# CERN IT GPU Update

Ricardo Rocha, CERN IT

Compute & Accelerator Forum, June 12th 2024

https://indico.cern.ch/event/1329690/

# Reminder

https://clouddocs.web.cern.ch/gpu/index.html

GPU Availability

Bare metal, VMs, Batch, Kubernetes Clusters

High Level Services: lxplus-gpu, GitLab runners, SWAN, ml.cern.ch

Requesting dedicated GPUs

Dedicated functional element: GPU Platform Consultancy

#GPU channel on IT-dep mattermost

# What's New

Long time since the last update

Since then

Nvidia A100 GPUs available, with partitioning capabilities (MIG)

Nvidia H100 GPUs, expected online August 2024

Updates to the GitLab CI runners

CERN IT ML Infrastructure Workshops, also covering GPU requirements

Benchmarks for GPU sharing capabilities

# Resources

Assignment to the different services depending on needs

| Card Type | Number of Cards | Notes |
| --- | --- | --- |
| T4 | 76 | |
| V100 & V100S | 40 | |
| A100 | 72 | |
| H100 | 52 | Available August 2024 |

# Resources

Dedicated resources via request to the [GPU Platform Consultancy FE](#)

# Resources

Dedicated resources via request to the [GPU Platform Consultancy FE](#)

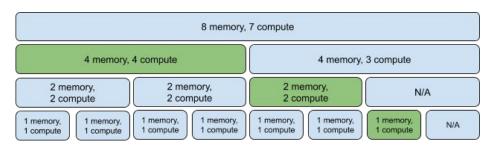| Flavor Name | GPU | RAM | vCPUs | Disk | Ephemeral | Comments |
|---|---|---|---|---|---|---|
| g1.xlarge | V100 | 16 GB | 4 | 56 GB | 96 GB | [^1], deprecated |
| g1.4xlarge | V100 (4x) | 64 GB | 16 | 80 GB | 528 GB | [^1] |
| g2.xlarge | T4 | 16 GB | 4 | 64 GB | 192 GB | [^1], deprecated |
| g2.5xlarge | T4 | 168 GB | 28 | 160 GB | 1200 GB | [^1] |
| g3.xlarge | V100S | 16 GB | 4 | 64 GB | 192 GB | [^1] |
| g3.4xlarge | V100S (4x) | 64 GB | 16 | 128 GB | 896 GB | [^1] |
| g4.p1.40g | A100 (1x) | 120 GB | 16 | 600 GB | - | [^1], AMD CPUs |
| g4.p2.40g | A100 (2x) | 240 GB | 32 | 1200 GB | - | [^1], AMD CPUs |
| g4.p4.40g | A100 (4x) | 480 GB | 64 | 2400 GB | - | [^1], AMD CPUs |

# Multi-Instance GPUs

https://www.nvidia.com/en-us/technologies/multi-instance-gpu/

Available for both Nvidia A100 and H100 GPUs

Physical partitioning of GPU cards, up to 7 times



H100 brings MIG v2 allowing partition reconfiguration without workload eviction

# GitLab CI GPU Runners

[Documentation](#)

Single flavor GPU runners, no differentiation for specific cards (at least for now)

CVMFS and EOS both available

```
job:
  tags:
    - k8s-gpu
  image: rochaporto/gpu_burn # overrides the default image.
  script:
    - nvidia-smi
    - cd /app
    - ./gpu_burn 120
```

```
19  Running on runner-zfzb5pgwc-project-184215-concurrent-3-amit007d via runners-k8s-gpu-866cb88495-rnp6c...
20  Getting source from Git repository                                                          00:01
21  Fetching changes with git depth set to 20...
22  Initialized empty Git repository in /builds/rbritoda/gpu-runner-test/.git/
23  Created fresh repository.
24  Checking out 04d4f54e as detached HEAD (ref is master)...
25  Skipping Git submodules setup
26  Executing "step_script" stage of the job script                                             01:06
27  $ nvidia-smi
28  Wed Jun 12 11:54:47 2024
29  +-----------------------------------------------------------------------------------------+
30  | NVIDIA-SMI 550.54.15              Driver Version: 550.54.15      CUDA Version: 12.4       |
31  |-----------------------------------------+------------------------+----------------------+
32  | GPU  Name               Persistence-M | Bus-Id        Disp.A | Volatile Uncorr. ECC |
33  | Fan  Temp   Perf          Pwr:Usage/Cap |         Memory-Usage | GPU-Util  Compute M. |
34  |                                         |                        |               MIG M. |
35  |=========================================+========================+======================|
36  |   0  NVIDIA A100-PCIE-40GB      On  | 00000000:00:06.0 Off |                   On |
37  | N/A   62C    P0             174W / 250W |                 N/A |     N/A       Default |
38  |                                         |                        |              Enabled |
39  +-----------------------------------------+------------------------+----------------------+
40  +-----------------------------------------------------------------------------------------+
41  | MIG devices:                                                                             |
42  +------------------+----------------------+-----------+-----------------------+
43  | GPU  GI  CI  MIG |          Memory-Usage |        Vol|        Shared         |
44  |      ID  ID  Dev |           BAR1-Usage | SM     Unc| CE ENC DEC OFA JPG    |
45  |                  |                       |        ECC|                       |
46  |==================+=======================+===========+=======================|
47  |   0    9   0   0 |           5MiB /  4864MiB | 14       0| 1   0   0   0   0 |
48  |                  |            0MiB /  8191MiB |           |                    |
49  +------------------+----------------------+-----------+-----------------------+
```

# CERN IT ML Infrastructure - Workshops

Report from 2nd workshop: https://indico.cern.ch/event/1358625/

## Action items

1. **Better advertising of available tools and resources in IT (and how to use them)**

   Request for a single entrypoint documenting access to GPUs and ML capable services, with recommendations

2. **Clarify reported bottlenecks on accessing storage**

   Unclear which backends were being referred, discussion on patterns accessing from public cloud

3. **AutoML / Hyper-parameter optimization seen as a crucial aspect**

   Integrated in ml.cern.ch, request for multi-GPU support in Batch (essential for upcoming years)

   Evaluate a central ML model repository

4. **Better coordination of access to accelerator resources**

   Sharing of GPUs between services, hurdles integrating online experiment resources (ALICE EPN, LHCb HLT1)

5. **Help with profiling and benchmarking of ML workloads**

   Deployment optimization requires expertise, IT can help. Opportunity to collaborate with industry partners

6. **Tackle needs for dedicated architectures and licensed software (including drivers)**

# GPU Sharing and Concurrency - Benchmark Analysis

Goal: Improve overall utilization of available hardware

Extensive benchmark analysis of different sharing and concurrency techniques

Part 1: Motivations and Use Cases

Part 2: Setup and configuration of GPU concurrency

Part 3: Benchmarking Use Cases

Part 4: Time slicing results

Part 5: MIG results

Part 6: ML results

Ongoing… part 7 with results for multi-process service (MPS)

# Ongoing work

Implement feedback from the ML workshops (action items)

Ongoing effort to improve overall GPU usage

  Improved sharing and flexibility (re)allocating resources to services

Integration with public cloud resources

  Potentially a larger set of options for accelerators, on-demand

  Ongoing technical work, many details still to be defined

# Q & A