



**TECHNION**  
Israel Institute  
of Technology



# Tensor-Parallelism with Partially Synchronized Activations

Itay Lamprecht, Asaf Karnieli, Yair Hanani, Niv Giladi and Daniel Soudry

The Thirty-Ninth Annual Conference on Neural Information Processing Systems  
(NeurIPS 2025)



# Background and Motivation

# Background and Motivation

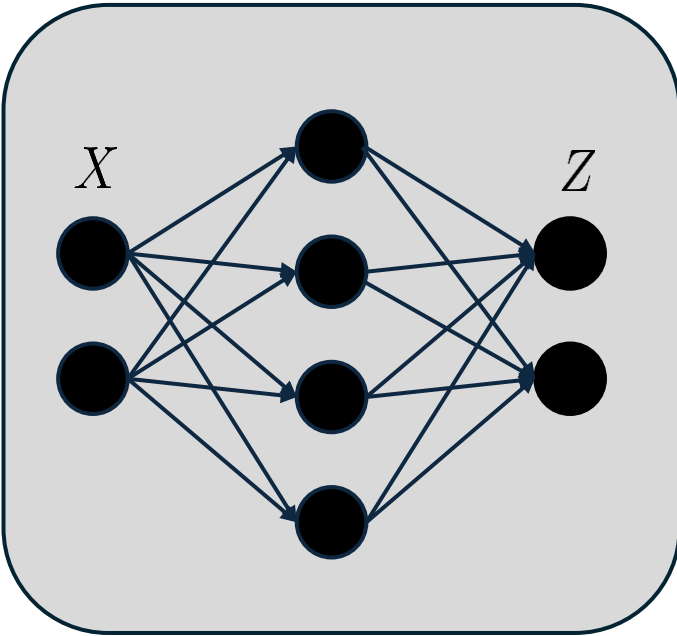
- Parallelism for Deep Neural Networks
  - Data Parallelism
  - Context Parallelism
  - Pipeline Parallelism
  - Tensor Parallelism

# Background and Motivation

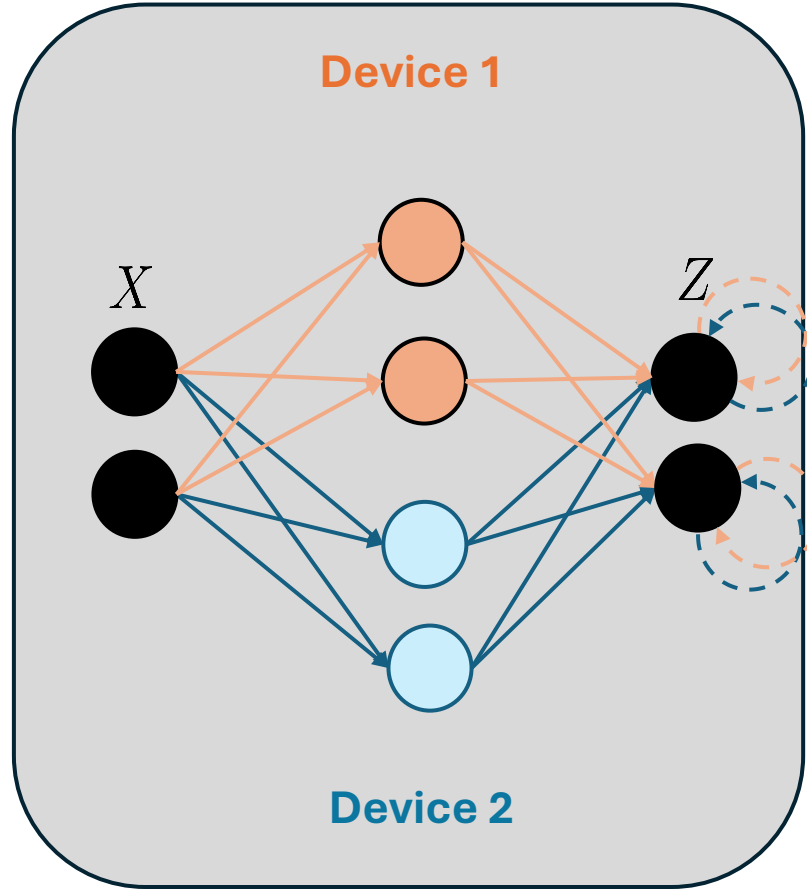
- Parallelism for Deep Neural Networks
  - Data Parallelism
  - Context Parallelism
  - Pipeline Parallelism
  - Tensor Parallelism

# Tensor-Parallelism

Original Model



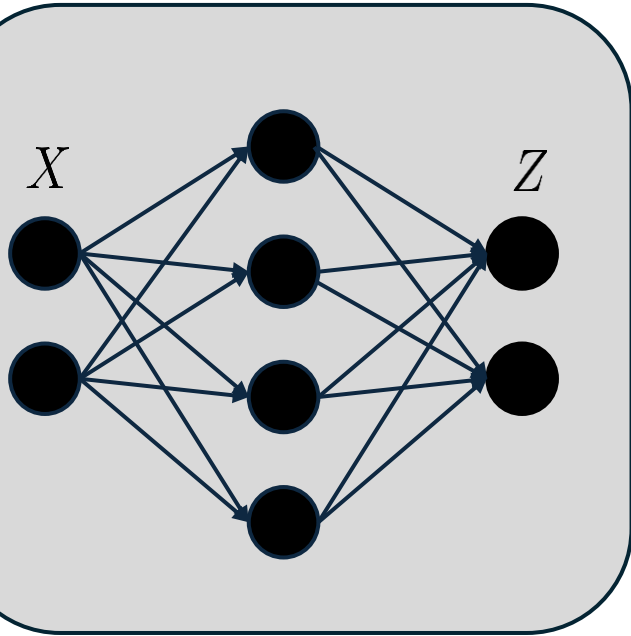
Tensor-Parallel



----- Values sent to device  
—— Values calculated on device

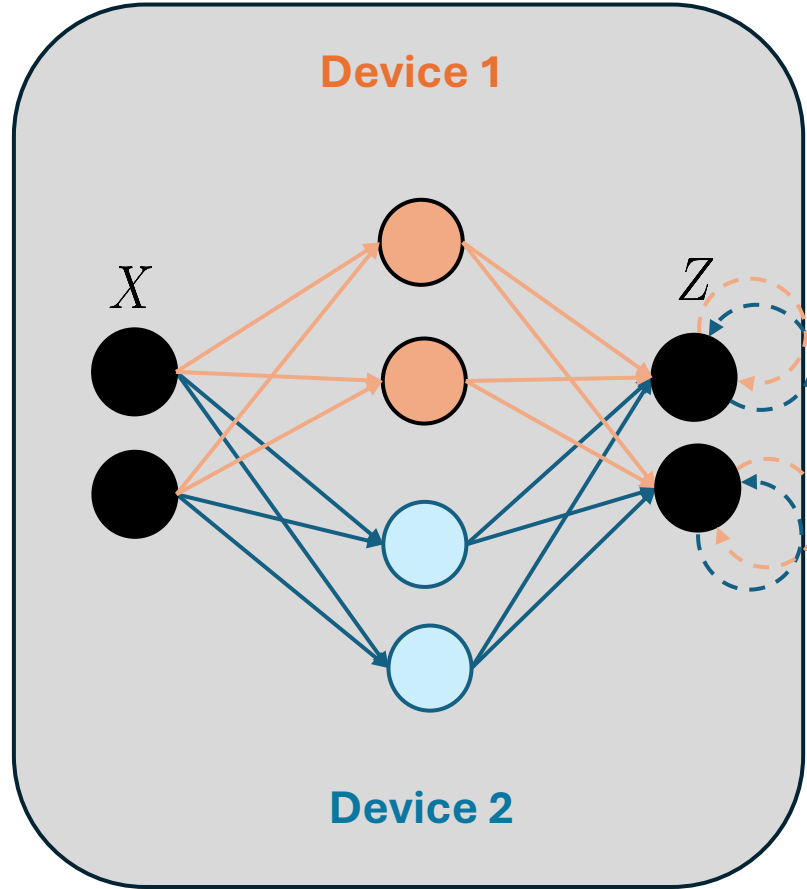
# Tensor-Parallelism

Original Model



----- Values sent to device  
——— Values calculated on device

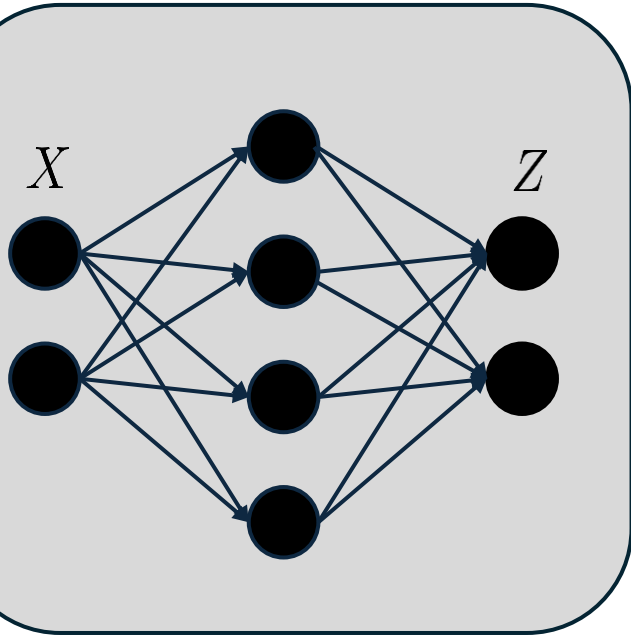
Tensor-Parallel



Memory Efficient

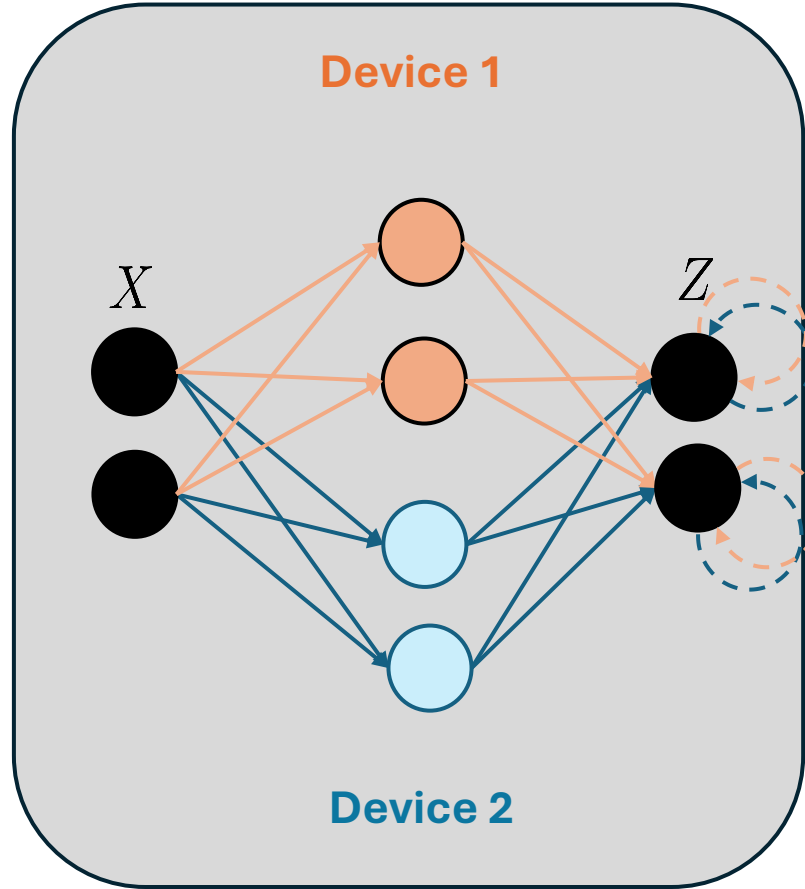
# Tensor-Parallelism

Original Model



----- Values sent to device  
——— Values calculated on device

Tensor-Parallel



Memory Efficient

Costly Communication  
On Critical Path

# Recent works

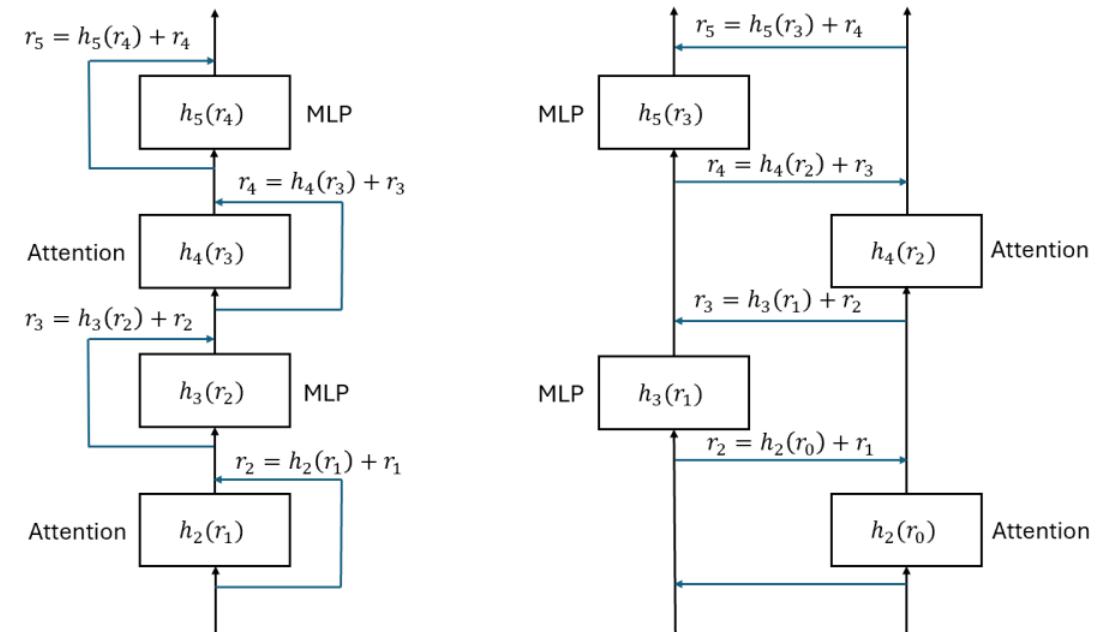
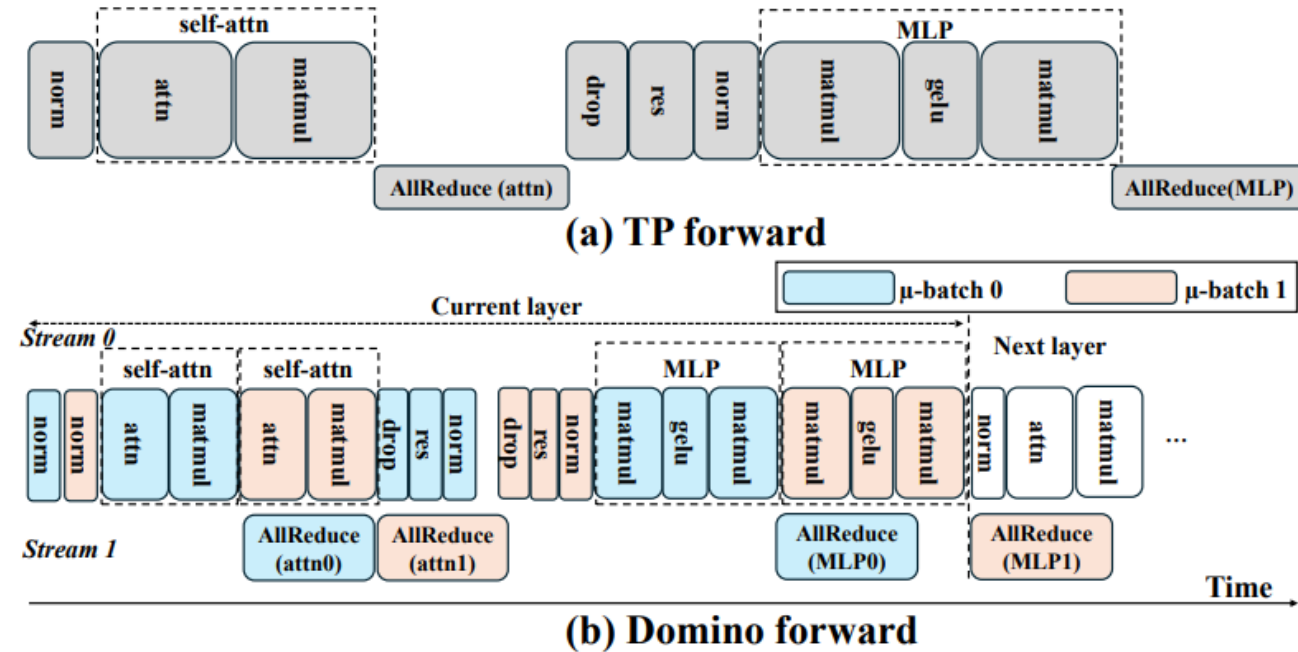
- Pipelining Communication and Computation
- Asynchronous Training
- Compression methods
- Inference time methods



# Recent works

- Pipelining Communication and Computation
- Asynchronous Training
- Compression methods
- Inference time methods

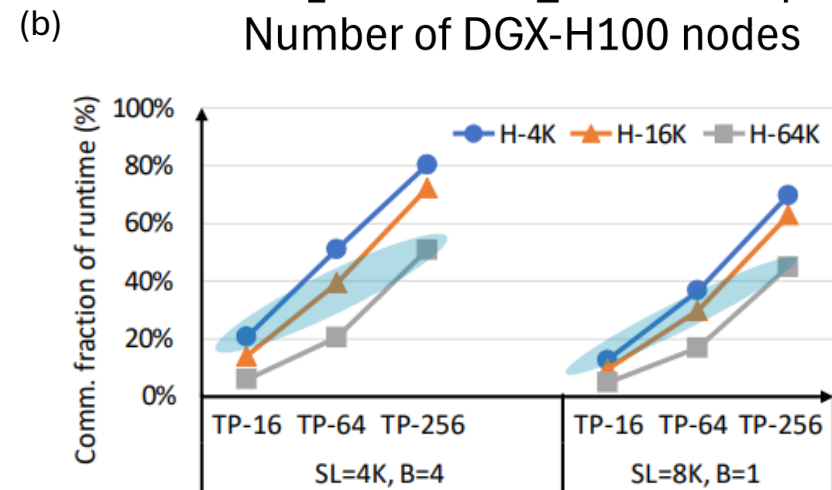
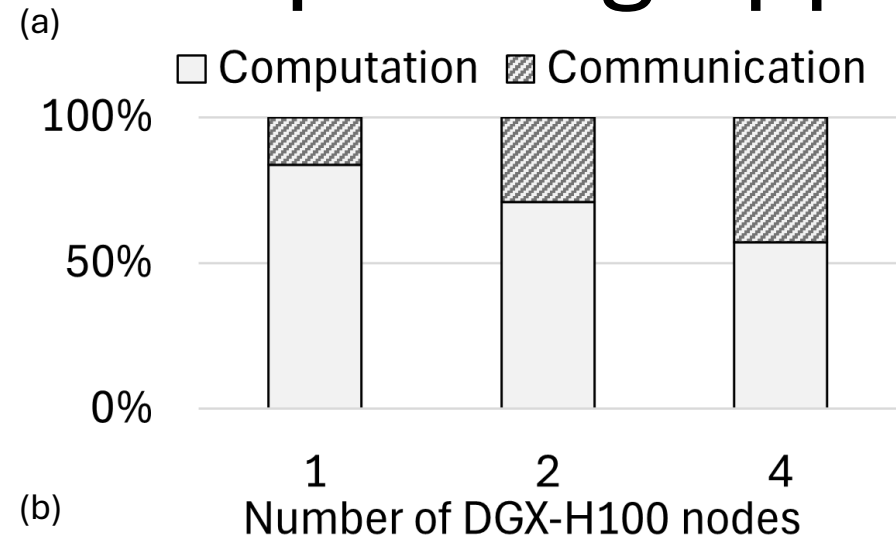
# Pipelining Approaches



**Ladder-residual: parallelism-aware architecture for accelerating large model inference with communication overlapping.** Muru Zhang, Mayank Mishra, Zhongzhu Zhou, William Brandon, Jue Wang, Yoon Kim, Jonathan Ragan-Kelley, Shuaiwen Leon Song, Ben Athiwaratkun, and Tri Dao. 2025.

**Domino: Eliminating communication in llm training via generic tensor slicing and overlapping** Guanhua Wang, Chengming Zhang, Zheyu Shen, Ang Li, and Olatunji Ruwase. 2024.

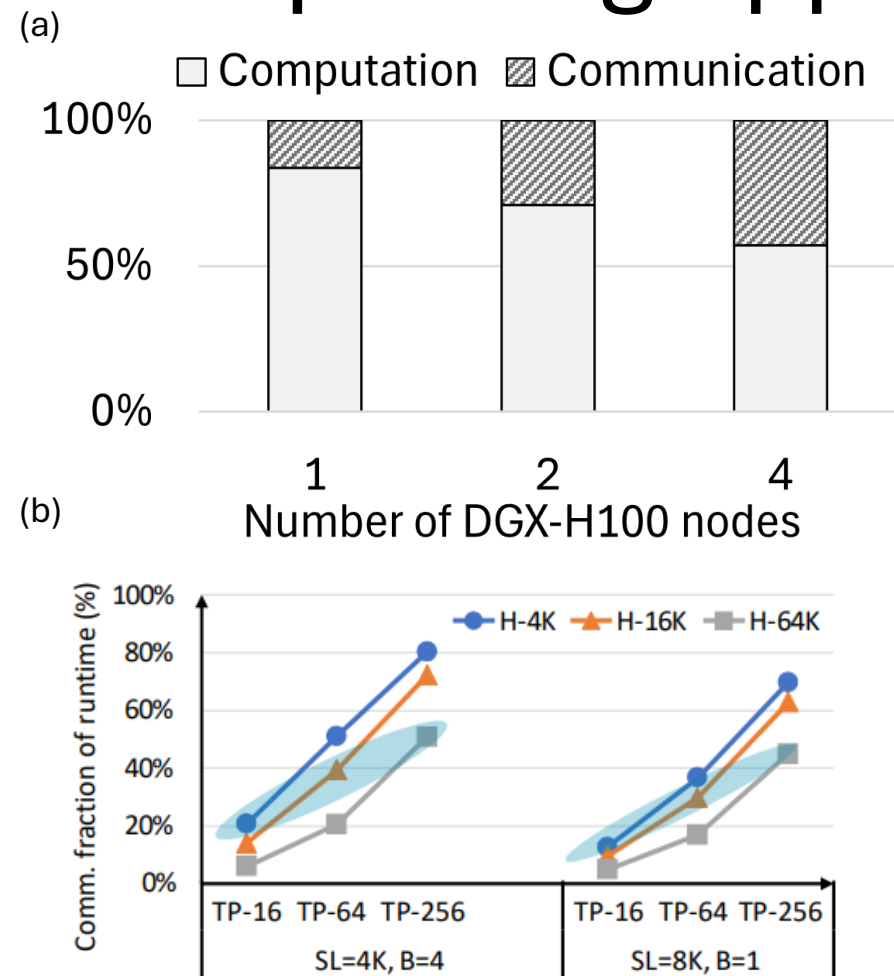
# Pipelining Approaches



(a) **Domino: Eliminating communication in llm training via generic tensor slicing and overlapping**  
Guanhua Wang, Chengming Zhang, Zheyu Shen, Ang Li, and Olatunji Ruwase. 2024.

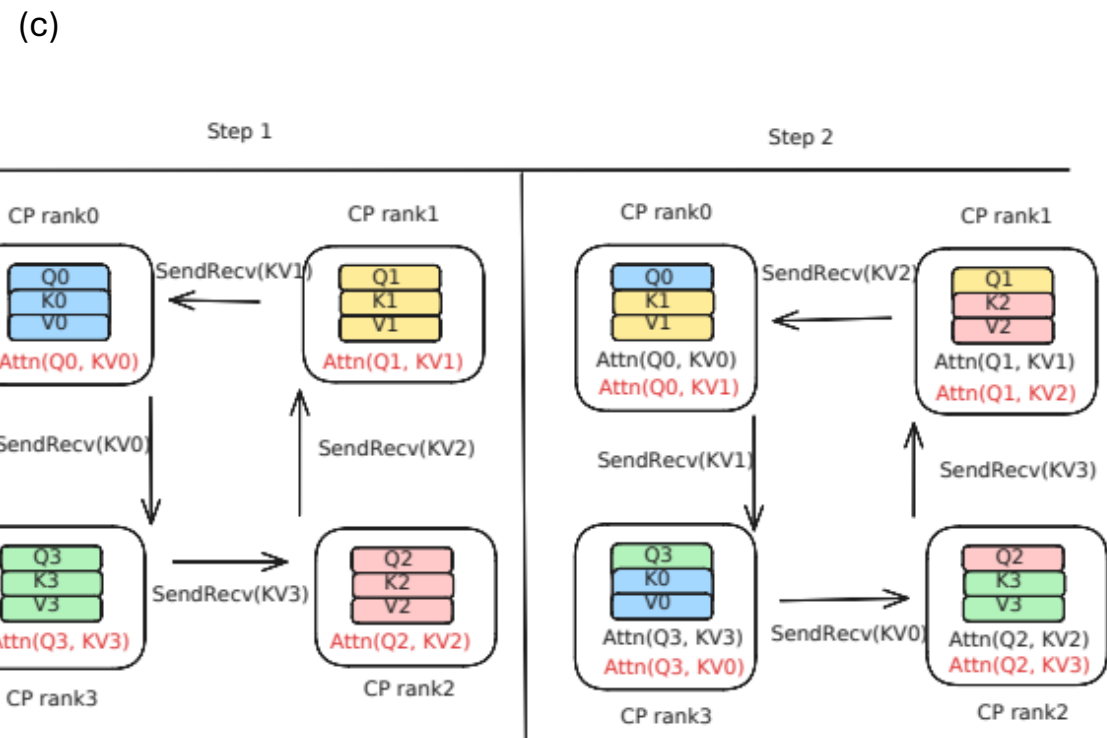
(b) **Computation vs. communication scaling for future transformers on future hardware.**  
Suchita Pati, Shaizeen Aga, Mahzabeen Islam, Nuwan Jayasena, and Matthew D. Sinclair. 2023

# Pipelining Approaches



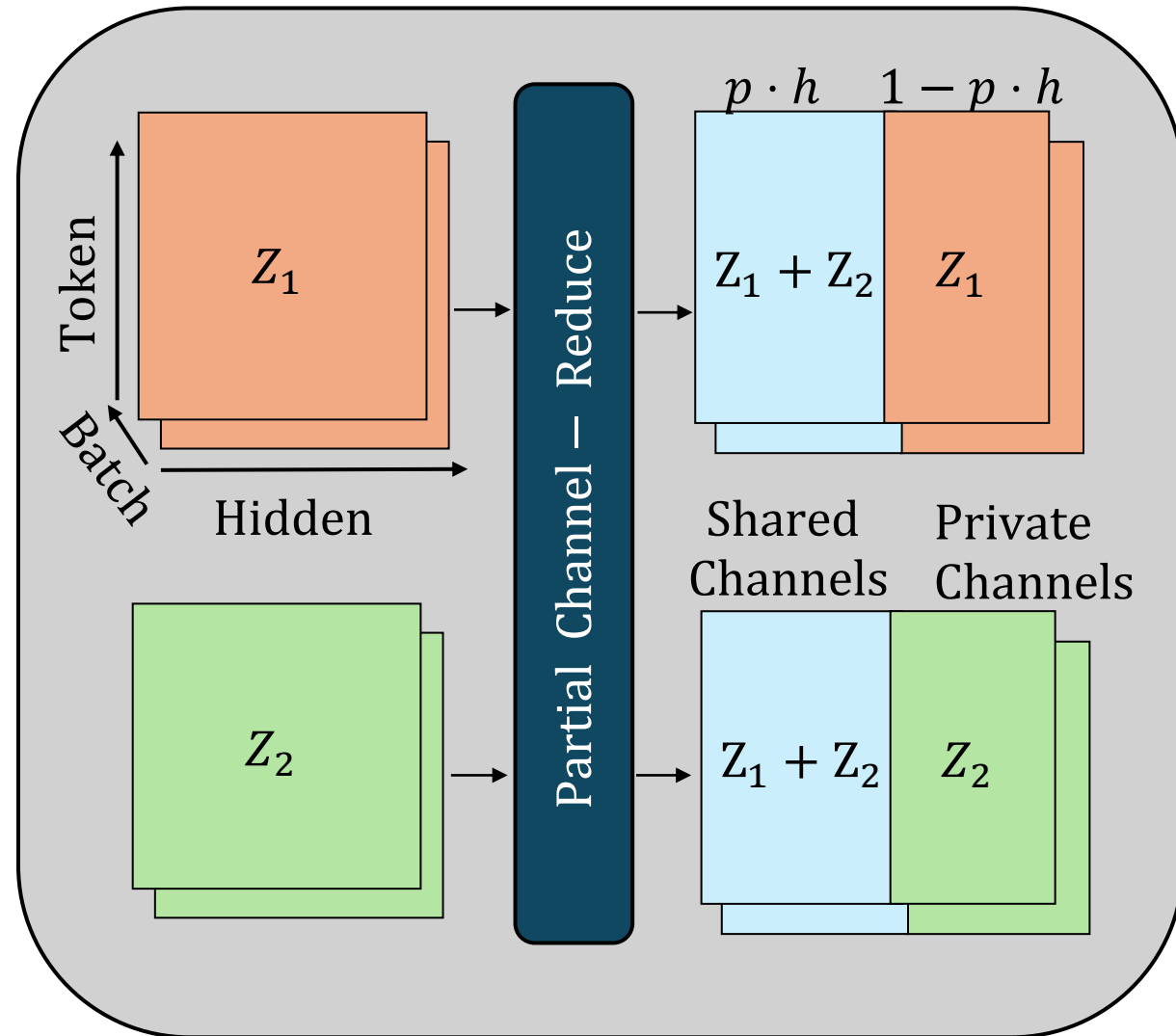
(a) **Domino: Eliminating communication in llm training via generic tensor slicing and overlapping**  
Guanhua Wang, Chengming Zhang, Zheyu Shen, Ang Li, and Olatunji Ruwase. 2024.

(b) **Computation vs. communication scaling for future transformers on future hardware.**  
Suchita Pati, Shaizeen Aga, Mahzabeen Islam, Nuwan Jayasena, and Matthew D. Sinclair. 2023

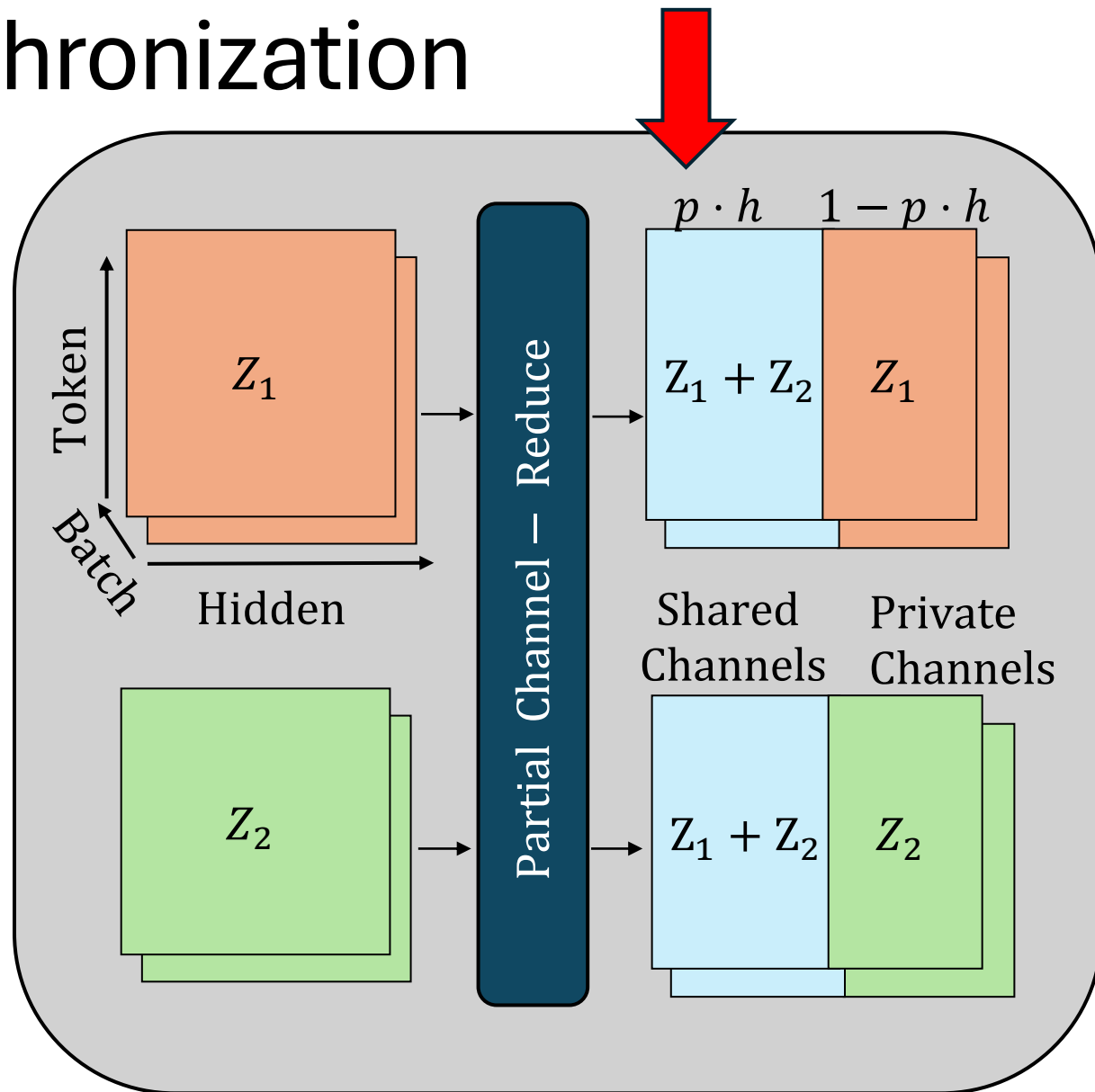


(c) **Context parallelism for scalable million-token**  
Amy Yang, Jingyi Yang, Aya Ibrahim, Xinfeng Xie, Bangsheng Tang, Grigory Sizov, Jeremy Reizenstein, Jongsoo Park, and Jianyu Huang. 2025

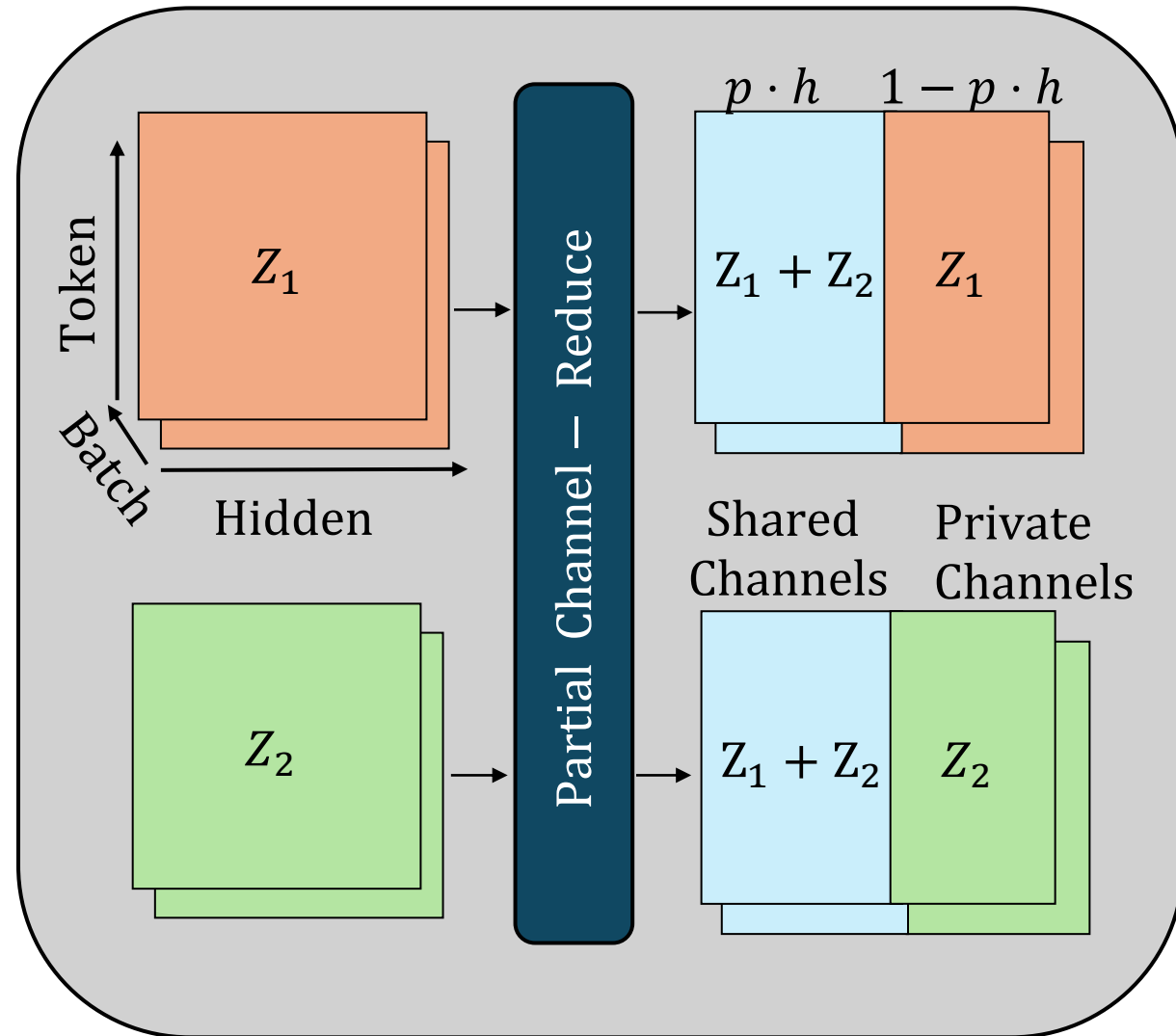
# Partial Synchronization



# Partial Synchronization

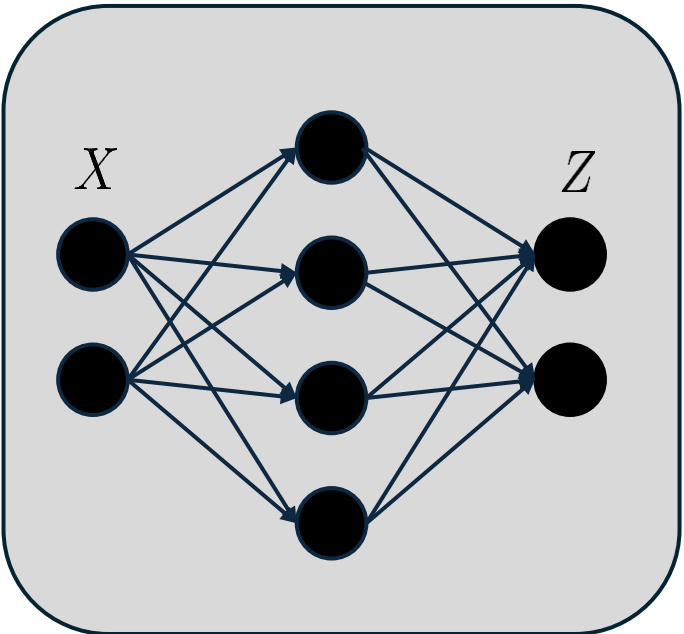


# Partial Synchronization



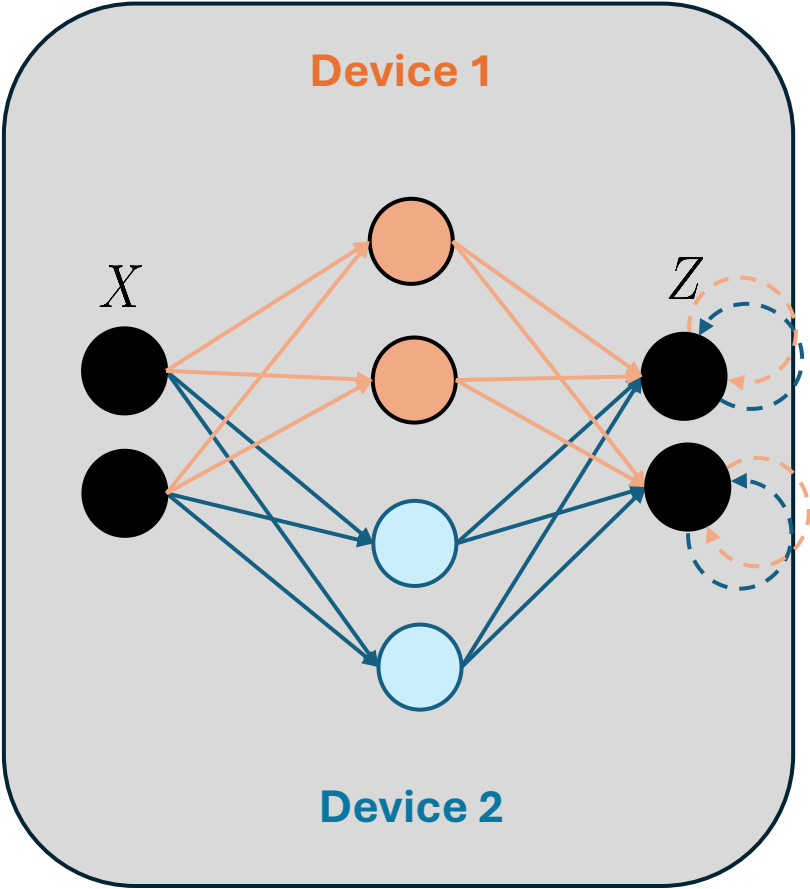
# CAAT-Net

Original Model

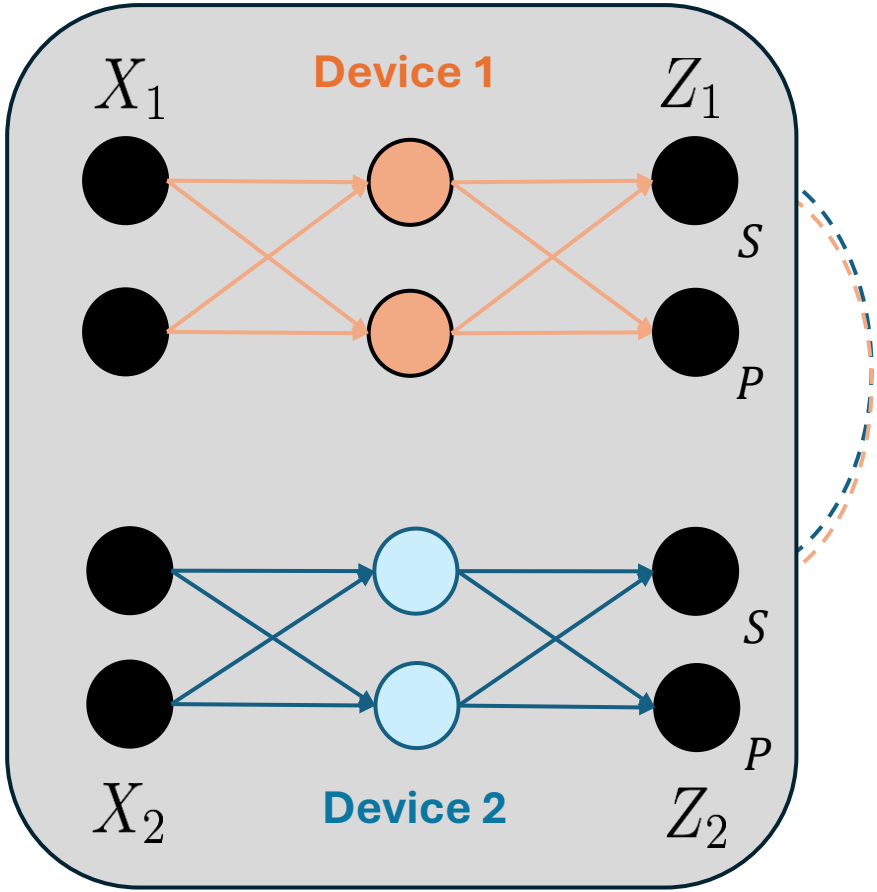


----- Values sent to device  
——— Values calculated on device

Tensor Parallel

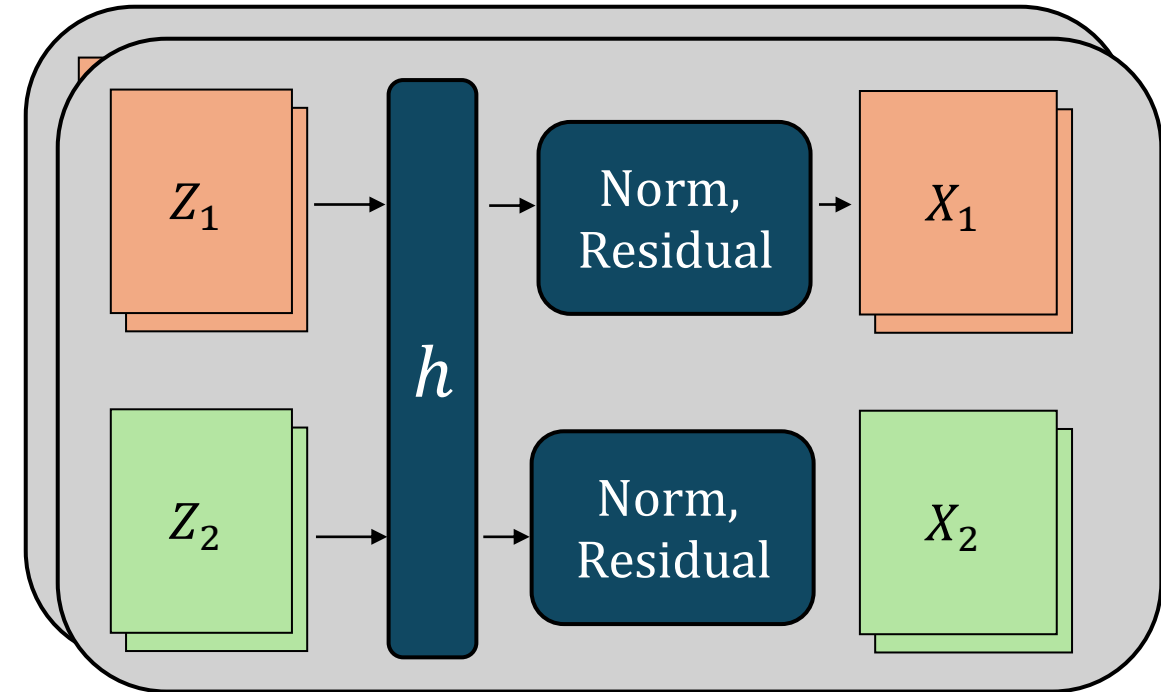
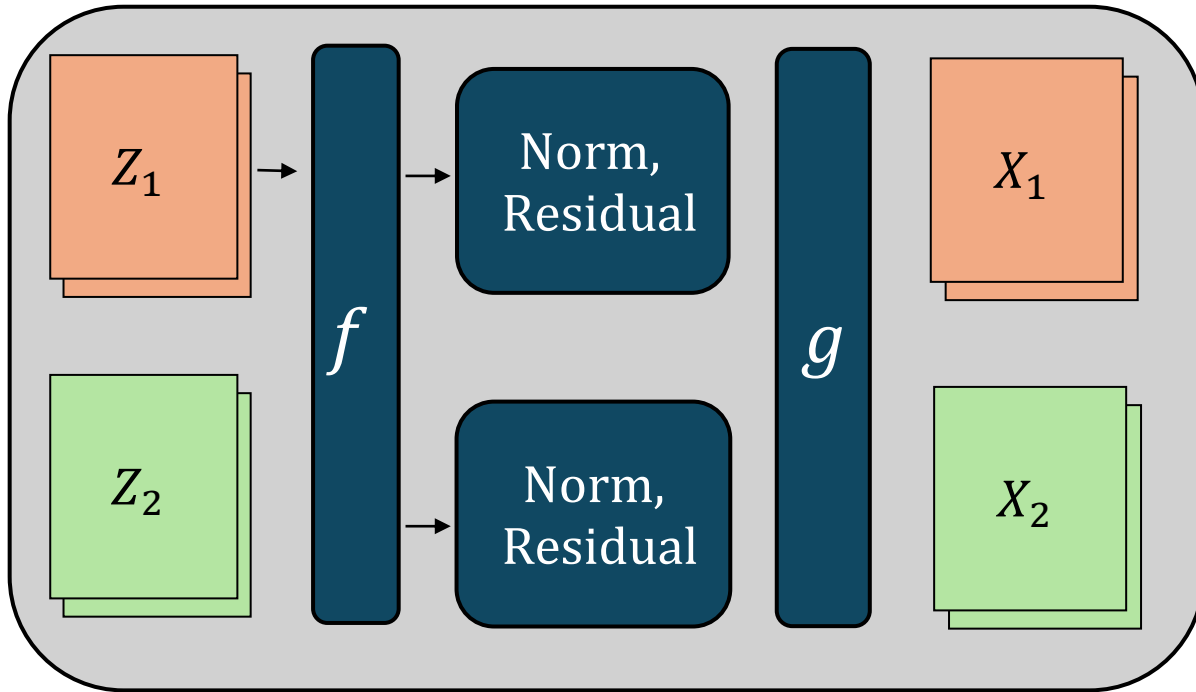


CAAT-Net





# Partial Synchronization – Implementation



# Results and Evaluation

Table 1: CAAT-Net vs baseline: Zero-shot accuracy after pretraining, 7B parameter models, with  $p = 0.5$  and tensor-parallel 8.

Model	LAMBADA (acc)	Hellaswag (acc)	WinoGrande (acc)	PIQA (acc)
Baseline	<b>61.34 <math>\pm</math> 0.68</b>	45.85 $\pm$ 0.50	61.48 $\pm$ 1.37	<b>72.91 <math>\pm</math> 1.06</b>
CAAT-Net	61.05 $\pm$ 0.68	<b>46.10 <math>\pm</math> 0.50</b>	<b>62.19 <math>\pm</math> 1.36</b>	72.86 $\pm$ 1.04
	OpenBookQA (acc)	BOOL-Q (acc)	WikiText (ppl)	Validation Loss
Baseline	<b>26.60 <math>\pm</math> 1.98</b>	<b>64.89 <math>\pm</math> 0.83</b>	12.51	1.01
CAAT-Net	24.00 $\pm$ 1.87	62.51 $\pm$ 0.85	<b>12.46</b>	<b>1.00</b>

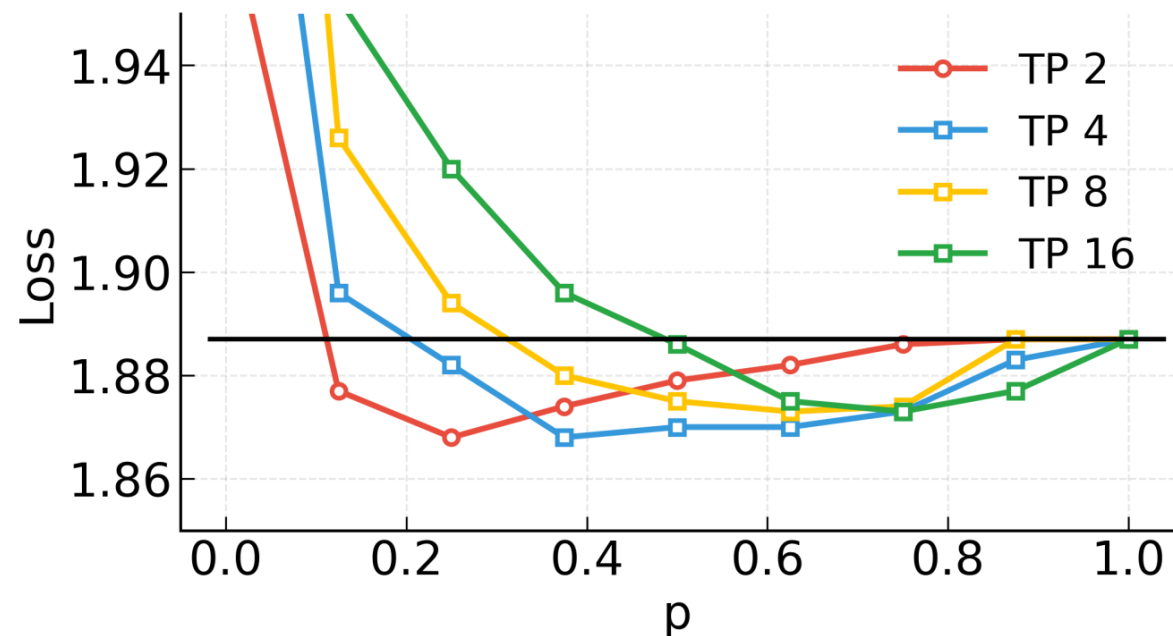
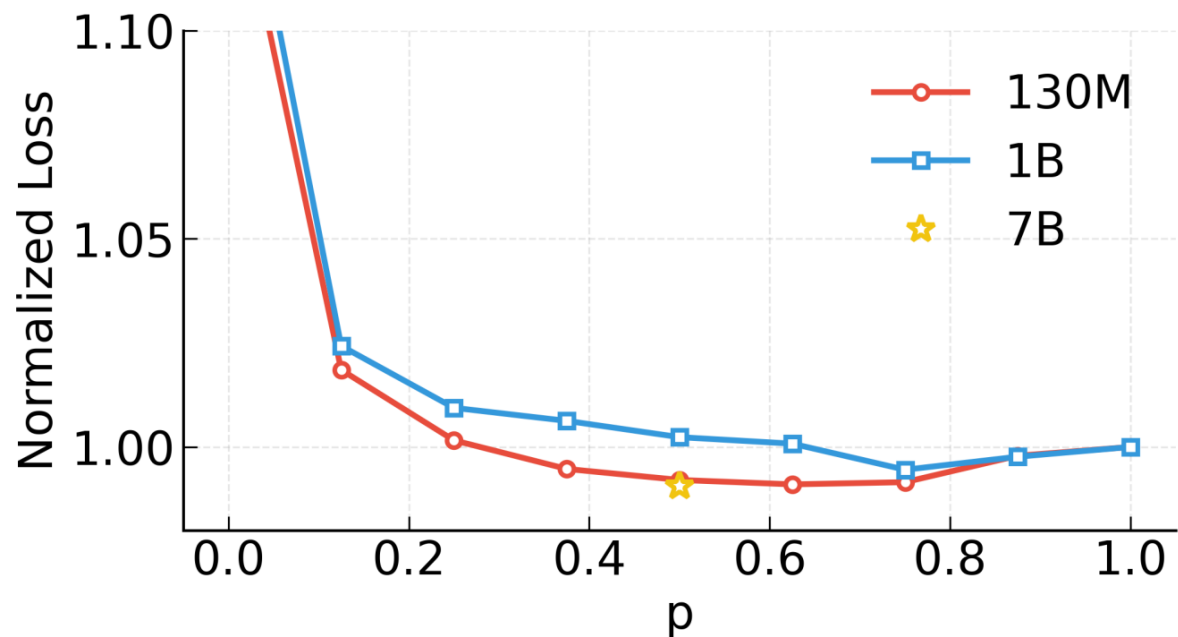
# Results and Evaluation

Table 1: CAAT-Net vs baseline: Zero-shot accuracy after pretraining, 7B parameter models, with  $p = 0.5$  and tensor-parallel 8.

Model	LAMBADA (acc)	Hellaswag (acc)	WinoGrande (acc)	PIQA (acc)
Baseline	<b>61.34 <math>\pm</math> 0.68</b>	45.85 $\pm$ 0.50	61.48 $\pm$ 1.37	<b>72.91 <math>\pm</math> 1.06</b>
CAAT-Net	61.05 $\pm$ 0.68	<b>46.10 <math>\pm</math> 0.50</b>	<b>62.19 <math>\pm</math> 1.36</b>	72.86 $\pm$ 1.04

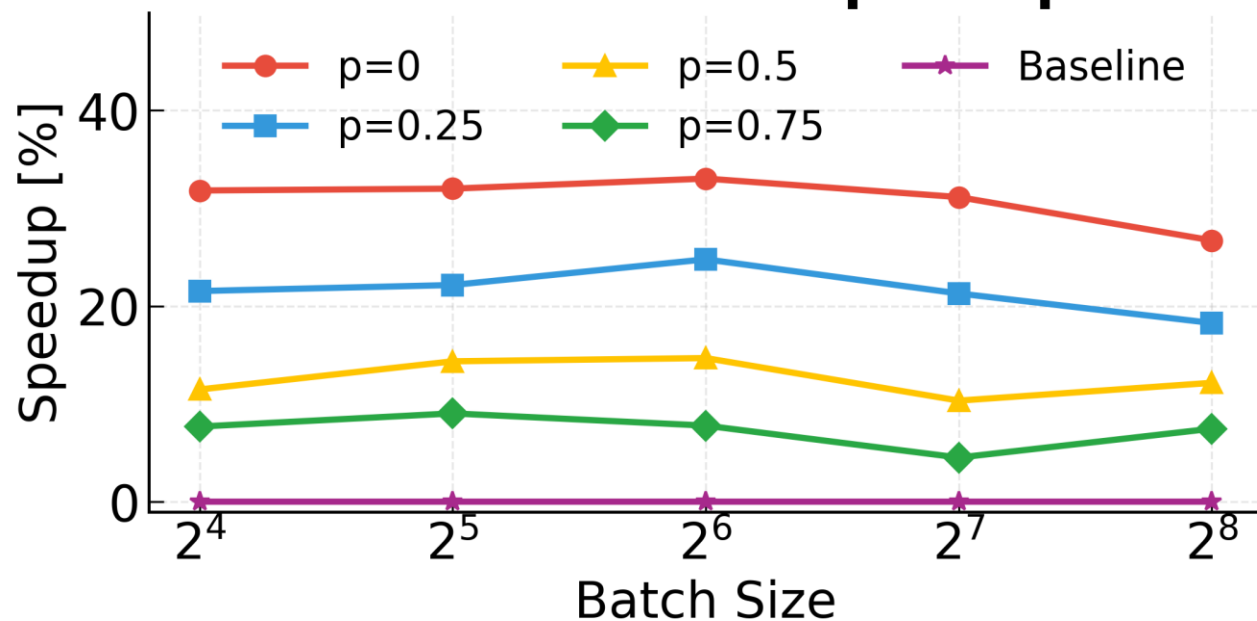
  

	OpenBookQA (acc)	BOOL-Q (acc)	WikiText (ppl)	Validation Loss
Baseline	<b>26.60 <math>\pm</math> 1.98</b>	<b>64.89 <math>\pm</math> 0.83</b>	12.51	1.01
CAAT-Net	24.00 $\pm$ 1.87	62.51 $\pm$ 0.85	<b>12.46</b>	<b>1.00</b>

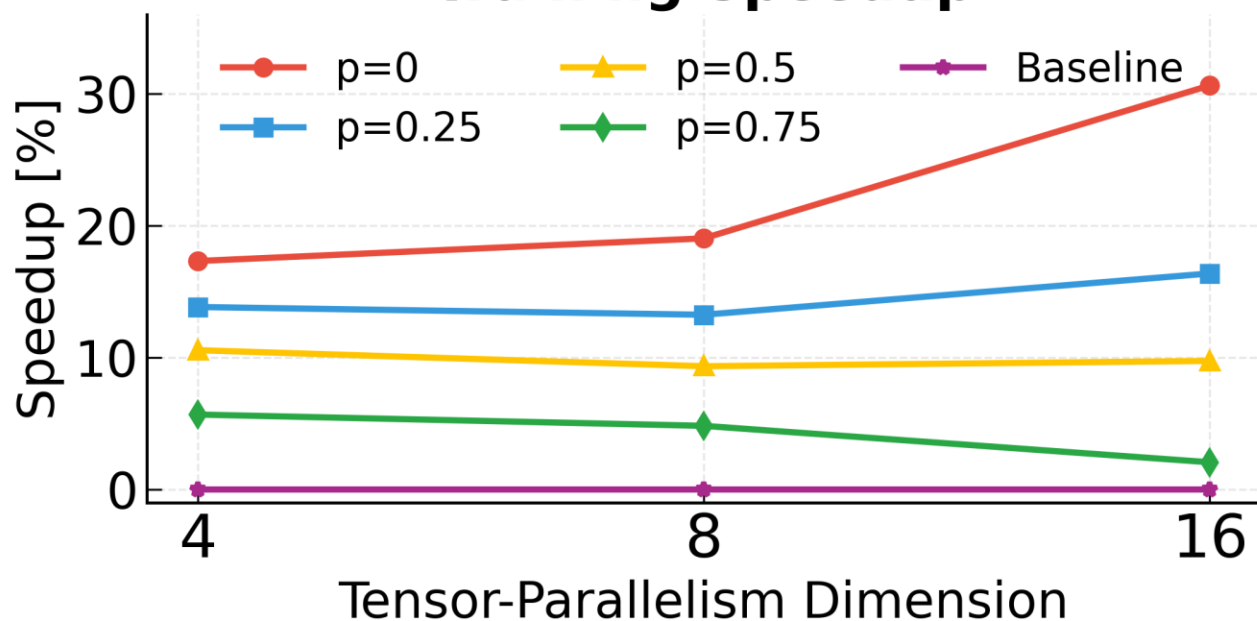


# Results and Evaluation

## TP 8 Inference Speedup



## Training Speedup





**TECHNION**  
Israel Institute  
of Technology

**intel**®



# Thank You!



The research of DS was funded by the European Union (ERC, A-B-C-Deep, 101039436).