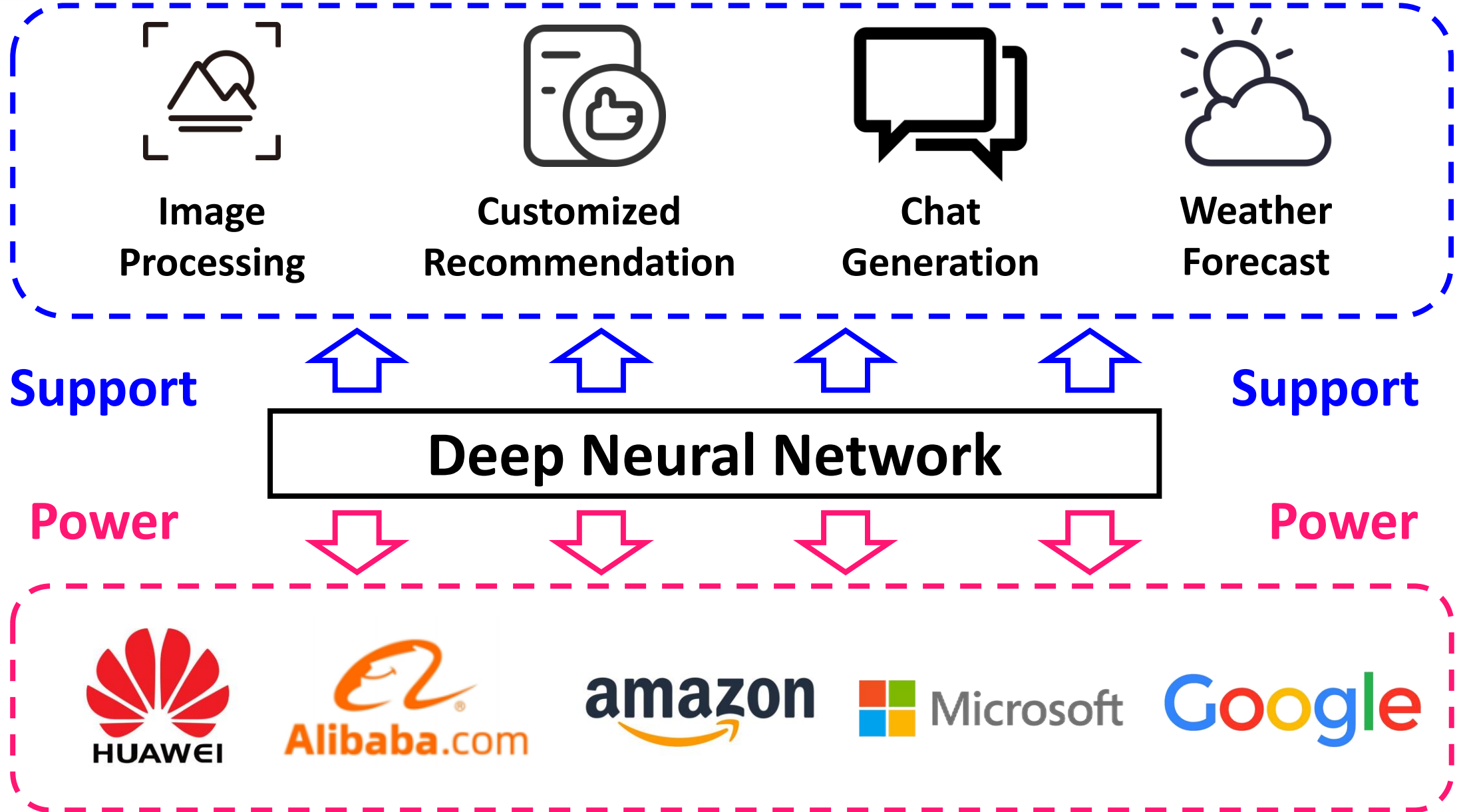


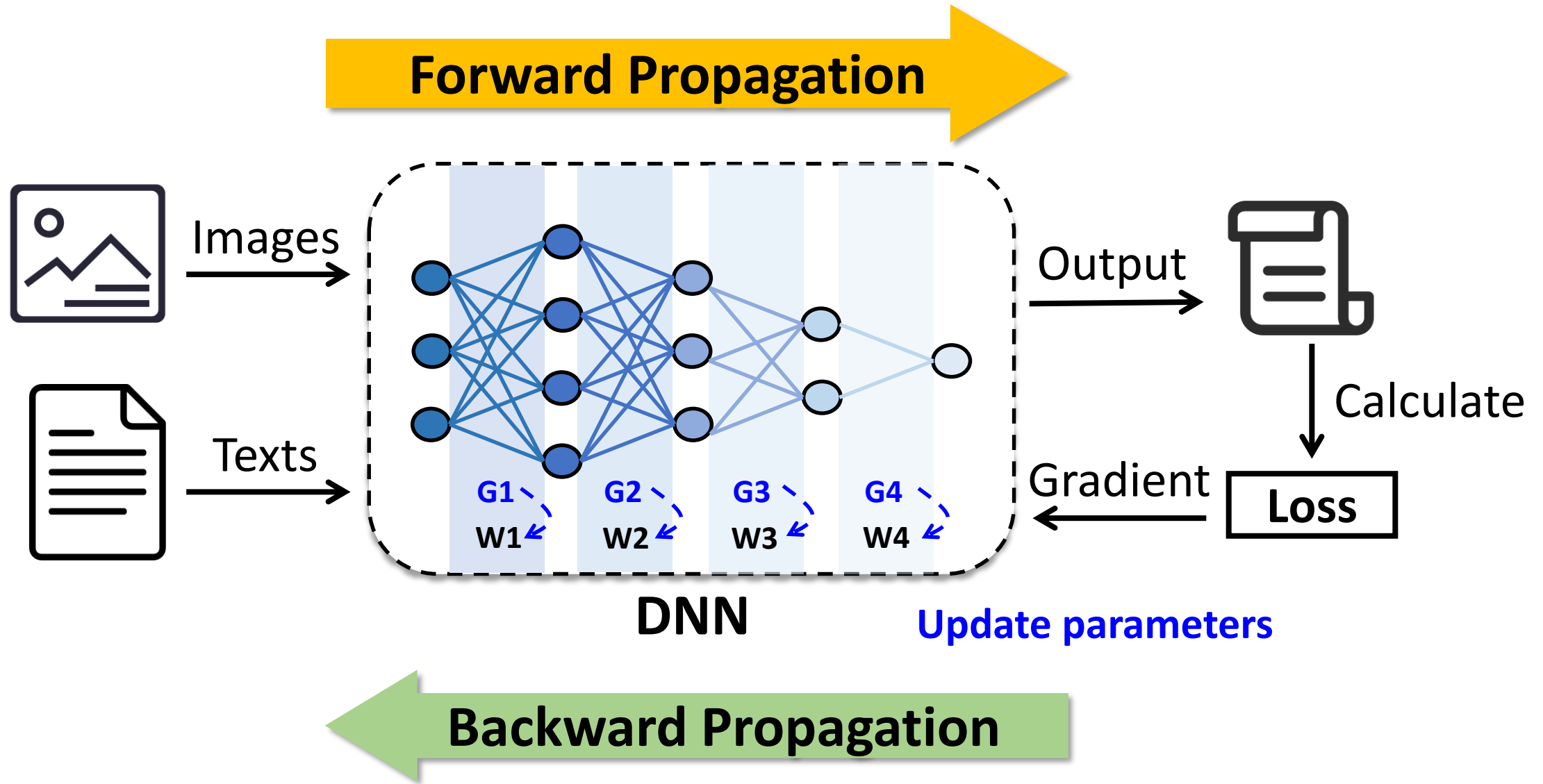
# A Cost-Efficient Failure-Tolerant Scheme for Distributed DNN Training

Menglei Chen, Yu Hua, Rong Bai, Jianming Huang  
*Huazhong University of Science and Technology, China*

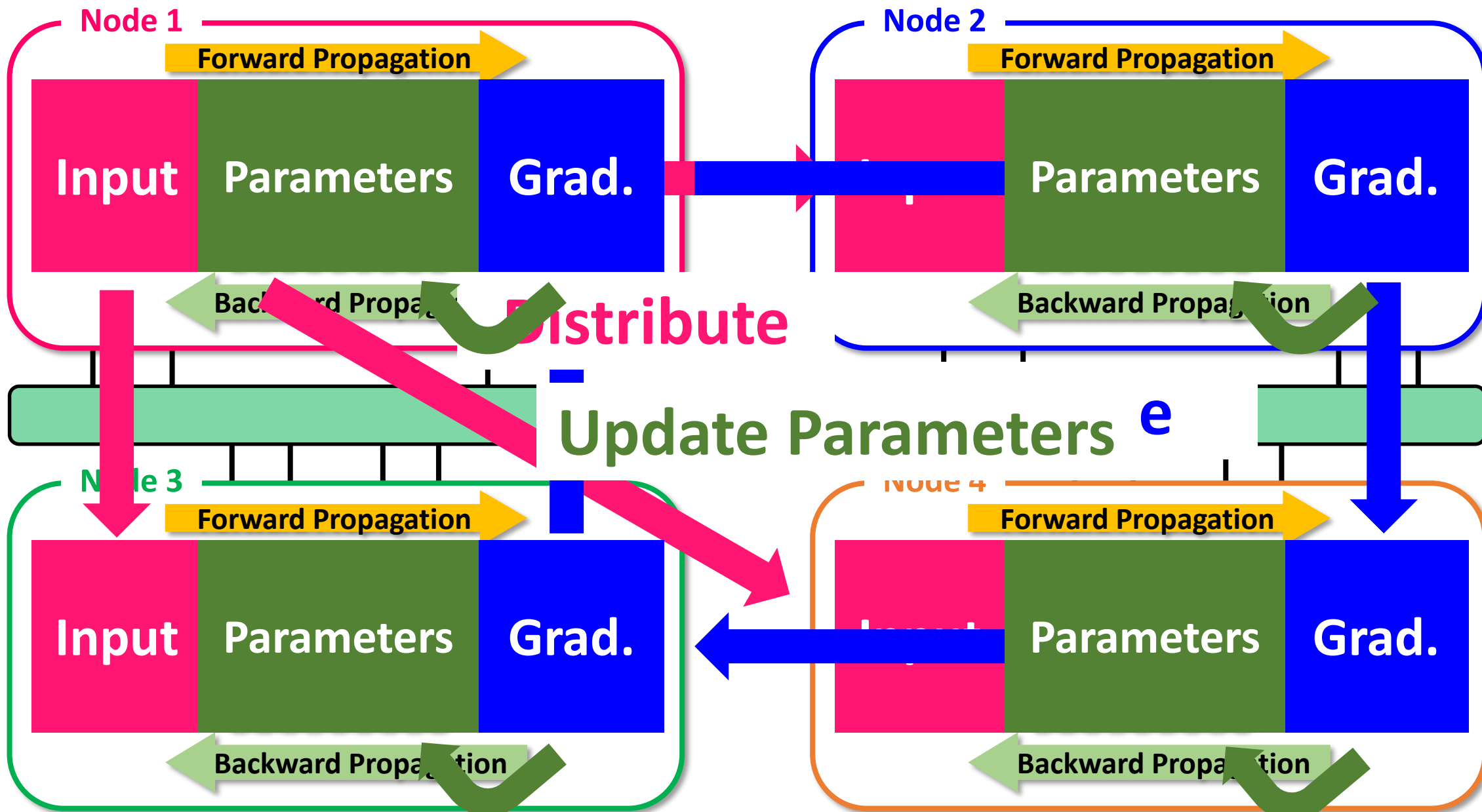
# Deep Neural Network (DNN)



# DNN Training

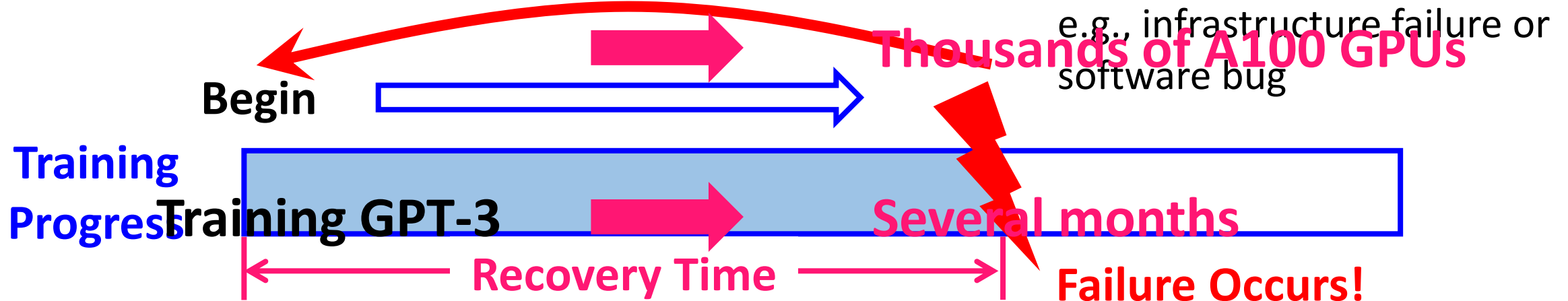


# Distributed DNN Training

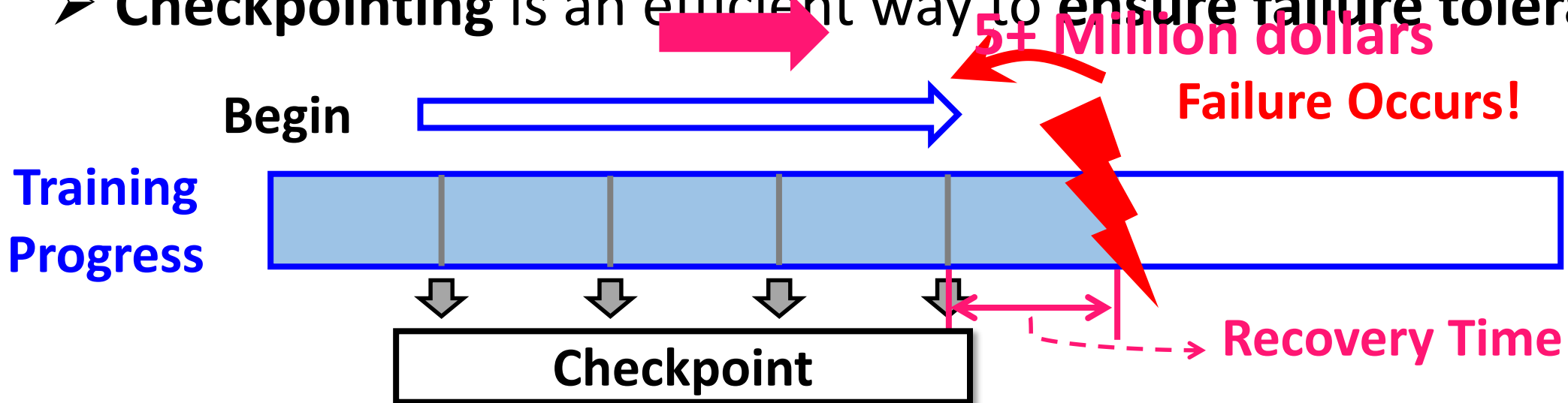


# The importance of Failure Tolerance

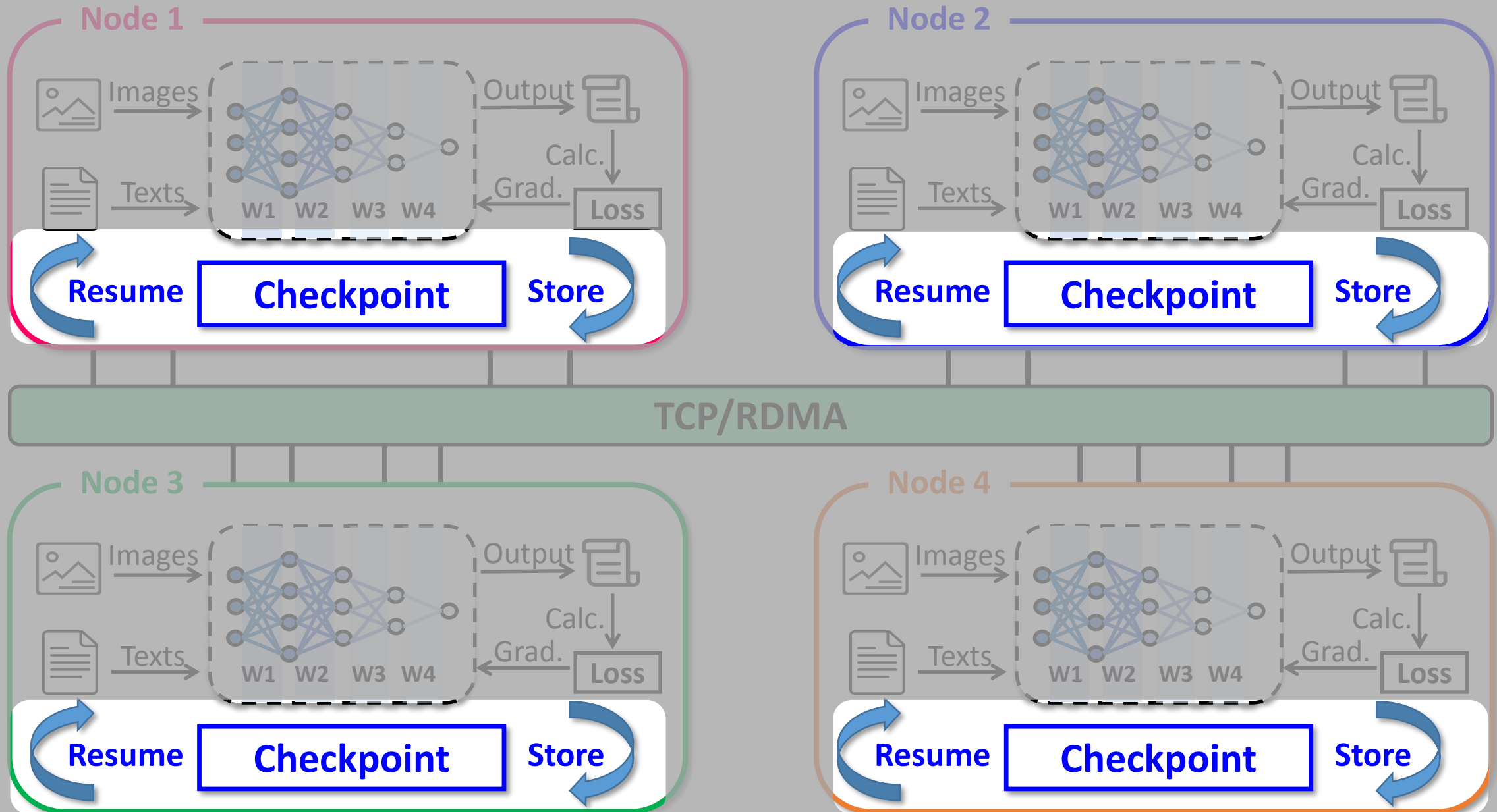
- DNN training is **time-consuming** and **expensive**



- **Checkpointing** is an efficient way to ensure **failure tolerance**



# Checkpointing in Distributed DNN Training



# The Need of Frequent Checkpointing

- Failures are common in large-scale GPU clusters
  - The mean time between failures is low to **a few minutes**
- Frequent job switches in the preemptive GPU cluster scheduling
  - The interval between two switches may be only **a few seconds**



**Frequent Checkpointing**



**High Runtime Overhead**

# Existing Checkpointing Schemes are Inefficient

- **Synchronous checkpointing**<sup>[1]</sup>
  - Introduce **severe training stall**
  - Suffer from **high runtime overhead**
- **Asynchronous checkpointing**<sup>[2-4]</sup>
  - Two-phase checkpointing
  - Pipeline the checkpointing with computation
  - Sub-optimal due to **monolithic** checkpointing process
  - **Fail to** fully pipeline checkpointing with communication



# Persistent Memory (PM)

- Intel Optane PM
- Samsung Memory-Semantic CXL (Compute Express Link) SSD



OR



**Byte-addressable**

**Fine-grained Persistence**

**Near-DRAM performance**



# Our Design

---

**LightCheck**: A cost-efficient checkpointing scheme for distributed DNN training

➤ Asynchronous layer-wise checkpointing

- Fine-grained pipelining
- Communication-aware

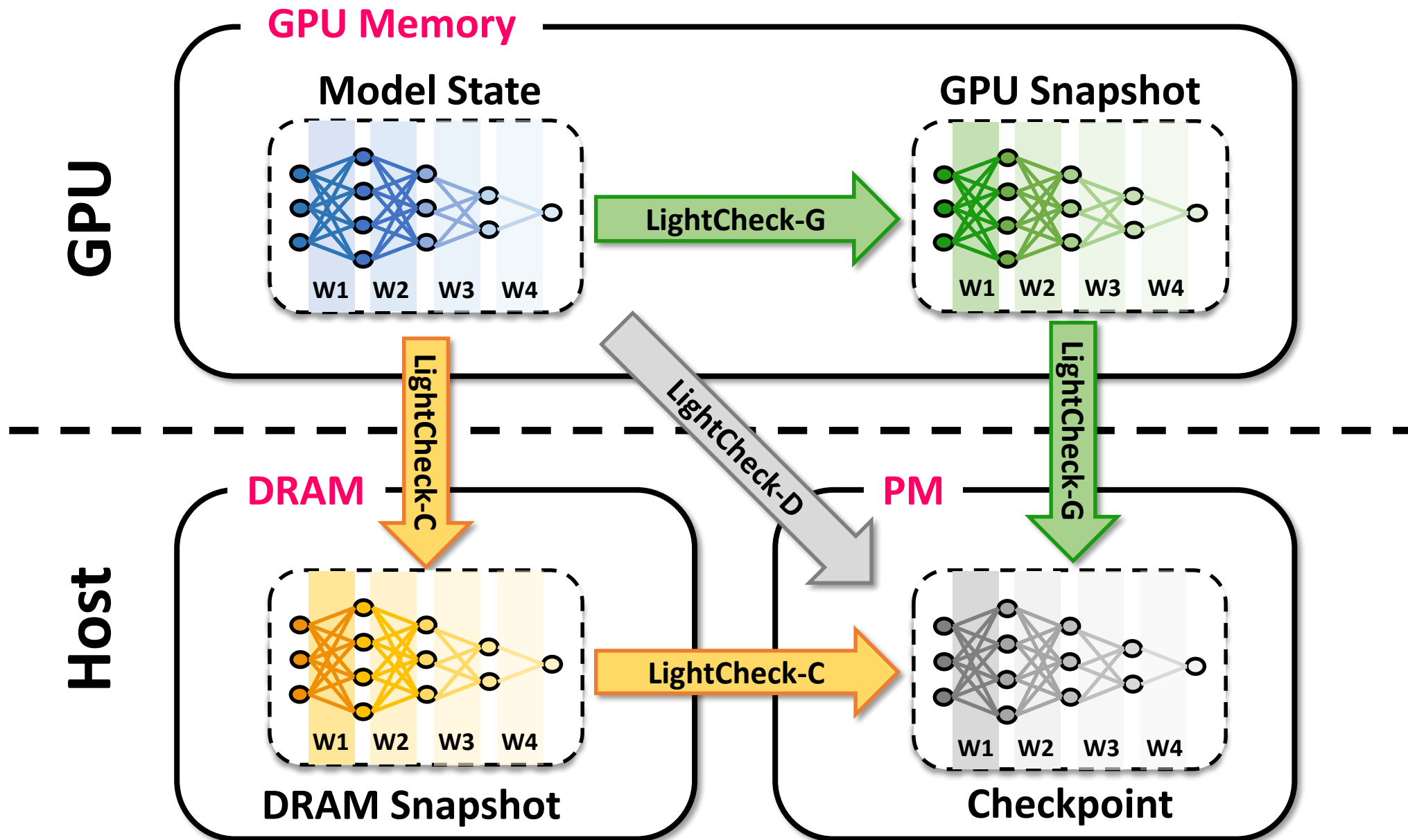
➡ **Minimizing training stalls**

➤ Efficient persistent memory management

- Direct access
- Metadata-aware

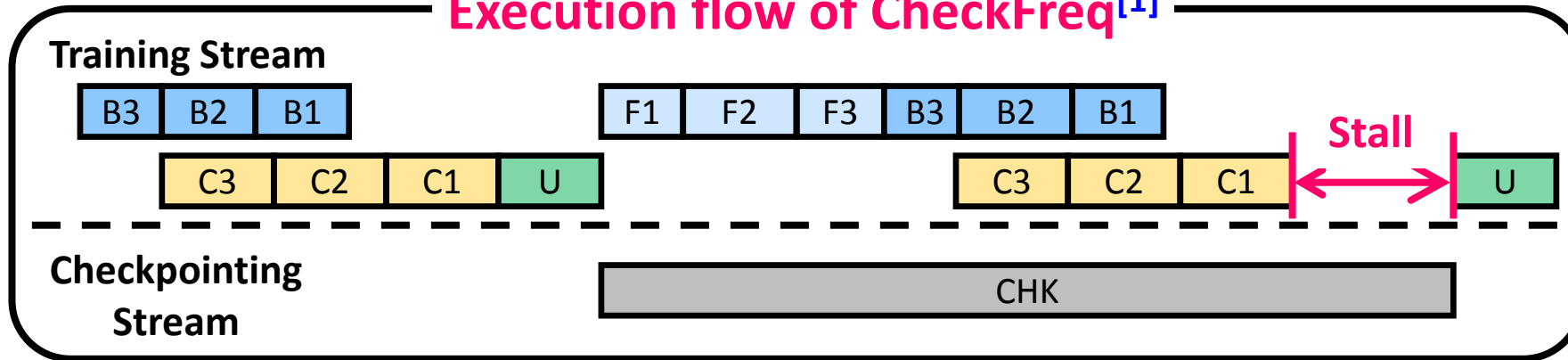
➡ **Fully exploiting persistent memory**

# Checkpointing Strategies



# Asynchronous Layer-wise Checkpointing

## Execution flow of CheckFreq<sup>[1]</sup>



B Backward Propagation

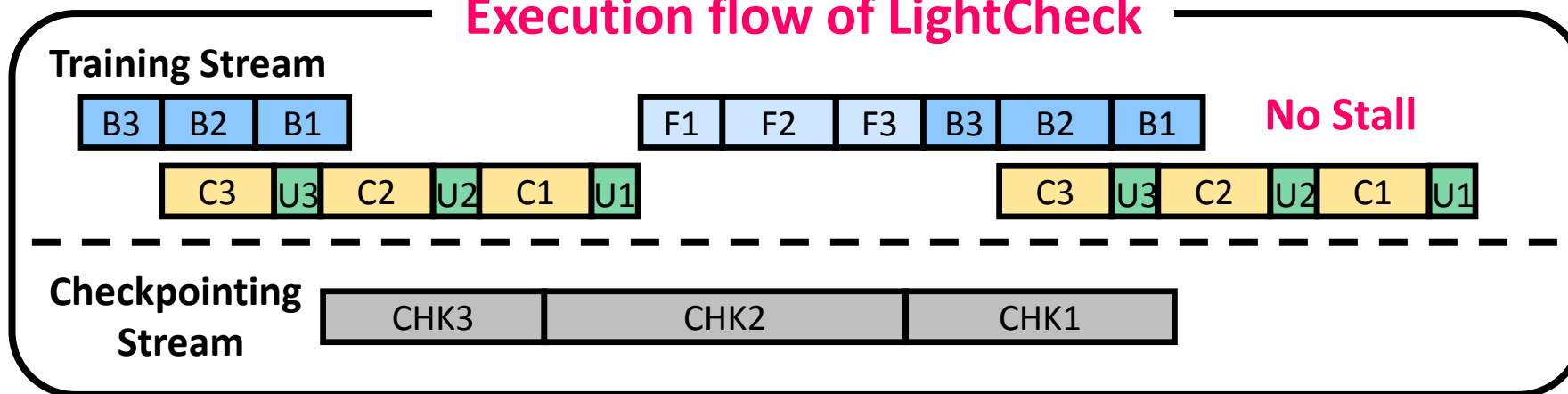
F Forward Propagation

C Communication

U Update Parameters

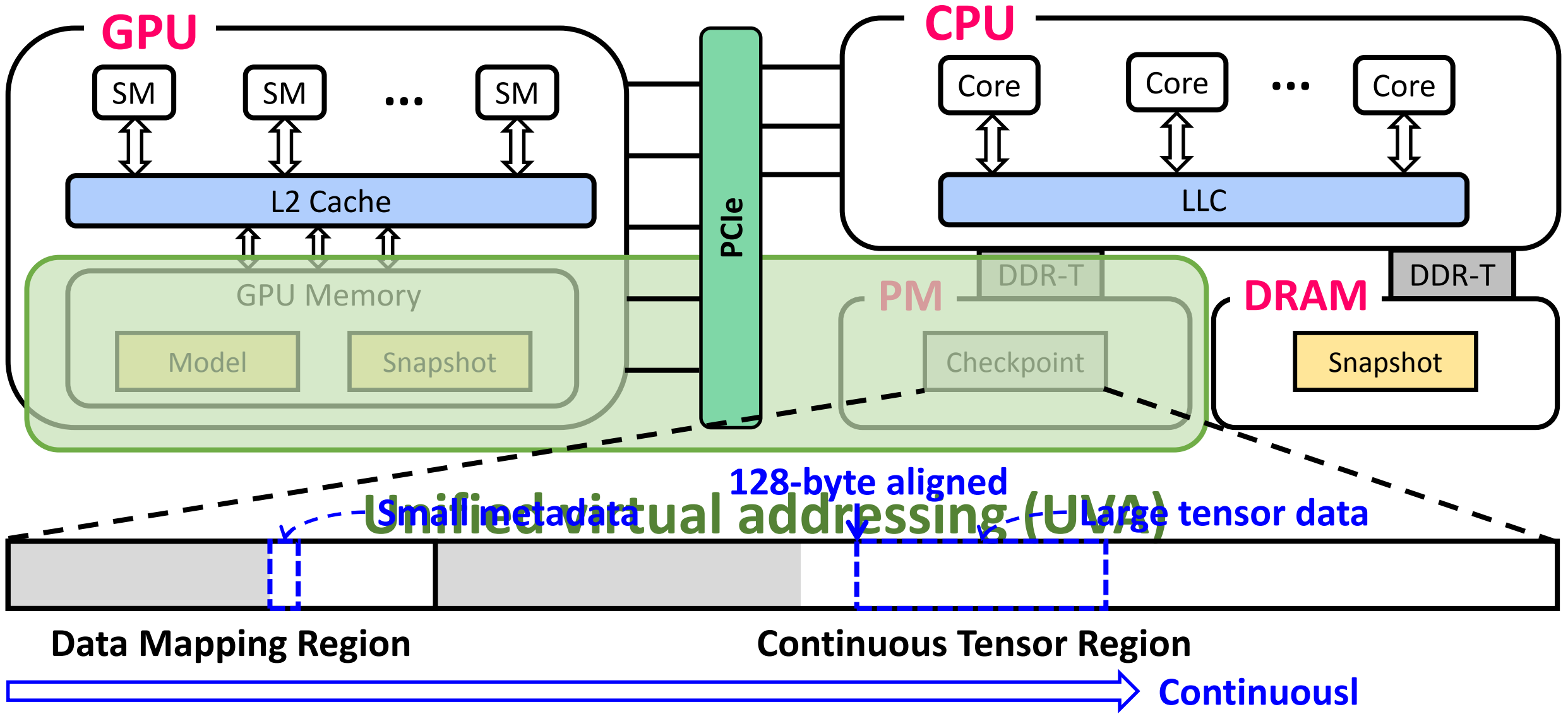
CHK Checkpointing

## Execution flow of LightCheck



[1] J. Mohan, A. Phanishayee, and V. Chidambaram, "Checkfreq: Frequent, fine-grained dnn checkpointing," in FAST, 2021

# Efficient persistent memory management



# Evaluation

## ➤ Platform

- Three nodes connected via 100 Gbps Mellanox InfiniBand switch

## ➤ DNN Models

- ResNet-18, VGG-16, Inception-V3, AlexNet, GPT-2, BERT

## ➤ Comparisons

- CheckFreq, Pytorch

Sever Configuration

Machine	CPU	GPU	Memory	Storage	Network
3 nodes	Intel Xeon Gold 6230R, 26 cores	1 Tesla V100, 16GB	192GB DRAM, 6 X 128GB Intel Optane PM Modules	3.6TB HDD	100Gbps Mellanox InfiniBand Switch

# Checkpointing Frequency

- Limit runtime overhead within 5%

Models	Checkpoint Size (MB)	Number of Iterations					
		LightCheck-G	LightCheck-C	LightCheck-D	LightCheck-disk	CheckFreq	torch.save
ResNet-18	90	1	1	1	7	20	102
VGG-16	1,056	6	6	6	64	146	904
Inception-V3	183	14	14	14	30	40	118
AlexNet	467	8	8	8	95	164	1,084
GPT-2	1,508	6	6	6	46	100	682
BERT	4,004	10	10	10	82	200	1,100

LightCheck can achieve  
with mode

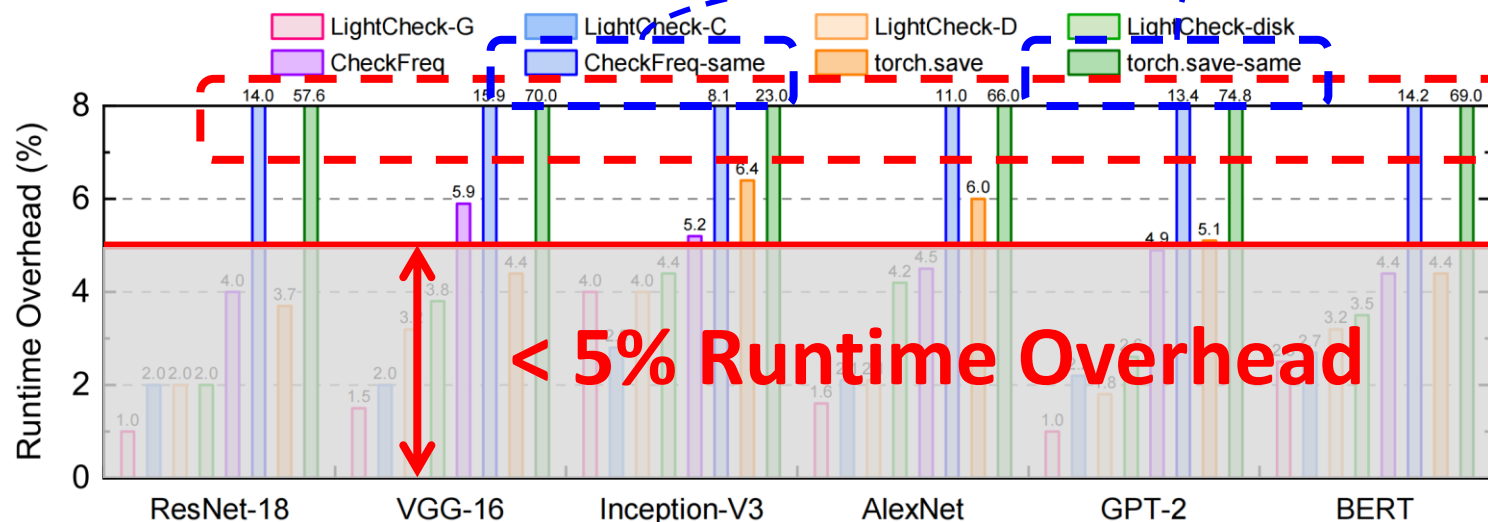
Asynchronous layer-wise checkpointing  
reduces the runtime overhead

to 10x to 2x

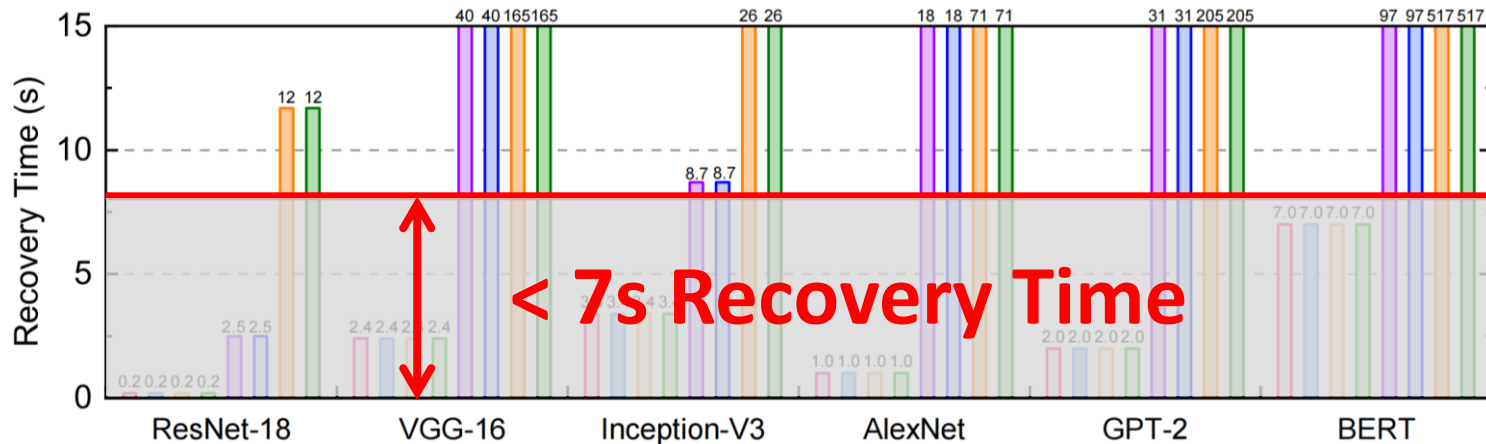
# Overall Performance

➤ With the aboved checkpointing frequency

with the **same** checkpointing frequency as LightCheck



< 5% Runtime Overhead



< 7s Recovery Time

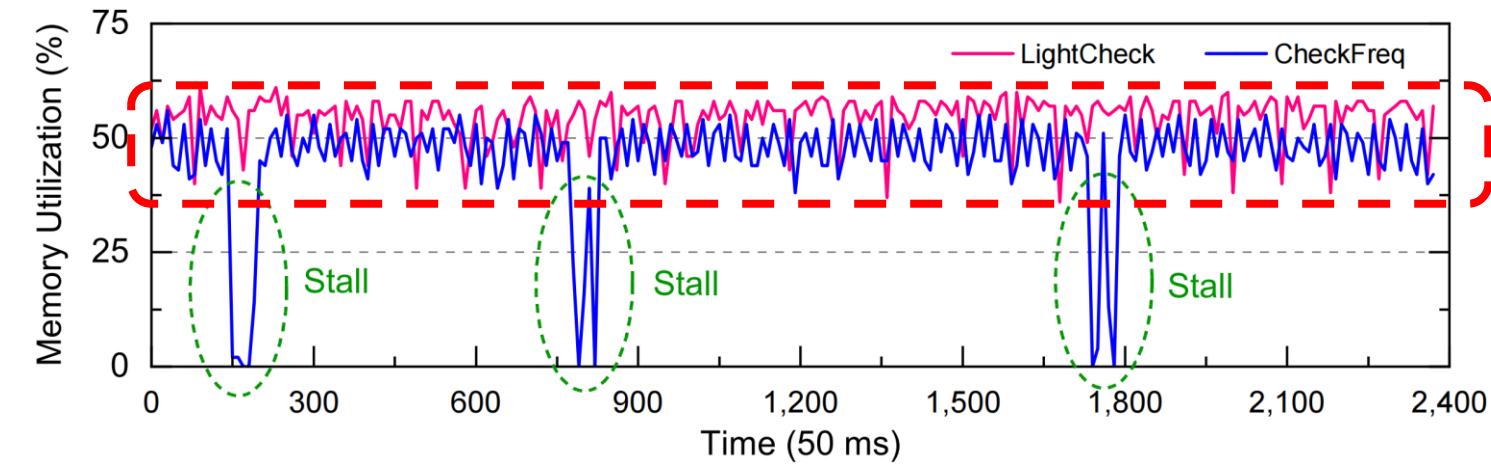
Incurring **high** runtime overhead when performing **frequent** checkpointing

LightCheck provides **lower** recovery time and overhead than existing schemes

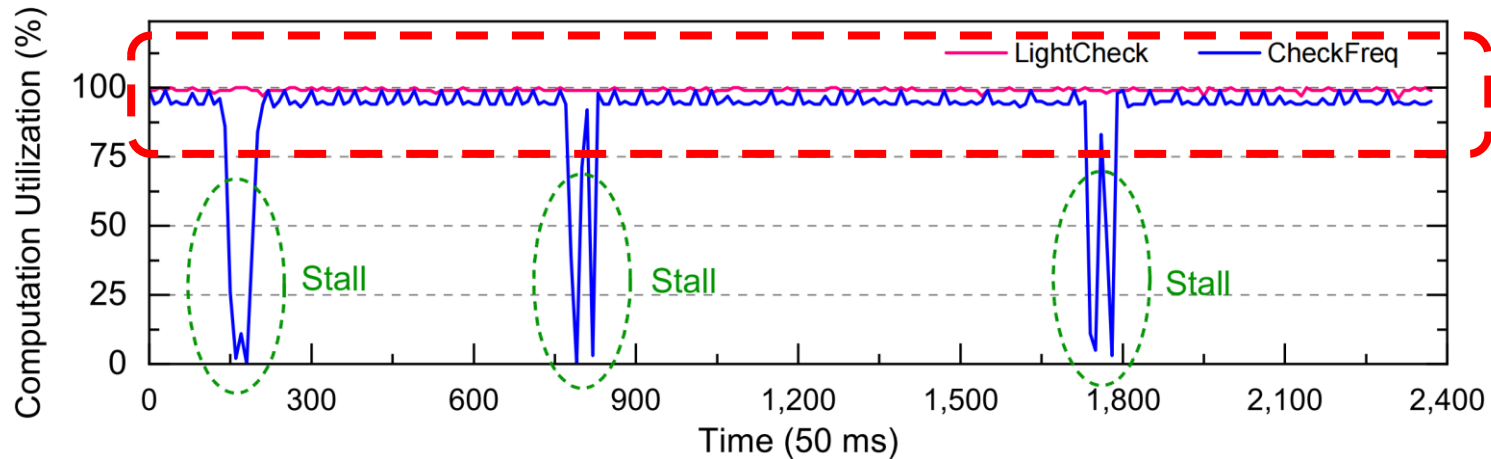


# GPU Utilization

- Record the GPU utilization every 50 ms, VGG-16



LightCheck **eliminates** training stall by leveraging find-grained pipelining





# Conclusion

---

- **LightCheck: A cost-efficient checkpointing scheme for DNN training**
  - Asynchronous layer-wise checkpointing
  - Efficient persistent memory management
- More evaluation results and analysis are in the paper
- Available at: <https://github.com/LighT-chenml/LightCheck.git>

***Thank you! Q&A***