

The background features abstract, overlapping geometric shapes in various shades of green, ranging from light lime to dark forest green. These shapes are primarily located on the left and right sides of the slide, framing the central text. The overall aesthetic is clean and modern.

# ZeRO: Memory Optimizations Toward Training Trillion Parameter Models

# Background

- ▶ DNN model computation involves difference types of data:
- ▶ Activation
- ▶ Parameter
- ▶ Gradient
- ▶ Optimizer

# Background

- ▶ Different parallelisms partitions different types of data

	Partitioned Data Type	Benefits	Drawbacks
Data Parallelism	Activation	Low communication volume	High memory volume
Model Parallelism	Parameter Gradient Optimizer	Low memory volume	High communication volume Small compute granularity
Pipeline Parallelism	Parameter Gradient Optimizer	Low communication volume	Imbalance memory volume Bubbles Latency cannot scale down

# Background

- ▶ As models become larger, existing parallelisms are ineffective
- ▶ Data parallelism: size of data that cannot be partitioned exceeds the GPU memory
- ▶ Model parallelism: large models requires many GPUs and fine-grained partitioning, causing excessive communication overhead
- ▶ Pipeline parallelism: each layer is still too large; and with deeper layers, the devices handling the first few layers must hold more gradients/optimizers for the ongoing batches.

# Challenge

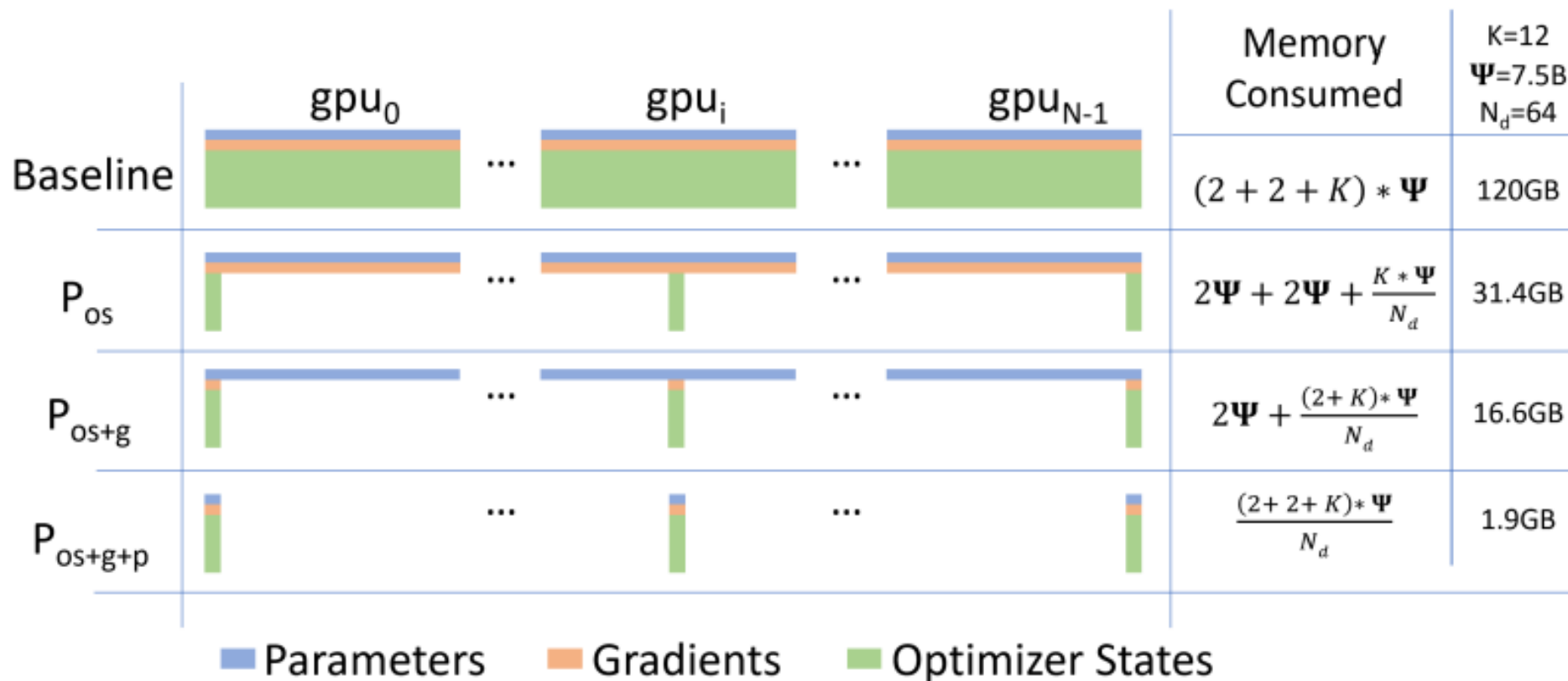
- ▶ Reduce communication overhead
- ▶ Partition all types of data to save memory in a balanced way
- ▶ Preserve large operator size per device for TensorCore utilization

# Observation

- ▶ In data/model parallelism, all types of data may be duplicated.
- ▶ In DNN model executions, most data stays idle when other layers are be computed.
- ▶ Thus, this data can be scatted into multiple devices and gathered when they are demanded

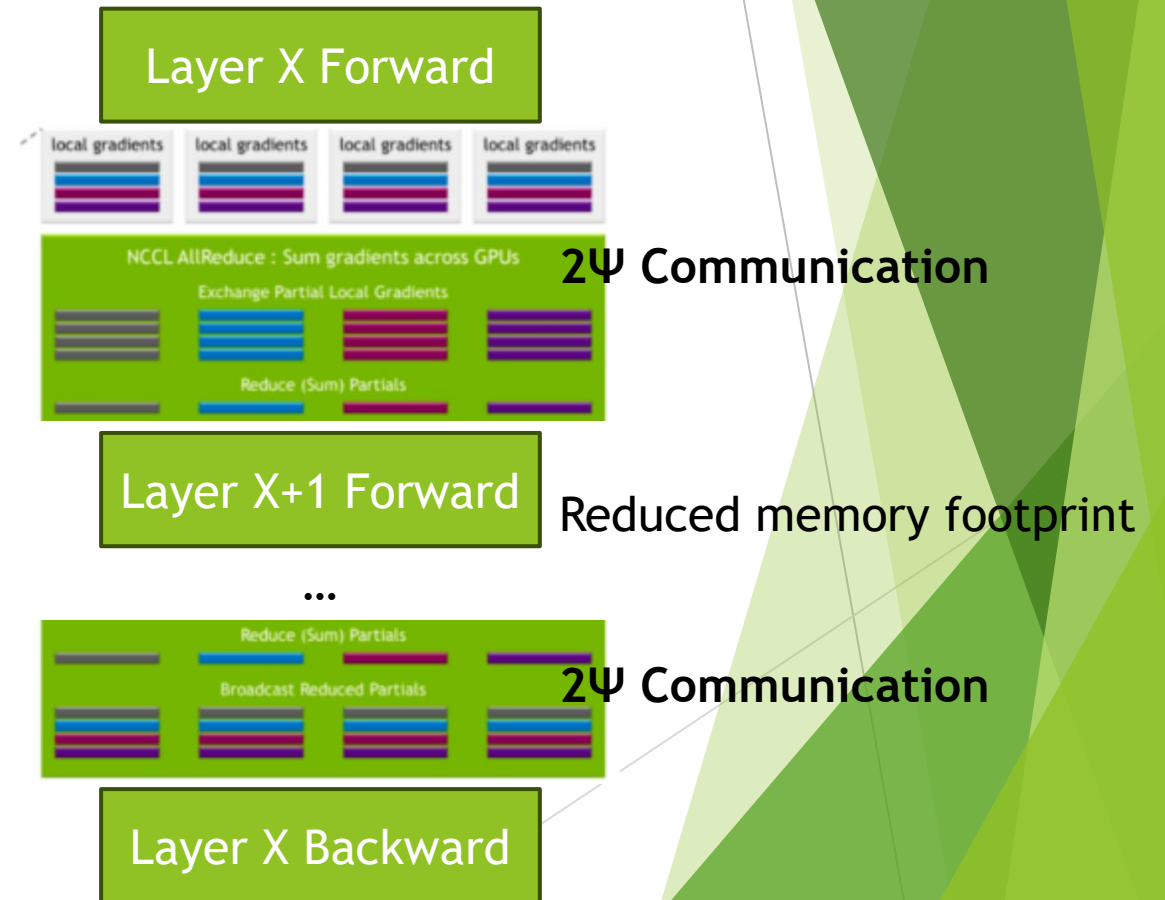
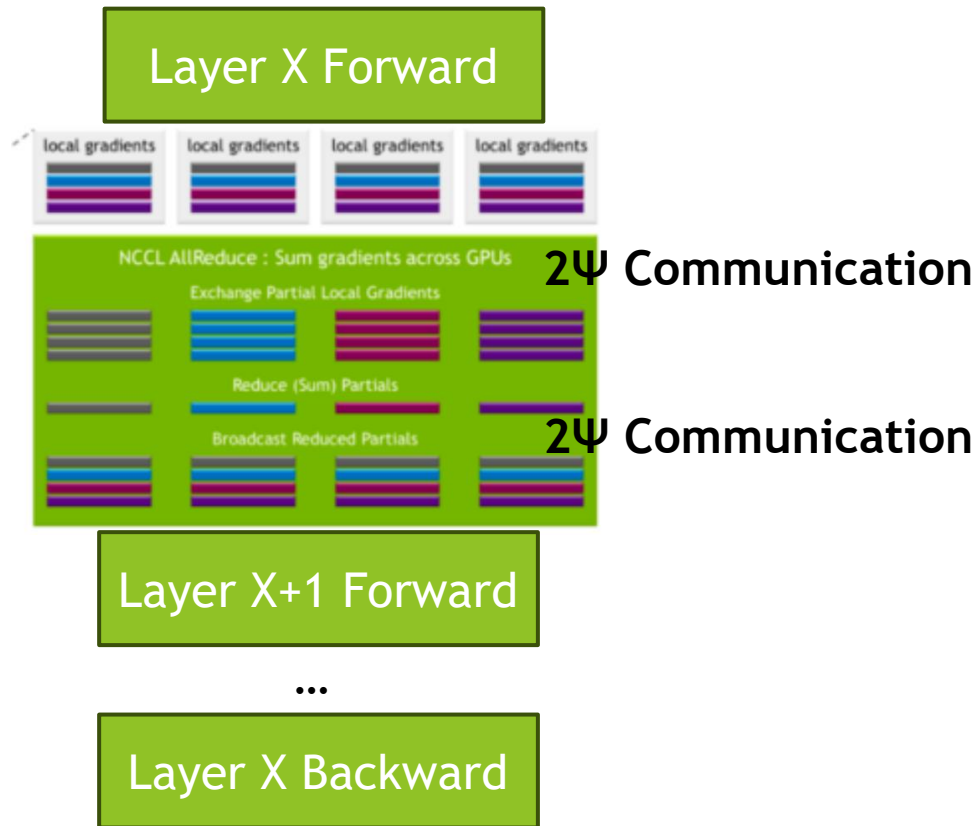
# ZeRO-DP

$\Psi$  : model size (number of parameters)  
 $K$  : memory multiplier of optimizer states,  
 $N_d$  : DP degree (number of devices)  
 $2$  : size of FP16



# ZeRO-DP (Grad / OS)

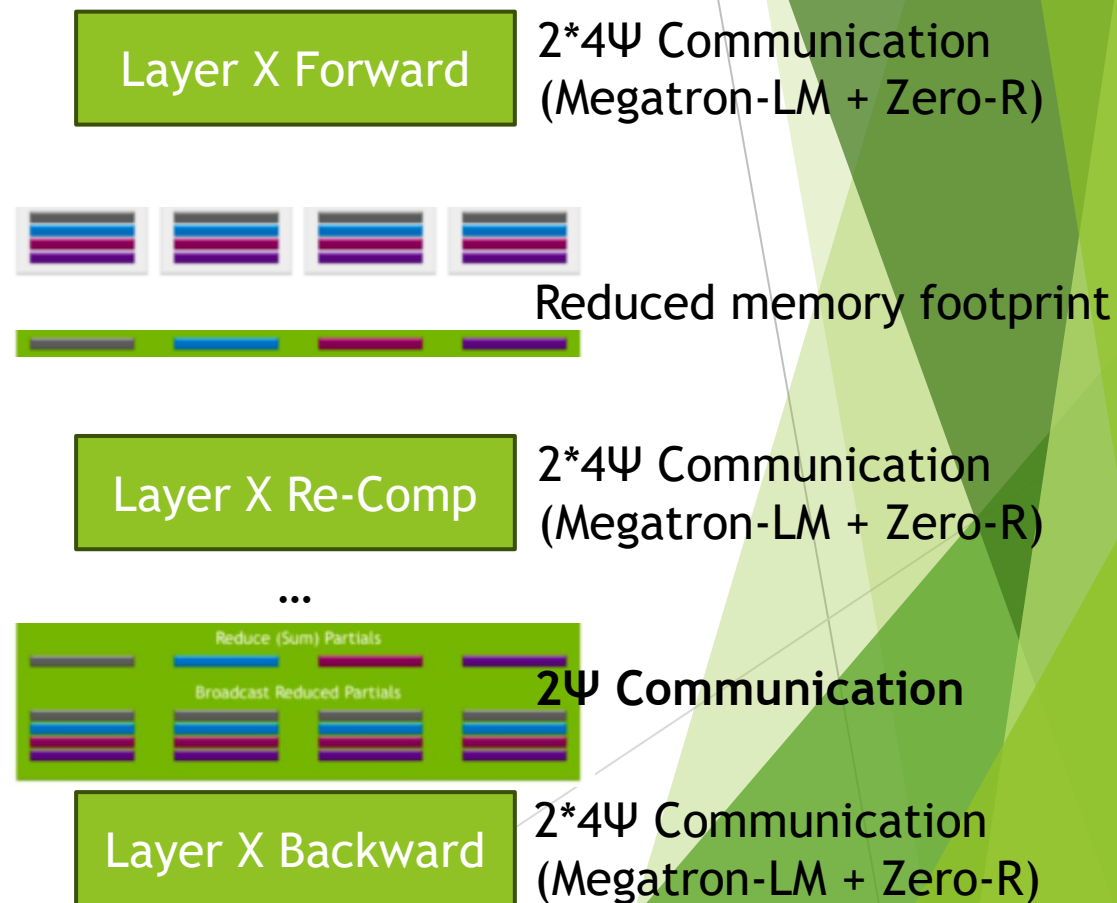
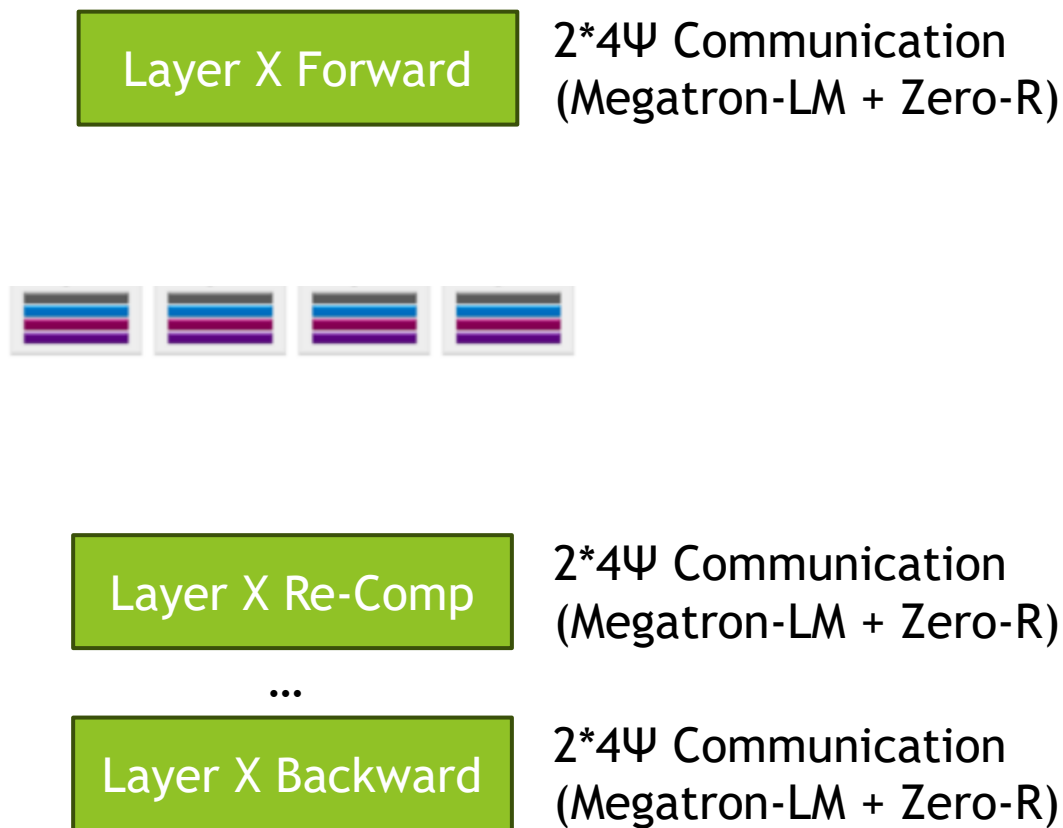
$\Psi$  : model size (number of parameters)  
 $K$  : memory multiplier of optimizer states,  
 $N_d$  : DP degree (number of devices)  
 $2$  : size of FP16





# ZeRO-R / DP (Para)

$\Psi$  : model size (number of parameters)  
 $K$  : memory multiplier of optimizer states,  
 $N_d$  : DP degree (number of devices)  
 $2$  : size of FP16



# Data Management

## ▶ Defragmentation

- ▶ Short lived memory: discarded activation
- ▶ Long lived memory: checkpointed activation, model states
- ▶ Placing short-lived and long-lived memory together causes fragmentation
- ▶ Solution: reserve contiguous memory for long lived memory

## ▶ Buffer size

- ▶ To enable efficient reduction, Nvidia fuses all parameters into one buffer
- ▶ Buffer of a large model can be too large for GPU memory
- ▶ Solution: partition the buffer into fixed-sized buffers, which still have high efficiency in reduction

# Evaluation

- Performance scales with number of GPUs

