

AI at Scale: HPC Primitives in the Cloud-Native Era

Ricardo Rocha

CERN IT, CNCF TOC, CNCF End User TAB

1. Mission of the Cloud Native Computing Foundation.

The Foundation's mission is to make cloud native computing ubiquitous. The CNCF Cloud Native Definition v1.0 says:

Cloud native technologies empower organizations to build and run scalable applications in modern, dynamic environments such as public, private, and hybrid clouds. Containers, service meshes, microservices, immutable infrastructure, and declarative APIs exemplify this approach.

These techniques enable loosely coupled systems that are resilient, manageable, and observable. Combined with robust automation, they allow engineers to make high-impact changes frequently and predictably with minimal toil.

The Cloud Native Computing Foundation seeks to drive adoption of this paradigm by fostering and sustaining an ecosystem of open source, vendor-neutral projects. We democratize state-of-the-art patterns to make these innovations accessible for everyone.



Immediate popularity and traction

For traditional, stateless services (micro services)

Later on also for stateful services (databases, in memory, ...)

Jobs existed as a concept but not as we would expect them

Immediate popularity and traction

For traditional

Later on also

Jobs existed :

```
apiVersion: batch/v1
kind: Job
metadata:
  name: job-backoff-limit-per-index-example
spec:
  completions: 10
  parallelism: 3
  completionMode: Indexed # required for the feature
  backoffLimitPerIndex: 1 # maximal number of failures per index
  maxFailedIndexes: 5 # maximal number of failed indexes before terminating the Job execution
  template:
    spec:
      restartPolicy: Never # required for the feature
      containers:
      - name: example
        image: python
        command:
          # The jobs fails as there is at least one failed index
          # (all even indexes fail in here), yet all indexes
          # are executed as maxFailedIndexes is not exceeded.
          - python3
          - -c
          - |
            import os, sys
            print("Hello world")
            if int(os.environ.get("JOB_COMPLETION_INDEX")) % 2 == 0:
              sys.exit(1)
```

Immediate popularity and traction

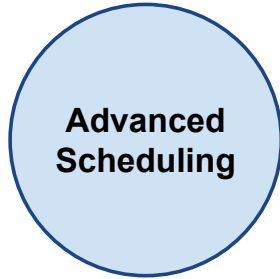
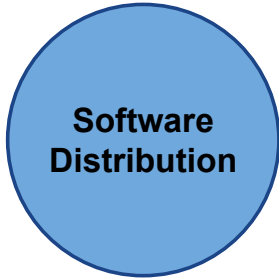
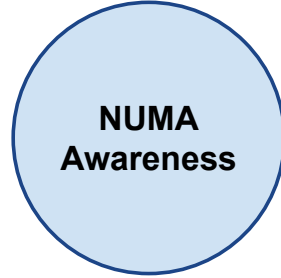
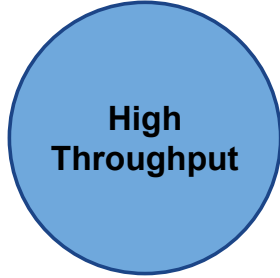
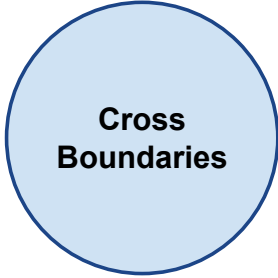
For traditional, stateless services (micro services)

Later on also for stateful services (databases, memory, ...)

Jobs existed as a concept but not as we would expect them

And lacking concepts for queues, quotas, gang scheduling, fair sharing, ...

Meet GenAI!



Where are we today?

Next Generation Triggers

Our research



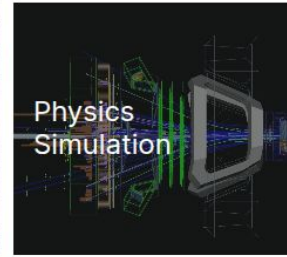
ATLAS and CMS experts in NextGen develop new workflows and data processing techniques to increase the sensitivity of future particle physics triggers



NextGen develops and benchmarks software and techniques to exploit accelerated computing architectures keeping in mind cost and energy efficiency



NextGen investigates the use of AI technologies to improve the experiments' physics impact leveraging massive data throughput pipelines and large-scale models



Theoretical physicists and software engineers in NextGen look at new event generators techniques to improve simulation and detection of exotic signatures



NextGen works with the High-Energy Physics community and computer science experts to provide new generations of scientists with the right skills for tomorrow's challenges

<https://nextgentriggers.web.cern.ch/>

Description	Nodes	Resources	Cores	RAM	Network
Nvidia H100 188GB NVL NVLink	12	8x GPUs	192	1.6TB	100G
Nvidia H100 80GB SXM w/IB	6	4x GPUs	64	1.6TB	100G + 4x400G IB
Nvidia L40S 48GB w/ RoCEv2	7	4x GPUs	128	768GB	200G
AMD MI300X HBM3 192GB w/RoCEv2	2	8x GPUs	128	768GB	200G
AMD Radeon PRO W7900 48GB w/RoCEv2	6	4x GPUs	128	768GB	200G

Pos	Device	Measured Power	Computed Power
36	P6A01	5.27 / 24.58 kW	22.58 / 24.58 kW
35	P6A02	7.99 / 24.58 kW	24.68 / 24.58 kW
34	P6A03	7.70 / 24.58 kW	25.92 / 24.58 kW
33	P6A04	7.58 / 24.58 kW	24.68 / 24.58 kW
32	P6A05	7.20 / 24.58 kW	23.88 / 24.58 kW



ContainerSSH
Launch containers on demand



Model Management

SSH

Notebook

VSCode

kubectl

GitLab CI

GitHub Actions

Interactive

MPI

Training



Serving

Shared Resources



AI/HPC workloads on Kubernetes

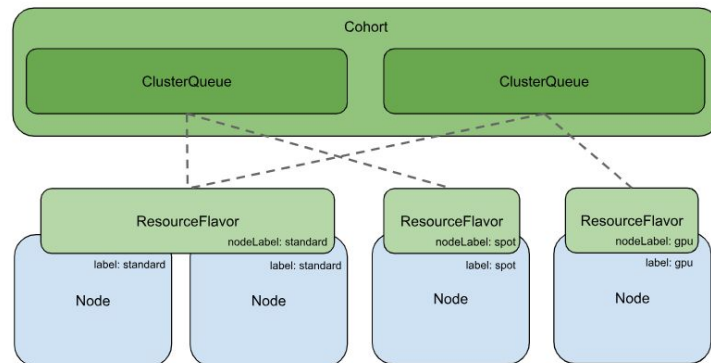
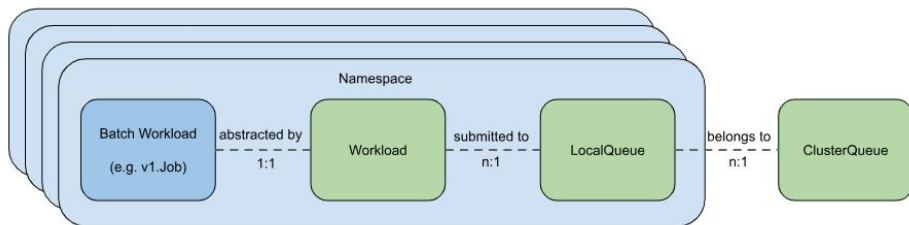
CPU Pinning, NUMA awareness

```
0[| 8.7%] 8[100.0%] 16[100.0%] 24[100.0%] 32[| 2.0%] 40[ 0.0%] 48[ 0.0%] 56[ 0.0%] 64[| 2.7%] 72[100.0%] 80[100.0%] 88[100.0%] 96[| 0.7%] 104[| 2.0%] 112[| 1.3%] 120[| 0.7%]
1[100.0%] 9[100.0%] 17[100.0%] 25[100.0%] 33[| 0.7%] 41[ 0.0%] 49[| 0.7%] 57[| 2.0%] 65[100.0%] 73[100.0%] 81[100.0%] 89[100.0%] 97[| 0.0%] 105[| 0.7%] 113[| 0.0%] 121[| 0.7%]
2[100.0%] 10[100.0%] 18[100.0%] 26[100.0%] 34[ 0.0%] 42[| 0.7%] 50[| 0.7%] 58[| 0.7%] 66[100.0%] 74[100.0%] 82[100.0%] 90[100.0%] 98[| 1.3%] 106[| 2.0%] 114[| 0.7%] 122[| 0.7%]
3[100.0%] 11[100.0%] 19[100.0%] 27[100.0%] 35[| 0.7%] 43[ 0.0%] 51[| 1.3%] 59[| 0.7%] 67[100.0%] 75[100.0%] 83[100.0%] 91[100.0%] 99[| 1.3%] 107[| 0.7%] 115[| 0.0%] 123[| 0.7%]
4[100.0%] 12[100.0%] 20[100.0%] 28[100.0%] 36[ 0.0%] 44[| 0.7%] 52[| 2.0%] 60[ 0.0%] 68[100.0%] 76[100.0%] 84[100.0%] 92[100.0%] 100[| 0.7%] 108[| 0.7%] 116[| 0.0%] 124[ 0.0%]
5[100.0%] 13[100.0%] 21[100.0%] 29[100.0%] 37[ 0.0%] 45[ 0.0%] 53[ 0.0%] 61[| 2.0%] 69[100.0%] 77[100.0%] 85[100.0%] 93[100.0%] 101[| 0.7%] 109[| 1.3%] 117[| 4.6%] 125[| 2.0%]
6[100.0%] 14[100.0%] 22[100.0%] 30[100.0%] 38[| 0.7%] 46[| 1.3%] 54[| 1.3%] 62[| 1.4%] 70[100.0%] 78[100.0%] 86[100.0%] 94[100.0%] 102[| 0.7%] 110[ 0.0%] 118[ 0.0%] 126[| 0.7%]
7[100.0%] 15[100.0%] 23[100.0%] 31[| 4.2%] 39[ 0.0%] 47[| 0.7%] 55[| 2.0%] 63[| 2.0%] 71[100.0%] 79[100.0%] 87[100.0%] 95[| 2.8%] 103[| 0.7%] 111[| 0.7%] 119[ 0.0%] 127[| 2.0%]
Mem[|||||] 23.2G/1008G Tasks: 134, 0 thr, 0 kthr; 0 running
Swp[ 0K/0K] Load average: 21.02 6.15 2.93
Uptime: 62 days, 22:17:35
```

AI/HPC workloads on Kubernetes

CPU Pinning, NUMA awareness

Kueue as a built-in component for advanced scheduling primitives

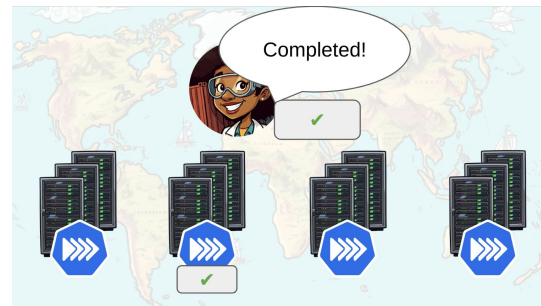
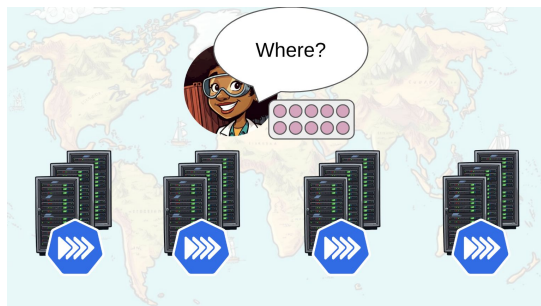


AI/HPC workloads on Kubernetes

CPU Pinning, NUMA awareness

Kueue as a built-in component for advanced scheduling primitives

Support for multi-cluster, multiple administrative domains



AI/HPC workloads on Kubernetes

CPU Pinning, NUMA awareness

Kueue as a built-in component for advanced scheduling primitives

Support for multi-cluster, multiple administrative domains

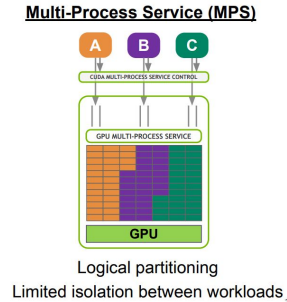
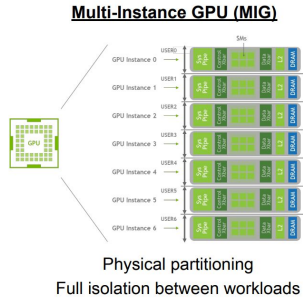
Support for Infiniband, RoCEv2

What's coming next?

Dynamic Resource Allocation

Previous to DRA resources are made available in a static way

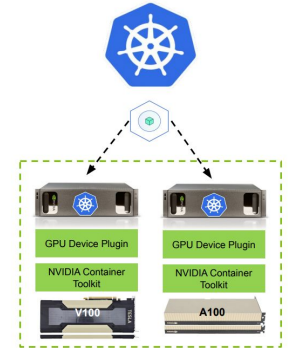
Strong demand for a more dynamic resource allocation (partitioning, sharing)



...

```
apiVersion: v1
kind: Pod
metadata:
  name: gpu-example
spec:
  containers:
    - name: gpu-example
      image: nvidia/cuda
      resources:
        limits:
          nvidia.com/gpu: 1
```

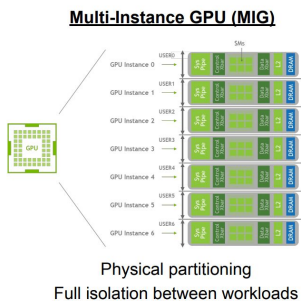
App GPU



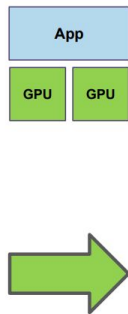
Dynamic Resource Allocation

Resources are made available in a static way

Strong demand for a more dynamic resource allocation (partitioning, sharing)



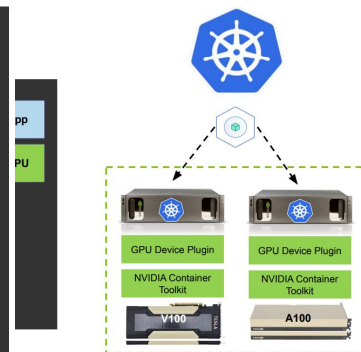
```
apiVersion: v1
kind: Pod
metadata:
  name: gpu-example
spec:
  containers:
  - name: ctr
    image: nvidia/cuda
    resources:
      limits:
        nvidia.com/gpu: 2
```



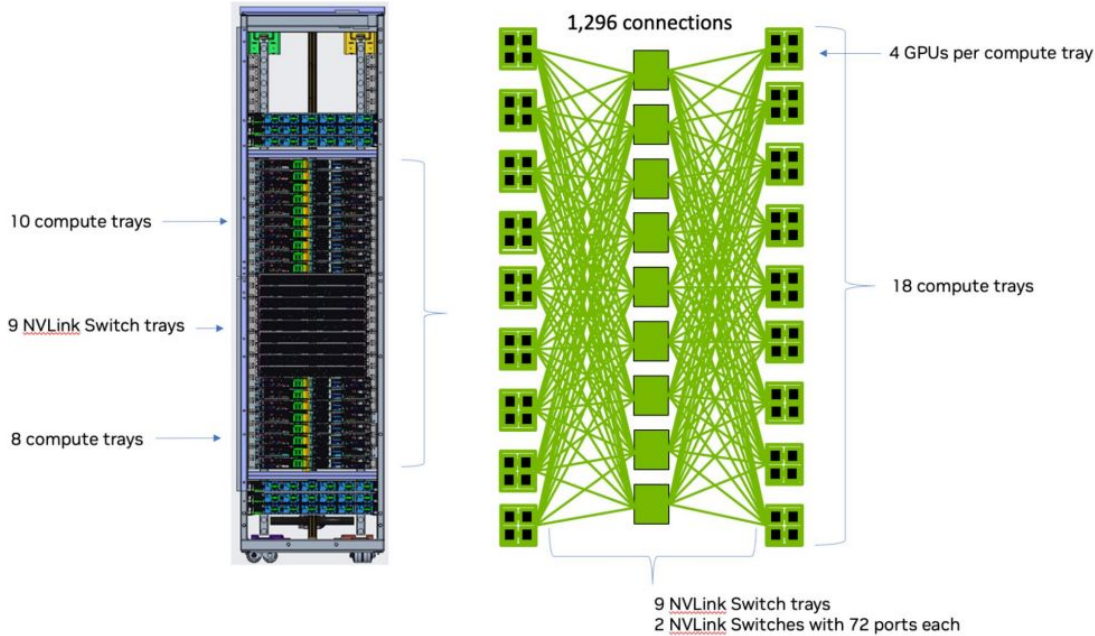
```
apiVersion: resource.k8s.io/v1alpha3
kind: ResourceClaimTemplate
metadata:
  name: unique-gpu
spec:
  devices:
    requests:
      - name: gpu
        deviceClassName: gpu.nvidia.com
---
apiVersion: v1
kind: Pod
metadata:
  name: gpu-example
spec:
  containers:
  - name: ctr
    resources:
      claims:
        - name: gpu0
        - name: gpu1
  resourceClaims:
  - name: gpu0
    resourceClaimTemplateName: unique-gpu
  - name: gpu1
    resourceClaimTemplateName: unique-gpu
```

Defines the "template" for creating a ResourceClaim

Associated with the DRA Driver and installed by the cluster admin



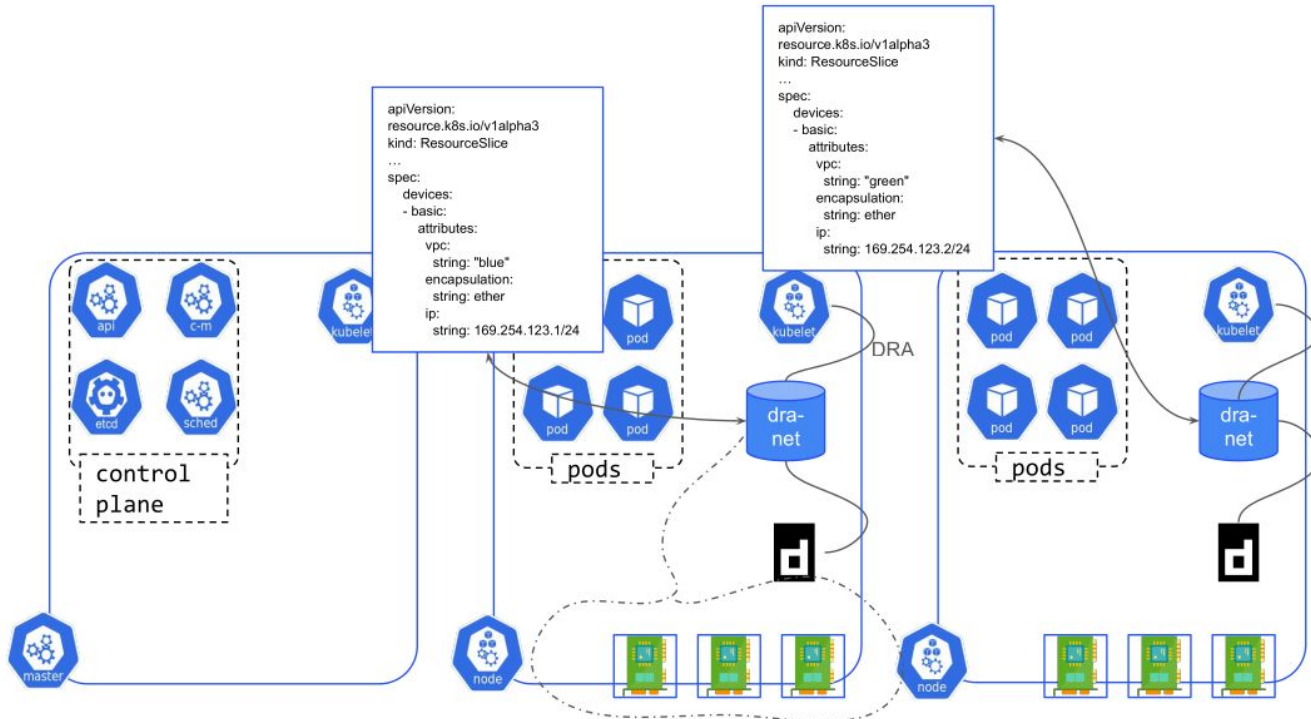
NVIDIA GB200 NVL72



```
apiVersion: resource.nvidia.com/v1beta1
kind: ComputeDomain
metadata:
  name: compute-domain
spec:
  numNodes: 18
  channel:
    resourceClaimTemplate:
      name: compute-domain-channel
```

```
apiVersion: apps/v1
kind: Deployment
metadata:
  name: test-workload
spec:
  replicas: 18
  template:
    spec:
      containers:
        resources:
          limits:
            nvidia.com/gpu: 4
          claims:
            - name: channel
      resourceClaims:
        - name: channel
          resourceClaimTemplateName: compute-domain-channel
```

Dynamic Resource Allocation... networking



Q & A