# 2022 GTC

# GTC 2022

**All Sessions (955)**   Mon, Mar 21 (64)   Tue, Mar 22 (163)   Wed, Mar 23 (299)   Thu, Mar 24 (200)

| Category | session count |
|---|---|
| Accelerated computing & Dev Tools | 78 |
| **Autonomous Machine (Robotics)** | 50 |
| Autonomous Vehicles | 38 |
| Computer Vision | 55 |
| Conversation AI/NLP | 47 |
| Data Center | 81 |
| Graphics | 48 |
| Healthcare | 69 |
| HPC | 84 |
| IoT/5G/Edge | 31 |

☑ Graphics - AI Applications, Art

☑ Graphics - Animation / VFX / Virtual Production

☑ Graphics - Production Rendering and Ray Tracing
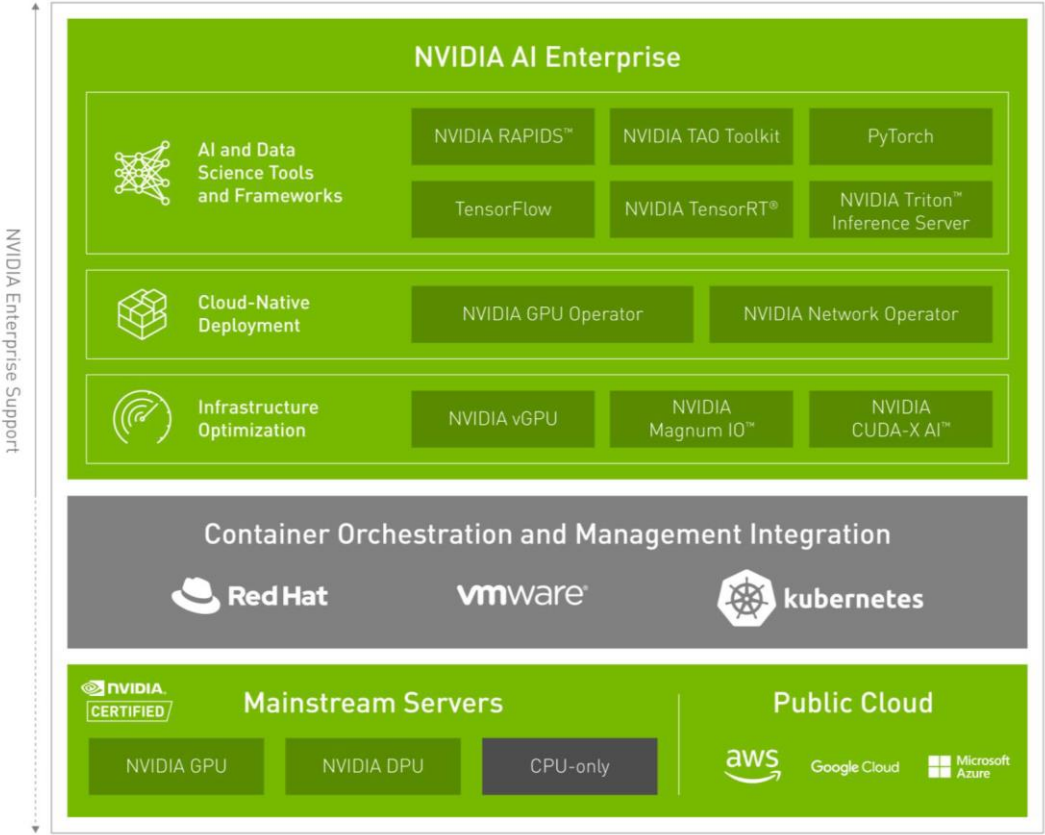
☑ Graphics - Real-Time Rendering and Ray Tracing

☑ Healthcare – Drug Discovery, Genomics

☑ Healthcare – Medical imaging

☑ Healthcare – Smart Hospitals & Instruments

☑ HPC - Astronomy / Astrophysics

☑ HPC - Climate / Weather / Ocean Modeling

☑ HPC - Computational Chemistry and Materials Science

☑ HPC - Computational Fluid Dynamics

☑ HPC - Computational Physics

☑ HPC - Quantum Computing

☑ HPC - Scientific Visualization
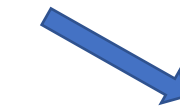
☑ HPC - Supercomputing

# New Products (H/W + S/W)

- 1、**H100 GPU**：采用台积电 4N 工艺，拥有 800 亿个晶体管，实现了首个 GPU 机密计算，相比 A100，FP8 性能提升 6 倍，FP16、TF32、FP64 性能各提升 3 倍。

- 2、**全新 NVLink Switch 系统**：高度可扩展，支持 256 块 H100 GPU 互联。

- 3、**融合加速器 H100 CNX**：耦合 H100 GPU 与 ConnectX-7 和以太网智能网卡，可为 I／O 密集型应用提供更强劲的性能。

- 4、**DGX H100**：配备 8 块 H100 GPU，总计有 6400 亿个晶体管，在全新的 FP8 精度下 AI 性能比上一代高 6 倍，可提供 900GB／s 的带宽。

- 5、**DGX SuperPOD**：最多由 32 个 DGX H100 组成，AI 算力可达 1EFLOPS。

- 6、**Eos 超级计算机**：全球运行速度最快的 AI 超级计算机，配备 576 台 DGX H100 系统，FP8 算力达到 18EFLOPS，FP64 算力达到 275PFLOPS。

- 7、**Grace CPU 超级芯片**：由两个 CPU 芯片组成，采用最新 Armv9 架构，拥有 144 个 CPU 核心和 1TB／s 的内存带宽，将于 2023 年上半年供货。

- 8、**为定制芯片集成开放 NVLink**：采用先进封装技术，与英伟达芯片上的 PCIe Gen 5 相比，能源效率高 25 倍，面积效率高 90 倍。英伟达还将支持通用小芯片互连传输通道 UCIe 标准。

- 9、**CUDA-X**：60 多个针对 CUDA-X 的一系列库、工具和技术的更新。

- 10、**Riva 2.0**：对话式 AI 服务 Riva 全面发行，2.0 版本支持识别 7 种语言，可将神经文本转换为不同性别发声的语音。

- 11、**Merlin 1.0**：可帮助企业快速构建、部署和扩展先进的 AI 推荐系统。

- 12、**Sionna**：一款用于 6G 通信研究的 AI 框架。

- 13、**OVX 与 OVX SuperPod**：面向工业数字孪生的数据中心级服务器和超级集群。

- 14、**Spectrum-4**：全球首个 400Gbps 端到端网络平台，交换吞吐量比前几代产品高出 4 倍，达到 51.2Tbps。

- 15、**Omniverse Cloud**：支持协作者们随时随地实现远程实时协同工作。

- 16、**DRIVE Hyperion 9**：汽车参考设计，拥有 14 个摄像头、9 个雷达、3 个激光雷达和 20 个超声传感器，总体传感器数量是上一代的两倍。

- 17、**DRIVE Map**：多模态地图引擎，包含摄像头、激光雷达和雷达的数据，同时兼顾安全性。

- 18、**Clara Holoscan MGX**：可供医疗设备行业在边缘开发和部署实时 AI 应用的计算平台，AI 算力可达每秒 254~610 万亿次运算。

- 19、**Isaac for AMR**：提供自主移动机器人系统参考设计。
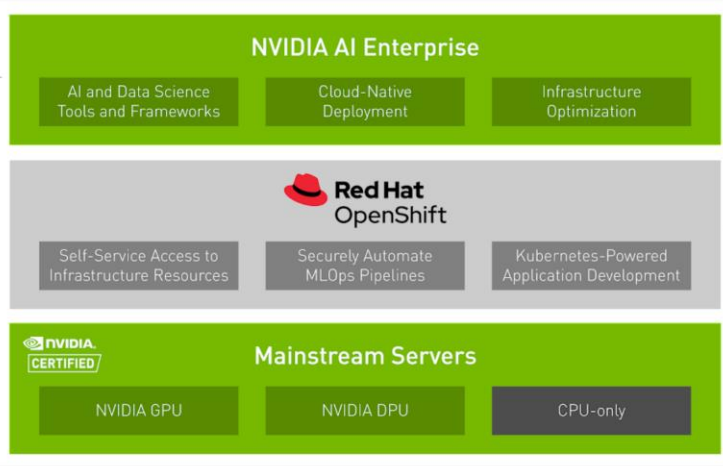
- 20、**Jetson AGX Orin 开发者套件**：在边缘实现服务器级的 AI 性能。
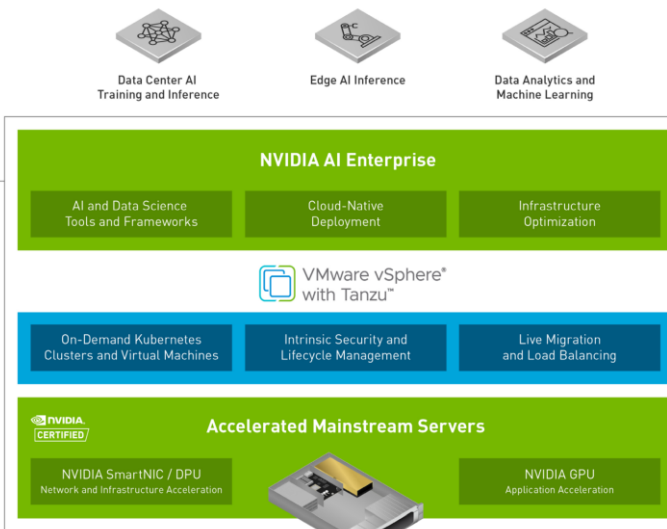
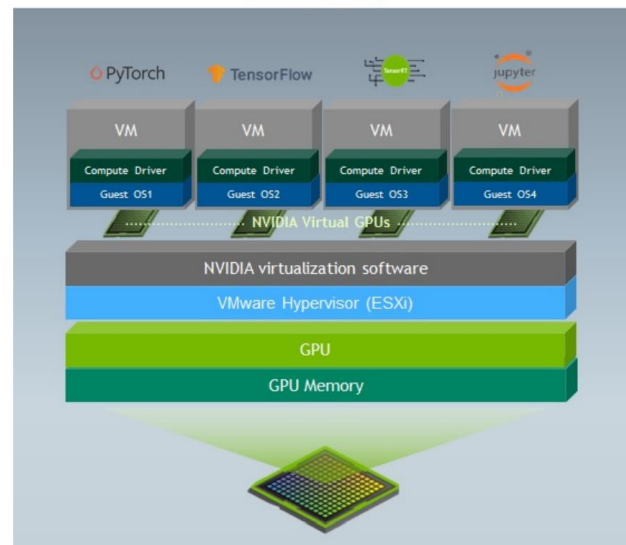# NVIDIA AI ENTERPRISE SOFTWARE SUITE

# S41838: Tuning Virtualized GPUs for Optimal Performance on ML/AI Workloads

## vGPU vs. MIG with vGPU with A100

### vGPU
- Allow up to **10 VMs** with vGPU per GPU
- Memory is equally partitioned among vGPUs
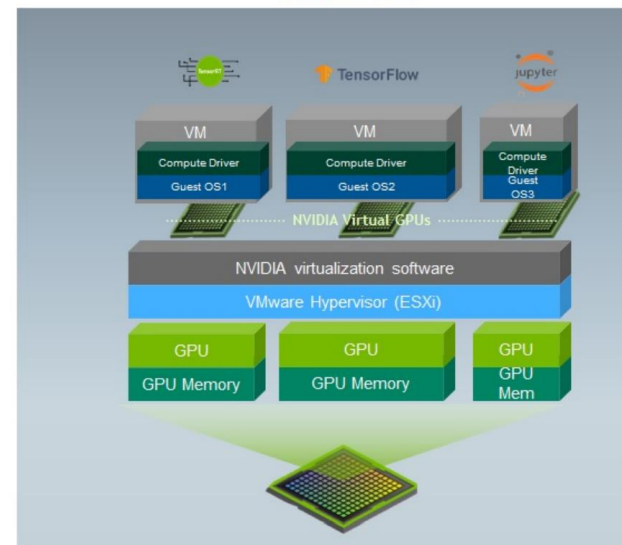- Compute is shared based on time slicing (Best Effort, Fixed Share, Equal Share)

### MIG with vGPU
- Allow up to **7 VMs** per GPU
- Memory & compute are partitioned among vGPUs as below
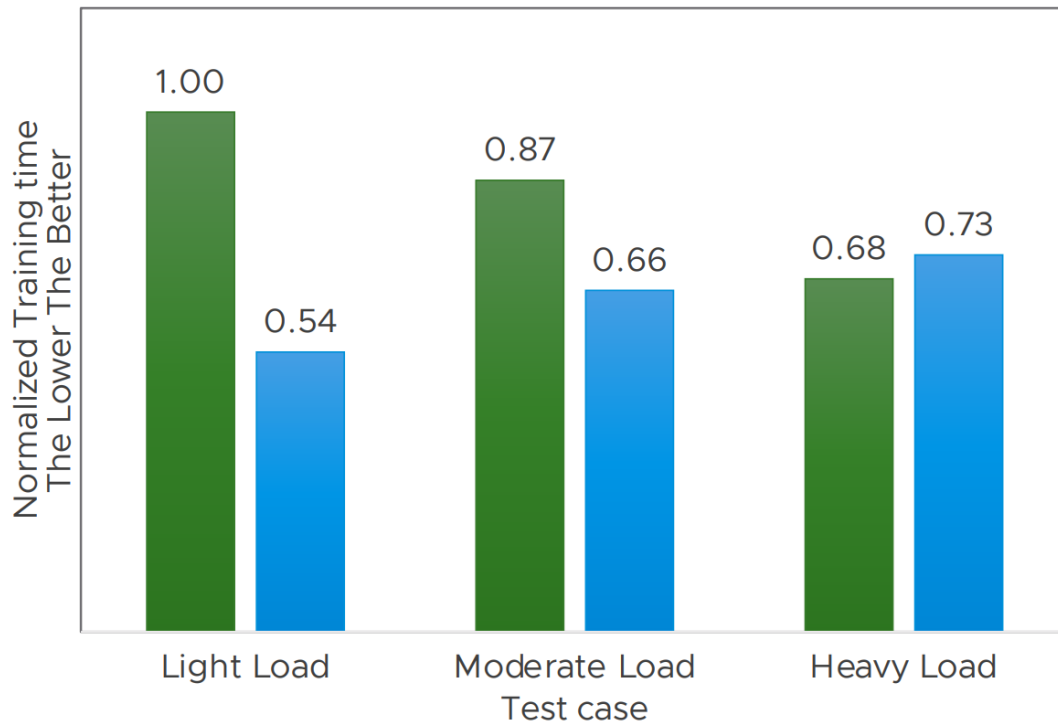
# vGPU vs. MIG with vGPU: Sizing ML Training Workload

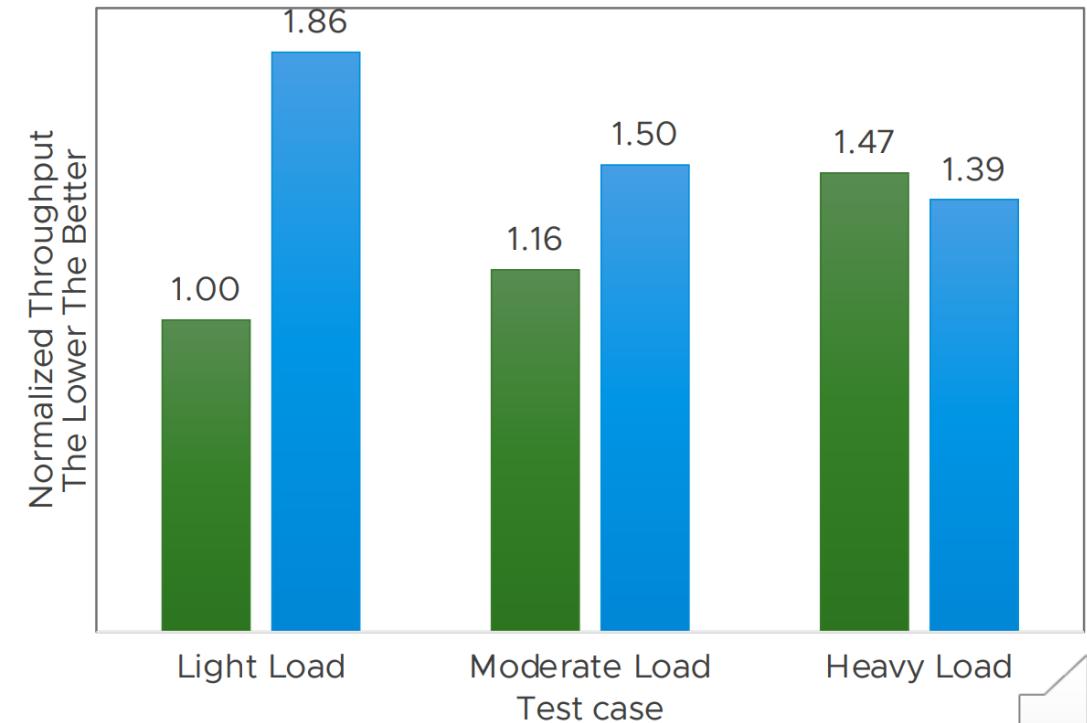- Light & moderate load → MIG better
- Heavy load → vGPU lightly better


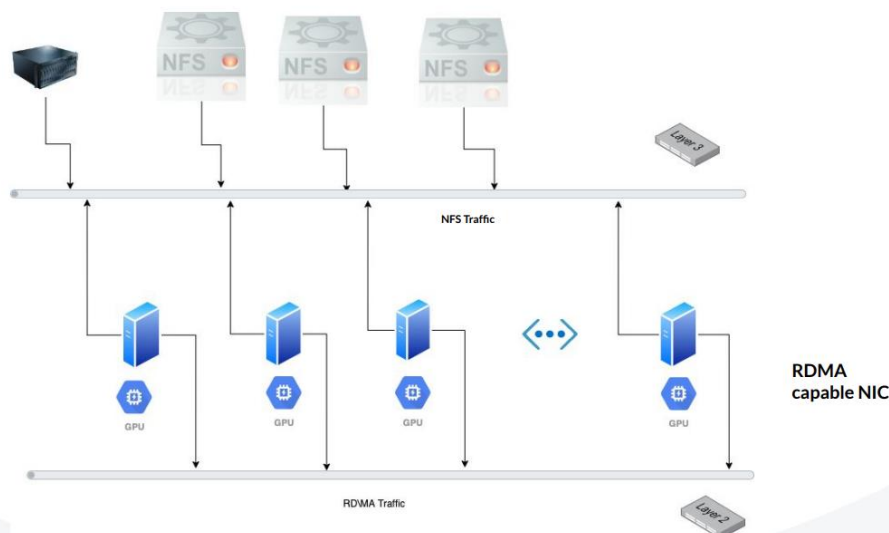
Normalized Training Time

Normalized Combined Throughput

# Scale and Accelerate the Distributed Model Training in Kubernetes Cluster [S42498]

In order to orchestrate deep learning workloads that scale across multiple GPUs and nodes, Kubernetes offers a compelling solution. With Kubernetes and Kubeflow PyTorch, we can easily schedule and track a distributed training job on single-GPU multi-node, multi-GPU single-node, and multi-GPU multi-nodes in a shared GPU resource pool. To accelerate deep learning training at Zoom, we enable RDMA, RoCE to bypass the CPU kernel and offload the TCP/IP protocol. We apply this technology in Kubernetes with SRIOV by NVIDIA Network Operator in a heterogenous GPUs cluster with four GPU servers and eight GPU servers. By combining NVIDIA NCCL, Apex, and PyTorch PowerSGD, we can reach a near-linear performance increase as the GPU number and worker node increases.

Jack Jin, Lead Machine Learning Infra Engineer, Zoom Video

Kaixing Wu, Machine Learning Engineer, Zoom

## Networking



RDMA capable NIC

# ☆ Nvidia Network Operator

https://github.com/Mellanox/network-operator/tree/master/deployment/network-operator

- The **NVIDIA Network Operator** loads the required drivers, libraries, device plugins, and CNIs on any cluster node with an network interface.

  - sriovNetworkOperator.
  - Nfd
  - nvPeerDriver
  - secondaryNetwork
  - cniPlugins
  - Multus
  - ipamPlugin

- The NVIDIA Network Operator uses Kubernetes CRD and the Operator Framework to provision the host software needed for enabling accelerated networking.

  Install with helm

## Kubernetes SR-IOV

The **Single Root I/O Virtualization** (SR-IOV) can segment a compliant network device on the host node as a physical function (PF), into multiple virtual functions (VFs).

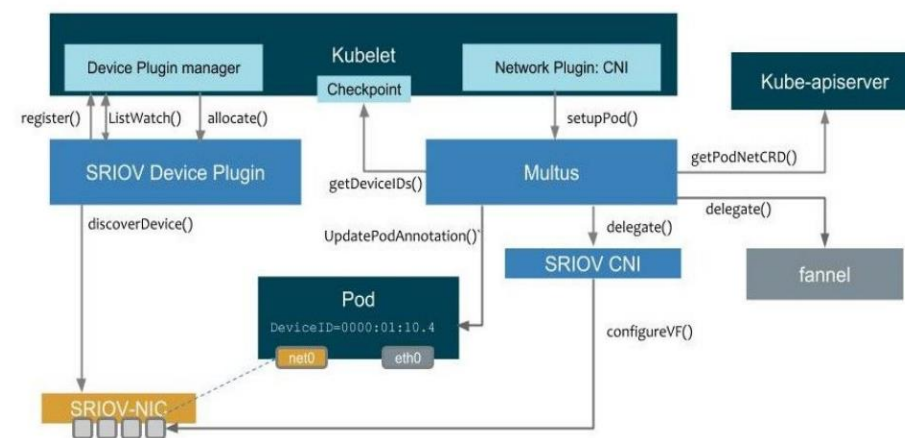Allow to use RDMA-enabled VF as a normal network device.

➢ Network configuration daemon
➢ Network device plug-in
➢ CNI plug-in

Configuration:
- SR-IOV network device.
Creating SR-IOV network node policy with node selector

- SR-IOV Ethernet network attachment.
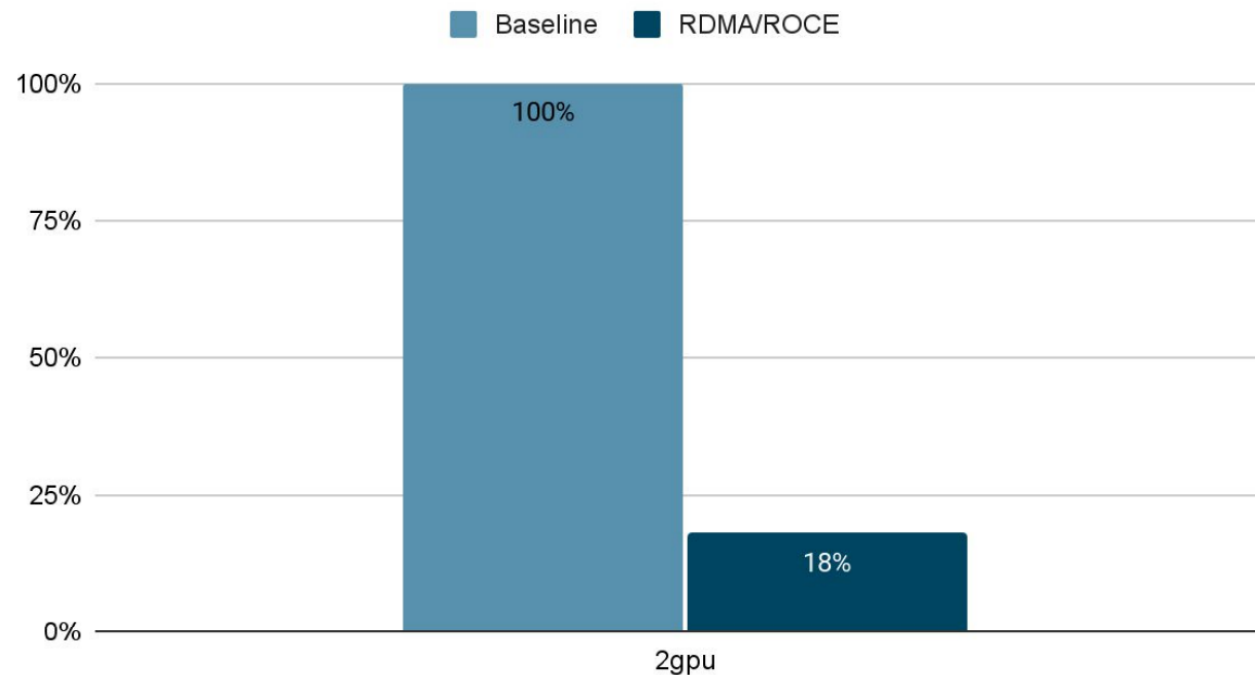Creating ethernet network device by defining an SriovNetwork object.

# Scale and Accelerate the Distributed Model Training in Kubernetes Cluster [S42498]

**Task 2**: Hardware update to enable faster multi-node communication. With RDMA/ROCE enabled, the multi-node 2gpu training is around 5x faster.



Baseline and RDMA/ROCE

# Fluid: build data orchestration in K8s, kubecon 2021

- https://www.youtube.com/watch?v=Pjt8v4GYvRQ



2.Find the cacheable Node

client

1. Create Pod

Kubernetes Scheduler

Fluid Scheduler

3.Order the nodes by the cache capabilities

Fluid Runtime Service

4.Start pod in N1

```
apiVersion: v1
kind: Pod
metadata:
  name: resnet50
spec:
  containers:
      - name: train
        image: resnet50
        volumeMounts:
          - mountPath: /data
            name: imagenet
  volumes:
    - name: imagenet
      persistentVolumeClaim:
        claimName: imagenet
```

Dataset Name

Pod    Alluxio    10G Cached    N1

Alluxio    5G Cached    N2

N3

9