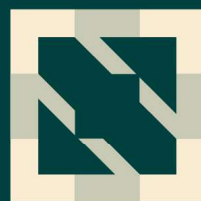


KubeCon



CloudNativeCon

S OPEN SOURCE SUMMIT

China 2023



KubeCon



CloudNativeCon



OPEN SOURCE SUMMIT

China 2023

Practice of Building Large-Scale AI Training Cluster Based on Kubernetes and RoCEv2

Dekui WANG, IEI

Agenda

1. Background and Challenges

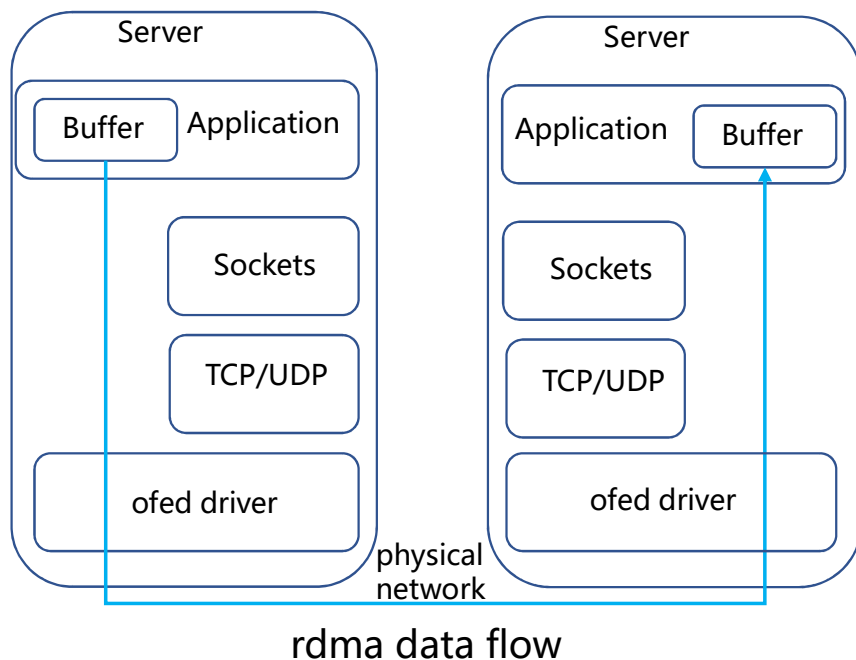
2. RoCEv2 solution

3. Practice and Tests

Network issues in AI training infrastructure

1. large scale training jobs require multiple nodes
2. network congestion problems
3. the difference between Infiniband and RoCE
4. GPU node with multiple GPU cards and multiple rdma cards

RDMA (Remote Direct Memory Access)



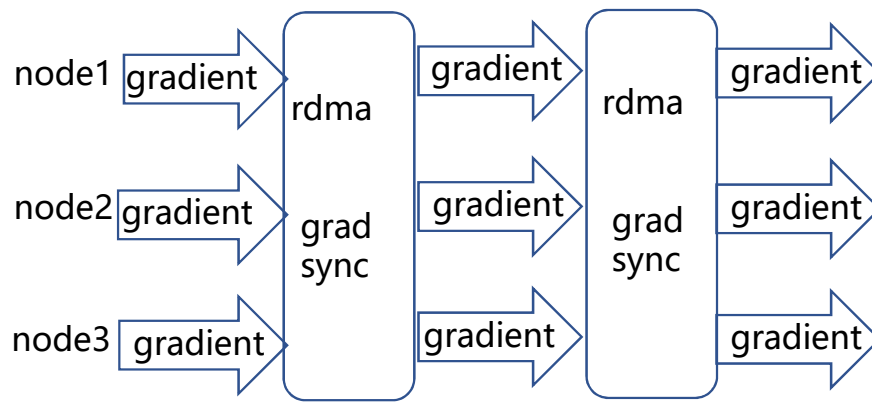
	Infiniband	RoCEv2
End-to-end delay	2us	5us
Flow Control Mechanism	Credit-based flow control mechanism	PFC/ECN, DCQN
Forwarding Mode	Forwarding based on Local ID	IP-based Forwarding
Load Balancing Mode	Packet-by-Packet Adaptive Routing	ECMP Routing
Recovery	Self-Healing Interconnect Enhancement for Intelligent Datacenters	Route Convergence
Network Configuration	Zero configuration through UFM	Manual Configuration

ref: <https://www.naddod.com/blog/infiniband-vs-roce-v2-which-is-best-network-architecture-for-ai-computing-center>

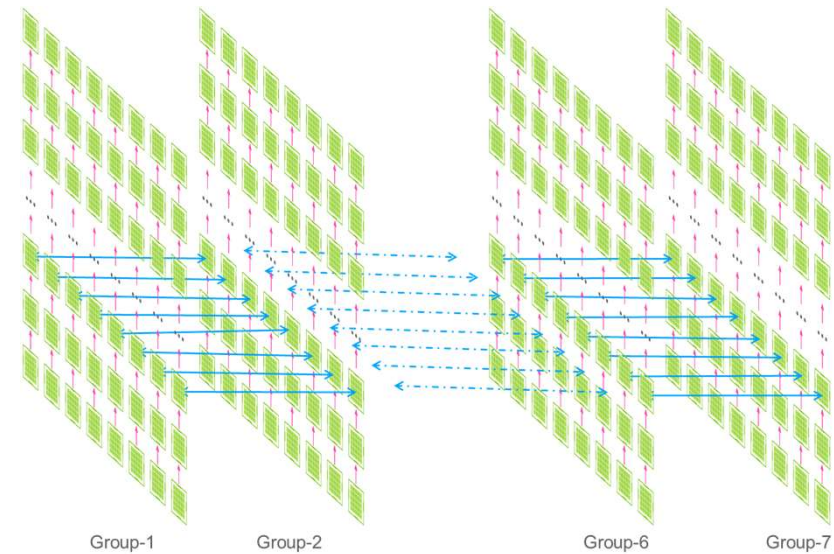
Infiniband VS RoCEv2

- Infiniband, not compatible with existing Ethernet devices, require all Infiniband devices
- RoCE, compatible with existing Ethernet devices. RoCEv2 is implemented based on udp
- iWARP, based on tcp, need many memory resources , implemented based on tcp

RDMA network requirements



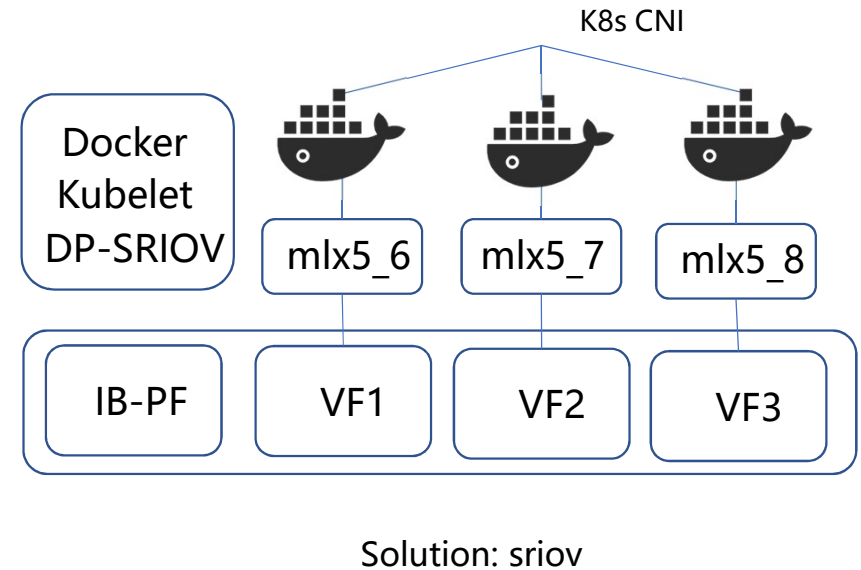
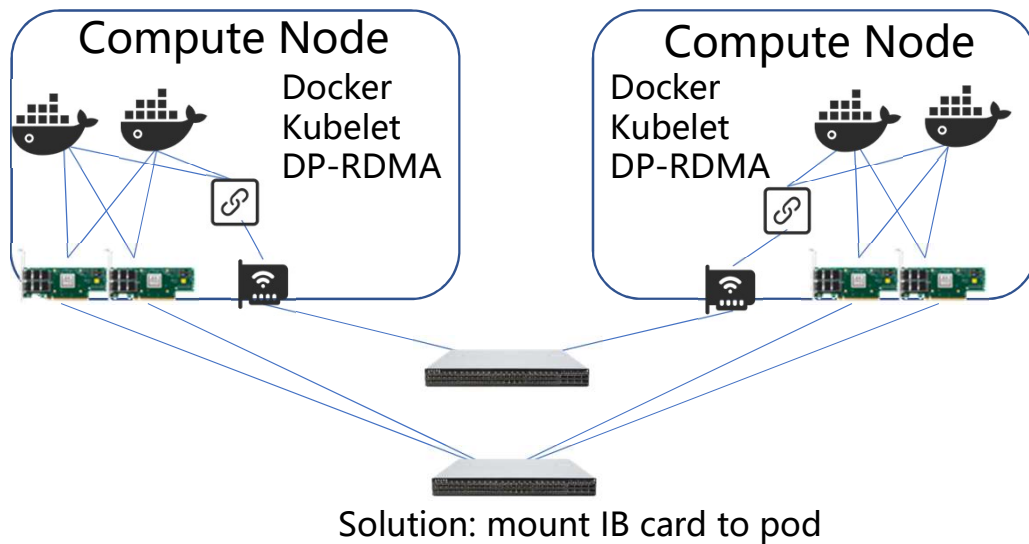
gradient sync based on rdma



3D parallel computing based on rdma

- data parallelism: network between nodes should have high bandwidth and low latency
- model parallelism and pipeline parallelism require rdma network
- gpt-3 with 128 A100 nodes
 - pipeline parallelism bandwidth between nodes is 12GB/s, with 0.1GB data per communication, 0.16s each time
 - data parallelism bandwidth between nodes is 27.4 GB/s, with 44GB data per communication, 32s each time

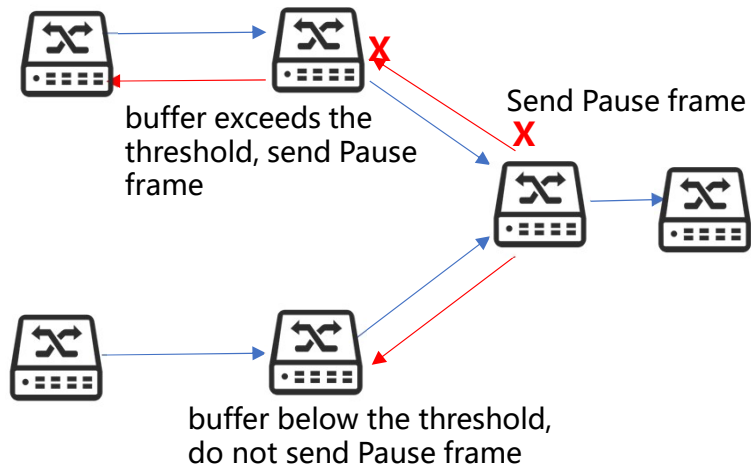
AI cluster with IB



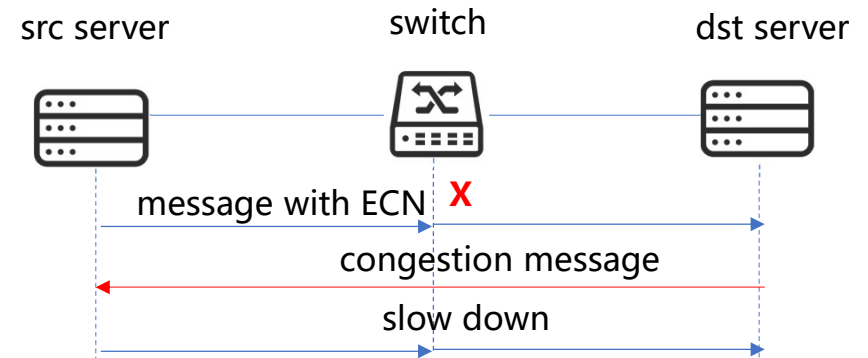
- the metadata exchange problem during RDMA communication
- IB communication between nodes is based on OpenSM, LID, UFM
- clusters with over 10000 nodes

ref: <https://github.com/Mellanox/k8s-rdma-shared-dev-plugin.git>
<https://docs.nvidia.com/networking/pages/releaseview.action?pagelD=18481842>

PFC+ECN



Priority-based Flow Control



Explicit Congestion Notification

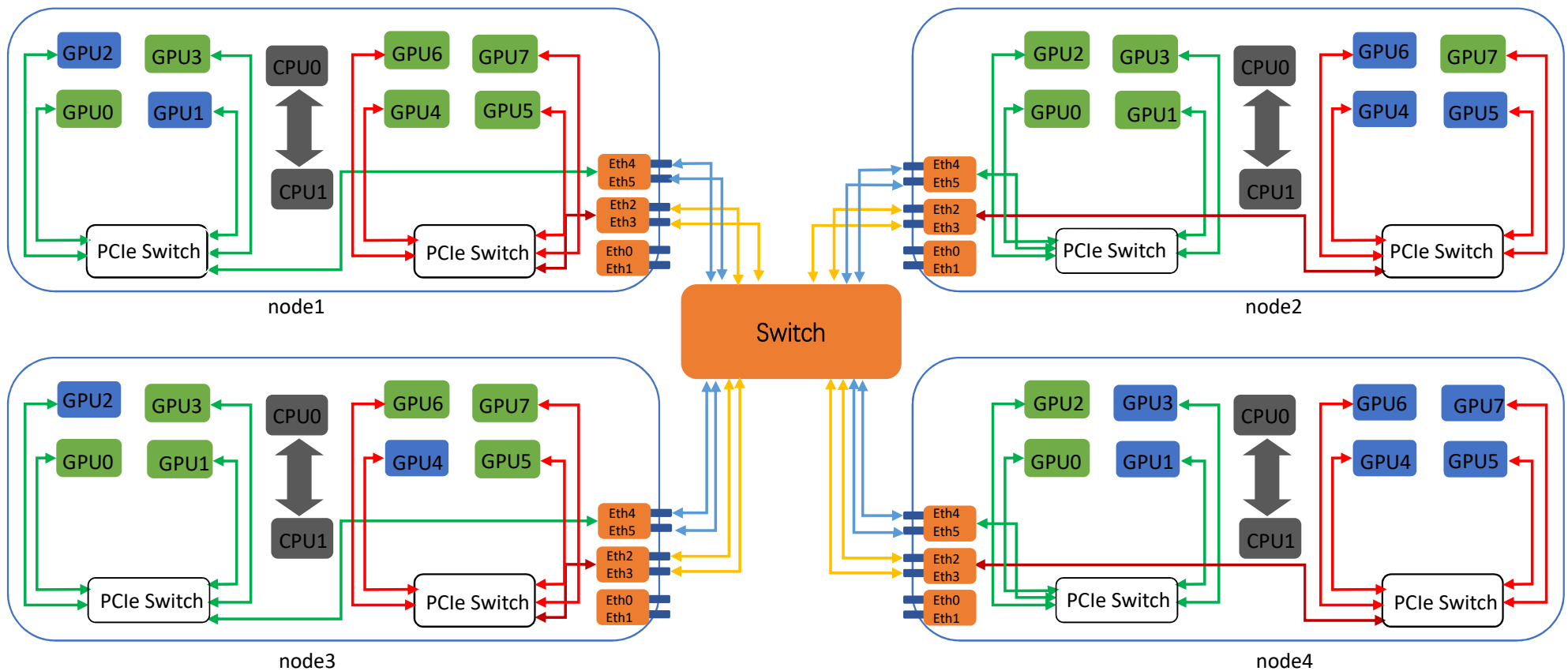
1.Switch Configuration

- PFC: data link layer, based on the packet priority and queue priority
- ECN: network layer, based on the identification bit in the data packet header

2. Host Configuration

- Linux、OFED Driver

GPU fragmentation



- cluster GPU resources are fragmented , idle GPU cards are disordered
- GPU fragmentation affects the used RoCE network card in multi-node training tasks

Agenda

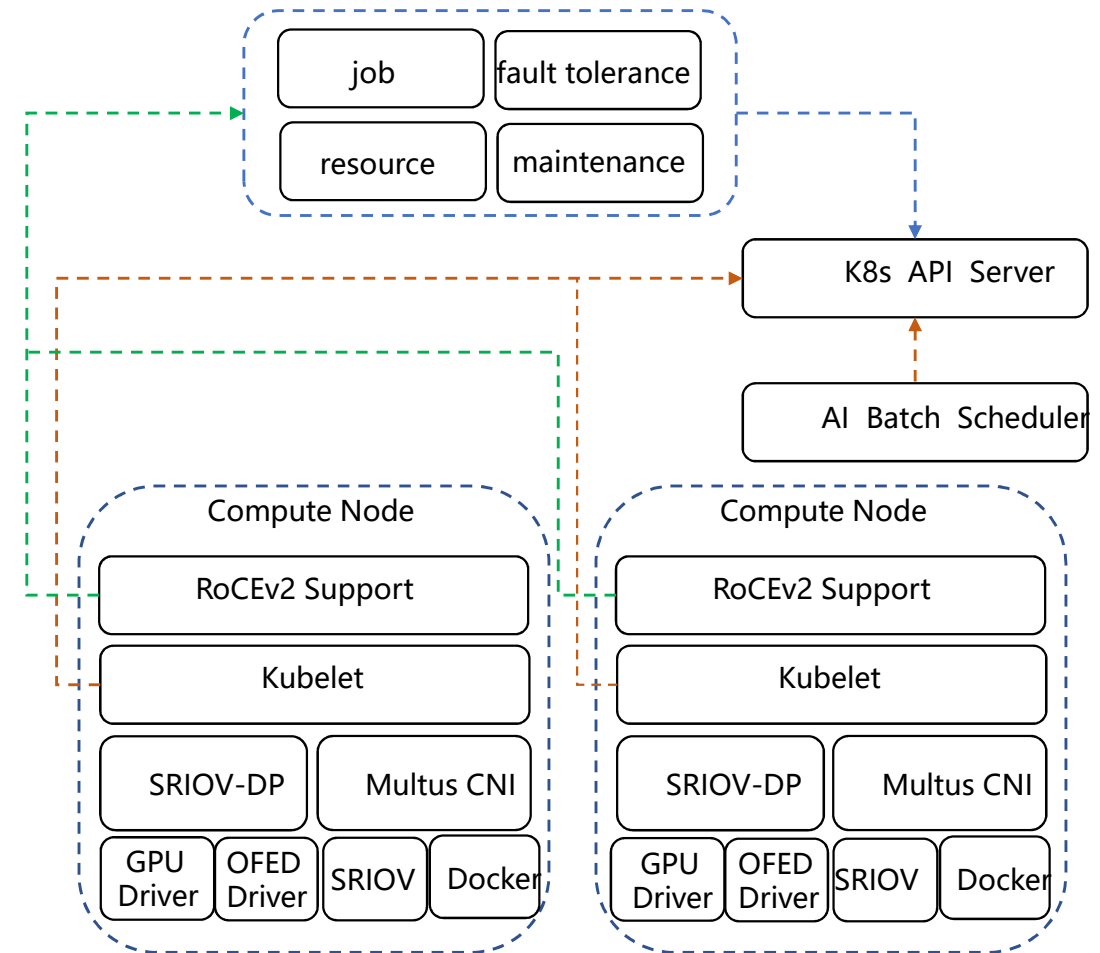
1. Background and Challenges

2. RoCEv2 solution

3. Practice and Tests

Software architecture for RoCEv2

- resource management
 - allocation and configuration for RoCE resources
 - PF/VF network traffic monitor, alarm, job fault tolerance
 - resource scheduling of computing nodes with different network types
- network management
 - business network based on Calico, multiple VF as the computing network
 - cross subnet management , route management
- components
 - sriov-dp and multus-cni , support multiple VF,
 - represent multiple VF network cards with only one K8s resources

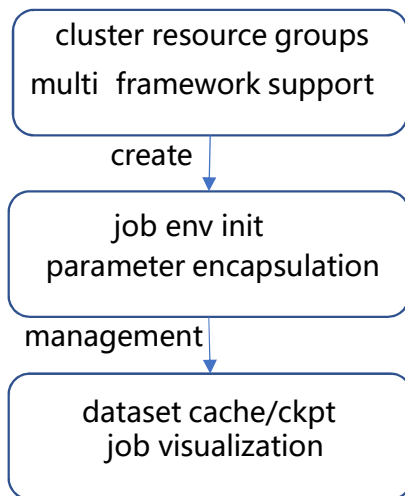


AI job management

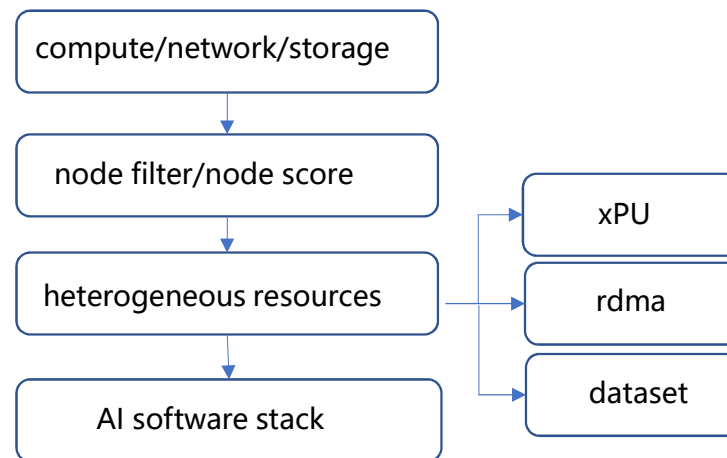
framework



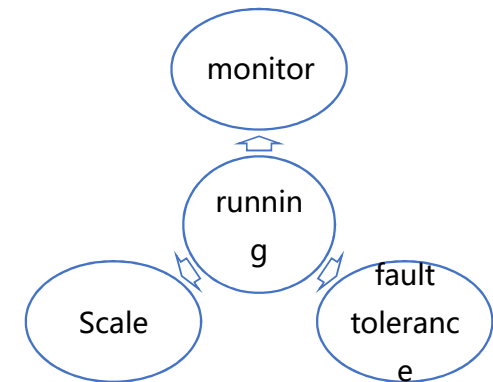
job management



resource scheduling

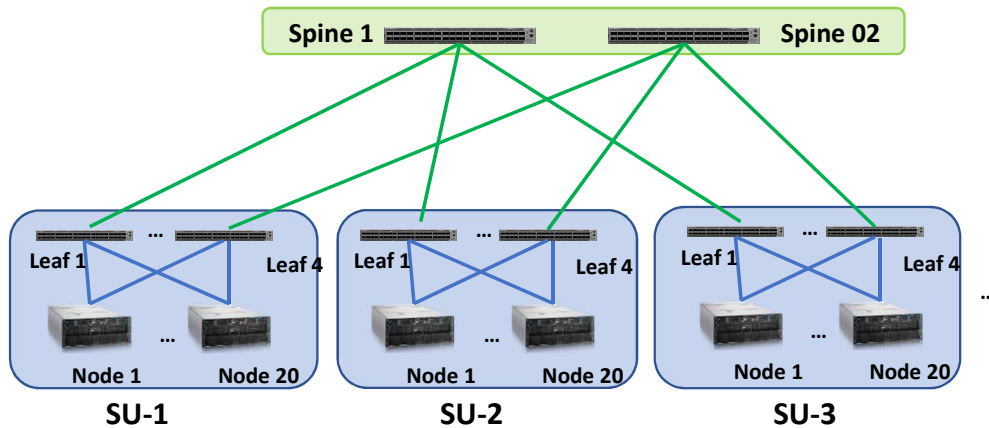


job maintenance

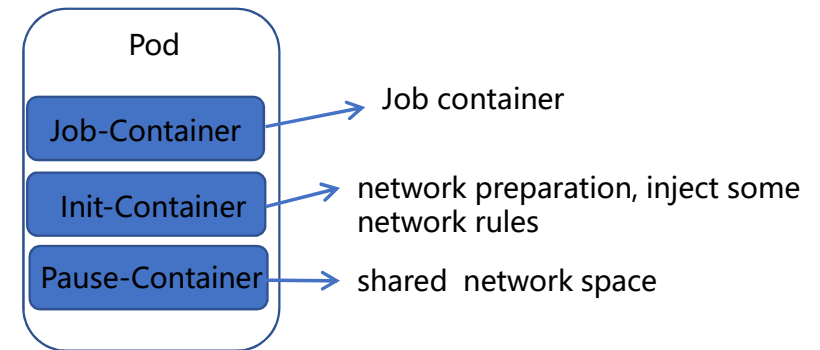


training task management on AI training platform

Physical architecture for RoCEv2



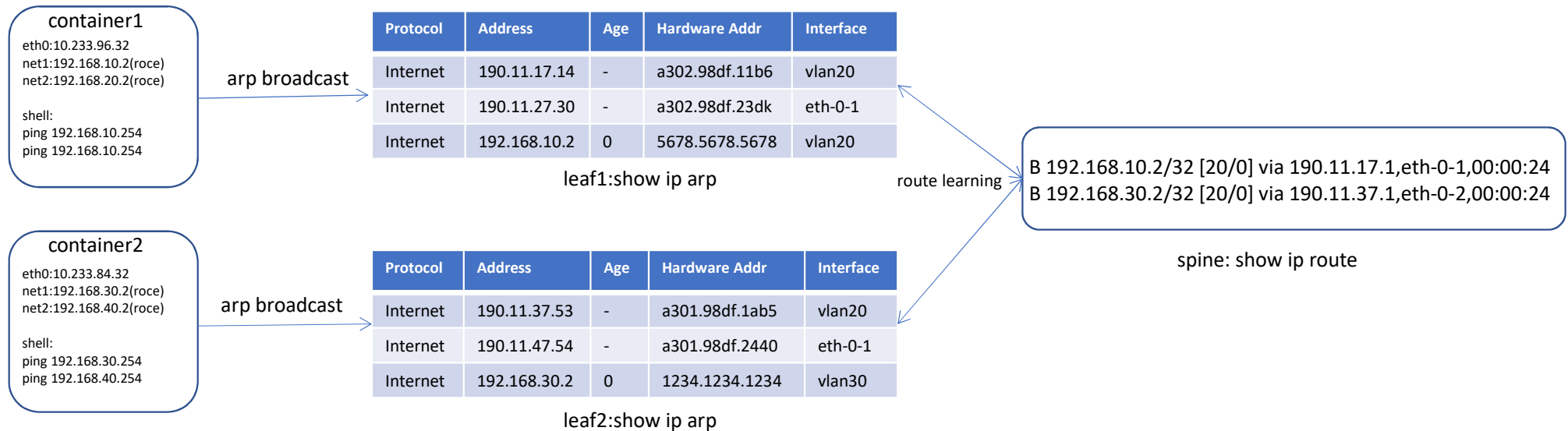
physical network topo



network preparation

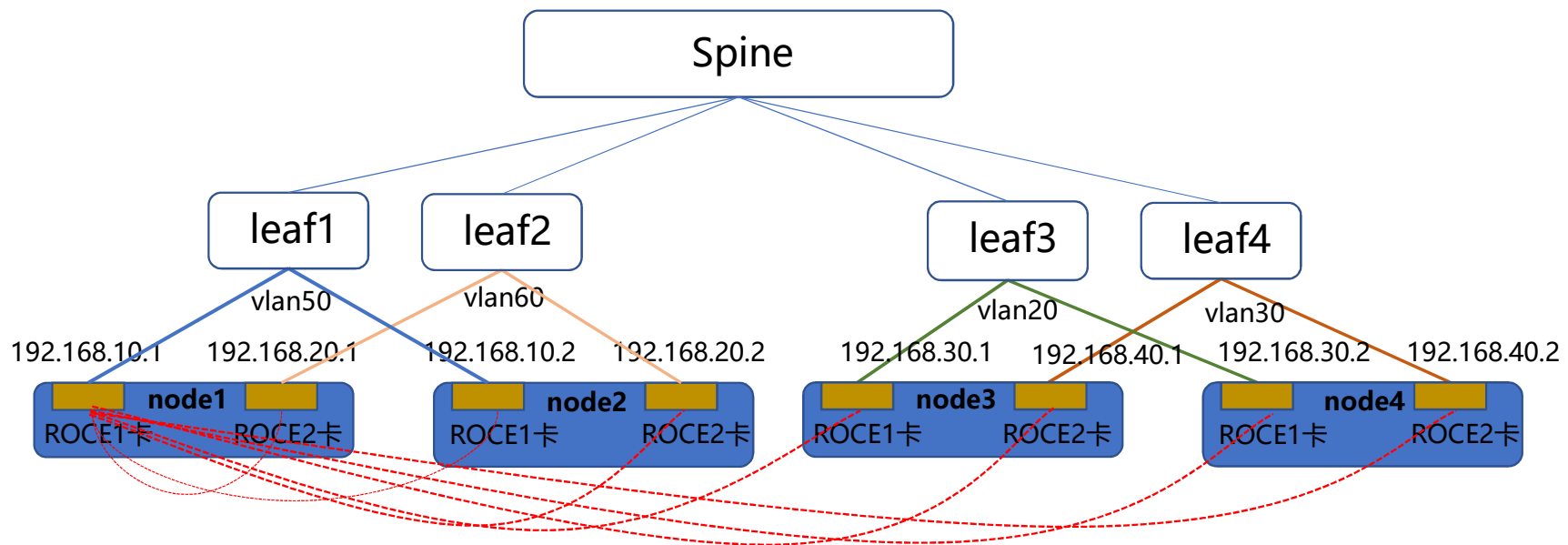
- spine-leaf network, horizontal scaling support for spine switch and leaf switch
- different RoCE cards with different vlan
- configure subnet information at switch, VF using physical subnet and routing

Network simulation in container



- container starts quickly, switches can not update the arp table before training tasks start
- ip reused by another container, but the arp table of the switch was not refreshed immediately
- adjust the aging time of arp table
- adjust the arp table capacity of switch

RoCE nic with P2P communication



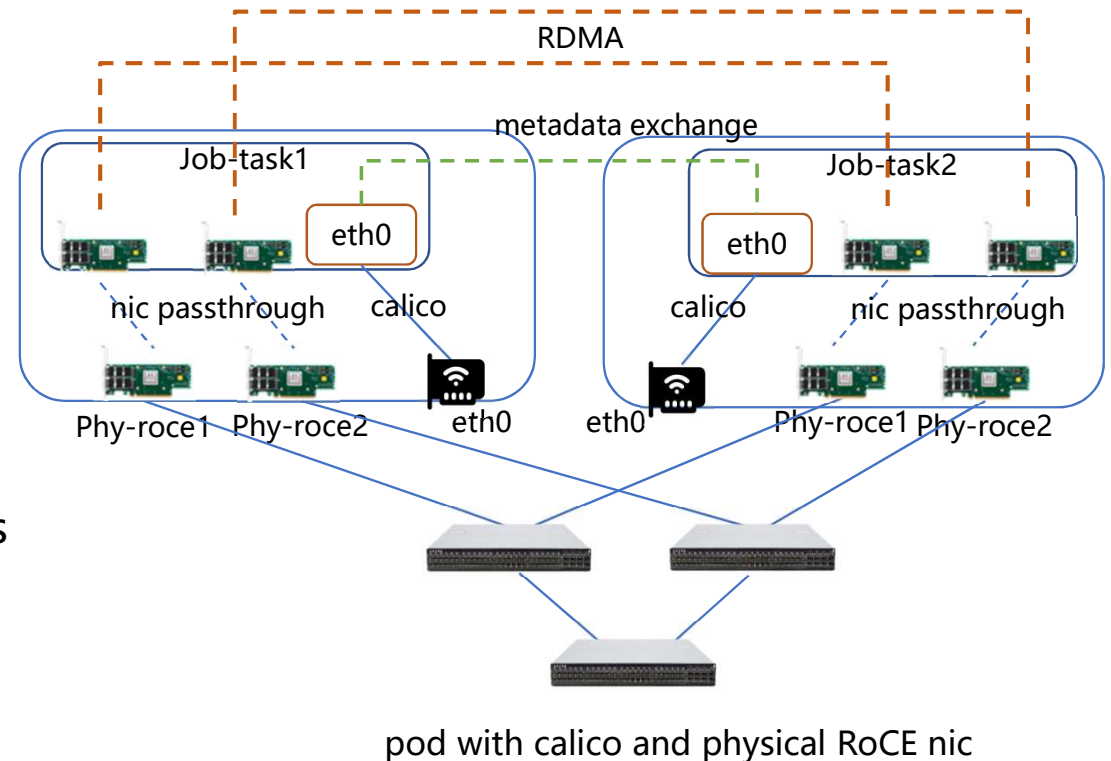
the VF of the RoCE1 network card at node1 can communicate with any VF of any node in the cluster

Other considerations

- GPU p2p exception when sriov virtualization enabled
- network traffic sharing problem between multiple vf
- roce gid index problem when using macvlan
- all vf and pf of the node can be recognized in container
- the maximum number of VFs for RoCE network cards

RoCEv2 for large model training

- large model training scenarios, all GPU of the node will be used by one pod
- calico network for metadata exchange , using physical RoCE nic in pod
- multus cni,sriov-dp support RoCE PF
- large model training jobs use the characteristics of nccl, such as pxn



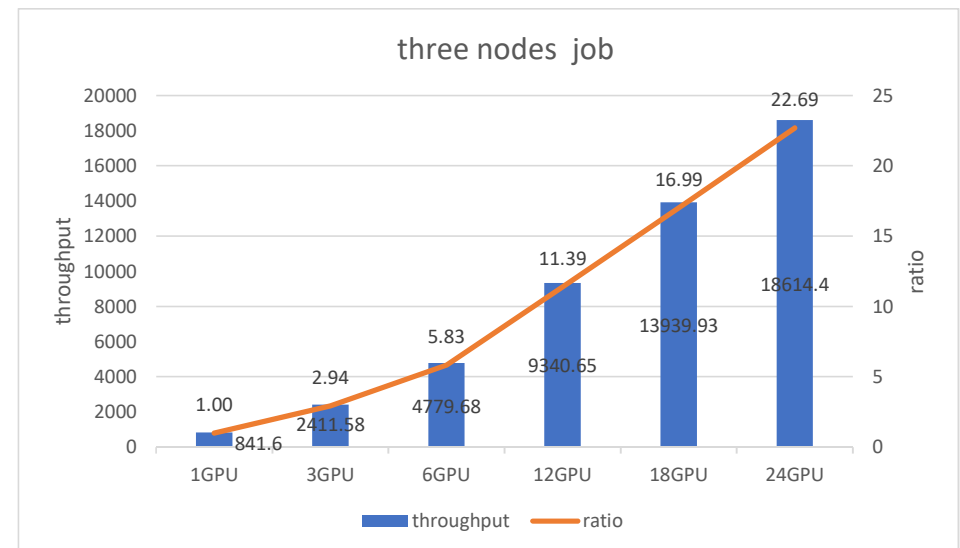
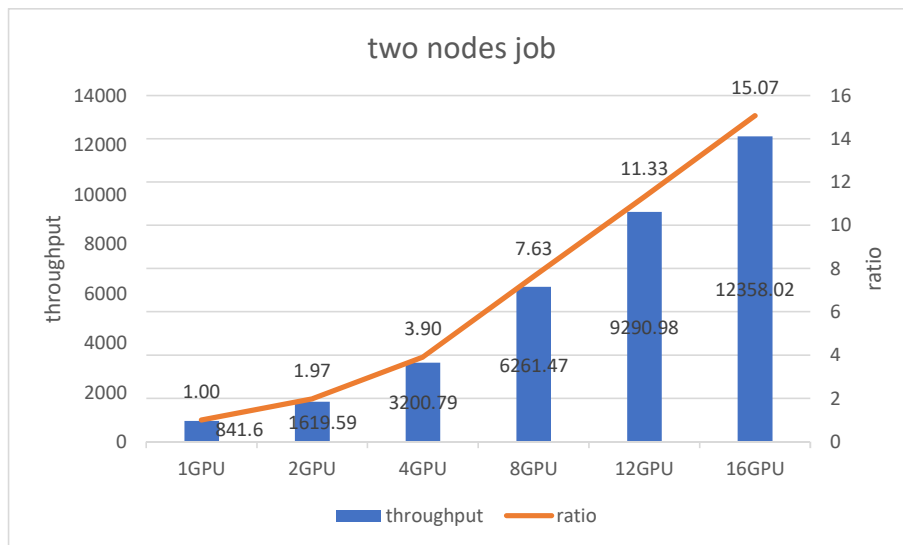
Agenda

1. Background and Challenges

2. RoCEv2 solution

3. Practice and Tests

Job tests



GPU server info:

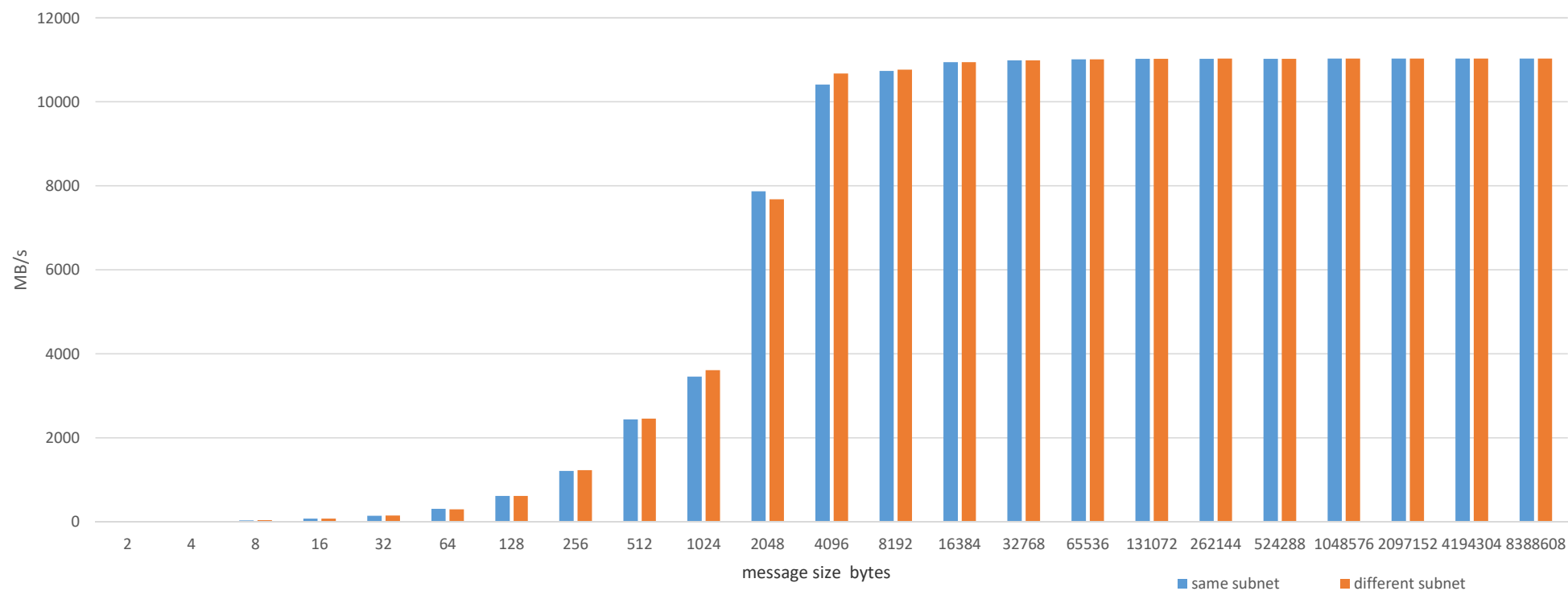
NF5468M5
CPU: Intel(R) Xeon(R) Gold 6230R CPU @ 2.10GHz
GPU: A100-PCIE-40GB
IB: Mellanox Technologies MT27800 Family 100Gb
GPU driver: 450.102.04
IB Driver: 5.4-1.0.3.0

software info:

CUDA: 11.0
NCCL: 2.12.6
Tensorflow: 1.15.3+nv
Tensorflow-cnn-
benchmark,imagenet(synthetic),resnet50,bs=256,iter=500

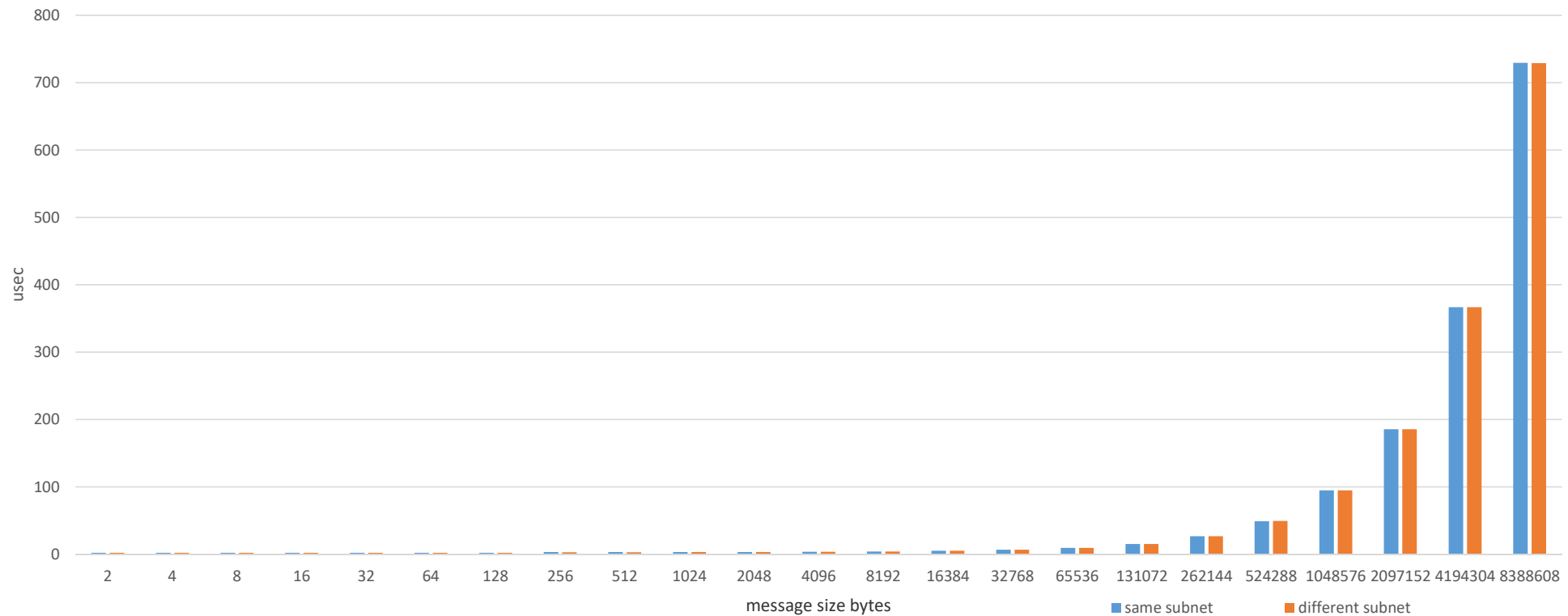
Bandwidth test between containers

configure routing rules in container

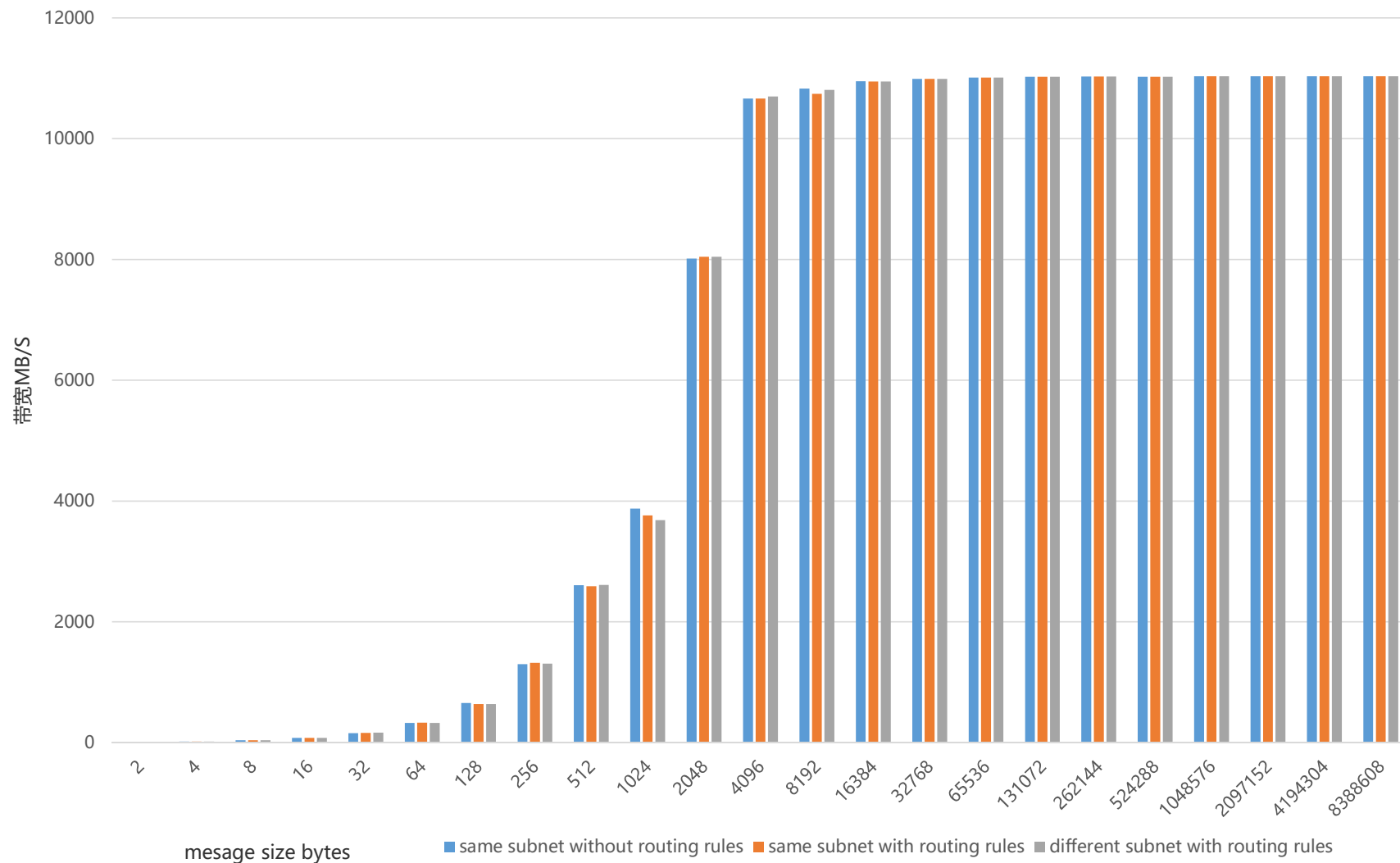


Latency test between containers

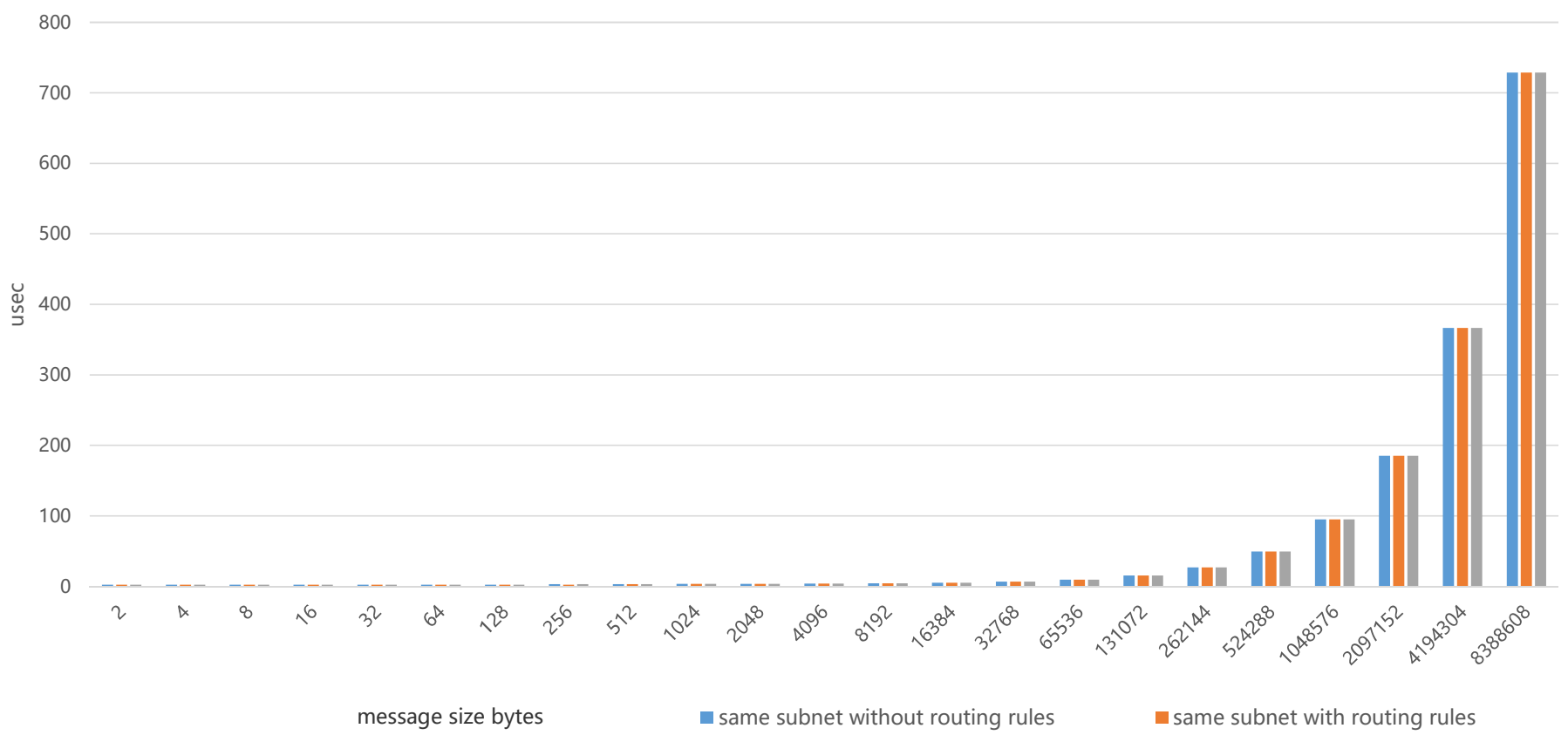
configure routing rules in container



Bandwidth test between hosts



atency test between hosts





KubeCon



CloudNativeCon



OPEN SOURCE SUMMIT

China 2023

AIStation Platform

Chao WANG, IEI

Gen.AI Trend

LLM, Multimodal, ...

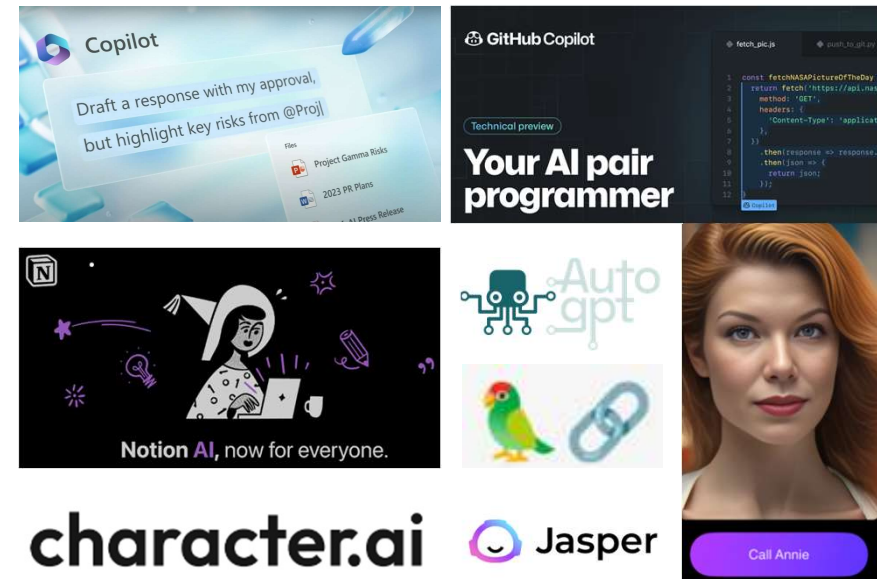
Open Source



Closed Source



Gen.AI Applications



LLM Training Challenges

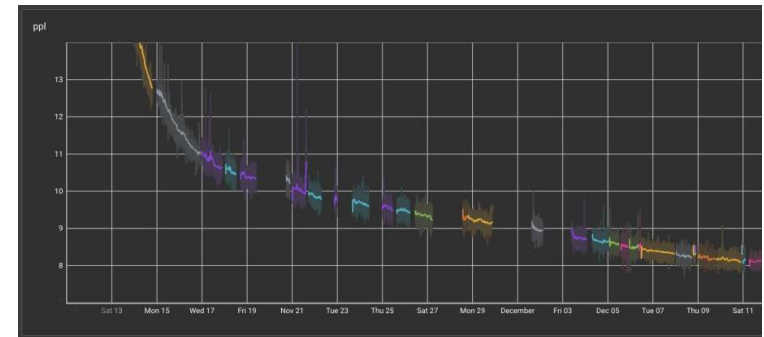
LLM Training

- 1000+ GPU Cluster
- PB data collecting, cleaning, etc.
- Training optimization
- GPU malfunction
- Unstable loss

Issues from customers

- CUDA initialization failure
- Poor NCCL performance
- GPU direct RDMA malfunction
- RoCE network
- Distributed training task
- GPU cluster performance optimization
-

- Meta OPT-175B top three record long runs of the experiment these past two weeks, lasting 1.5, 2.8, and 2 days each.

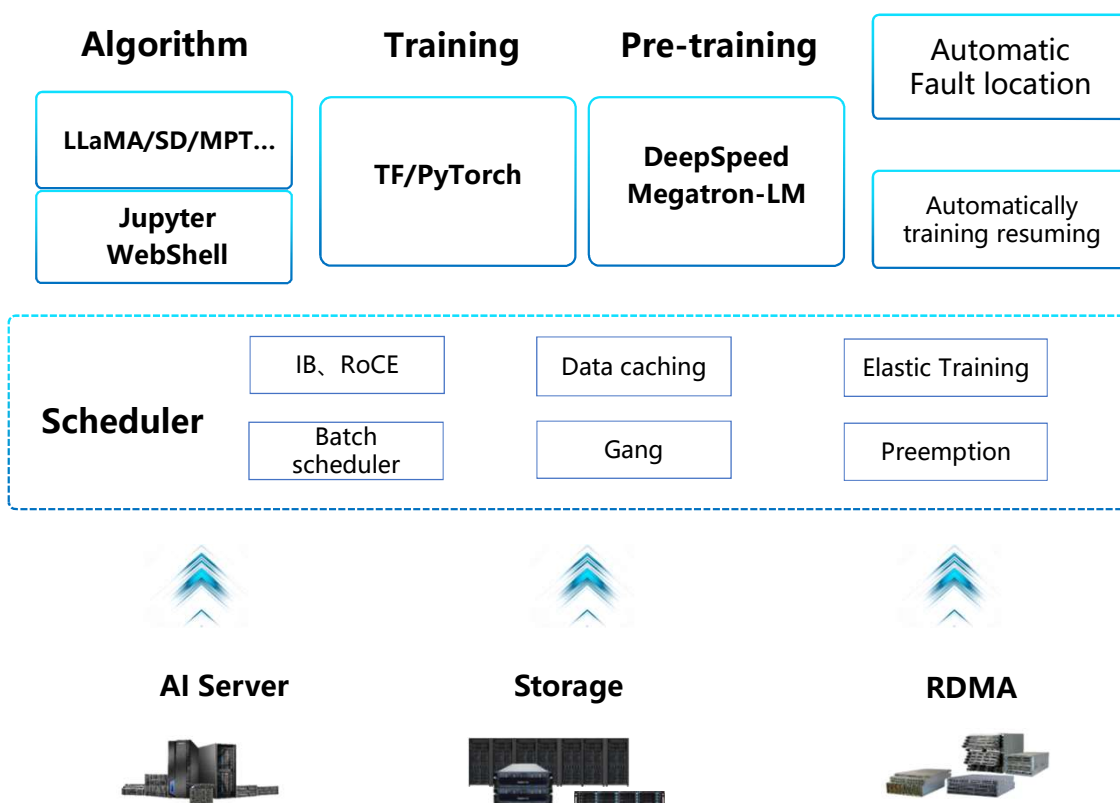


Analysis

2021-12-01 8:39am ET: 12.38 True Adam with Lower LR
2021-12-01 2:21am ET: [Stephen oncall] Run 12.37 [WARNING: See 2021-12-02 17:16 ET: Debrief on why it may not be SGD]
2021-11-30 7:24am ET: [Stephen oncall] Run 12.37 Manual request of 12.36 [WARNING: See 2021-12-02 17:16 ET: Debrief on why it may not be SGD]
2021-11-30 7:24am ET: [Stephen oncall]
2021-11-30 10:10am PT: 12.36 restart from 37k SGD mimicking [WARNING: See 2021-12-02 17:16 ET: Debrief on why it may not be SGD]
2021-11-30 8:00am PT: 12.35 restart from 37k SGD mimicking [WARNING: See 2021-12-02 17:16 ET: Debrief on why it may not be SGD]
2021-11-30 9:00am ET: 12.34 request
2021-11-29 7:43am PT [Susan]: 12.34 restart
2021-11-30 7:43am PT [Susan]: 12.33 request
2021-11-28 6:34am ET [Stephen]: 12.33
2021-11-28 5:52am ET [Stephen]: 12.32
2021-11-28 12:28am ET [Stephen]: 12.31
2021-11-28 10:09am ET [Stephen]: 12.30
2021-11-28 9:41am ET [Stephen]: 12.29
2021-11-28 3:20am ET [Stephen]: 12.28
2021-11-28 1:50am ET [Stephen]: 12.27
2021-11-27 11:38 ET [Stephen]: Run 12.26
2021-11-27 8:10am PT: Run 12.25 [Susan restart]
2021-11-27 10:59am PT: Run 12.24 [Mye rerunning job, but AFK rest of day]
2021-11-26 8:47am ET [Stephen managing cluster]
2021-11-25 8:53am ET [Susan]: Run 12.23
2021-11-25 11:35am ET [Mye]: Run 12.22
2021-11-25 11:20am ET: Run 12.21 [request]
2021-11-24 11:18am ET [Susan]: Run 12.21
2021-11-24 10:40pm ET [Susan]: Run 12.20
2021-11-24 3:30pm ET [Susan]: Run 12.19
2021-11-24 2:10pm ET [Susan]: Run 12.18
2021-11-24 1:00pm ET [Susan]: Run 12.17

Run 4
Run 3
[Undated] Run 2
2021-10-20: Run 1
Kitchen sink: Analysis of Exp ...
Description of experiments:
Exp 21 - 23
Exp 23, 24 and 25 (Drop out ...
Exp 26 and 27: Adding in Bo...
Exp 29: Manually cleaned up...
Learned Embeddings
Exp 27 vs 28
Oncall Debugging
Philosophy
Responsibilities
Onboarding & Gotchas
Questions
Run is stuck in loop of "lowe...
Remember to document y...
WPS has dropped a lot

Distributed Training Adoption



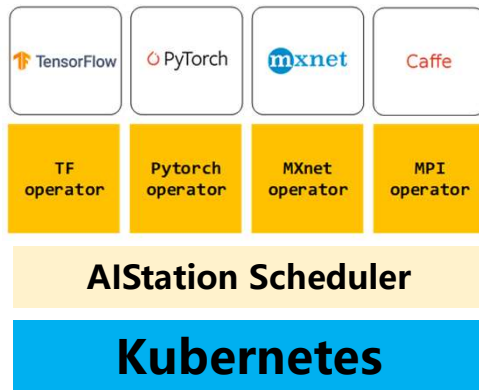
Fast training abnormal locating, automatic training resuming

- ✓ Fast chip, network abnormal locating and fault pause process to hold global training.
- ✓ Calculating standby node and automatic elastic replacement.
- ✓ Health node CheckPoint reading and automatic training resuming.

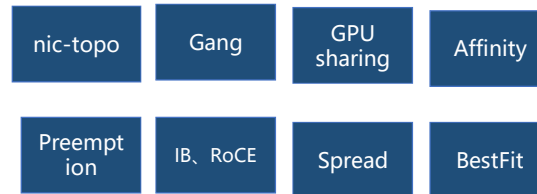
Simplify network adaptation, flexible and efficient

- ✓ Compatible with IB, RoCE and other complex cluster networking environment.
- ✓ Flexible resource matching for large-scale training scenario.
- ✓ Automatic fault tolerance to ensure long-term model training efficient and stable.

Distributed Training Optimization



AIStation scheduling strategy



Convenient way

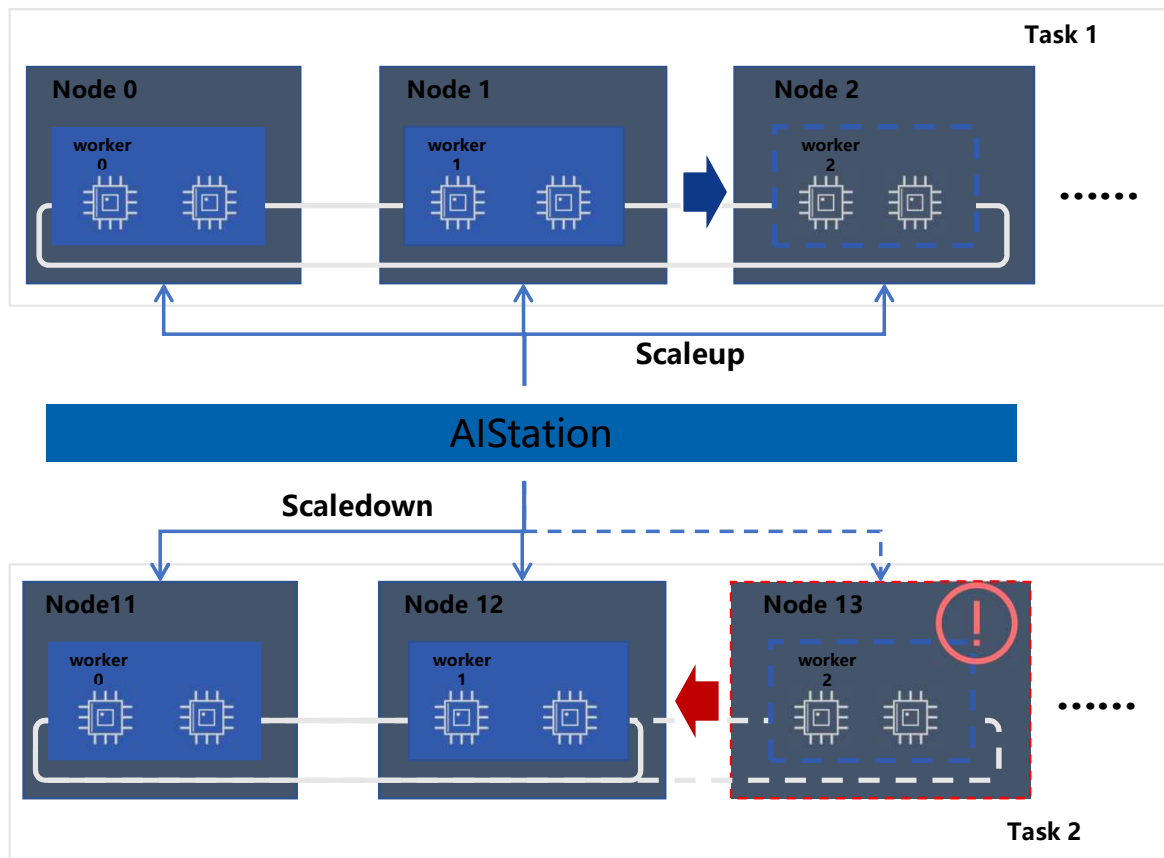
- ✓ Simple configuration, one-click launch distributed job.
- ✓ Large model training scenarios, quick start and support Megatron-LM, DeepSpeed, etc.

Tensorflow	Pytorch	Mxnet	Paddle
ParameterServer	Distributed Data Parallel Training (DDP)	Data Parallel (Server-worker-scheduler)	ParameterServer
Mirrored			
MultiWorker Mirrored	Collective Communication		Collective
CentralStorage			
MPI	MPI	MPI	MPI

Professional optimization

- ✓ Operator optimization, support tensorflow, pytorch, mxnet, caffe, paddle native distributed training and MPI mode.
- ✓ Optimized distributed scheduling strategy to achieve a fast computing resources allocation, and automatic distributed training process launch.

Elastic and Dynamic Resource Usage



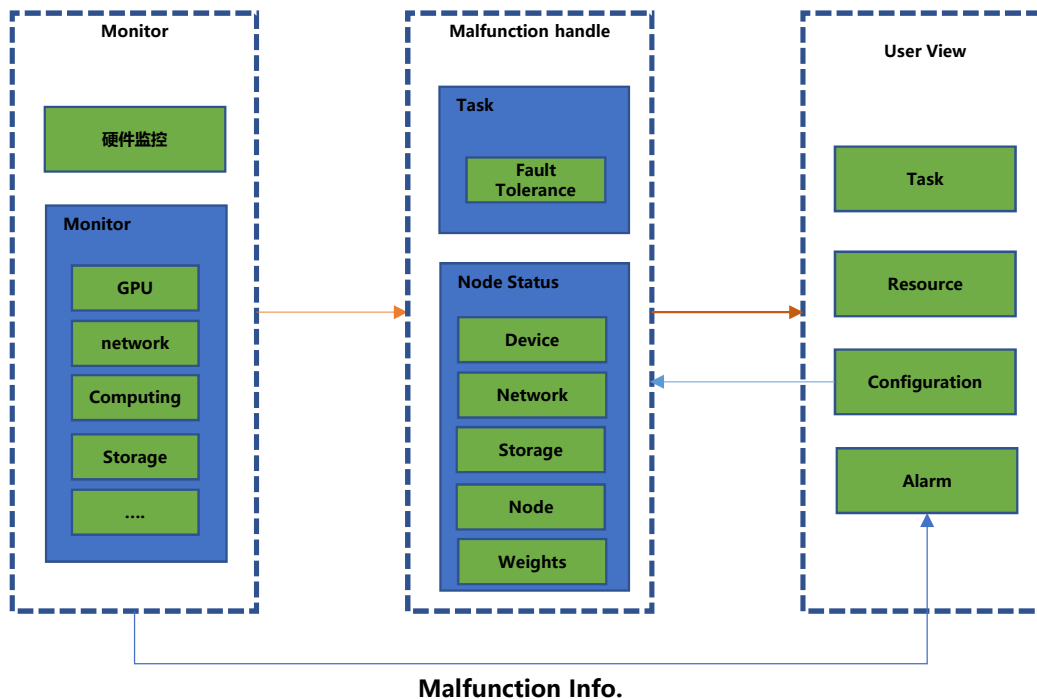
Elastic training

- ✓ Elastic training mode, dynamic resource allocation on demand.
- ✓ Maximizing the use of GPUs to achieve a high utilization.

Comprehensive guarantee for large-scale training

- ✓ Simplify the resources evaluation strategy, dynamic adjustment training resources.
- ✓ Timeliness and reliability of the huge scale of training.
- ✓ Training anomaly recovery with limited resources, automatic fault awareness and self-healing elastic resource usage.

Automatic Fault Tolerant Processing



Basic function

- When training mission abort, such as worker exit, master exit.
- Fault tolerance: GPU node displacement process, task to restart.

For elastic training task

- Master malfunction: resubmit training task
- Worker malfunction: handle by framework
- NIC fault tolerance: handle by framework
- GPU malfunction: replace the abnormal mode and restart the task



KubeCon



CloudNativeCon



OPEN SOURCE SUMMIT

China 2023

Thanks