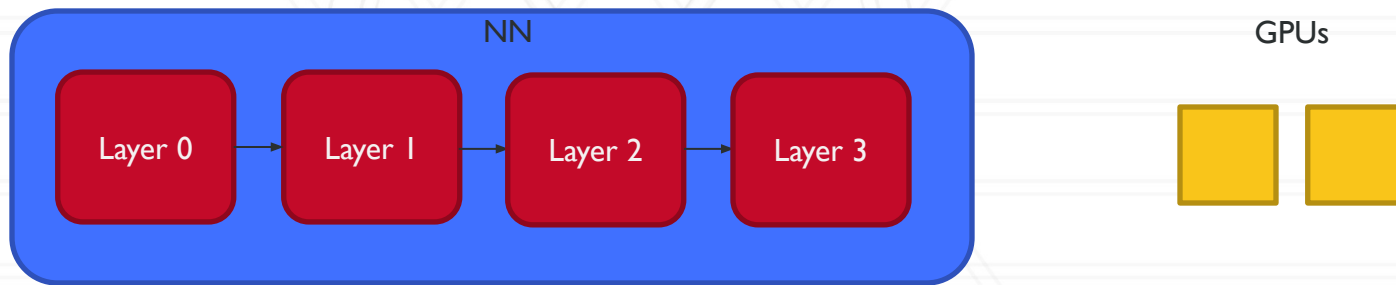# Distributed Training of Deep Neural Networks

Abhinav Bhatele, Siddharth Singh, Daniel Nichols
Department of Computer Science

P-S
S-G PARALLEL SOFTWARE
AND SYSTEMS GROUP

UNIVERSITY OF
MARYLAND

# Inter-layer parallelism

- Map contiguous subsets of layers to GPUs.

# Inter-layer parallelism

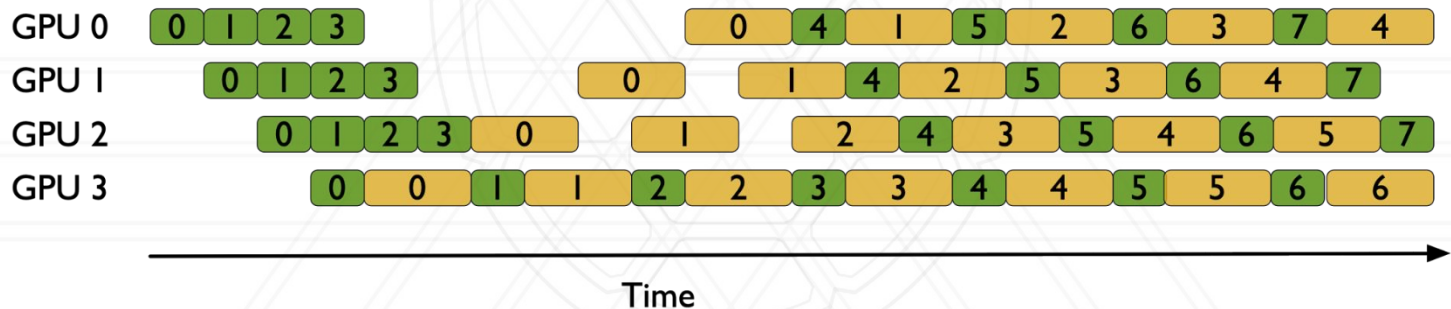| Layer 0 | Layer 1 | | Layer 2 | Layer 3 |

Point-to-point communication of activations and their gradients

# Inter-layer 'parallelism' ?

- Break batch intro micro-batches.
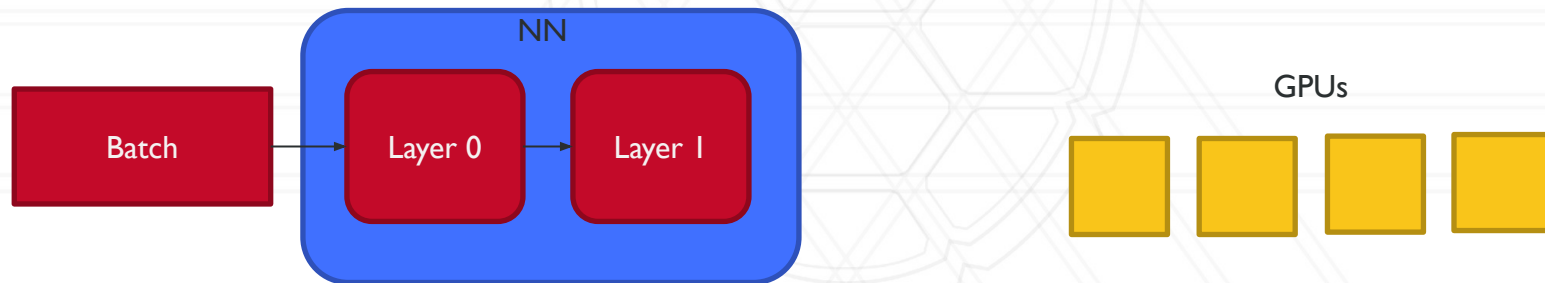- Process micro-batches in a pipelined fashion.

# Running the code
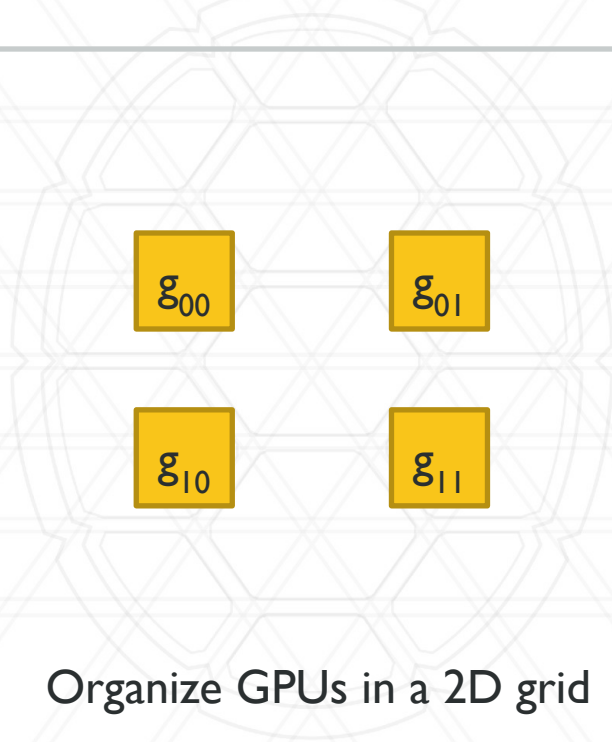
- **Code –** `train_axonn_inter_layer.py`

```
cd session_4_inter_layer_parallelism
sbatch --reservation=isc2024 run.sh
```

# Hybrid parallelism

- AxoNN can combine inter/intra-layer parallelism with data parallelism too.
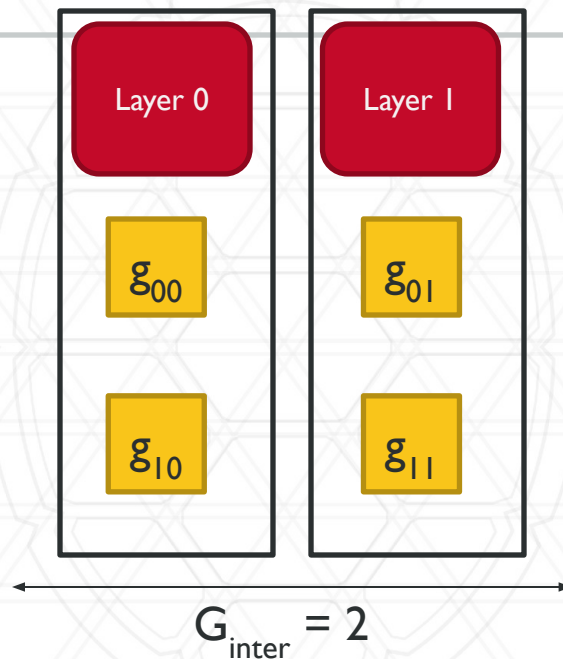
- Let us discuss inter-layer + data parallelism.

# Inter-Layer+Data Parallelism



$g_{00}$  $g_{01}$

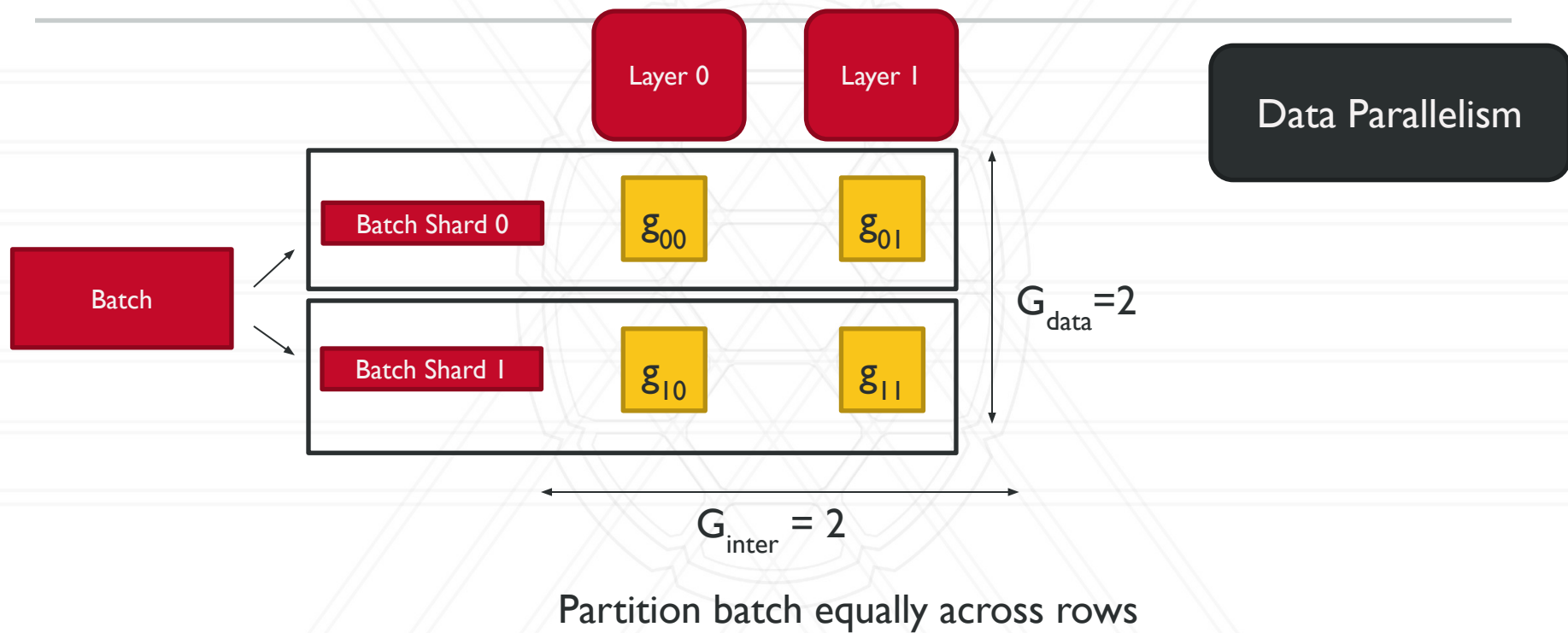$g_{10}$  $g_{11}$

Organize GPUs in a 2D grid

# Inter-Layer+Data Parallelism

Layer 0

Layer 1

$g_{00}$

$g_{01}$

$g_{10}$

$g_{11}$

Inter-Layer
Parallelism

$G_{inter} = 2$

Partition layers equally across columns

# Inter-Layer+Data Parallelism

Layer 0    Layer 1

Data Parallelism

Batch Shard 0    $g_{00}$    $g_{01}$

Batch

Batch Shard 1    $g_{10}$    $g_{11}$

$G_{data} = 2$

$G_{inter} = 2$

Partition batch equally across rows

# Running the code

- **Code –** `train_axonn_inter_layer.py`

```
cd session_4_inter_layer_parallelism
HYBRID_PARR=true sbatch --reservation=isc2024 run.sh
```