2g.10gb

Hardware level sharing - MIG

Advantages Disadvantages

Hardware isolation allows processes to run securely in parallel and not influence each other

Monitoring and telemetry data available at partition level

Only available for Ampere, Hopper, and Blackwell architecture

Reconfiguring the partition layout requires all running processes to be evicted

Allows partitioning based on use cases, making the solution flexible

* Potential loss of available memory depending on chosen profile layout

* Not a risk if the partitioning layout is chosen in an informed way after careful consideration.

We established that GPU sharing increases overall usage.

But how do we share in the best way?

But how do we share in the best way?

- 1. Summarize **in your team** all workloads that need GPUs. Run them on one cluster and collocate them using time-slicing and MIG.
- 2. Single point of GPU Access across multiple teams

- 1. GPUs are always in-use
	- a. As soon as a GPU is released by an user, it is reassigned to another one requesting a GPU
- 2. People can get access to multiple types of GPUs, or even other accelerators (TPUs, IPUs) through public cloud.

GPU sharing tradeoffs

Benchmarked script:

- Simulation script that generates collision events. [Find more](https://kubernetes.docs.cern.ch/blog/2023/03/20/efficient-access-to-shared-gpu-resources-part-3/#compute-intensive-particle-simulations)
- Built with Xsuite (Suite of python packages for multiparticle simulations for particle accelerators)
- Very heavy on GPU usage
- Low on memory accesses
- Low on CPU-GPU communication

Environment:

- NVIDIA A100 40GB PCIe GPU
- Kubernetes version 1.22
- Cuda version utilized: 11.6
- Driver Version: 470.129.06

Time-slicing Performance Analysis

The GPU **context switching** caused a ~38% performance loss

Time-slicing Performance Analysis

There is no additional performance loss when sharing the GPU between more processes (4, 8, and even more).

MIG Performance Analysis

MIG Performance Analysis

The theoretical loss of **9.25%** is seen experimentally.

MIG Performance Analysis

The scaling between partitions converges to ideal values.

GPU Sharing Use Cases

¹ Independent workloads can trigger OOM errors between each other. Needs an external mechanism to control memory usage (similar to kubelet CPU memory checks)

Monitoring

<https://grafana.com/grafana/dashboards/18288-nvidia-gpu/>

...

DCGM_FI_PROF_PIPE_TENSOR_ACTIVE, gauge, Ratio of cycles the tensor (HMMA) pipe is active (in %). DCGM_FI_PROF_PIPE_FP64_ACTIVE, gauge, Ratio of cycles the fp64 pipes are active (in %). DCGM_FI_PROF_PIPE_FP32_ACTIVE, gauge, Ratio of cycles the fp32 pipes are active (in %). DCGM_FI_PROF_PIPE_FP16_ACTIVE, gauge, Ratio of cycles the fp16 pipes are active (in %).

Profiling the A100 compute pipeline utilization

<https://docs.nvidia.com/datacenter/dcgm/2.4/dcgm-api/dcgm-api-field-ids.html>

GPU access using Kubeflow

\leftarrow New notebook

CPU / RAM @

\vee Advanced Options

GPUs

Find more:

- [https://ml.docs.cern.c](https://ml.docs.cern.ch/) [h/](https://ml.docs.cern.ch/)
- <https://ml.cern.ch/>

… and in the next ML session :)

Conclusions

- 1. It is easy to create a cluster with GPU nodes
	- a. The user is abstracted away from having to set any drivers
- 2. GPU sharing is useful to improve the overall GPU utilization, but it comes with performance tradeoffs
	- a. Sharing helps us to offer GPUs to more users
	- b. For use cases that can fully utilize the GPU, we need to consider allocating dedicated GPUs
- 3. Monitoring is very important
	- a. The current infrastructure is flexible enough to cater for various use cases
- 4. For ML workloads consider using Kubeflow

Thank you!

Bonus slides

Not a service, but very good to know

Data Formats

We've heard about fp32 and fp64, but what about ML?

Data formats

Data formats

Data formats

Data formats

Data formats people starting using half precision (fp16).

> Since fp16 has small range (5 bits), Trainings with fp16 only, very often results in numerical errors (underflow, overflow).

But ML needs range!

This is how mixed precision was invented:

- What is numerical stable \rightarrow done in fp16
- For the rest \rightarrow there is a copy stored in fp32

You get the speed, while avoiding numerical errors, but memory consumption is bigger (V100, T4).

Data formats

* ML needs range

This is how BFloat16 was created (available for Ampere and newer)

The range of fp32, the precision less than fp16

* ML doesn't need a lot of precision!

Using low precision works because ML tries to minimize the error. As long as it can represent the number and move towards a smaller error - low precision will work and bring speedup.

Low precision can be a problem for regression tasks (when you want to get a real number as output) - then think about using another format.

Data formats

tf32 is a solution to use when you still want the speedup of bf16, but with a bigger precision.

Things to keep in mind

- This is an NVIDIA solution (amd still not very clear what is the current state)
- not existing on older GPUs

Exercises Time

- Access the exercises:
- [https://indico.cern.ch/event/1441237/contributions/6073469/attachments/2952019/52113](https://indico.cern.ch/event/1441237/contributions/6073469/attachments/2952019/5211317/exercises) [17/exercises](https://indico.cern.ch/event/1441237/contributions/6073469/attachments/2952019/5211317/exercises)
- Group 1
	- Shahzaib Aamir
	- Berk Balci
	- Nayana Bangaru
	- Gábor Bíró
	- Abhishek Bohare
	- Marco Buonsante

- Group 2
	- Gianluca De Bonis
	- Elena de la Fuente Garcia
	- Jesse Geens
	- Ediz Genc
	- Panagiotis Georgopoulos
	- Daniel Goncalves Portovedo
- Group 3
	- Panagiotis Gkonis
	- Dmytro Gruzdo
	- Hannes Jakob Hansen
	- Musa Kaymaz
	- Idriss Larbi
	- Manuel Ramirez Garcia
- Group 4
	- Rimsky Alejandro Rojas Caballero
	- Jonathan Samuel
	- Nikolaos Smyrnioudis
	- Lorenzo Valentini
	- Jan Hauke Voss
	- Julian Weick

- Group 5
	- Joao Ramiro
	- Lorenzo Ventura Vagliano
	- Pascal Egner
	- Jakub Jelinek
	- Luis Pelaez Bover
	- Franciska-Leonora Toeroek

