



# Using Multiple GPUs for Distributed Deep Learning on Neuron

---

**Soonwook Hwang**

December 2022

Korea Institute of Science and Technology Information

# Some Motivational thoughts on large-scale DL/LM on top of national supercomputing facilities

---

- **KISTI-6 /w ~600PF is coming in 2 years**
  - ✓ It is expected that several thousands of GPUs(?) are available
- **Is large-scale LM (Language Model) training a exclusive task that can be conducted by big tech companies (e.g., Google, Meta, MS) running data center facilities?**
- **Why is it that large-scale LM R&D is so hard in Korea?**
  - ✓ Lack of computing resources
  - ✓ Lack of datasets
  - ✓ Lack of skills??
- **What can KISTI do in contributing to large-scale LM R&D in Korea?**
  - ✓ KISTI is uniquely positioned to running a genenal-purpose national supercomputing facility in Korea
- **Is KISTI's supercomputing facility easy to access for users to do large-scal distributed deep learning R&D?**
- **What are the most significant barriers that prevent users from having access to KISTI supercomputing facilities in conducting large-scale distributed training?**
  - ✓ Is it because of the traditioinal batch-scheduling based service?
  - ✓ ??

# Objectives

---

- **Guide users to run his/her DL codes using multiple GPU nodes on Neuron (KISTI GPU Cluster)**
- **Introduce how to set up a virtual environment**
  - install Conda and create/set up his/her virtual environment
  - run his/her codes interactively on Neuron
- **Introduce how to run Distributed DL training on Neuron**
  - Horovod
  - Singularity Container
  - ...

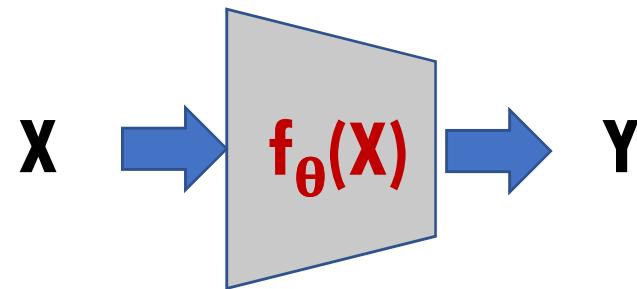
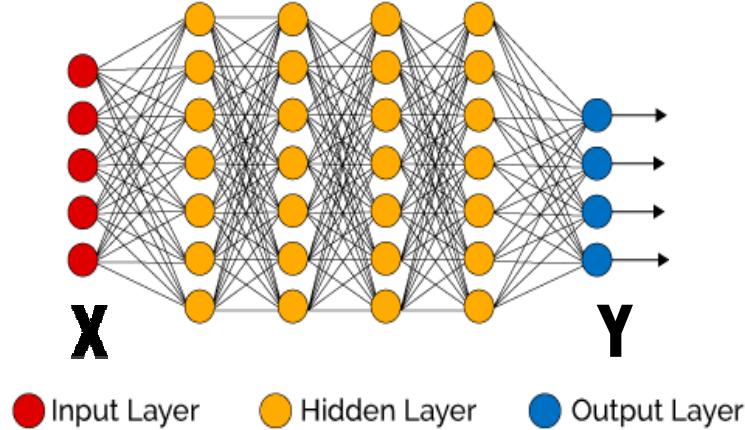
# Agenda

---

- **Why distributed training?**
  - **Distributed DL Approaches**
    - Data Parallel
      - Parameter Server(s)
      - Ring AllReduce
    - Model Parallel
      - Pipeline parallel
      - Tensor parallel
    - Multi-dimensional (3D) Parallel
  - **KISTI GPU Cluster: Neuron**
  - **Hands-on Exercises**
    - Conda Virtual Environment
    - Distributed Data Parallel(DDP) using Horovod
    - DDP using Singularity Container
-

# Deep Learning in a slide

Courtesy of Prof. Sanjeev Arora



$\theta$  : parameters

(X, Y) : (Input, Label) pair of  
Training data

$$\ell(\theta, x, y)$$

Loss function (how well net output **matched true label**  $y$  on point  $x$ ) Can be  $\| \cdot \|_2$ , cross-entropy....

Objective  $\operatorname{argmin}_{\theta} E_i[\ell(\theta, x_i, y_i)]$

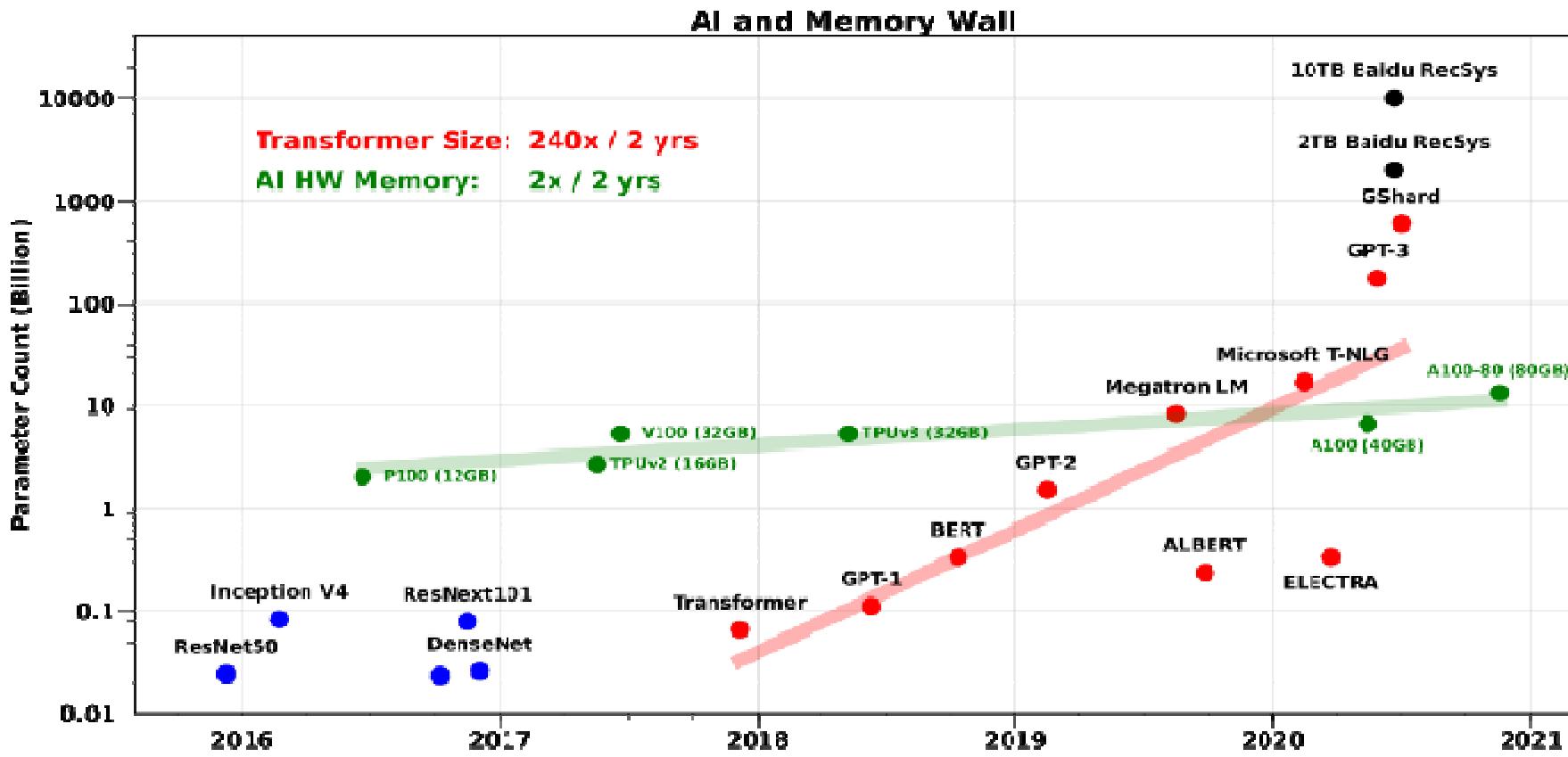
Gradient Descent  $\theta^{(t+1)} \leftarrow \theta^{(t)} - \eta \nabla_{\theta} (E_i[\ell(\theta^{(t)}, x_i, y_i)])$

Stochastic GD: Estimate  $\nabla$  via small sample of training data.

# Why Distributed Training?

# Why Distributed Training

- Model size is growing too big to fit
  - Transformer('17): 465M, GPT-3('20):175B



<https://medium.com/riselab/ai- and- memory- wall- 2cb4265cb0b8>

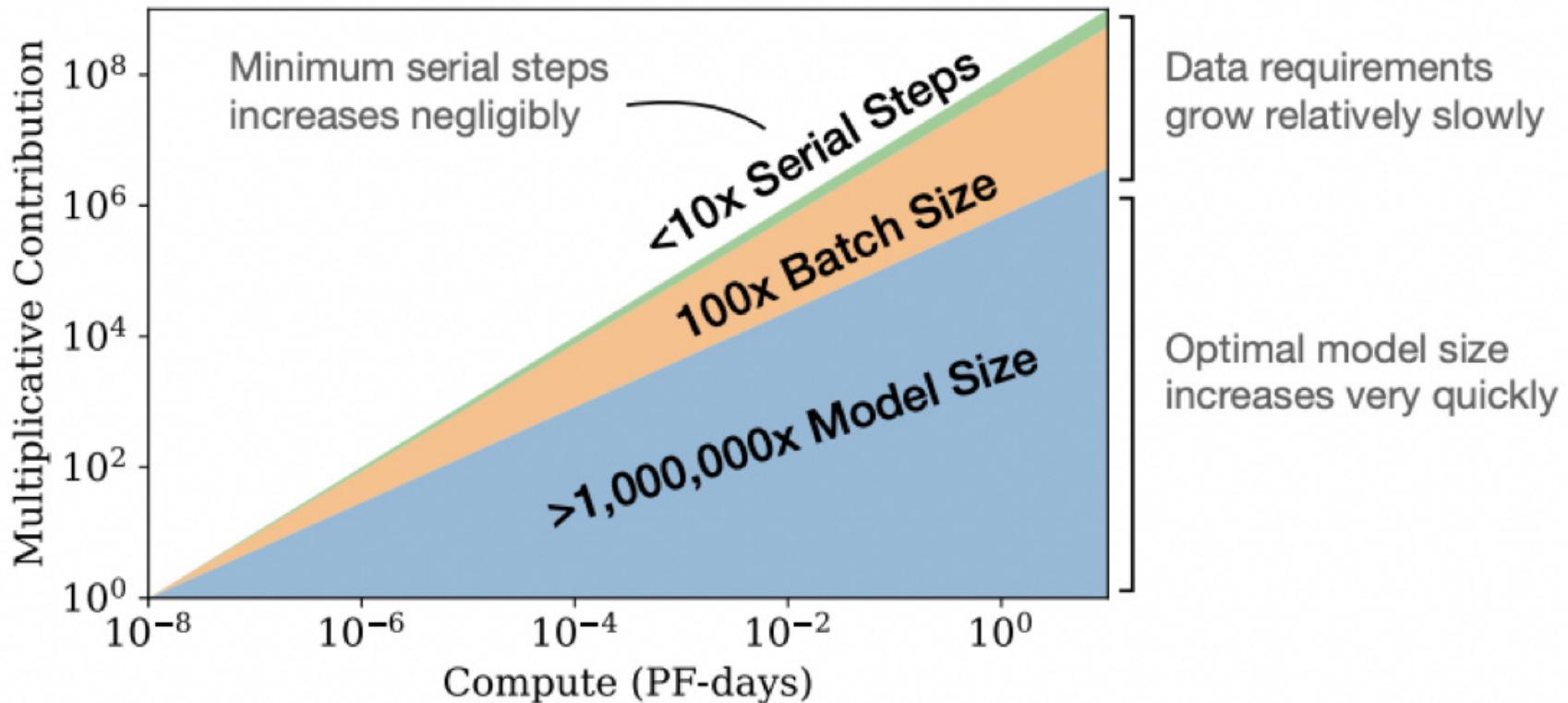
# Why Distributed Training

---

- **With More compute power and memory using multiple devices (GPU/TPU/CPU)**
  - Enable training large model
  - Speed up model training and shorten training time
  - Shorter training time allows you to do more experiments to reach your modeling goal

# Scaling Laws

- Larger models are significantly more sample-efficient
  - ✓ optimally compute-efficient training involves training very large models on a relatively modest amount of data and stopping significantly before convergence.

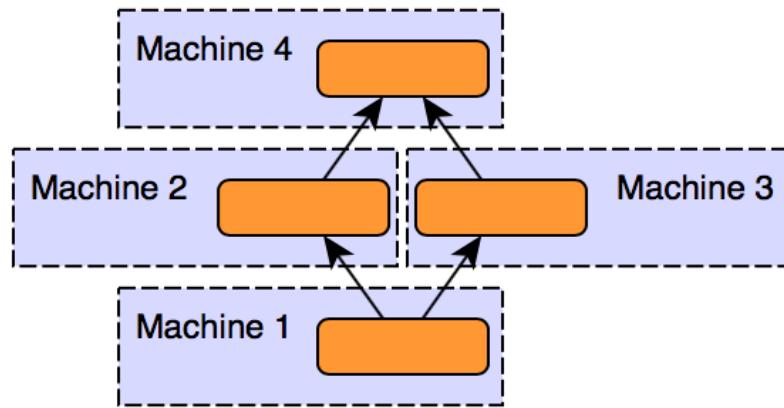


Scaling Laws for Neural Language Models  
<https://arxiv.org/pdf/2001.08361.pdf>

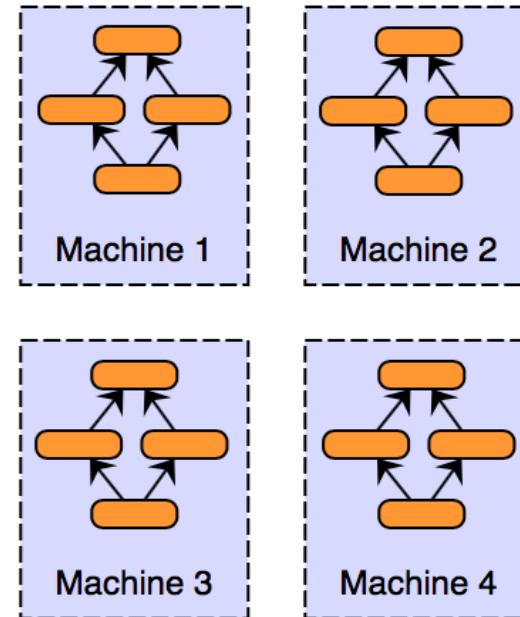
# Distributed DL Approaches

# Distributed DL approaches

Model Parallelism



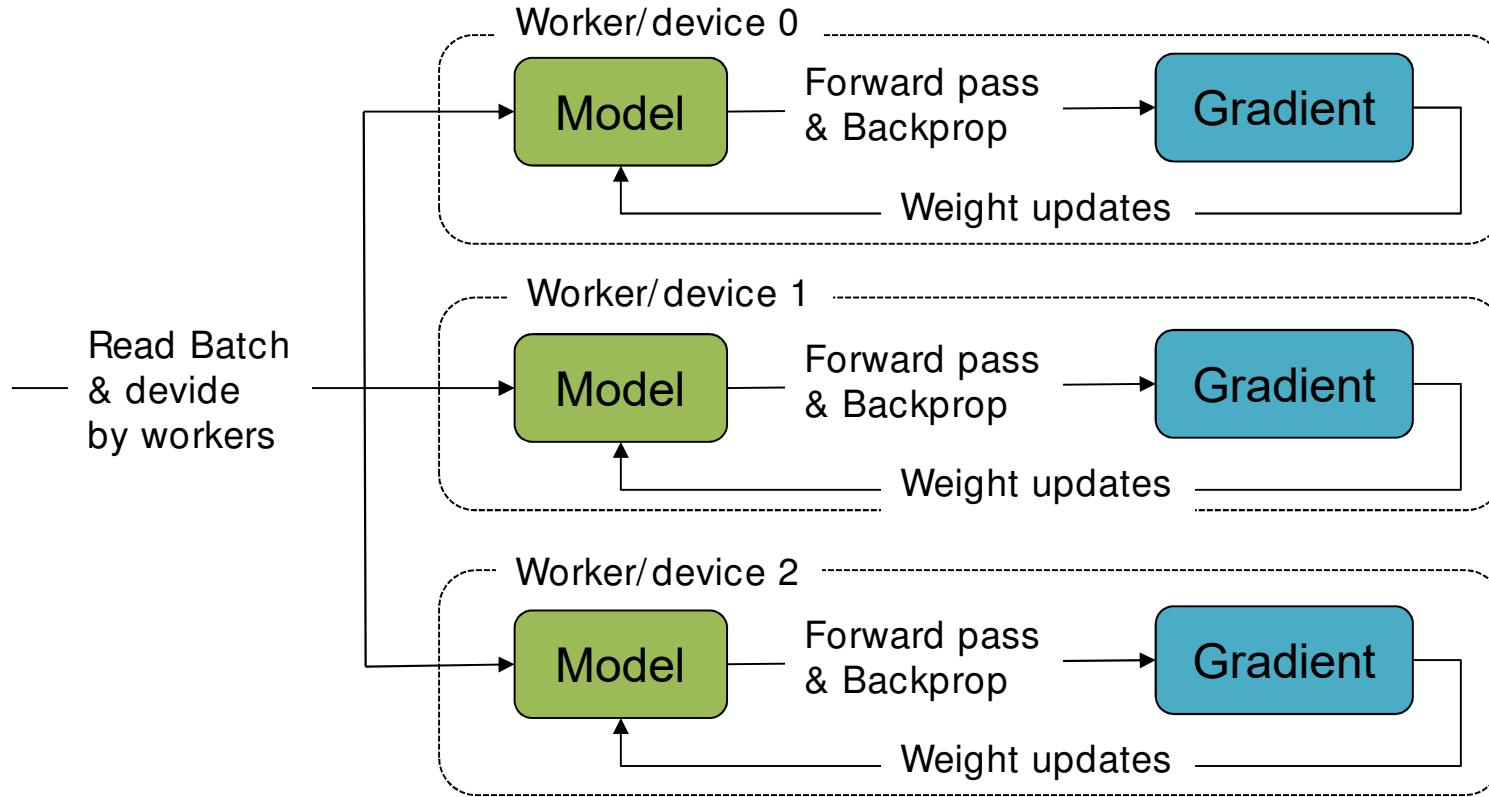
Data Parallelism



- Different parts of **model** running on multiple GPUs
- Model is too large, which cannot fit in a single device

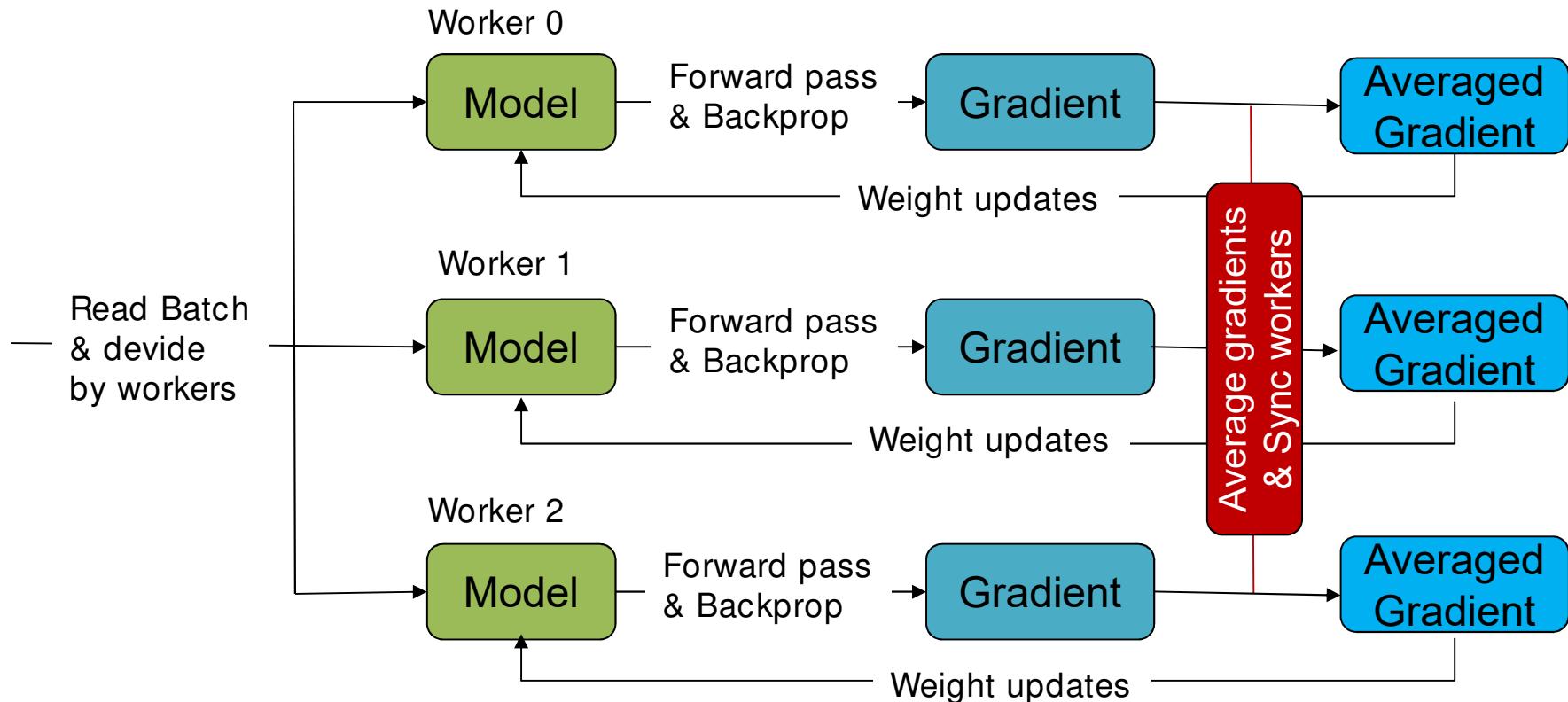
- Different parts of **data** running on multiple GPUs
- Data is too large, which need to be processed in parallel

# Data Parallel

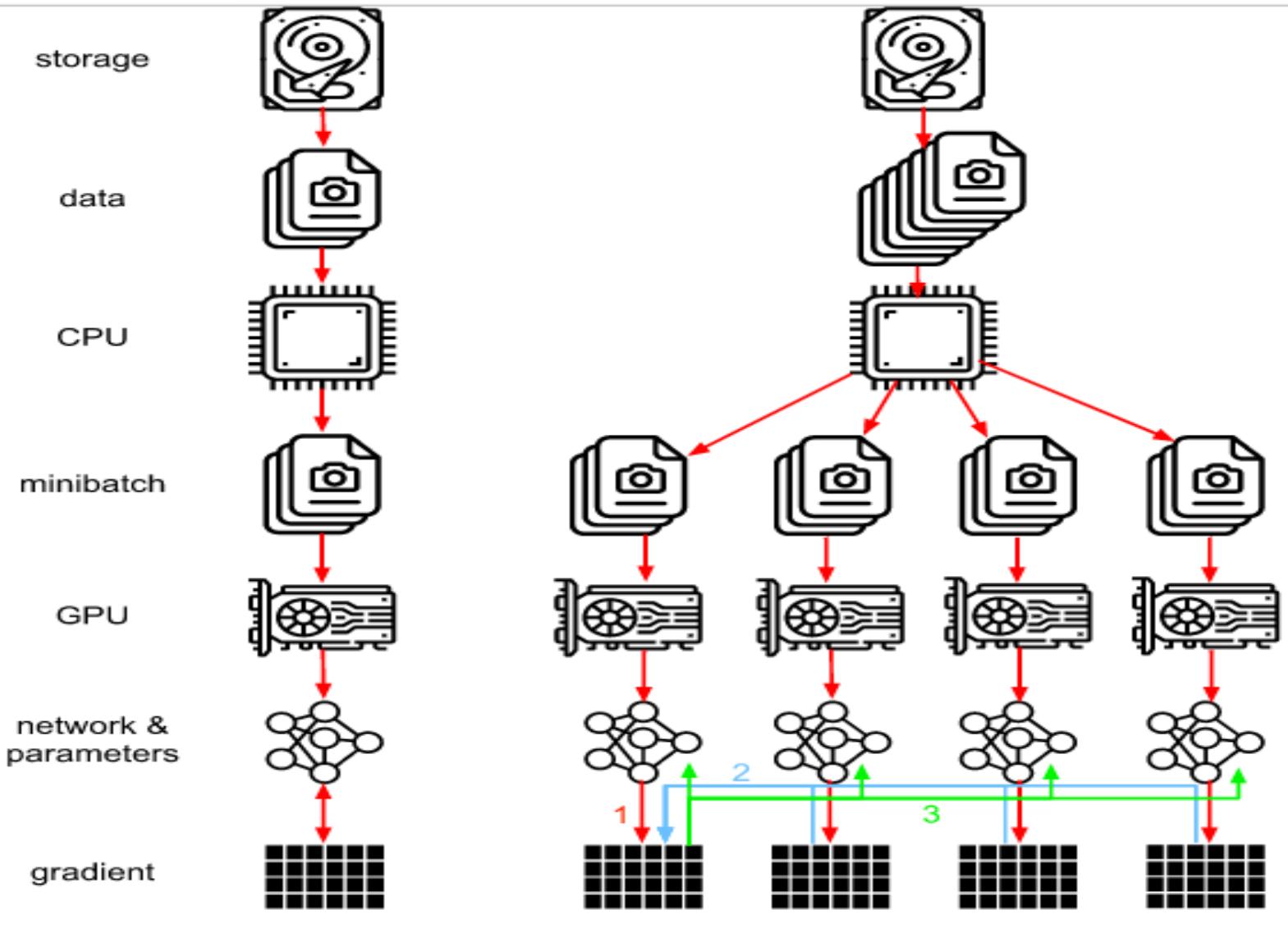


- **How to divide/distribute batches to workers**
- **When to synchronize workers**
- **How to maintain consistency in states**

# Data Parallel



# Data Parallel



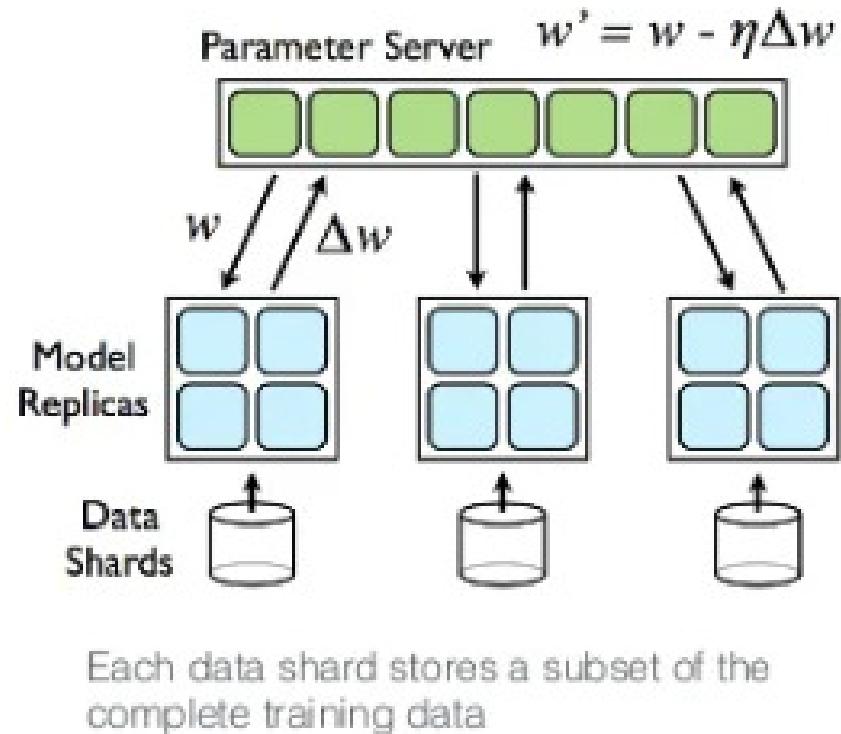
# How to synchronize model parameters

---

- **Parameter Server (centralized)**
  - Synchronous
  - Asynchronous
- **Ring AllReduce (decentralized)**

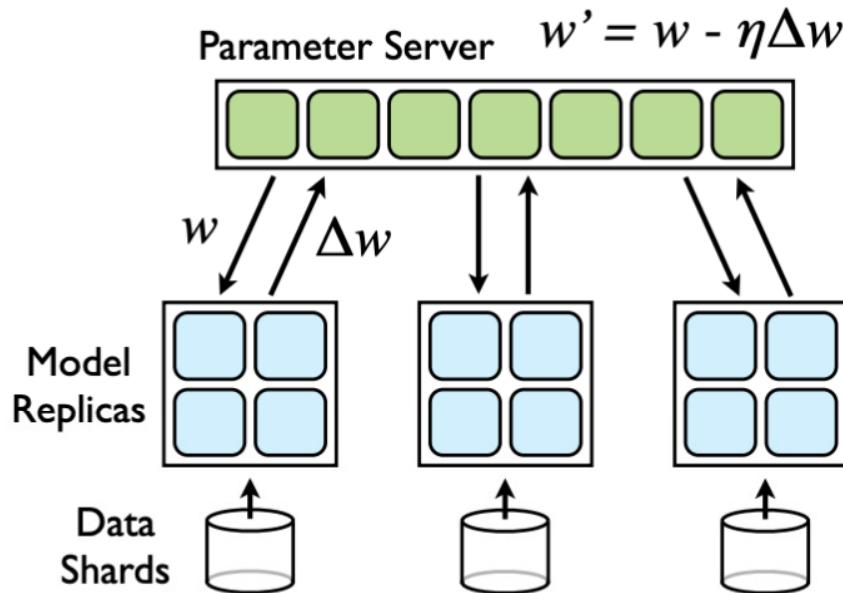
# Parameter Server (PS)

1. Split the training data into **shards** and assign a model replica to each data shard
2. For each model replica, fetch the parameters from the centralized **sharded parameter server**
3. Gradients are computed per model and pushed back to the parameter server



# Synchronous vs. Asynchronous PS

**Downpour SGD:**  
Online Asynchronous Stochastic  
Gradient Descent



**Sandblaser L-BFGS:**  
Batch Distributed Parameter Storage  
and Manipulation

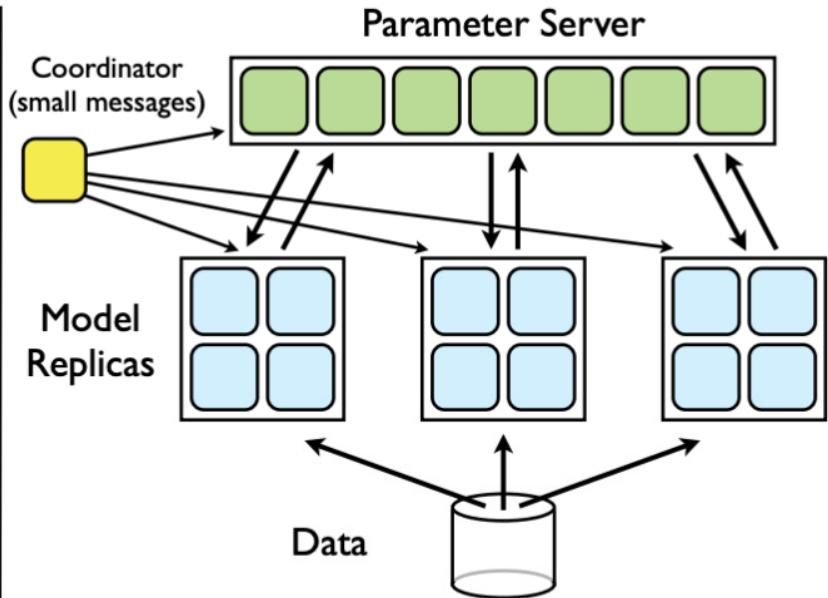
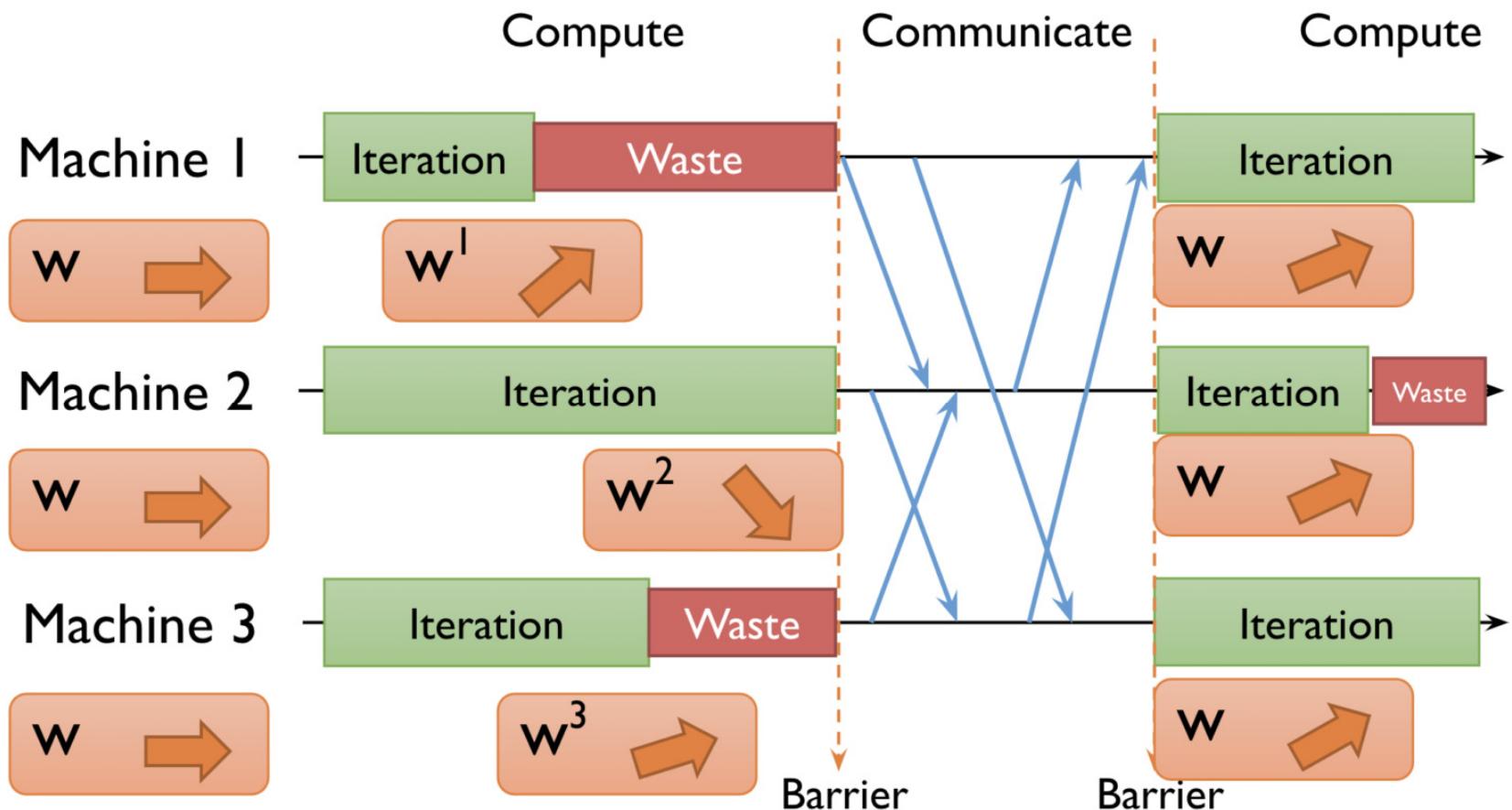


Figure 2: Left: Downpour SGD. Model replicas asynchronously fetch parameters  $w$  and push gradients  $\Delta w$  to the parameter server. Right: Sandblaster L-BFGS. A single ‘coordinator’ sends small messages to replicas and the parameter server to orchestrate batch optimization.

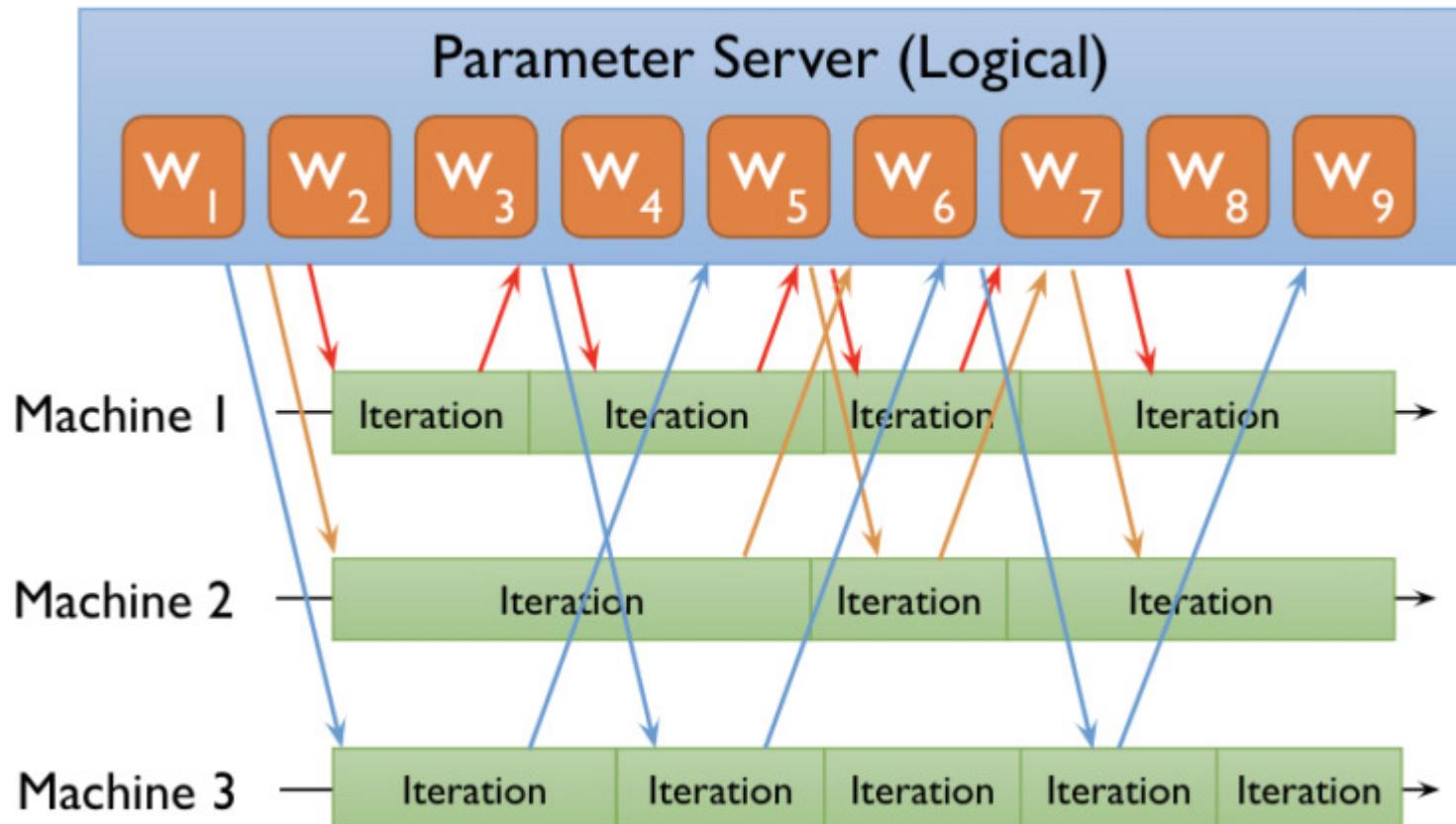
**Figure from Large Scale Distributed Deep Networks**

[https://static.googleusercontent.com/media/research.google.com/ko//archive/large\\_deep\\_networks\\_nips2012.pdf](https://static.googleusercontent.com/media/research.google.com/ko//archive/large_deep_networks_nips2012.pdf)

# Synchronous PS



# Asynchronous PS



# Ring AllReduce

- Synchronized w/o parameter server(s)

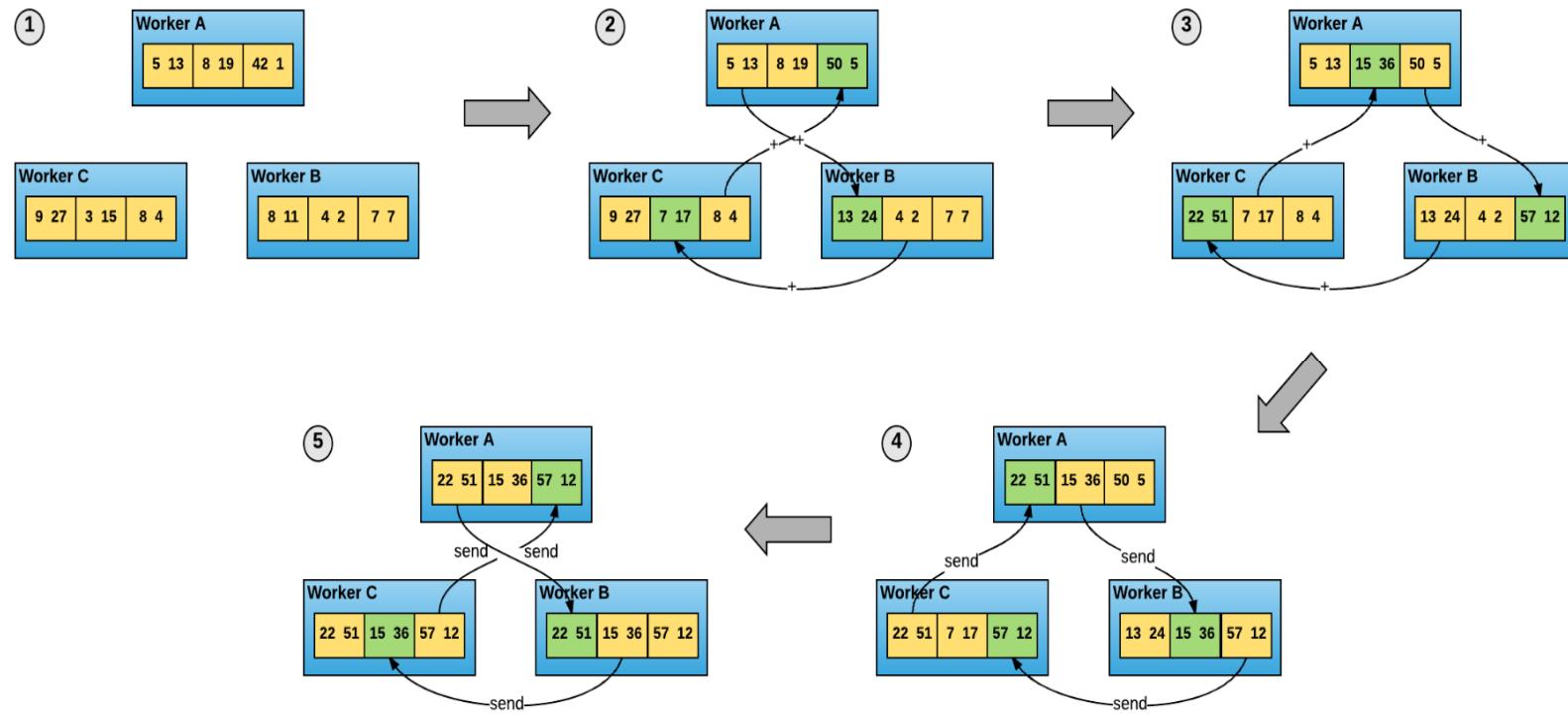


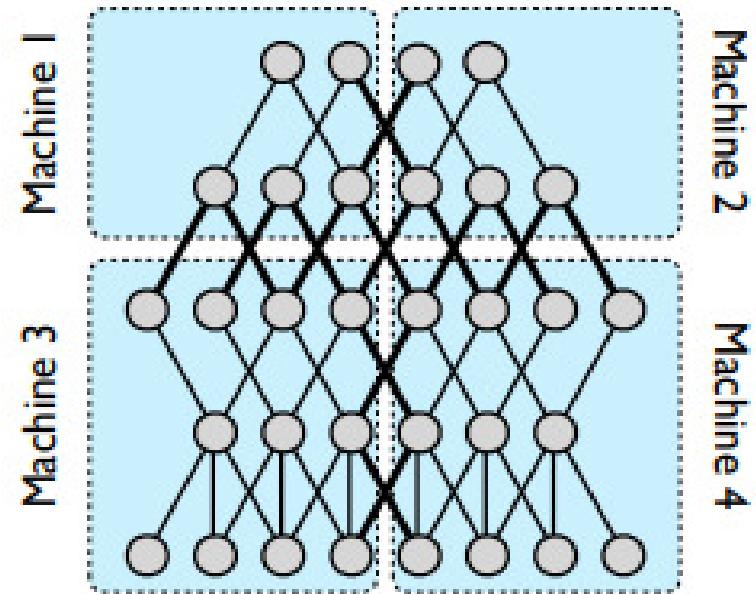
Figure 4: The ring-allreduce algorithm allows worker nodes to average gradients and disperse them to all nodes without the need for a parameter server.

**Horovod: fast and easy distributed deep learning in Tensorflow**

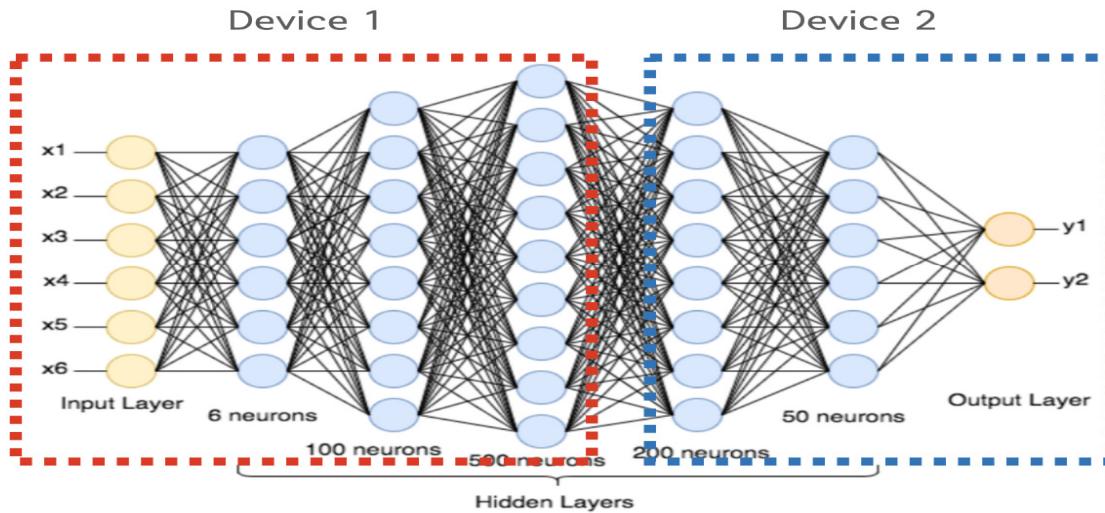
<https://arxiv.org/pdf/1802.05799.pdf>

# Model Parallel

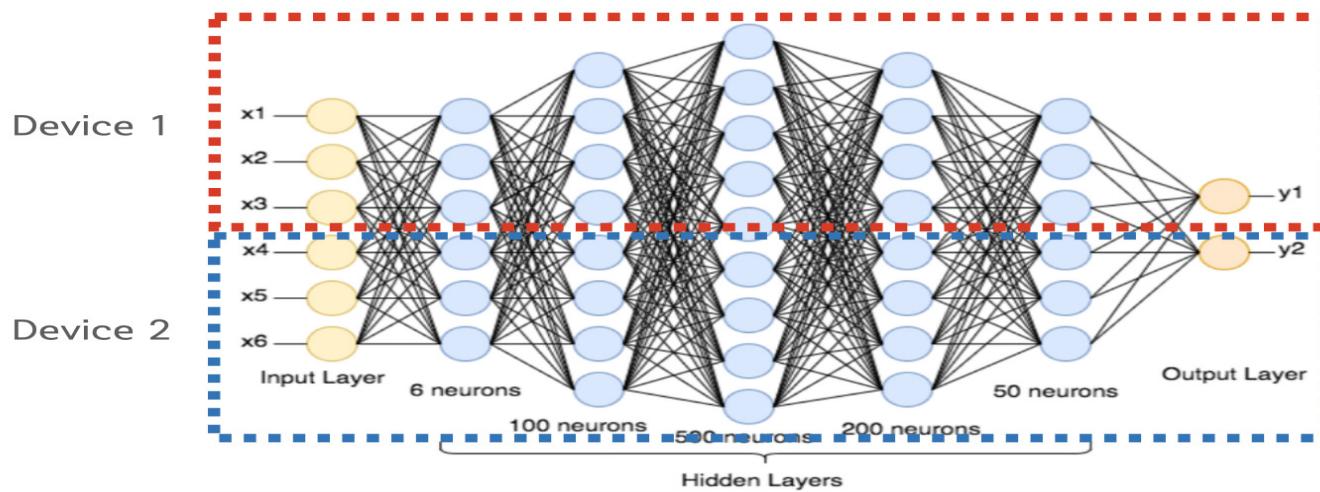
- **Inter-layer Model Parallel**
  - pipeline parallel
- **Intra-layer Model Parallel**
  - Tensor parallel
- **Multi-dimensional Parallel**
  - 3D parallel



# Inter/Intra-layer Model Parallel

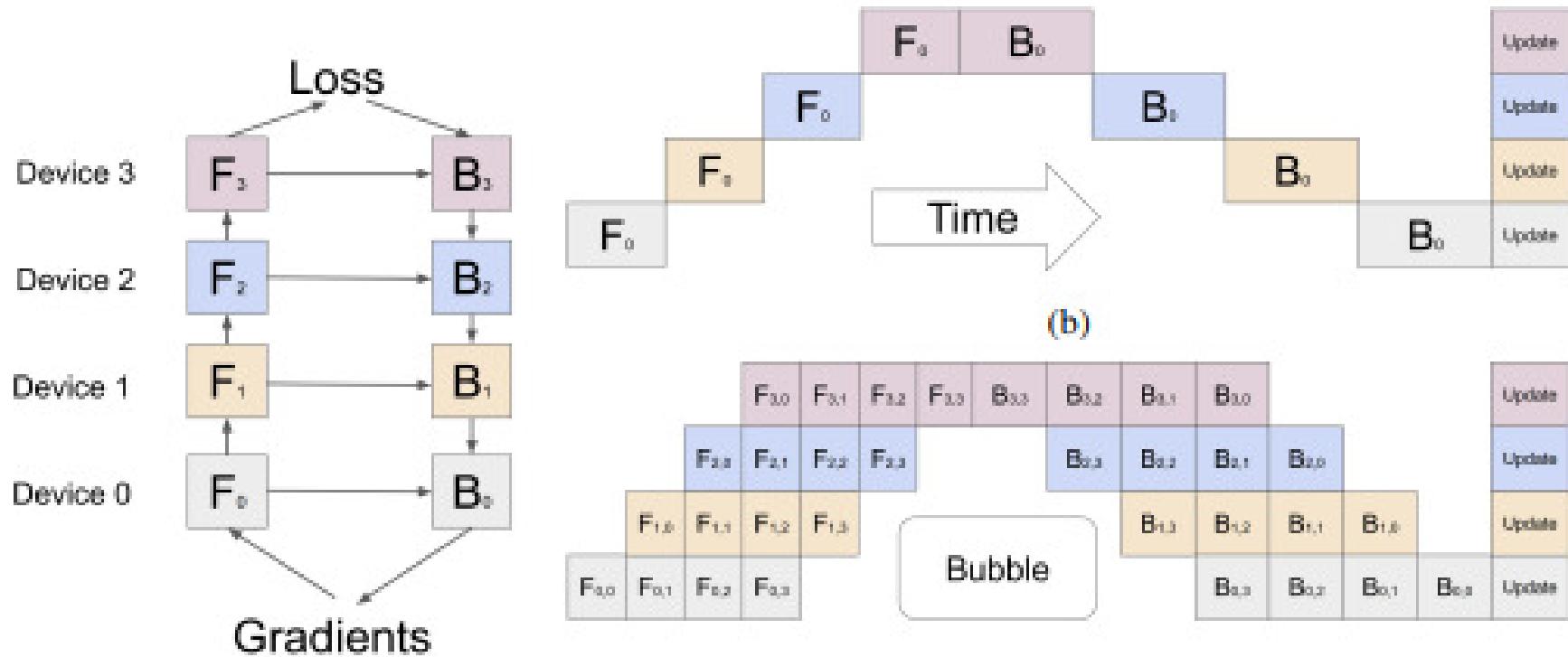


**Inter- layer  
model parallel**



**Intra- layer  
model parallel**

# Pipeline Model Parallel



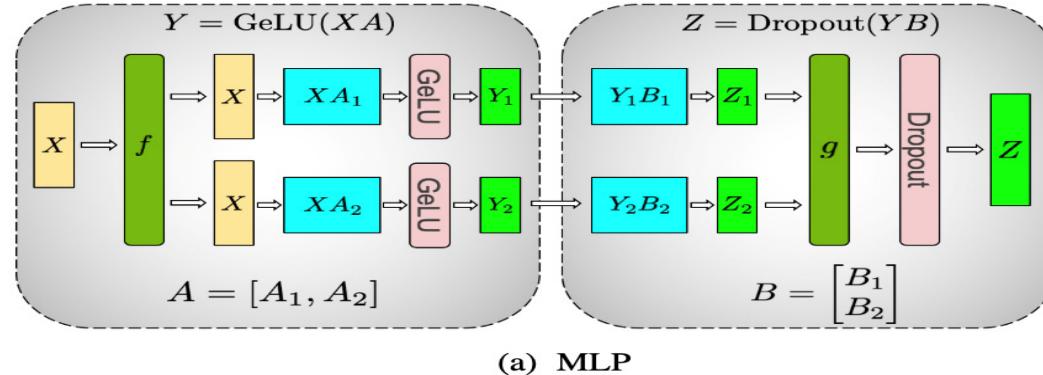
**GPipe: Easy Scaling with Micro- Batch Pipeline Parallelism**

<https://arxiv.org/pdf/1811.06965.pdf>

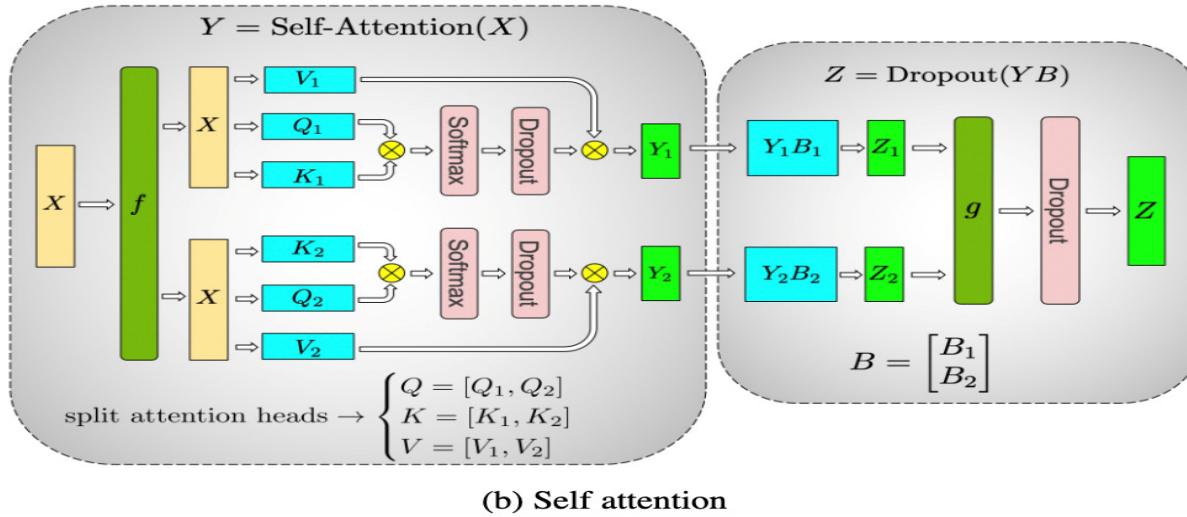
**PipeDream: Fast and Efficient Pipeline Parallel DNN Training**

<https://arxiv.org/pdf/1806.03377.pdf>

# Tensor Parallel



(a) MLP

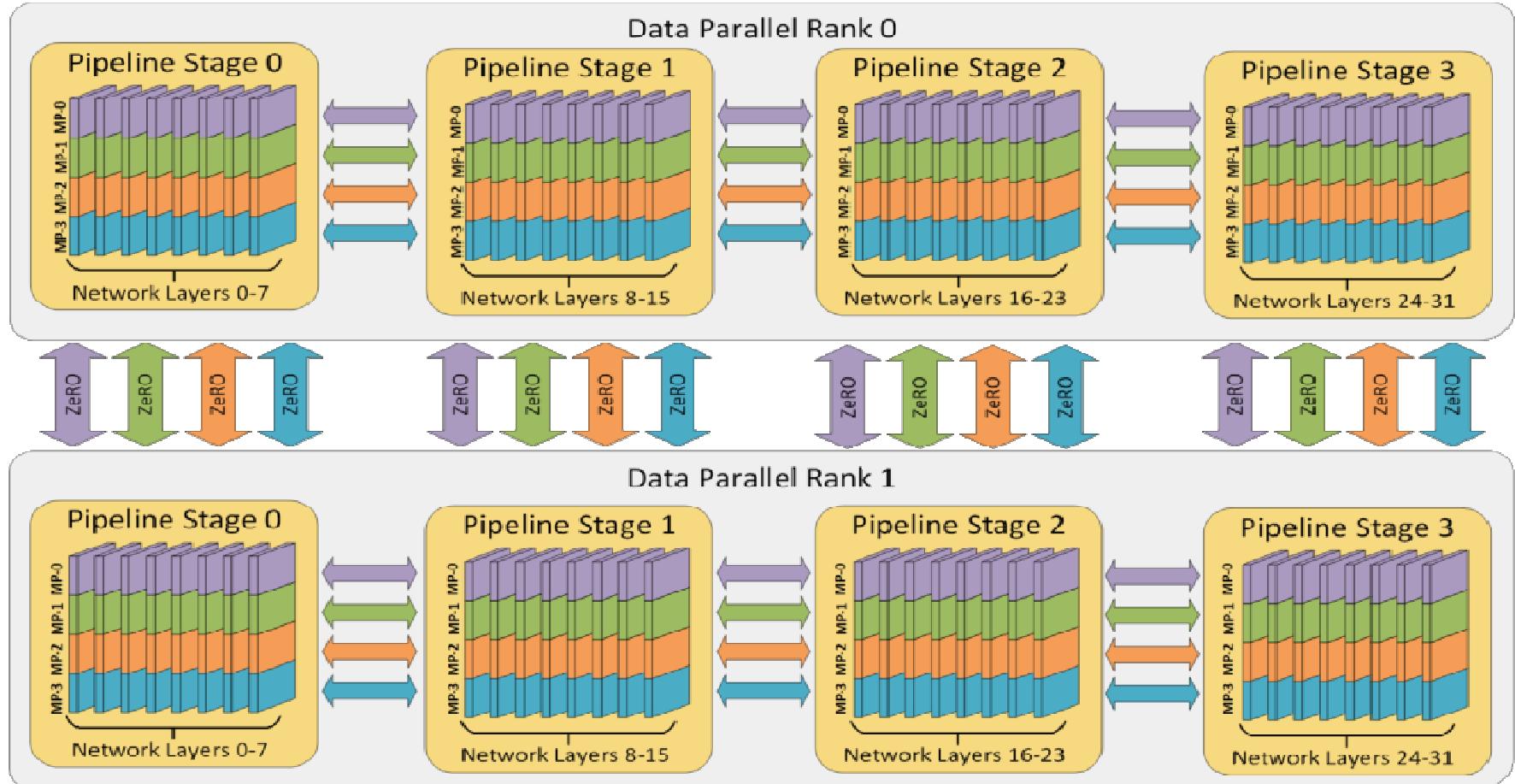


(b) Self attention

**Megatron- LM: Training Multi- Billion Parameter Language Models Using Model Parallelism**  
<https://arxiv.org/pdf/1909.08053.pdf>

# 3D Parallel

- Data Parallel + Pipeline Parallel + Tensor Parallel

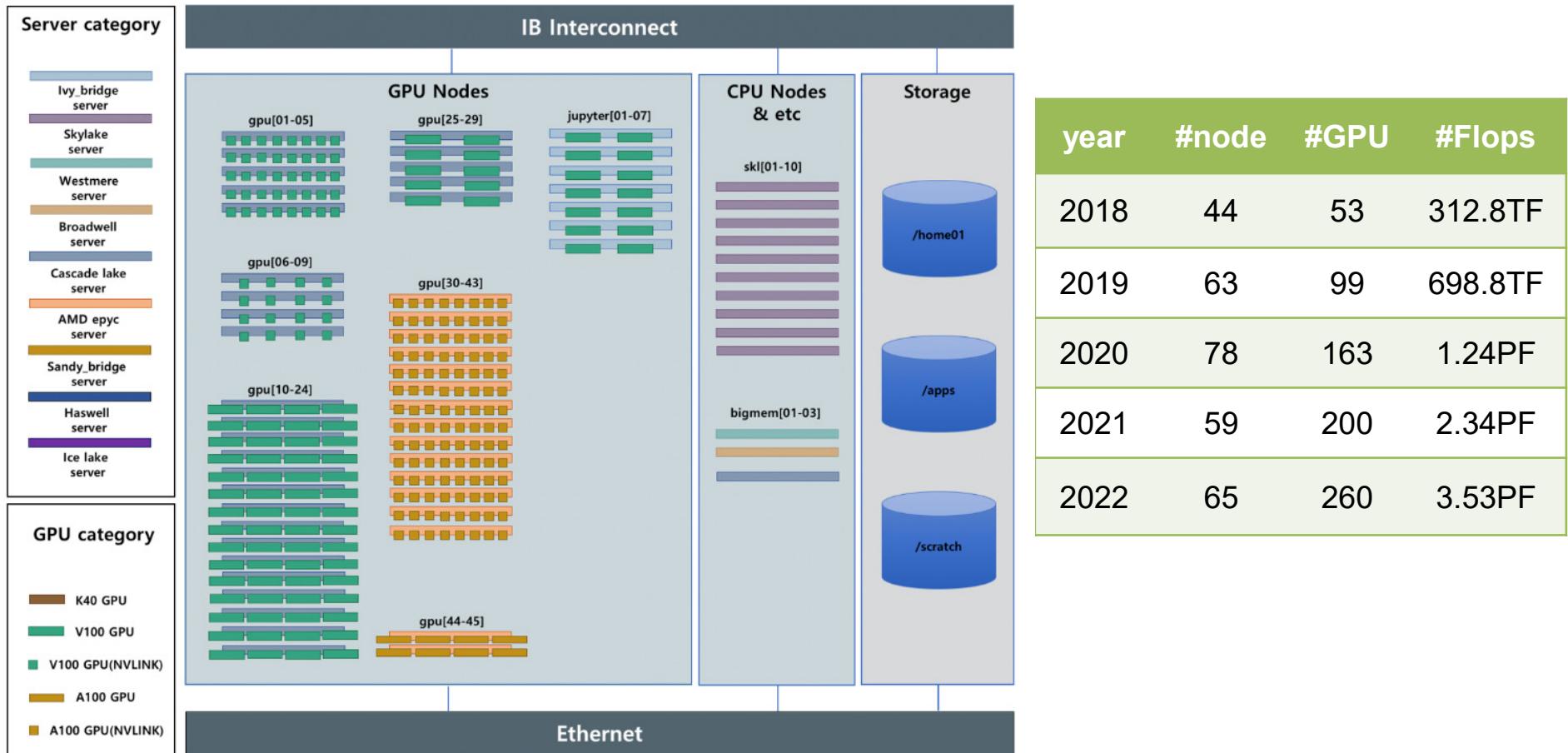


**DeepSpeed: Extreme- scale model training for everyone**

<https://www.microsoft.com/en-us/research/blog/deepspeed-extreme-scale-model-training-for-everyone>

# KISTI GPU Cluster: Neuron

# KISTI GPU Cluster: Neuron



# Slurm Queues on Neuron

Queue Name (CPU_GPU_GPU#)	#Node	#Total CPU Core	#Job Submission Limit per User	#GPU allocation Limit per User	
<b>cas_v100nv_8</b>	<b>5</b>	<b>160</b>	<b>2</b>	<b>40</b>	<b>V100 (NVlink) 8ea</b>
<b>cas_v100nv_4</b>	<b>4</b>	<b>160</b>	<b>2</b>	<b>16</b>	<b>V100 (NVlink) 4ea</b>
<b>cas_v100_4</b>	<b>15</b>	<b>600</b>	<b>4</b>	<b>40</b>	<b>V100 4ea</b>
<b>cas_v100_2</b>	<b>5</b>	<b>160</b>	<b>2</b>	<b>10</b>	<b>V100 2ea</b>
<b>amd_a100nv_8</b>	<b>14</b>	<b>868</b>	<b>4</b>	<b>64</b>	<b>A100 (Nvlink) 8ea</b>
<b>amd_a100_4</b>	<b>2</b>	<b>128</b>	<b>1</b>	<b>8</b>	<b>A100 4ea</b>
<b>amd_a100_2</b>	<b>2</b>	<b>128</b>	<b>1</b>	<b>4</b>	<b>A100 2ea</b>
<b>skl</b>	<b>10</b>	<b>360</b>	<b>2</b>	<b>-</b>	
<b>bigmem</b>	<b>3</b>	<b>120</b>	<b>1</b>	<b>-</b>	

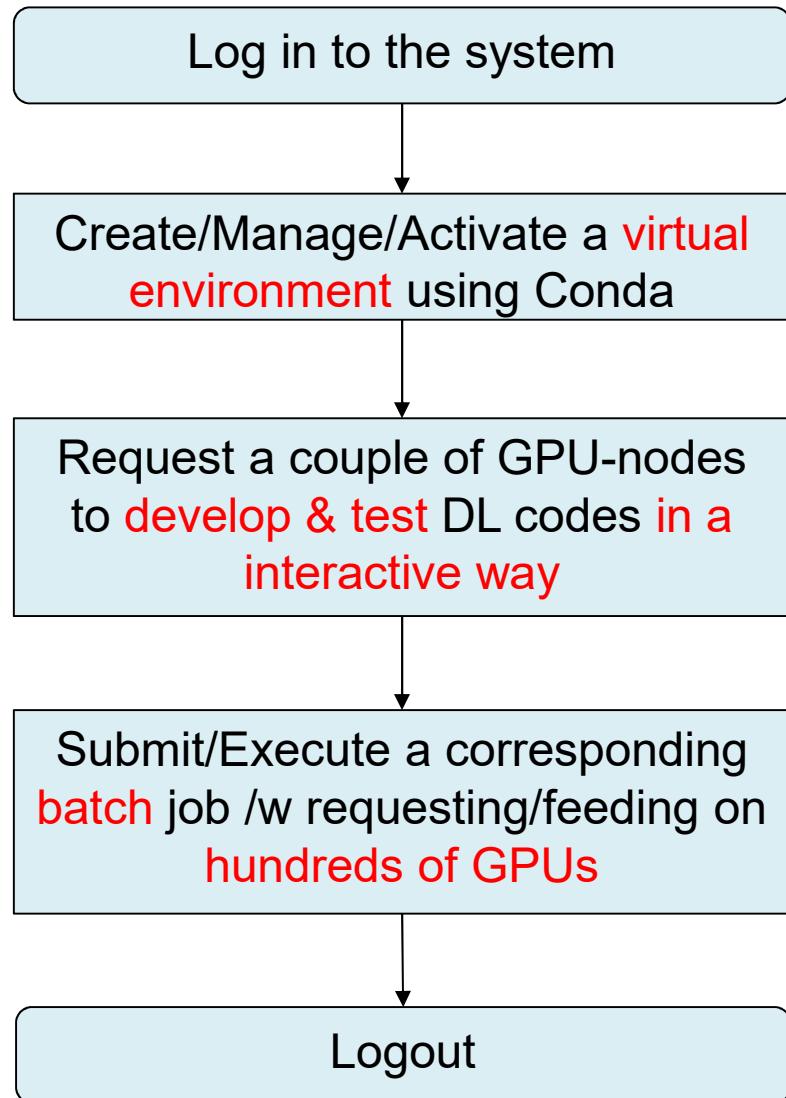
# Slurm Queues & available GPUs on Neuron

```
$ sinfo
PARTITION AVAIL TIMELIMIT NODES STATE NODELIST
jupyter    up 2-00:00:00   3 mix jupyter[02-04]
jupyter    up 2-00:00:00   1 idle jupyter01
cas_v100nv_8 up 2-00:00:00   1 mix gpu01
cas_v100nv_8 up 2-00:00:00   4 alloc gpu[02-05]
cas_v100nv_4 up 2-00:00:00   1 mix gpu09
cas_v100nv_4 up 2-00:00:00   2 alloc gpu[07-08]
cas_v100_4   up 2-00:00:00   2 mix gpu[13,17]
cas_v100_4   up 2-00:00:00  10 alloc gpu[10-12,18-24]
cas_v100_4   up 2-00:00:00   2 idle gpu[14-15]
cas_v100_2   up 2-00:00:00   1 mix gpu25
cas_v100_2   up 2-00:00:00   1 alloc gpu26
amd_a100nv_8 up 2-00:00:00   2 mix gpu[36-37]
amd_a100nv_8 up 2-00:00:00   6 alloc gpu[30,32-33,39-41]
amd_a100nv_8 up 2-00:00:00   1 idle gpu42
amd_a100_4   up 2-00:00:00   1 mix gpu44
amd_a100_4   up 2-00:00:00   1 alloc gpu45
skl        up 2-00:00:00   8 idle skl[01-06,08-09]
bigmem     up 2-00:00:00   2 idle bigmem[01-02]
exclusive   up infinite   1 mix gpu06
scidebert   up infinite   1 mix gpu35
scidebert   up infinite   1 alloc gpu34
new_service up infinite   4 down* bigmem03,jupyter[05-06],skl07
maintenance up infinite   6 idle gpu[16,29,31,38,43],jupyter07
```

- **Slurm Quick Start User Guide**
  - <https://slurm.schedmd.com/quickstart.html>
- **KISTI User Guide**
  - <https://www.ksc.re.kr/gsjw/jcs/hd>

# Distributed Training Workflow on Supercomputer

# Distributed Training Workflow on Neuron



```
## create/activate a virtual environment
$ conda create -n pt_env python=3.7
$ conda activate pt_env
$ (pt_env) conda install pytorch
$ (pt_env) python train.py
```

```
## request 2 nodes for interactive job
$ salloc --nodes=2 --time=8:00:00 --
gres=gpu:4 # available GPU-nodes allocated

## run& test ML/DL codes interactively
$ (pt_env) srun -n 8 python train_ddp.py
```

```
## submit a batch job script requesting 20
nodes with 8 GPUs each node
$ (pt_env) sbatch train_ddp_script.sh

## monitor/check the job status
$ (pt_env) squeue
```

# Conda Virtual Environment

# Conda Virtual Environment

## ▪ Download & Install Conda on /scratch/userID directory

```
## Miniconda
$ wget https://repo.anaconda.com/miniconda/Miniconda3-latest-Linux-x86_64.sh
$ chmod 755 Miniconda3-latest-Linux-x86_64.sh
$ ./Miniconda3-latest-Linux-x86_64.sh
```

- ✓ Type your Conda installation directory to “/scratch/userID/miniconda3”
  - ✓ Conda default installation directory: \$HOME/miniconda3
- ✓ Type Conda init: “yes”, which will add conda init scripts to ~/.bashrc

```
$ source ~/.bashrc # set conda path
$ conda config --set auto_activate_base false
$ which conda
/scratch/$USER/miniconda3/condabin/conda
$ conda --version
conda 4.12.0
$ ls /scratch/$USER/miniconda3
./      conda-meta/  lib/      mkspecs/   qml/      ssl/
../      doc/       libexec/  phrasebooks/ resources/  translations/
bin/    envs/     LICENSE.txt  pkgs/     sbin/      var/
....
```

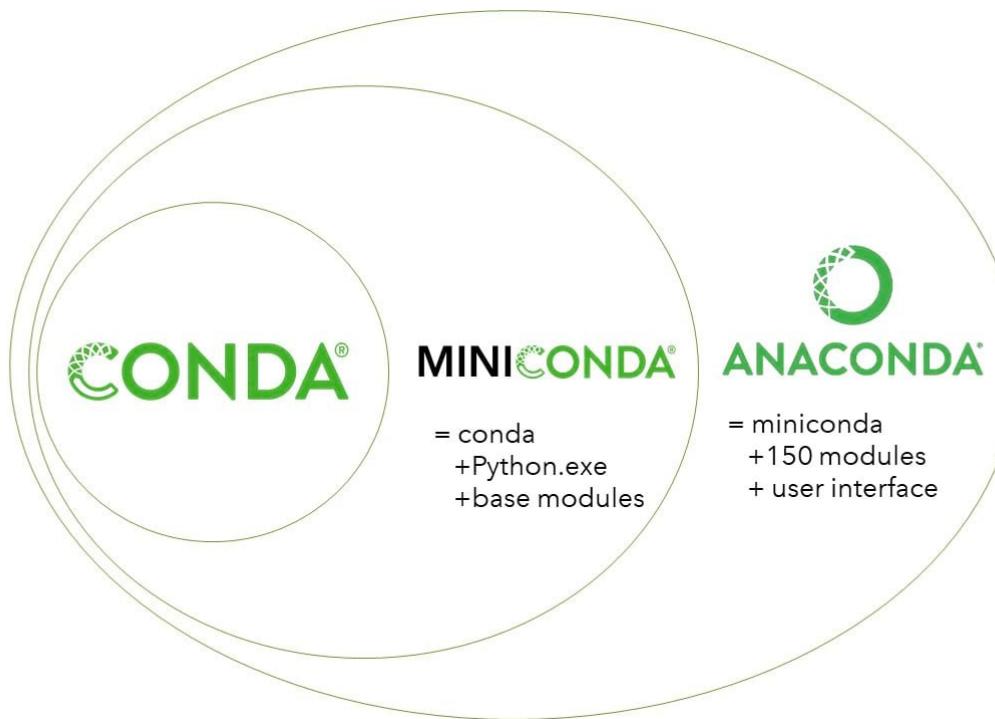
# Anaconda

---

- **Anaconda**
  - distribution of the Python and R programming languages for scientific computing
    - data science, machine learning applications, large-scale data processing, predictive analytics, etc
  - aims to simplify package management and deployment
- **Conda**
  - open source **package management** system and **environment management** system
  - runs on Windows, macOS, Linux and z/OS
  - quickly installs, runs and updates packages and their dependencies.
- **PIP**
  - package installer for Python
  - use pip to install packages from the Python Package Index and other indexes.

# Miniconda

- **a small, bootstrap version of Anaconda**
  - includes only conda, Python, the packages they depend on
- **a free minimal installer for conda**



# Anaconda vs Miniconda

---

- **Number of packages**
  - Anaconda comes with over 150 data science packages, whereas miniconda comes with only a handful
- **Interface**
  - Anaconda has a graphical user interface (GUI) called the Navigator, while miniconda has a command-line interface

```
### Anaconda Download
```

```
$ wget https://repo.anaconda.com/archive/Anaconda3-2022.10-Linux-x86_64.sh
```

```
### Miniconda Download
```

```
$ wget https://repo.anaconda.com/miniconda/Miniconda3-latest-Linux-x86_64.sh
```

# Conda command

clean	Remove unused packages and caches.
config	Modify configuration values in .condarc. This is modeled after the git config. Writes to the user .condarc file (/home01/userID/.condarc) by default.
create	Create a new conda environment from a list of specified packages.
help	Displays a list of available conda commands and their help strings.
info	Display information about current conda install.
init	Initialize conda for shell interaction. [Experimental]
install	Installs a list of packages into a specified conda environment.
list	List linked packages in a conda environment.
package	Low-level conda package utility. (EXPERIMENTAL)
remove	Remove a list of packages from a specified conda environment.
uninstall	Alias for conda remove.
run	Run an executable in a conda environment. [Experimental]
search	Search for packages and display associated information. The input is a MatchSpec, a query language for conda packages. See examples below.
update	Updates conda packages to the latest compatible version.
upgrade	Alias for conda update

# Horovod

# Horovod

---

- **distributed deep learning training framework for TensorFlow, Keras, PyTorch, and Apache MXNet**
  - developed by Uber
- **aims to make distributed deep learning fast and easy to use**

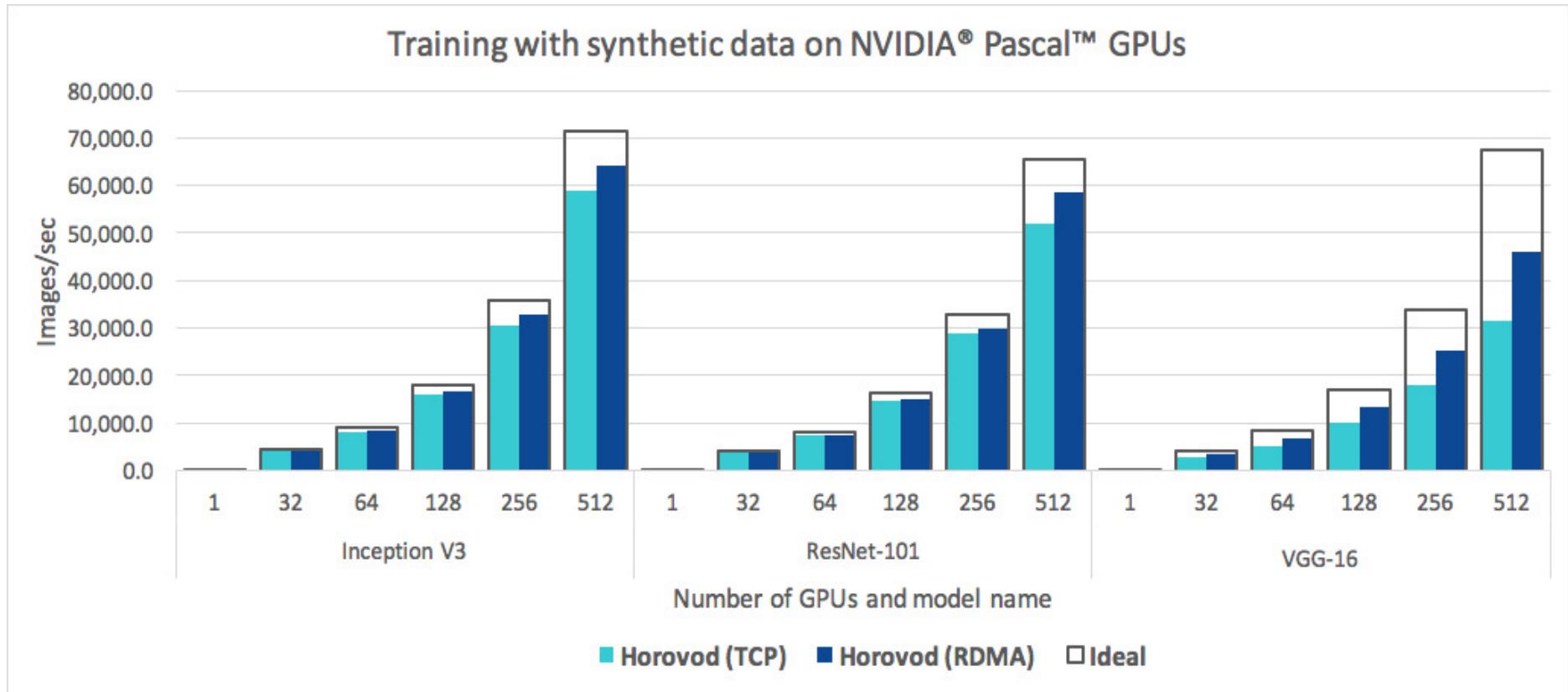
# Why Horovod?

---

- The primary motivation for this project is to make it easy to take a single-GPU training script and successfully scale it to train across many GPUs in parallel
  - ✓ How much modification does one have to make to a program to make it distributed, and how easy is it to run it?
  - ✓ How much faster would it run in distributed mode?
- Internally at Uber we found the MPI model to be much more straightforward and require far less code changes than previous solutions such as Distributed TensorFlow with parameter servers.

# Why Horovod?

- **easy to use and fast**
- **scalling with Horovod**



# Horovod Usage

---

- **5 steps/lines to be added in your code**
  - Initialize Horovod
  - Pin GPU to each worker
  - Wrap the optimizer
  - Syncroize state across workers
  - Checkpoint on the first worker

# Initialize Horovod

---

- **Tensorflow**
  - import horovod.tensorflow as hvd
  - hvd.init()
- **Kera**
  - import horovod.keras as hvd
  - hvd.init()
- **Pytorch**
  - import horovod.torch as hvd
  - hvd.init()

# Pin a GPU for each worker

---

- **Tensorflow**
  - `tf.config.experimental.set_visible_devices(gpus[hvd.local_rank()], 'GPU')`
- **Keras**
  - `tf.config.experimental.set_visible_devices(gpus[hvd.local_rank()], 'GPU')`
- **Pytorch**
  - `torch.cuda.set_device(hvd.local_rank())`

# Adjust learning rate and wrap the optimizer

---

- **Tensorflow**
  - `opt = tf.optimizers.Adam(0.01 * hvd.size())`
  - `opt = hvd.DistributedOptimizer(opt, ...)`
- **Keras**
  - `opt = keras.optimizers.Adadelta(0.01 * hvd.size())`
  - `opt = hvd.DistributedOptimizer(opt, ...)`
- **Pytorch**
  - `opt = optim.SGD(model.parameters(), 0.01 * hvd.size())`
  - `opt= hvd.DistributedOptimizer(opt, ...)`

# Synchronize state across workers

---

- **Tensorflow/Kera**
  - callbacks =  
[hvd.callbacks.BroadcastGlobalVariablesCallback(0)]
- **Pytorch**
  - hvd.broadcast\_parameters(model.state\_dict(),  
root\_rank=0)
  - hvd.broadcast\_optimizer\_state(optimizer, root\_rank=0)
- **Ensure all workers start with the same weights**

# Checkpoint on the first worker (rank 0)

---

- **Tensorflow/Keras**

- if hvd.rank() == 0:  
  callbacks.append(keras.callbacks.ModelCheckpoint(args.c  
  heckpoint\_format))

- **Pytorch**

- if hvd.rank() == 0:  
  state = {'model': model.state\_dict(),  
          'optimizer': optimizer.state\_dict(), }  
  torch.save(state, filepath)

# Pytorch Example

```
import torch
import horovod.torch as hvd

# Initialize Horovod
hvd.init()

# Horovod: pin GPU to local rank.
torch.cuda.set_device(hvd.local_rank())

# Build model.
model = Net()
model.cuda()
optimizer = optim.SGD(model.parameters())

# Horovod: wrap optimizer with DistributedOptimizer.
optimizer = hvd.DistributedOptimizer(
    optimizer,
    named_parameters=model.named_parameters())

# Horovod: broadcast parameters.
hvd.broadcast_parameters(
    model.state_dict(),
    root_rank=0)

for epoch in range(100):
    for batch_idx, (data, target) in enumerate(...):
        optimizer.zero_grad()
        output = model(data)
        loss = F.nll_loss(output, target)
        loss.backward()
        optimizer.step()
```

# Horovod execution command

---

- **MPI takes care of launching processes on all nodes**
- **Run on a 4-GPU machine**
  - \$ salloc ...
  - \$ mpirun -np 4 -H localhost:4 python train\_hvd.py
  - \$ horovodrun -np 4 -H localhost:4 python train\_hvd.py
  - \$ srun -n python train\_hvd.py
- **Run on 4 machines with 4-GPUs each**
  - \$ salloc ...
  - \$ mpirun/horovodrun -np 4 -H node1:4,node2:4,node3:4,node4:4 python train\_hvd.py
  - \$ srun -n 16 python train\_hvd.py

# Horovod Installation on Neuron

```
$ module load gcc/10.2.0 cuda/11.4 cudampi/openmpi-4.1.1 cmake/3.16.9
$ conda create -n horovod
$ conda activate horovod
$ conda install pytorch==1.12.0 torchvision==0.13.0 torchaudio==0.12.0 cudatoolkit=11.3 -c pytorch
$ pip install tensorflow-gpu==2.10.0
$ HOROVOD_GPU_OPERATIONS=NCCL HOROVOD_WITH_TENSORFLOW=1
HOROVOD_WITH_PYTORCH=1 \
HOROVOD_WITH_MPI=1 HOROVOD_WITH_GLOO=1 pip install --no-cache-dir horovod
$ horovodrun -cb
Horovod v0.26.1:
```

*Available Frameworks:*

- [X] TensorFlow
- [X] PyTorch
- [ ] MXNet

*Available Controllers:*

- [X] MPI
- [X] Gloo

*Available Tensor Operations:*

- [X] NCCL
- [ ] DDL
- [ ] CCL
- [X] MPI
- [X] Gloo

# Horovod execution on Neuron

1) request an allocation of available on neuron

```
$ salloc --partition=amd_a100nv_8 -J debug --nodes=2 --time=2:00:00 --gres=gpu:4 --comment=horovod
```

2) load modules

```
$ module load gcc/10.2.0 cuda/11.4 cudampi/openmpi-4.1.1 cmake/3.16.9
```

3) activate the horovod virtual environment

```
$ conda activate horovod
```

4) run horovod applications

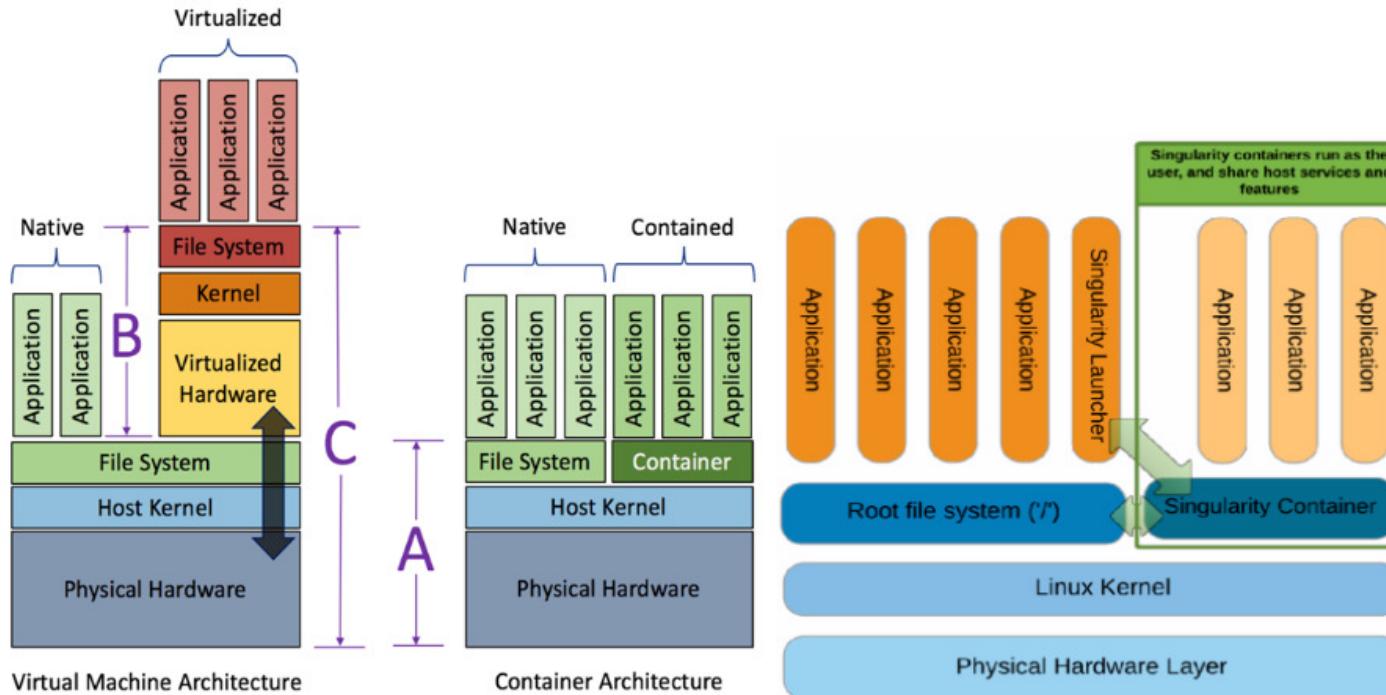
```
(horovod) $ mpirun/horovod -np 8 -H gpu#:4,gpu#:4 python train_hvd.py
```

```
(horovod) $ srun -n python train_hvd.py
```

# Singularity Container

# Singularity

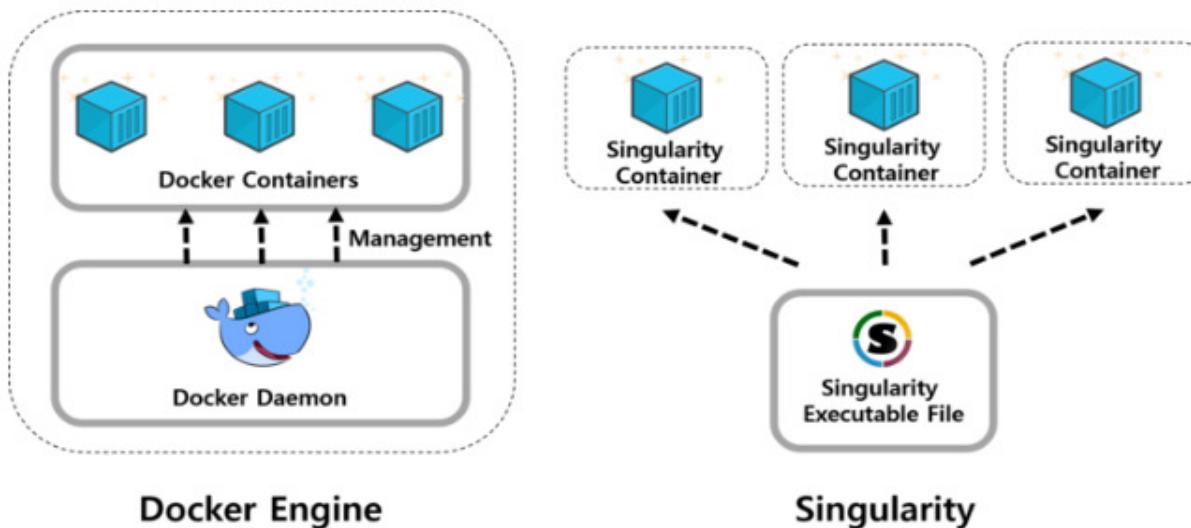
- a container platform for HPC



Container-based applications have *direct access* to the host kernel and hardware and, thus, are able to achieve similar performance to native applications. In contrast, VM-based applications only have *indirect access* via the guest OS and hypervisor, which creates a significant performance overhead.

# Why Singularity?

- A container platform for HPC
  - Each container is a single image file
  - No root owned daemon processes
  - Support shared/multi-tenant resource environment
  - Support HPC hardware
    - Infiniband, GPUs
  - Support HPC applications
    - MPI



# Running Horovod using Singularity on Neuron

- No bother to deal with conda & horovod
- Just allocate nodes using salloc and run singularity container. **That's it!!!**

```
$ salloc --partition=amd_a100nv_8 -J debug --nodes=2 --time=2:00:00 --  
gres=gpu:4 --comment=horovod  
  
$ srun -n 8 singularity run --nv  
/apps/applications/singularity_images/ngc/pytorch_22.03-hd-py3.sif python  
pytorch_imagenet_resnet50.py  
/usr/bin/rm: cannot remove '/usr/local/cuda/compat/lib': Read-only file system  
/usr/bin/rm: cannot remove '/usr/local/cuda/compat/lib': Read-only file system  
.....  
Train Epoch #10: 0% | 1/5005 [00:13<19:03:52, 13.72s/it, loss=2.29, accurTrain Epoch  
#10: 0% | 1/5005 [00:13<19:03:52, 13.72s/it, loss=2.22, accurTrain Epoch #10: 0% |  
2/5005 [00:13<19:03:38, 13.72s/it, loss=2.16, accurTrain Epoch #10: 0% | 3/5005  
[00:13<5:00:52, 3.61s/it, loss=2.16, accurTrain Epoch #10: 0% | 3/5005 [00:13<5:00:52,  
3.61s/it, loss=2.17, accurTrain Epoch #10: 0% | 4/5005 [00:13<5:00:48, 3.61s/it, loss=2.23,  
accura  
.....
```

# Singularity Usage on Neuron

---

- **Web site:** <https://www.ksc.re.kr/gsjw/jcs/hd>
- **DDL training job scripts directory**
  - /apps/applications/singularity\_images/examples
- **Singularity Container Images directory**
  - /apps/applications/singularity\_images/ngc
- **Pytorch examples directory**
  - Single node
    - /apps/applications/singularity\_images/examples/pytorch/resnet50v1.5
  - Multiple nodes
    - /apps/applications/singularity\_images/examples/horovod/examples/pytorch
- **Imagenet datasets directory**
  - Training data
    - /apps/applications/singularity\_images/imagenet/train
  - Validation data
    - /apps/applications/singularity\_images/imagenet/val

---

# Thank you

Contact: Soonwook Hwang <[hwang@kisti.re.kr](mailto:hwang@kisti.re.kr)>

[github.com/hwang2006/KISTI-DL-tutorial-using-horovod](https://github.com/hwang2006/KISTI-DL-tutorial-using-horovod)