



BUILDING FOR THE ROAD AHEAD

DETROIT 2022

Machine Learning Using Various GPU Technologies with Kubeflow

Jihye Choi, SAMSUNG SDS

Machine Learning Using Various GPU Technologies with Kubeflow



BUILDING FOR THE ROAD AHEAD

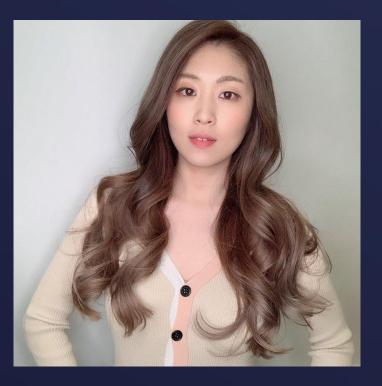
DETROIT 2022



BUILDING FOR THE ROAD AHEAD

DETROIT 2022

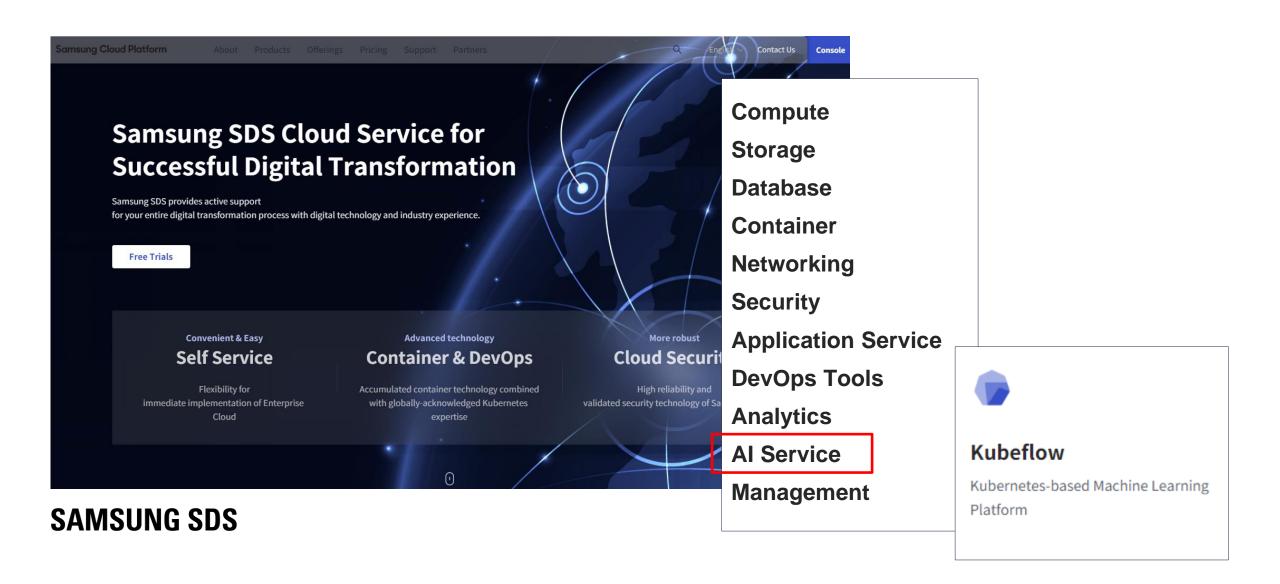
October 24-28, 2022



Jihye Choi
Cloud Architect
SAMSUNG SDS

Samsung Cloud Platform

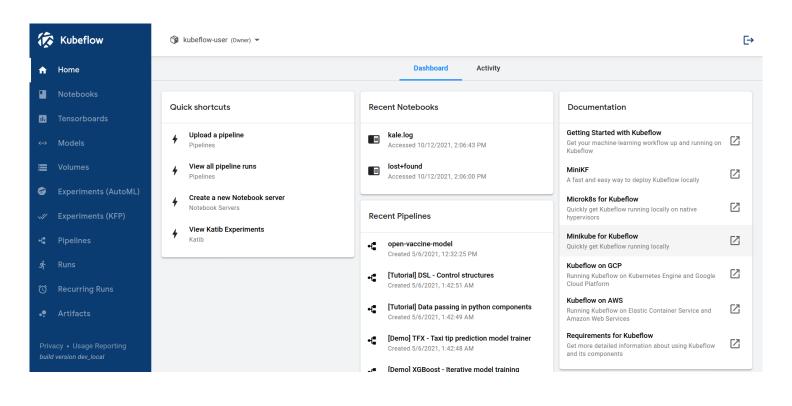




What is Kubeflow?











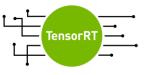














GPUs in Data Center VS Professional Desktop

	NVIDIA H100	NVIDIA A100	AMD MI250	
		op of		
GPU Memory	80GB	80GB	128GB	
Memory Bandwidth	2.0TB/s	1,935 GB/s	3,276 GB/s	
Interconnect	PCIe Gen5	PCIe Gen4	PCle Gen4	

GPUs in Data Center

NVIDIA RTX



4GB - 48GB

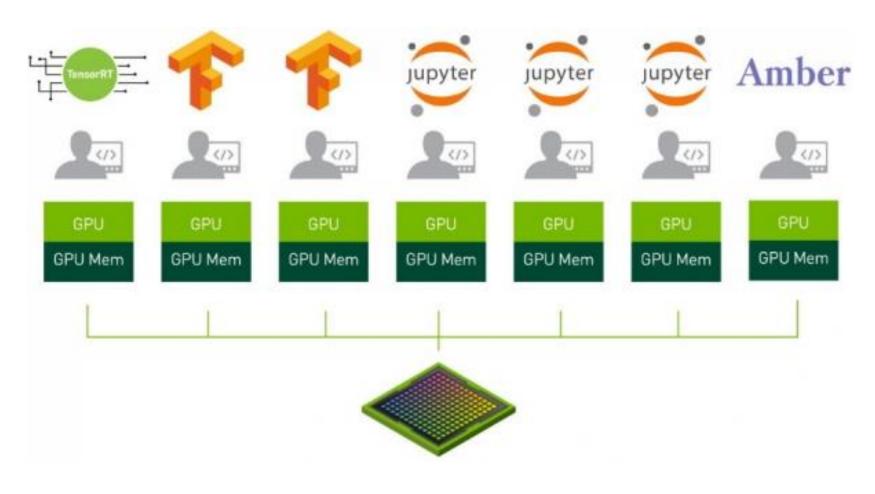
80 GB/s - 768 GB/s

PCIe Gen3 PCIe Gen4

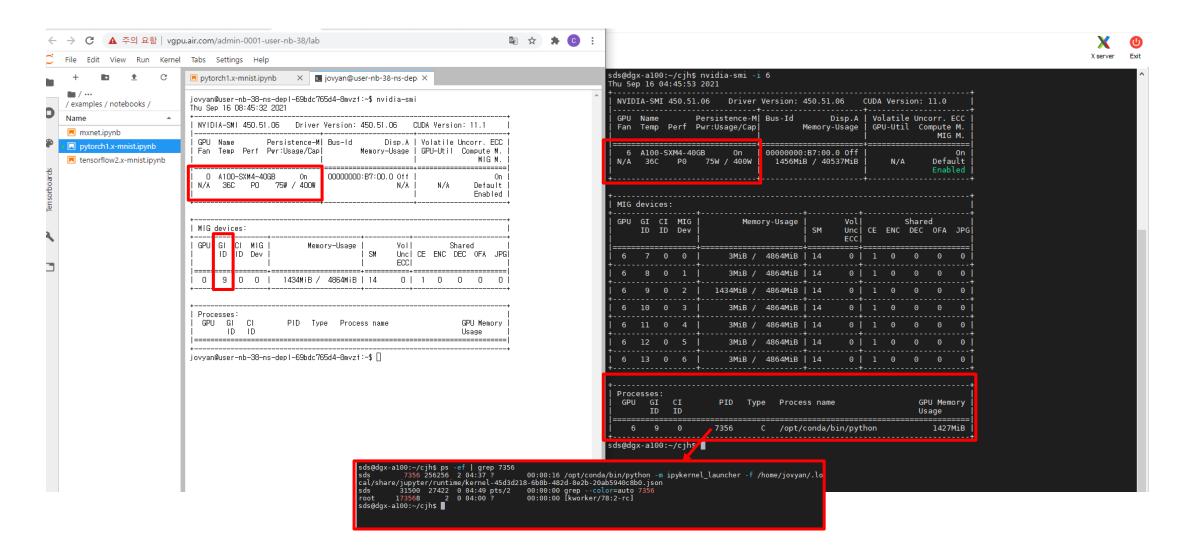
GPUs in Professional Desktop



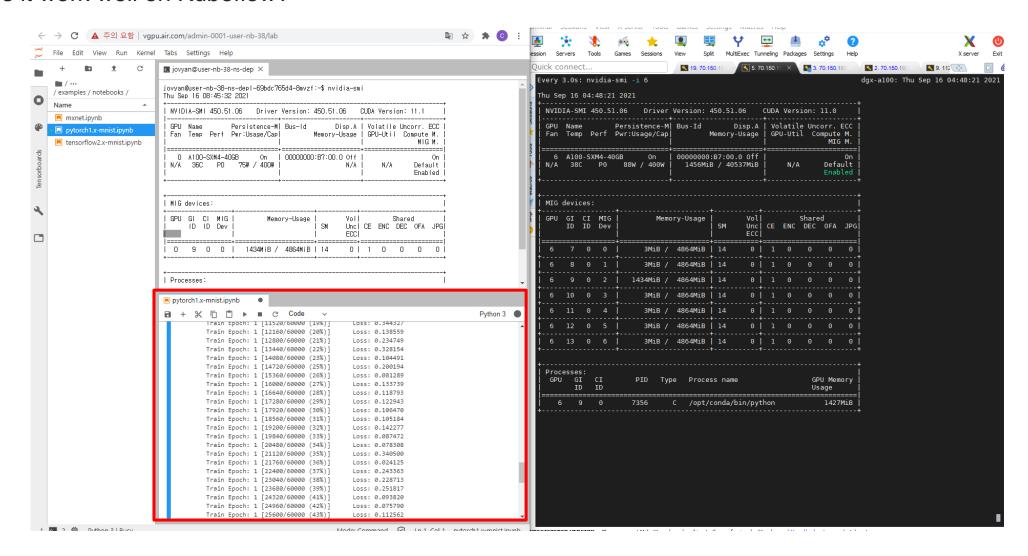
What is Multi-Instance GPU (MIG)?



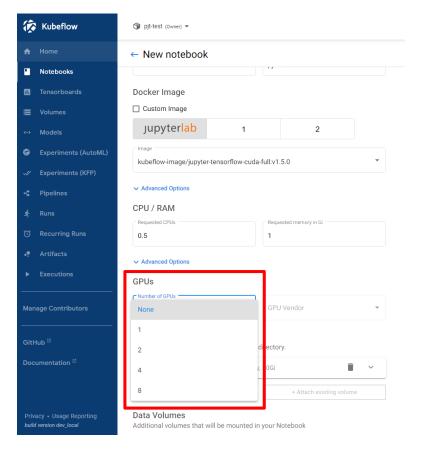




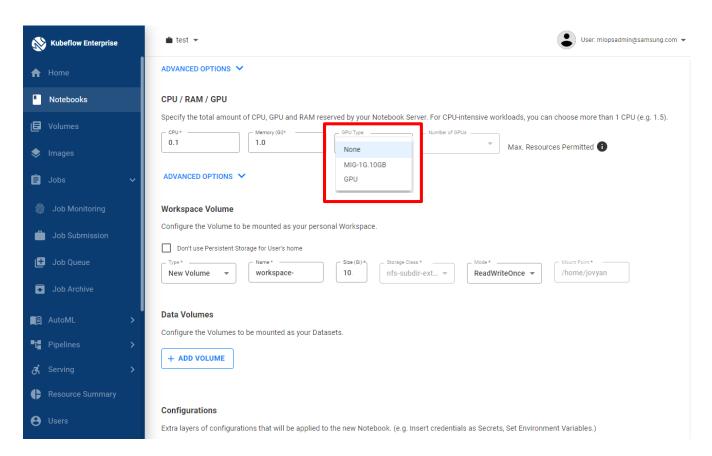






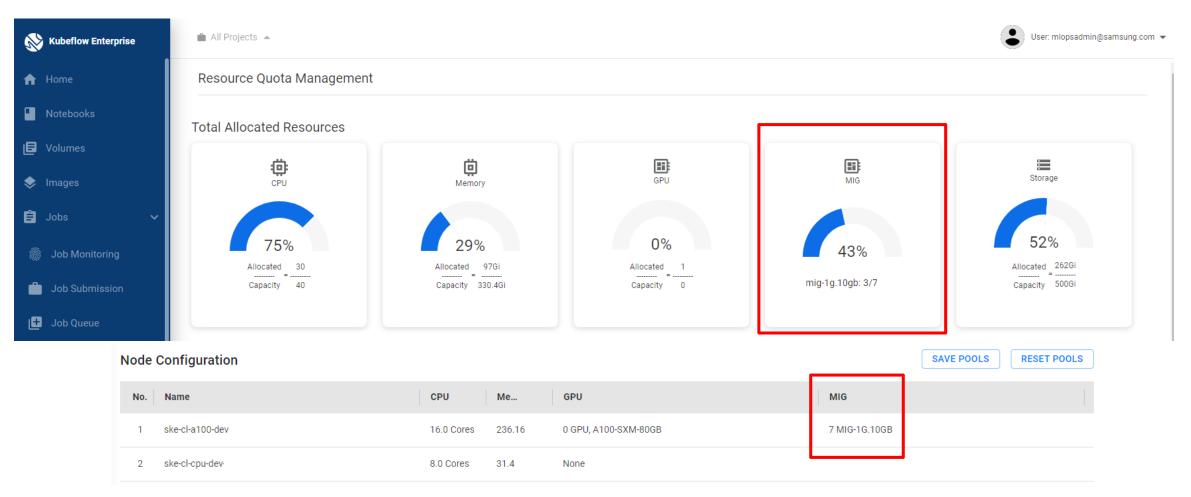


Open Source Kubeflow



SCP Kubeflow Service

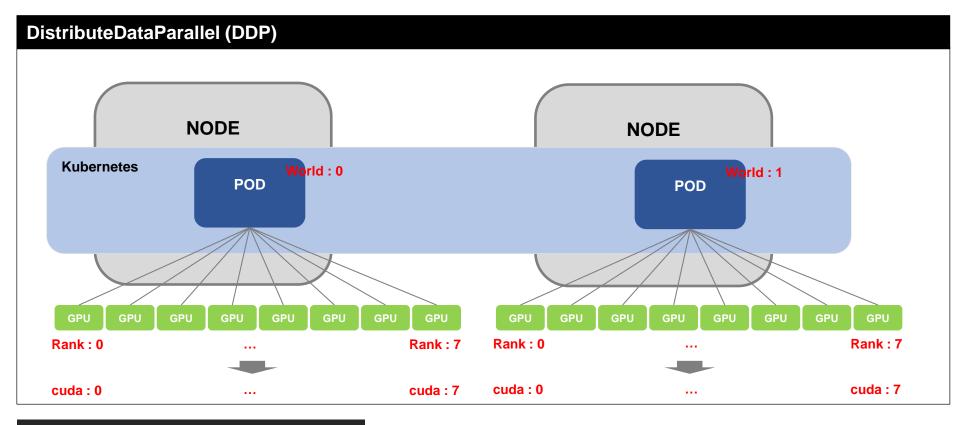




SCP Kubeflow Service



2. Is Distributed Training Feasible? (only one MIG device can be assigned per POD)



device = torch.device("cuda:{}".format(rank))

- World Size: Total Number of Process

- Rank: Name of Process



2. Is Distributed Training feasible? (only one MIG device can be assigned per POD)



<4 GPUs per POD, Total 2 PODs>

	2.0s: n							
	[A-SMI			river	Version:	470.57.02	CUDA Versi	on: 11.4
GPU Fan	Name Temp I	Perf	Persiste Pwr:Usaç		Bus-Id	Disp. Memory-Usag		Uncorr. ECC Compute M. MIG M.
0 N/A	NVIDIA 30C	A100 P0	-SXM 63W /			0:07:00.0 Of iB / 81251Mi		0 Default Disabled
1 N/A	NVIDIA 30C	A100 P0	-SXM 65W /			0:0B:00.0 Of iB / 81251Mi		0 Default Disabled
2 N/A	NVIDIA 30C	A100 P0	-SXM 66W /			0:48:00.0 Of iB / 81251Mi		0 Default Disabled
3 N/A	NVIDIA 31C	A100 P0	-SXM 69W /	Off 400W		0:4C:00.0 Of iB / 81251Mi		0 Default Disabled
4 N/A	NVIDIA 29C	A100 P0	-SXM 66W /	off 400W		0:88:00.0 Of iB / 81251Mi		0 Default Disabled
5 N/A	NVIDIA 30C	A100 P0	-SXM 67W /	off 400W		0:8B:00.0 Of iB / 81251Mi		0 Default Disabled
6 N/A	NVIDIA 30C	A100 P0	-SXM 68W /			0:C8:00.0 Of iB / 81251Mi		0 Default Disabled
7 N/A	NVIDIA 30C		-SXM 66W /			9:CB:00.0 Of iB / 81251Mi		0 Default Disabled
Proce GPU	esses: GI ID	CI ID	PIC	тур	oe Proce	ess name		GPU Memory Usage
9 1 2 3 4 5 6	N/A N/A N/A N/A N/A N/A N/A	N/A N/A N/A N/A N/A N/A N/A	4044499 4044500 4044501 4044502 4041834 4041840 4041842 4041843		C /opt, C /opt, C /opt, C /opt, C /opt, C /opt,	/conda/bin/p /conda/bin/p /conda/bin/p /conda/bin/p /conda/bin/p /conda/bin/p /conda/bin/p	ython ython ython ython ython ython	1009MiB 1009MiB 1009MiB 1009MiB 1009MiB 1009MiB 1009MiB



2. Is Distributed Training feasible? (only one MIG device can be assigned per POD)

```
pytorch-albert-4gpu-master-0
World Size: : 2
Node Rank : 0
Master Addr : localhost
Params : --do_train --fp16 --train_batch_size 32 --gradient_accumulation_steps 1 --num_train_epochs 1 --nccl_debug_INFO --nccl_debug_subsys_INIT --nccl_ib_disable 1 --nccl_ib_gdr_level 0
Setting OMP NUM THREADS environment variable for each process to be 1 in default, to avoid your system being overloaded, please further tune the variable for optimal performance in your appl
 ***************
Gradient Acculumation Steps:
NCCL_DEBUG_SUBSYS:
NCCL IB DISABLE:
NCCL_IB_GDR_LEVEL:
Model Name:
  onfig File:
                       /workspace/src/bert/albert_config_large_with_dropout.json
Train Batch Size:
Gradient Acculumation Steps: 1
NCCL DEBUĞ SUBSYS:
NCCL IB DISABLE:
NCCL_IB_GDR_LEVEL:
Model Name:
                       albert
 onfig File:
                       /workspace/src/bert/albert config large with dropout.json
Train Batch Size:
Gradient Acculumation Steps:
NCCL Debug:
NCCL DEBUG SUBSYS:
NCCL_IB_DISABLE:
NCCL_IB_GDR_LEVEL:
Model Name:
                                                                                                                              cuda: 1
Config File:
                       /workspace/src/bert/albert_config_large_with_dropout.json
                                                                                                                              cuda: 3
Train Batch Size:
Gradient Acculumation Steps:
NCCL Debug:
NCCL_DEBUG_SUBSYS:
                                                                                                                              cuda: 2
                       TNTT
NCCL_IB_DISABLE:
NCCL_IB_GDR_LEVEL:
Model Name:
Config File:
                                                                                                                              cuda: 0
                       /workspace/src/bert/albert_config_large_with_dropout.json
Init Process
Init Process
Init Proce
  _main__] device: cuda:1 n_gpu: 1, distributed training: True, 16-bits training: True
  aRNING: 0 tput directory /workspace/shared-dir/bert/results/squad/normal/ckpt_281181_2e-05_1 already exists and is not empty. ['albert_korquard_config.pt', 'albert_korquard_model_final.pt']
__main__| device: cuda:3 n_gpu: 1 distributed training: True, 16-bits training: True
           tput directory /workspace/s/shared-dir/bert/results/squad/normal/ckpt_281181_2e-05_1 already exists and is not empty. ['albert_korquard_config.pt', 'albert_korquard_model_final.pt']
device: cuda:2 n_gpu: 1, distributed training: True, 16-bits training: True
WARNING: 0 tput directory /workspace/shared-dir/bert/results/squad/normal/ckpt_281181_2e-05_1 already exists and is not empty. ['albert_korquard_config.pt', 'albert_korquard_model_final.pt' to device
  _main_ ] device: cuda:0 n_gpu: 1, distributed training: True, 16-bits training: True
  ARNING: Output directory /workspace/shared-dir/bert/results/squad/normal/ckpt_281181_2e-05_1 already exists and is not empty. ['albert_korquard_config.pt', 'albert_korquard_model_final.pt']
LOADED CHECKPOINT
```



2. Is Distributed Training feasible? (only one MIG device can be assigned per POD)

```
kind: PyTorchJob
metadata:
 name: pytorch-albert-mig2
spec:
 pytorchReplicaSpecs:
   Master:
     replicas: 1
     template:
        spec:
         containers:
           image:
           name: pytorch
           resources:
             limits:
               nvidia.com/mig-1g.10gb: 2
                                                     Total 4 MIGs
   Worker:
     replicas: 1
     restartPolicy: OnFailure
     template:
       metadata:
         annotations:
           sidecar.istio.io/inject: "false"
        spec:
         containers:
           image:
           name: pytorch
           resources:
             limits:
               nvidia.com/mig-1g.10gb: 2
```

```
root@ac-al00-10:~/cjh# kubectl logs -f pytorch-albert-mig2-master-0
World Size: : 2
Number of Gpus : 2
Node Kank : 0
Master Addr : localhost
Master Port: 23456
Params : --do train --fp16 --train batch size 1 --gradient accumulation steps 1 --num train
b_disable 1 --nccl_ib_gdr_level 0
Setting OMP_NUM_THREADS environment variable for each process to be 1 in default, to avoid
for optimal performance in your application as needed.
Train Batch Size:
Gradient Acculumation Steps: 1
NCCL Debug:
                    INFO
NCCL DEBUG SUBSYS: INIT
NCCL IB DISABLE:
NCCL IB GDR LEVEL:
Model Name:
                    /workspace/src/bert/albert_config_large_with_dropout.json
Confia File:
Traceback (most recent call last):
  File "/workspace/src/bert/benchmark.py", line 2250, in <module>
  File "/workspace/src/bert/benchmark.py", line 2207, in main
   torch.cuda.set device(args.local rank)
  File "/opt/conda/lib/python3.6/site-packages/torch/cuda/ init .py", line 263, in set dev
RuntimeError: CUDA error: invalid device ordinal
                                                            Error!
```

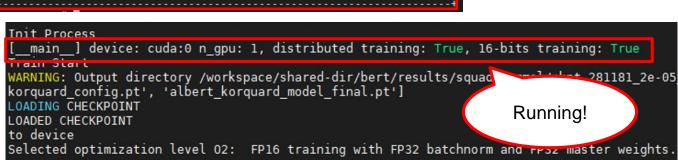
<2 MIGs per POD, Total 2 PODs>

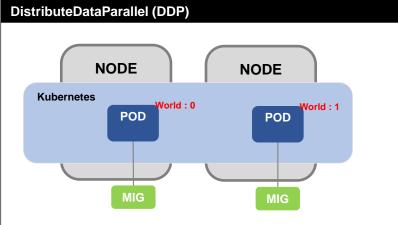


2. Is Distributed Training feasible? (only one MIG device can be assigned per POD)



<1 MIGs per POD, Total 2 PODs>



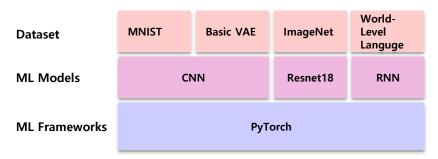


GPU Memory



3. Performance

□ Model



□ How to Test

- GPU : Sequentially 7 times on 1 GPU

- MIG: Parallelly 1 time on 7 MIGs (1g.10gb)

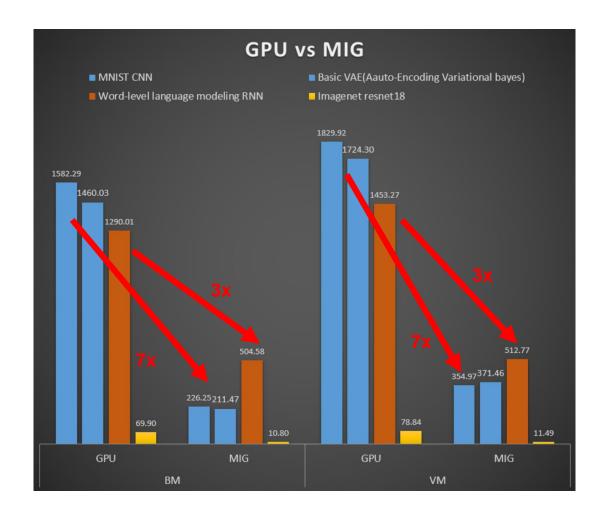
□ Result

- Comparing **Total Execution Time**

. CNN: 500-700%

. Resnet18: 700%

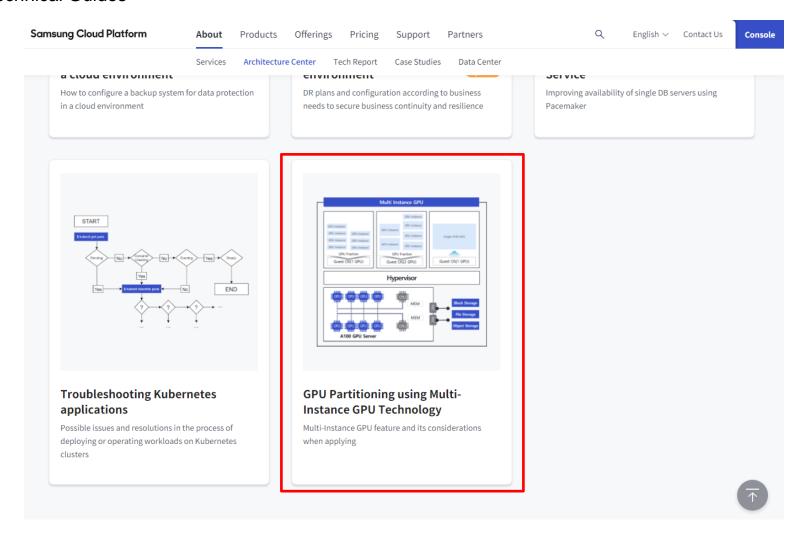
. RNN: 300%





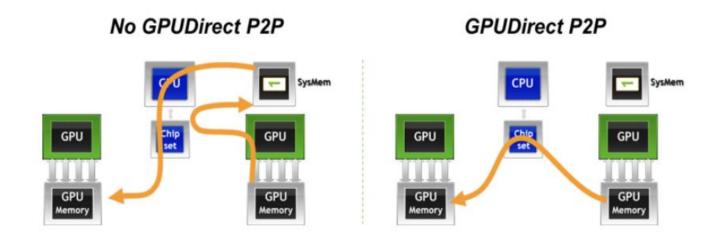
https://cloud.samsungsds.com/

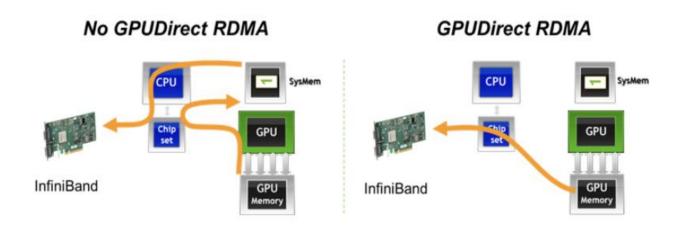
Architecture Center > Technical Guides





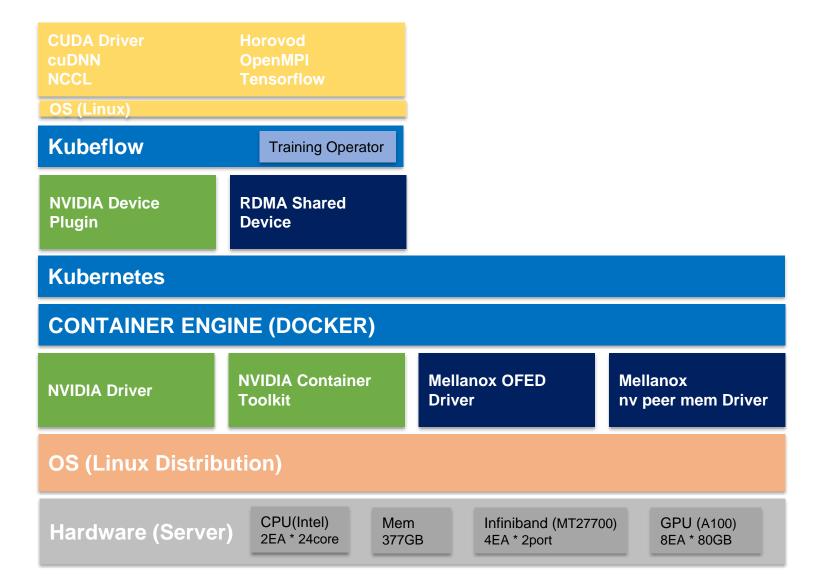
What is GPUDirect RDMA(Remote Direct Memory Access)?







Software Architecture





Examples

```
apiVersion: "kubeflow.org/v1"
kind: "PyTorchJob"
metadata:
  name: "pytorch-flower-rdma"
spec:
  pytorchReplicaSpecs:
    Master:
      replicas: 1
      restartPolicy: OnFailure
      template:
        metadata:
          annotations:
            sidecar.istio.io/inject: "false"
        spec:
                                                                                                               OS ENV
          containers:
            - name: pytorch
              image:
              command: ["python","/workspace/jovyan/src/launch.py","--dist","2","--nccl_ib_disable","0","--nccl_ib_gdr_level","1"]
             # Comment out the below resources to use the CPU.
              resources:
               limits:
                 nvidia.com/gpu: 8
                                                            Kubernetes Resource
                 rdma/hca_shared_devices_a: "1"_
              securityContext:
                                                             & SecurityContext
               capabilities:
                 add: [ "IPC_LOCK" ]
```

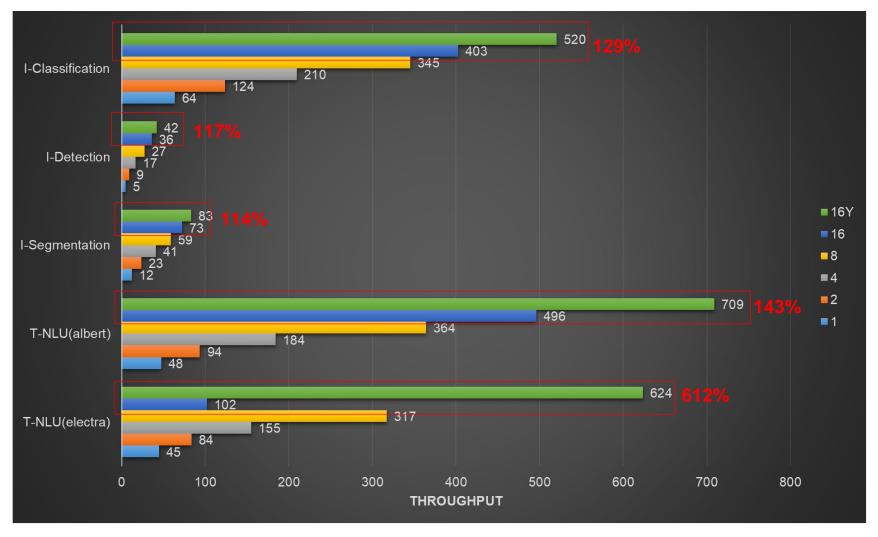


Examples

```
World Size:
Number of Gpus:
                   8
Dist World Size:
                 16 (World Size x Number of Gpus)
Node Rank:
NCCL_IB_DISABLE: 0
NCCL IB GDR LEVEL: 1
(중략)
INFO: ========== start =========
INFO: ## model_type
                          : segmentation
INFO: ## model_name
                           : unet
INFO: ## model_input : [512,512,3]
                           : A100 Graphics Device
INFO: ## gpu name
INFO: ## apus
                         : 16
INFO: ## epochs
INFO: ## batch size
                          : 128
pytorch-flower-rdma-master-0:149:149 [0] NCCL INFO Bootstrap: Using [0]eth0:40.244.64.19<0>
pytorch-flower-rdma-master-0:149:149 [0] NCCL INFO NET/Plugin: No plugin found (libnccl-net.so), using internal implementation
pytorch-flower-rdma-master-0:149:149 [0] NCCL INFO NCCL IB DISABLE set by environment to 0.
pytorch-flower-rdma-master-0:149:149 [0] NCCL INFO NET/IB: Using [0]mlx5_0:1/IB [1]mlx5_2:1/IB [2]mlx5_5:1/RoCE [3]mlx5_7:1/IB; OOB
eth0:40.244.64.19<0>
pytorch-flower-rdma-master-0:149:149 [0] NCCL INFO Using network IB
NCCL version 2.7.8+cuda11.1
(중략)
pytorch-flower-rdma-master-0:154:299 [5] NCCL INFO GPU Direct RDMA Enabled for GPU 8b000 / HCA 3 (distance 3 <= 3), read 0
pytorch-flower-rdma-master-0:153:300 [4] NCCL INFO Channel 01: 4[88000] -> 3[4c000] via P2P/IPC/read
pytorch-flower-rdma-master-0:156:296 [7] NCCL INFO GPU Direct RDMA Enabled for GPU cb000 / HCA 1 (distance 3 <= 3), read 1
pytorch-flower-rdma-master-0:154:299 [5] NCCL INFO Channel 01: 10[48000] -> 5[8b000] [receive] via NET/IB/3/GDRDMA
pytorch-flower-rdma-master-0:156:296 [7] NCCL INFO Channel 02 : 7[cb000] -> 12[88000] [send] via NET/IB/1/GDRDMA
(중략)
```



Performance





Please scan the QR Code above to leave feedback on this session



BUILDING FOR THE ROAD AHEAD

DETROIT 2022