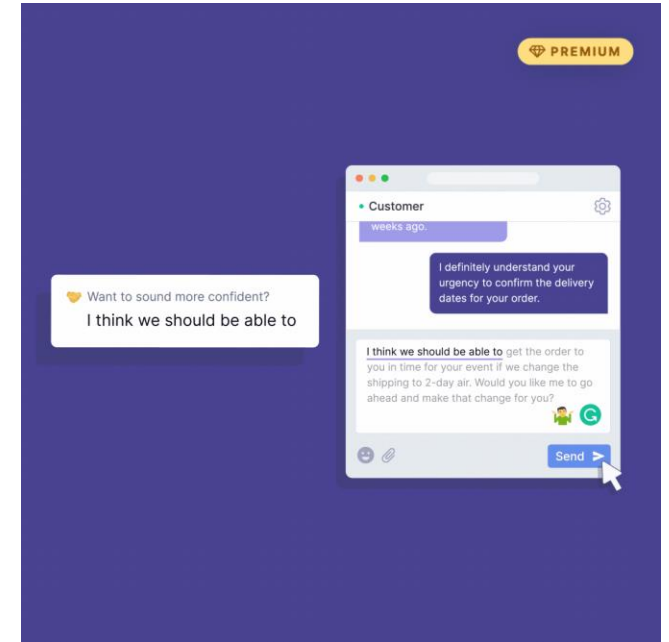# OptiReduce:
# Resilient and Tail-Optimal AllReduce for Distributed Deep Learning in the Cloud
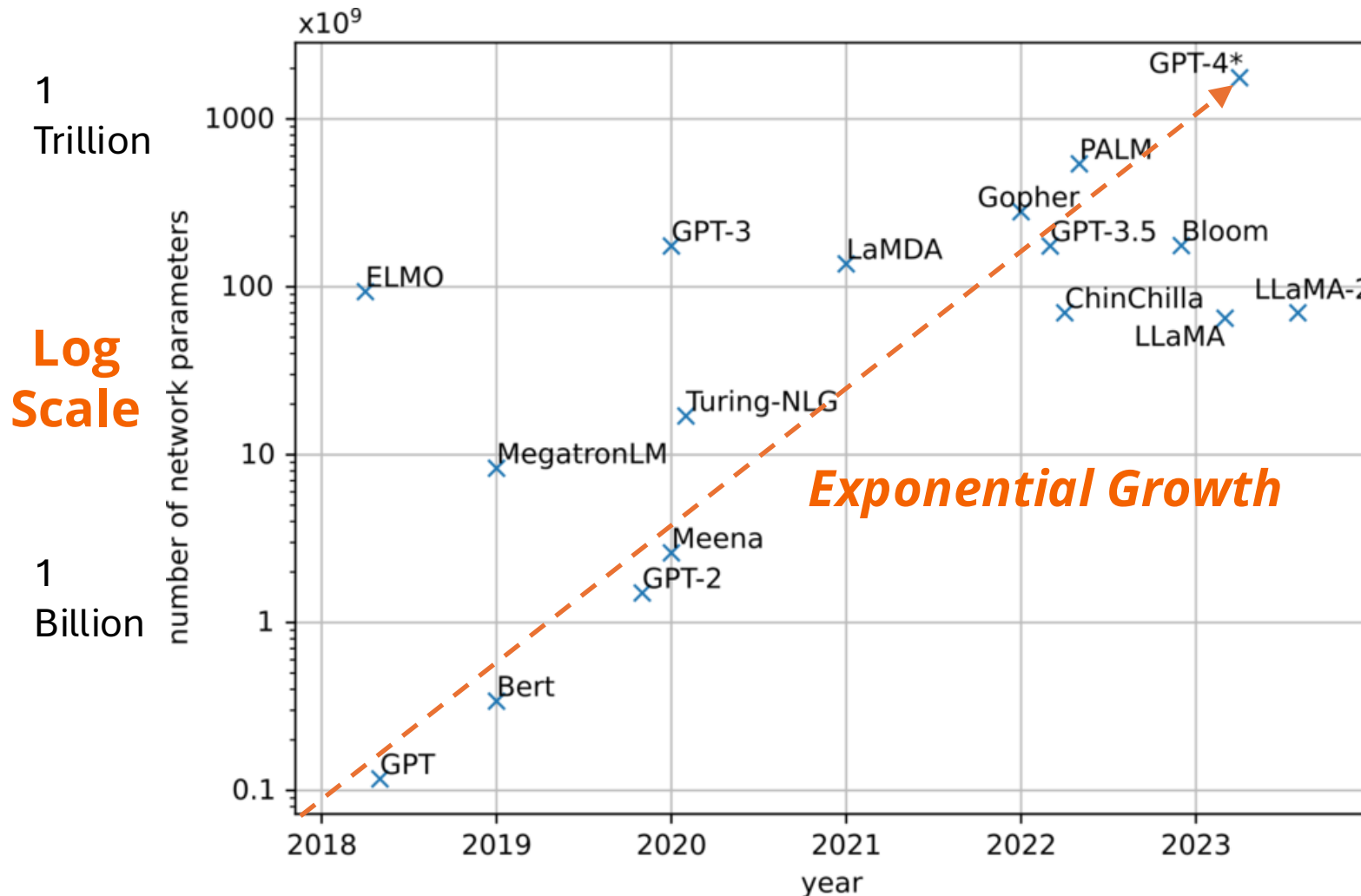
**Ertza Warraich**, Omer Shabtai, Khalid Manaa, Shay Vargaftik, Yonatan Piasetzky, Matty Kadosh, Lalith Suresh, Muhammad Shahbaz

# DNNs are Everywhere



ChatGPT



Midjourney



grammarly

# Growth of DNNs & Distributed Training



**Takes years to train!!!**

*Distribute training* over multiple GPUs: **Distributed Deep Learning (DDL)**

[1] Gerstmayr, Johannes et al. "Multibody Models Generated from Natural Language." Multibody System Dynamics '24
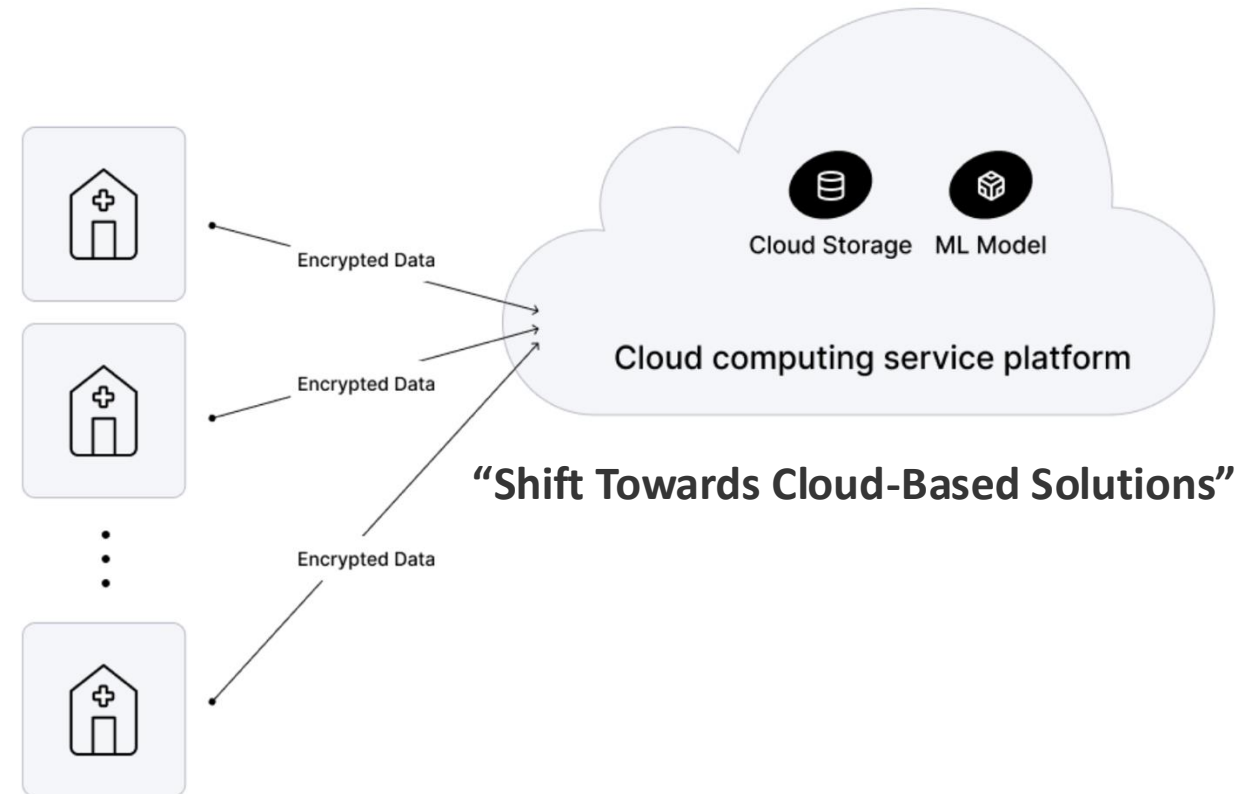
# Cloud for Distributed Training

Machine Learning in the Cloud: What Are the Benefits?

**mentormate**   Services ▾   Industries ▾   About ▾   Careers ▾

March 19, 2024

**The Power of Machine Learning in the Cloud: Transforming Business Operations**
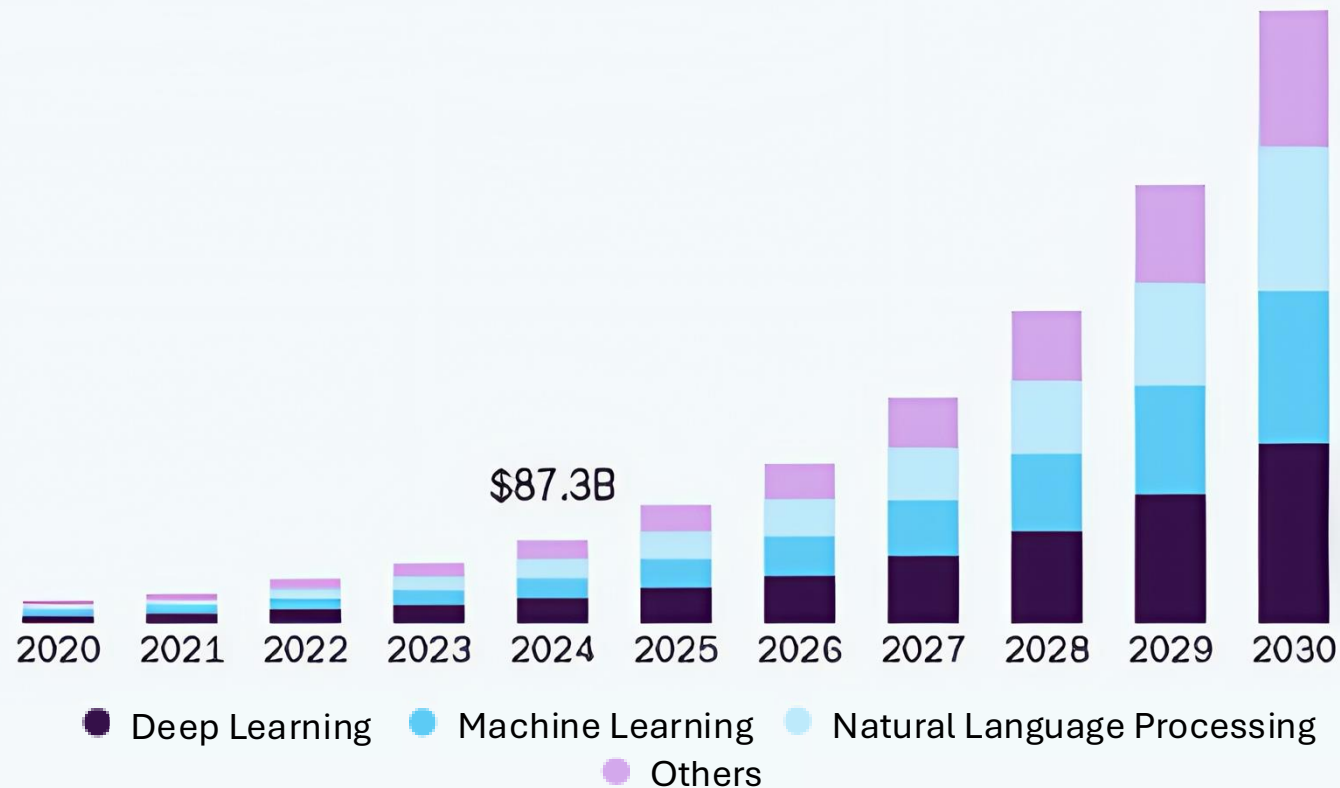
Encrypted Data

Encrypted Data

Encrypted Data

Cloud Storage   ML Model

Cloud computing service platform

**"Shift Towards Cloud-Based Solutions"**

[1] "https://encapture.com/the-rise-of-cloud-machine-learning"
[2] "https://redresscompliance.com/exploring-cloud-based-machine-learning-platforms"
[3] "https://symphony-solutions.com/insights/machine-learning-in-the-cloud-what-are-the-benefits"

# Cloud for Distributed Training

# Distributed Data Parallel (DDP)

# Distributed Data Parallel (DDP)



Node 2

Node 1

Gradients

AllReduce

Node 3

Node 4

*Synchronize and aggregate model gradients:*
**AllReduce Communication**
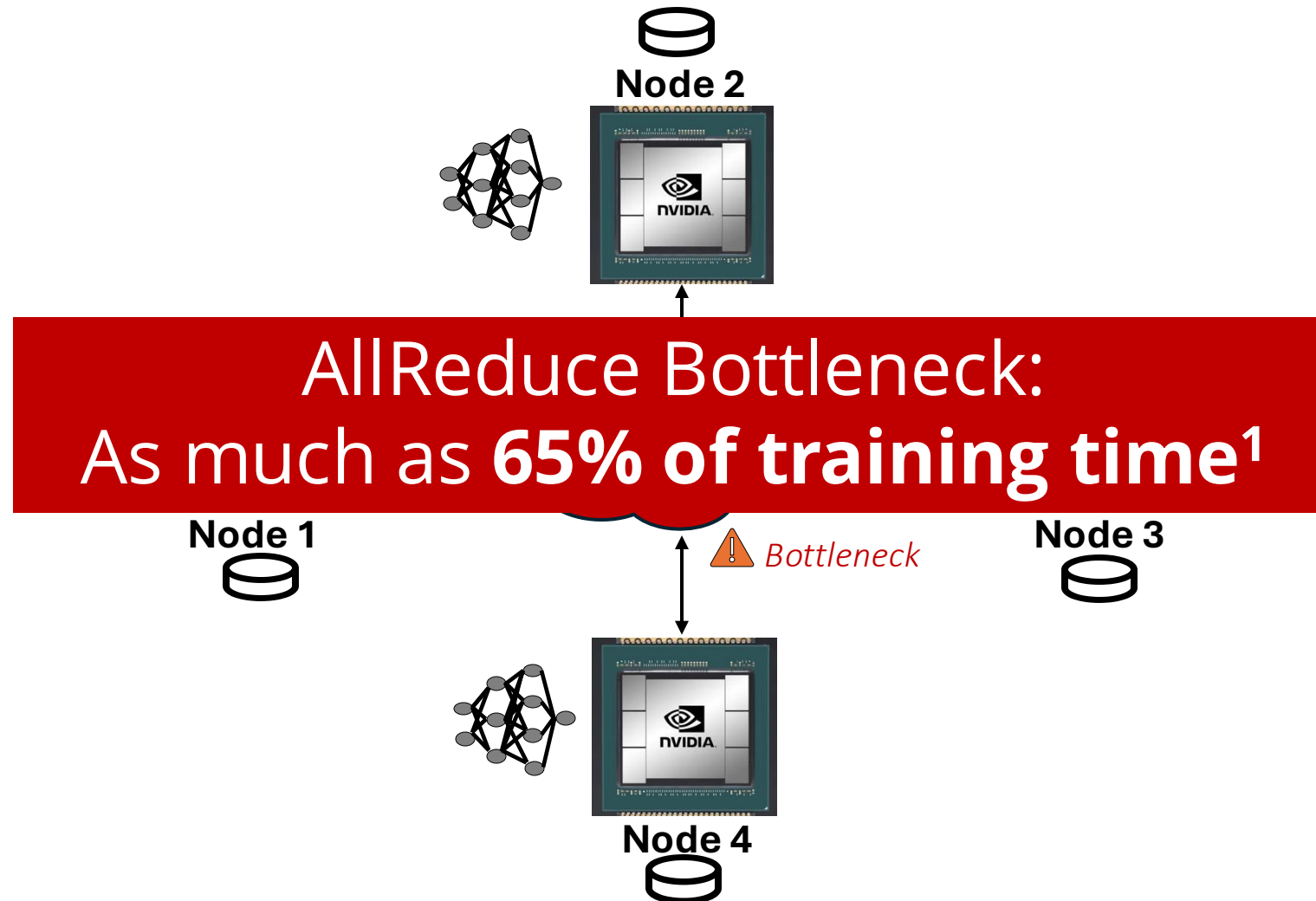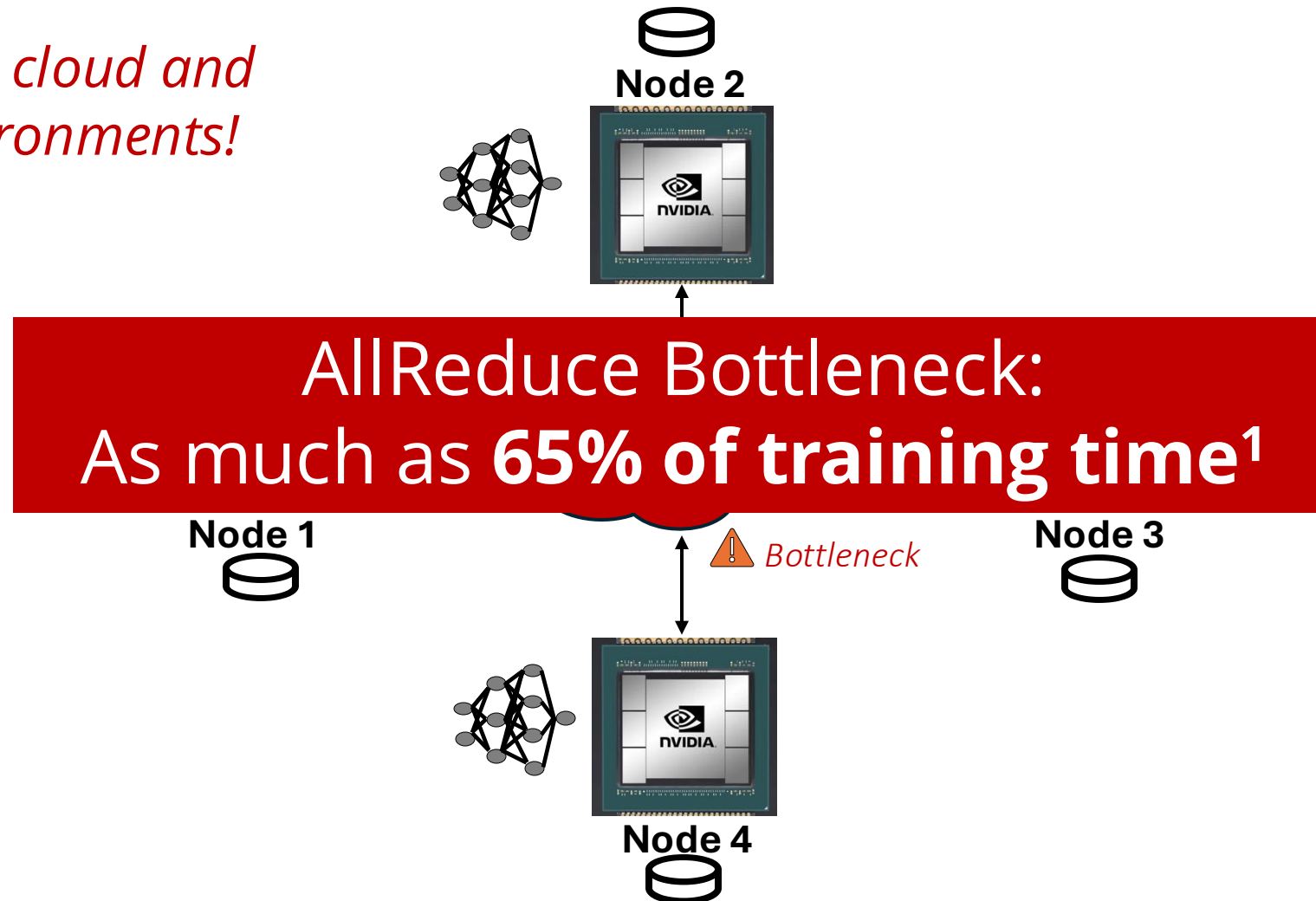
Dataset    DNN Model

8

# Distributed Data Parallel (DDP)

[1] Sapio, Amedeo, et al. "Scaling distributed machine learning with in-network aggregation." NSDI'21

# Distributed Data Parallel (DDP)
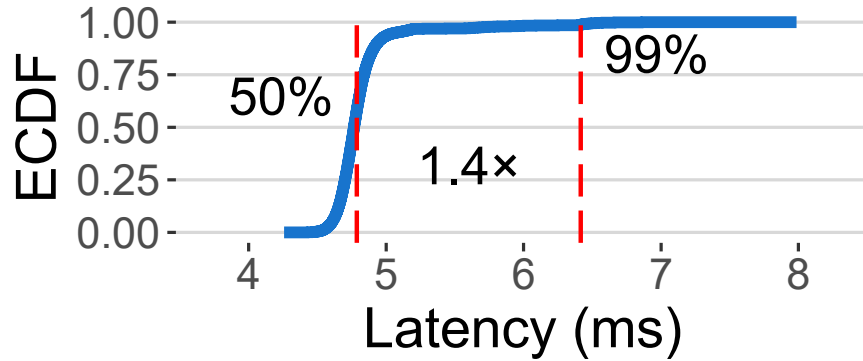


**Node 2**

## AllReduce Bottleneck:
## As much as **65% of training time**[1]

**Node 1**

⚠️ *Bottleneck*

**Node 3**

**Node 4**

[1] Sapio, Amedeo, et al. "Scaling distributed machine learning with in-network aggregation." NSDI'21

# Distributed Data Parallel (DDP)

*Worsened in cloud and shared environments!*

**Node 2**

**AllReduce Bottleneck:
As much as 65% of training time[1]**

**Node 1**

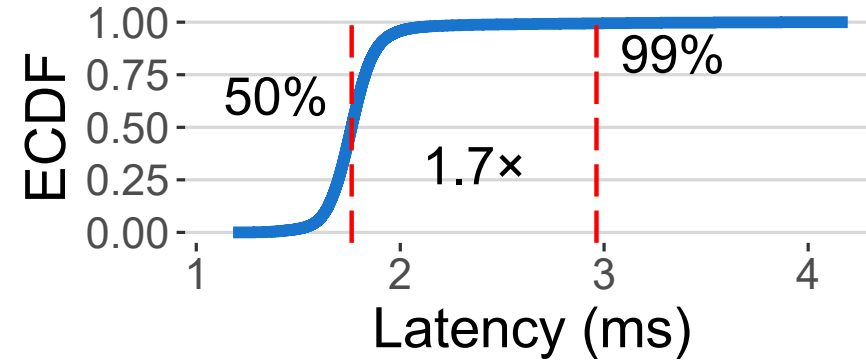⚠️ *Bottleneck*

**Node 3**

**Node 4**

[1] Sapio, Amedeo, et al. "Scaling distributed machine learning with in-network aggregation." NSDI'21
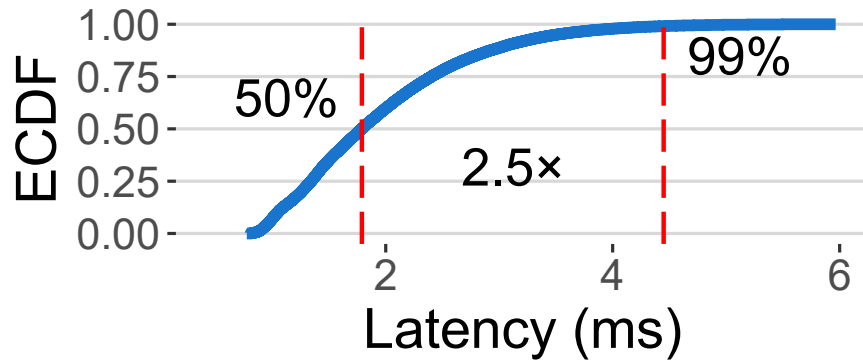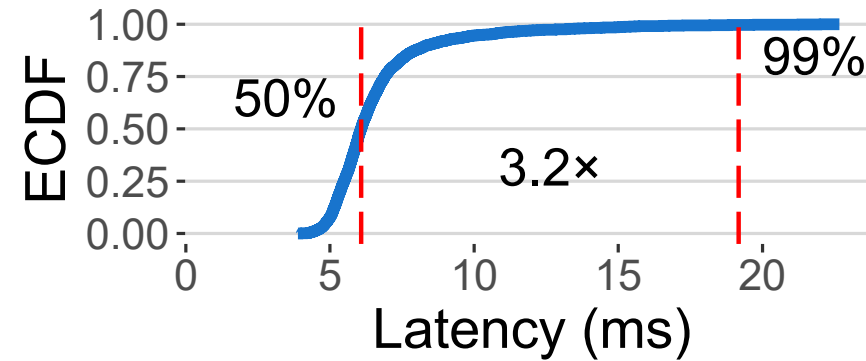
# Tail in the Cloud



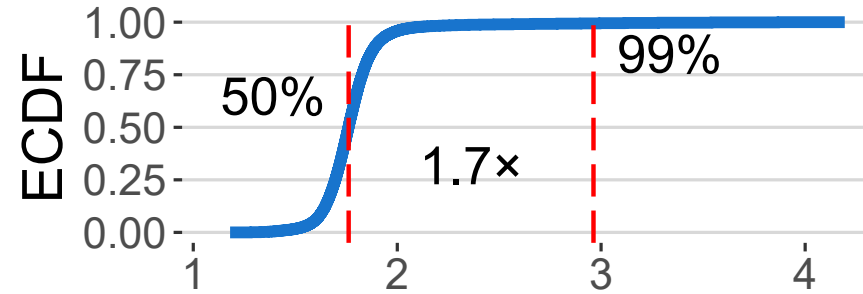(a) CloudLab

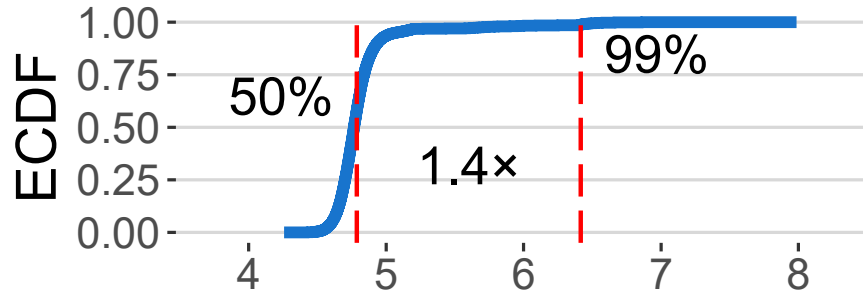(b) Hyperstack

(c) AWS EC2

(d) Runpod

# Tail in the Cloud



(c) AWS EC2
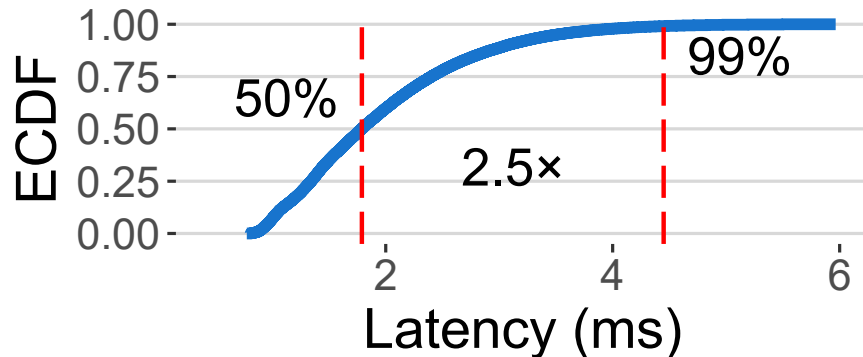
(d) Runpod
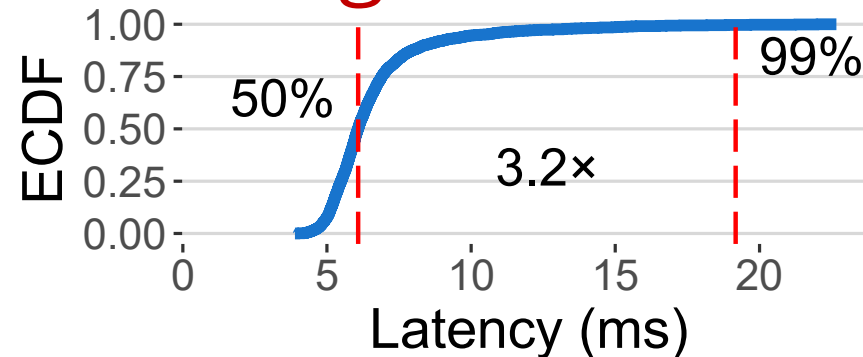
Tail: More than **3x higher** than median

GPUs remain idle during this time!

13

# How to Mitigate this Tail?



Node 2

Node 1

**Congested link**

AllReduce

**Network delays**

⚠️ *Bottleneck*

Node 3

Node 4

# How to Mitigate this Tail?



**Node 2**

**Slow**

Tail: Wait until the **last packet** arrives

Mitigate Tail: **Ignore slow arriving data**

**Node 4**

# How to Mitigate this Tail?

ML Models are Resilient
Against Gradient Loss:
**No need for 100%
reliability**

**Node 2**

*Any late arriving data is
considered lost and ignored*

**Slow**

Tail: Wait until the **last packet** arrives

Mitigate Tail: **Ignore slow arriving data**

**Node 4**

# Unreliable Communication



**Node 2**

**- No retransmissions or reliability guarantees**

*Congested link*

**Slow node**

**Node 1**

AllReduce

⚠ *Bottleneck*

**Node 3**

*Network delays*

**Node 4**

# Unreliable Communication



Node 2

**- No retransmissions or reliability guarantees**

**- A strict time-bound on the communication**

Congested ~~link~~

Slow ~~node~~

AllReduce

Node 1

Node 3

Network ~~delays~~

⚠️ *Bottleneck*

Node 4

# Unreliable Communication



Node 2

Node 1

AllReduce

Node 3

Node 4

- No retransmissions or reliability guarantees

- A strict time-bound on the communication

19

# Unreliable Communication

**Node 2**

- **No retransmissions or reliability guarantees**

- **A strict time-bound on the communication**

## ML Models are Resilient Against *Some* Gradient Loss!

**Node 1**

**Node 3**

**Node 4**

# Unreliable Communication

*Existing architectures
exacerbate the loss!*
→ *Accuracy degradation*

**Node 2**

**- No retransmissions or
reliability guarantees**

**- A strict time-bound on
the communication**

## ML Models are Resilient Against *Some* Gradient Loss!

**Node 1**

**Node 3**

**Node 4**

# UDP for Unreliable Transport

# UDP for Unreliable Transport

# UDP for Unreliable Transport

# UDP for Unreliable Transport



Local Gradient Data

Unreliable Link

**Node 1**

**Node 2**

**Node 3**

**Node 4**

Network

**Congestion**

*Fire and forget*

25

# UDP for Unreliable Transport



**Node 2**

**Node 1**

**Node 3**

**Node 4**

Network

Local Gradient Data

Unreliable Link

**Observation 1:**
Vanilla UDP fires away and leads to undue losses!

**Worsen Congestion**

**More Drops**

*Fire and forget*

# Timeouts to Mitigate Slow Workers

Local Gradient Data

Unreliable Link

**Node 2**

*How to set the timeout?*

Network

**Node 1**

**Node 3**

**Node 4**

Time-window

**X** Dropped

Time

# Timeouts to Mitigate Slow Workers

**Node 2**

**Node 1**

Network

**Node 3**

**Node 4**

Local Gradient Data

Unreliable Link

*How to set the timeout?*

1. Large timeout value

**Wasted time**

Time-window

Time

# Timeouts to Mitigate Slow Workers



**Node 2**

Local Gradient Data

Unreliable Link

*How to set the timeout?*
1. Large timeout value
2. Small timeout value

**Node 1**

Network

**Node 3**

<u>Dynamic</u>
<u>Conditions</u>

**Node 4**

**Excessive drops**

Time-window

**X** Dropped

Time

# Timeouts to Mitigate Slow Workers



Local Gradient Data

Unreliable Link

**Node 2**

Network

**Dynamic Conditions**

**Node 1**

**Node 3**

**Node 4**

**Observation 2:**
Static timeouts can add delays or increase loss!

**Excessive drops**

Time-window

**X** Dropped

Time

# Existing Architectures: Ring AR

# Existing Architectures: Ring AR

# Existing Architectures: Ring AR

Local Gradient Data

**Node 2**



$w$

**Node 1**

$x+y$

**Node 3**

**Node 4**

$z$

# Existing Architectures: Ring AR



Local Gradient Data

$w$

**Node 2**

**Node 1**

**Node 3**

**Node 4**

$x+y+z$

# Existing Architectures: Ring AR

# Existing Architectures: Ring AR

# Existing Architectures: Ring AR

# Existing Architectures: Ring AR
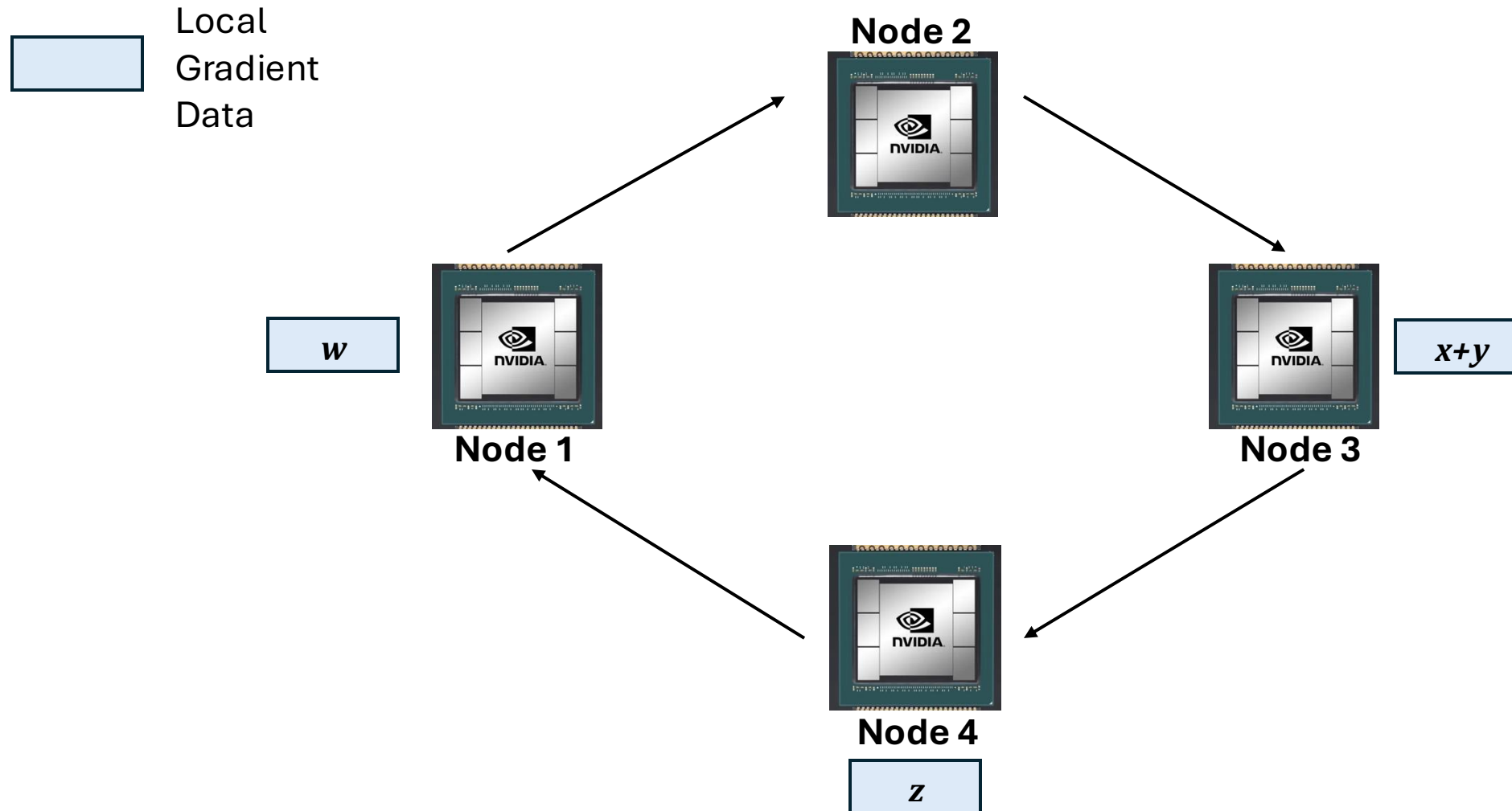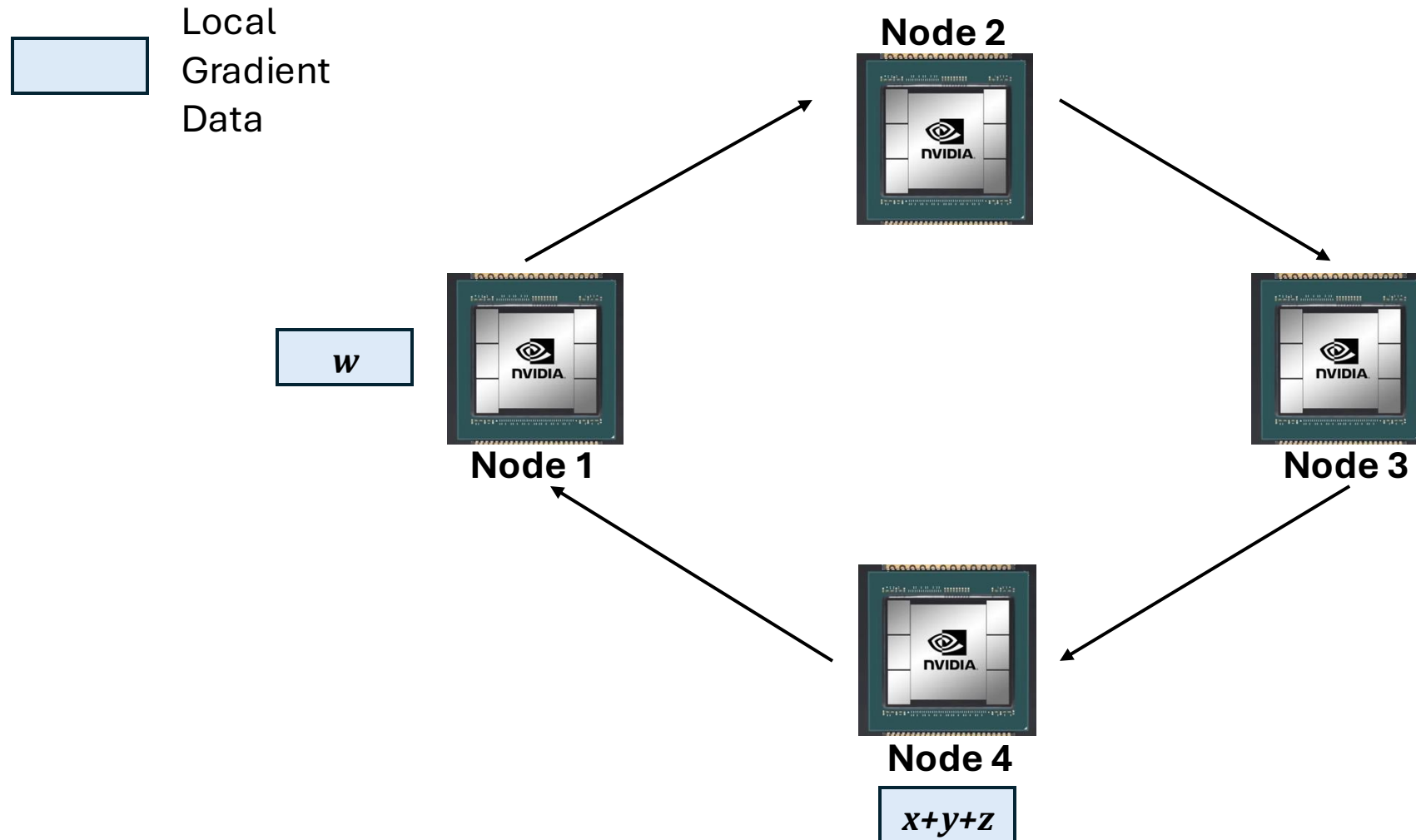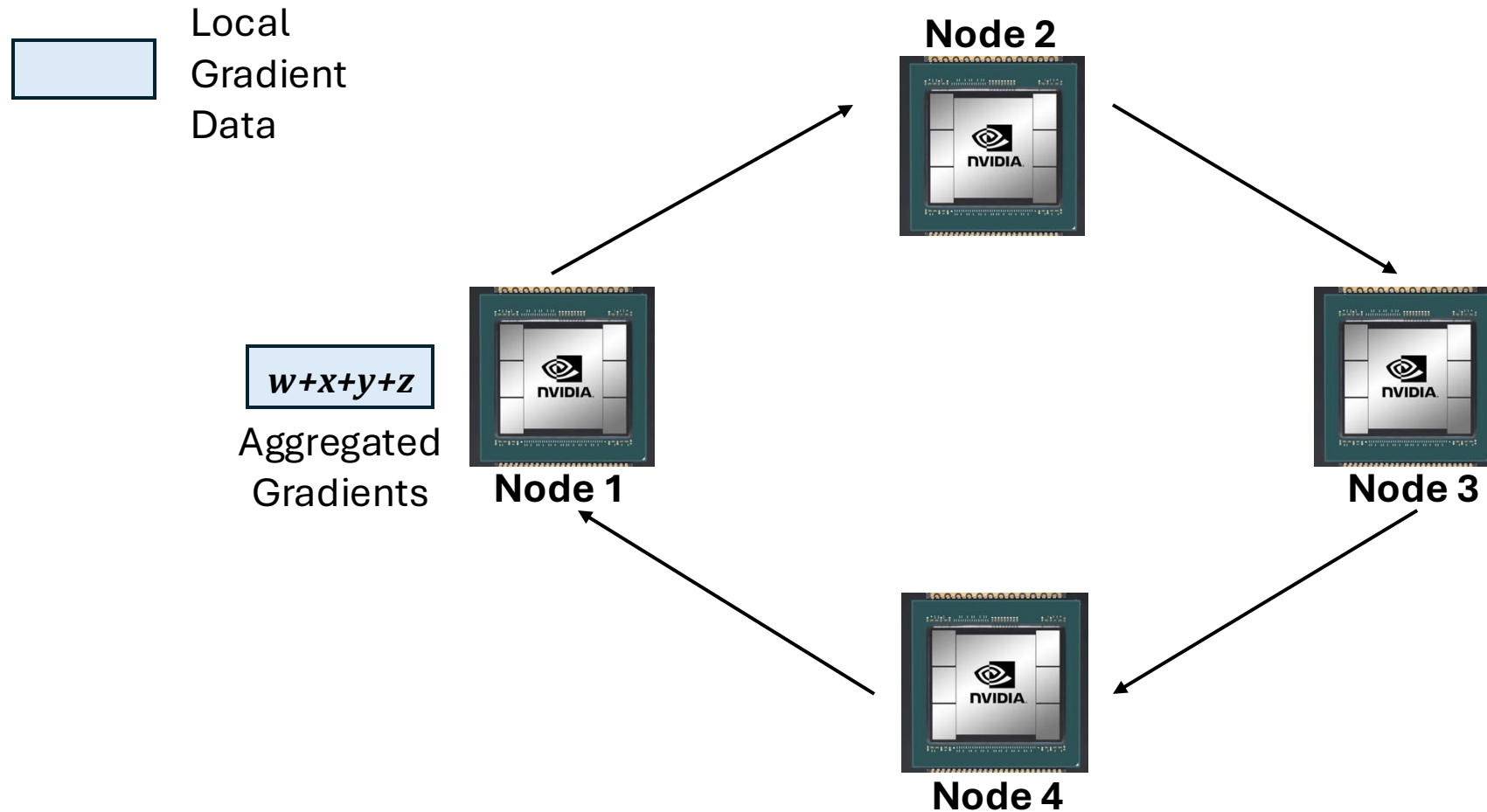
# Existing Architectures: Ring AR

# Existing Architectures: Ring AR

Local Gradient Data

Unreliable Link

**Node 2**

All the **aggregated data lost**, not just Node 4's gradients!

*x+y+z*

**Node 4**

*x+y+z*

# Existing Architectures: Ring AR



**Observation 3:**
Ring AllRreduce exacerbates the loss!

Local Gradient Data

Unreliable Link

**Node 2**

All the **aggregated data lost**, not just Node 4's gradients!

*x+y+z*

**Node 4**

*x+y+z*

# OptiReduce Design

**Observation 1:**
Vanilla UDP fires away and adds undue losses

→ *1. Unreliable Bounded Transport*
To capitalize & bound the loss

**Observation 2:**
Static timeouts can add delays or increase loss

→ *2. Adaptive Timeouts*
To handle dynamic delays

**Observation 3:**
Ring AllRreduce exacerbates the loss

→ *3. Transpose AllReduce*
To minimize loss

# OptiReduce Design

**Observation 1:**
Vanilla UDP fires away and adds undue losses

→

*1. Unreliable Bounded Transport*
To capitalize & bound the loss

Stat... ...ts
delays or increase loss ...ays

**100% reliability**

***Finish faster!***

**Observation 3:**
Ring AllRreduce exacerbates the loss

→

*3. Transpose AllReduce*
To minimize loss

43

# OptiReduce Design

**Observation 1:**
Vanilla UDP fires away and adds undue losses

➡

*1. Unreliable Bounded Transport*
To capitalize & bound the loss

**Observation 2:**
Static timeouts can add delays or increase loss

➡

*2. Adaptive Timeouts*
To handle dynamic delays

**Observation 1:**
Ring AllRreduce exacerbates the loss

➡

*1. Transpose AllReduce*
To minimize loss

44

# Unreliable Bounded Transport

*Reliability adds to the tail*          *Firing away exacerbates loss*

| Features | TCP | UDP |
|---|---|---|
| - Retransmissions | ✓ | ✗ |
| - In-order delivery | ✓ | ✗ |
| - Congestion Control | ✓ | ✗ |
| - Flow Control | ✓ | ✗ |

# Unreliable Bounded Transport

**Reliability adds to the tail**          **Firing away exacerbates loss**

| Features | TCP | UBT | UDP |
|---|---|---|---|
| - Retransmissions | ✓ | X | X |
| - In-order delivery | ✓ | X | X |
| - Congestion Control | ✓ | | X |
| - Flow Control | ✓ | | X |

🎯Direct placement of packets, enabled by *Offset* field!

# Unreliable Bounded Transport

**Reliability adds to the tail**   **Firing away exacerbates loss**

| Features | TCP | UBT | UDP |
|---|---|---|---|
| - Retransmissions | ✓ | ✗ | ✗ |
| - In-order delivery | ✓ | ✗ | ✗ |
| - Congestion Control | ✓ | ✓ | ✗ |
| - Flow Control | ✓ | ✓ | ✗ |

🎯 Direct placement of packets, enabled by *Offset* field!

⏱ Timely-inspired rate-control to throttle the sending!

# Unreliable Bounded Transport

*Effect of this approach:*
**Capitalizes and bounds the loss**

***Reliability adds to the tail***

| Features | TCP | UBT | UDP |
|---|---|---|---|
| - Retransmissions | ✓ | ✗ | ✗ |
| - In-order delivery | ✓ | ✗ | ✗ |
| - Congestion Control | ✓ | ✓ | ✗ |
| - Flow Control | ✓ | ✓ | ✗ |

🎯 Direct placement of packets, enabled by *Offset* field!

⏱ Timely-inspired rate-control to throttle the sending!

# Unreliable Bounded Transport

*Reliability adds to the tail*

*Effect of this approach:*
**Capitalizes and bounds the loss**

| Features | TCP | UBT | UDP |
|---|---|---|---|
| - Retransmissions | ✓ | ✗ | ✗ |
| - In-order delivery | ✓ | ✗ | ✗ |
| - Congestion Control | ✓ | ✓ | ✗ |
| - Flow Control | ✓ | ✓ | ✗ |

🎯 Direct placement of packets, enabled by *Offset* field!

Timely-inspired rate-control to throttle the sending!

Receiver-driven multicast, signaled using *Incast* field!

49

# OptiReduce Design

*Observation 1:*
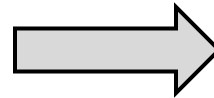Vanilla UDP fires away and adds undue losses

➡

*1. Unreliable Bounded Transport*
To capitalize & bound the loss

**Observation 2:**
Static timeouts can add delays or increase loss

➡

*2. Adaptive Timeouts*
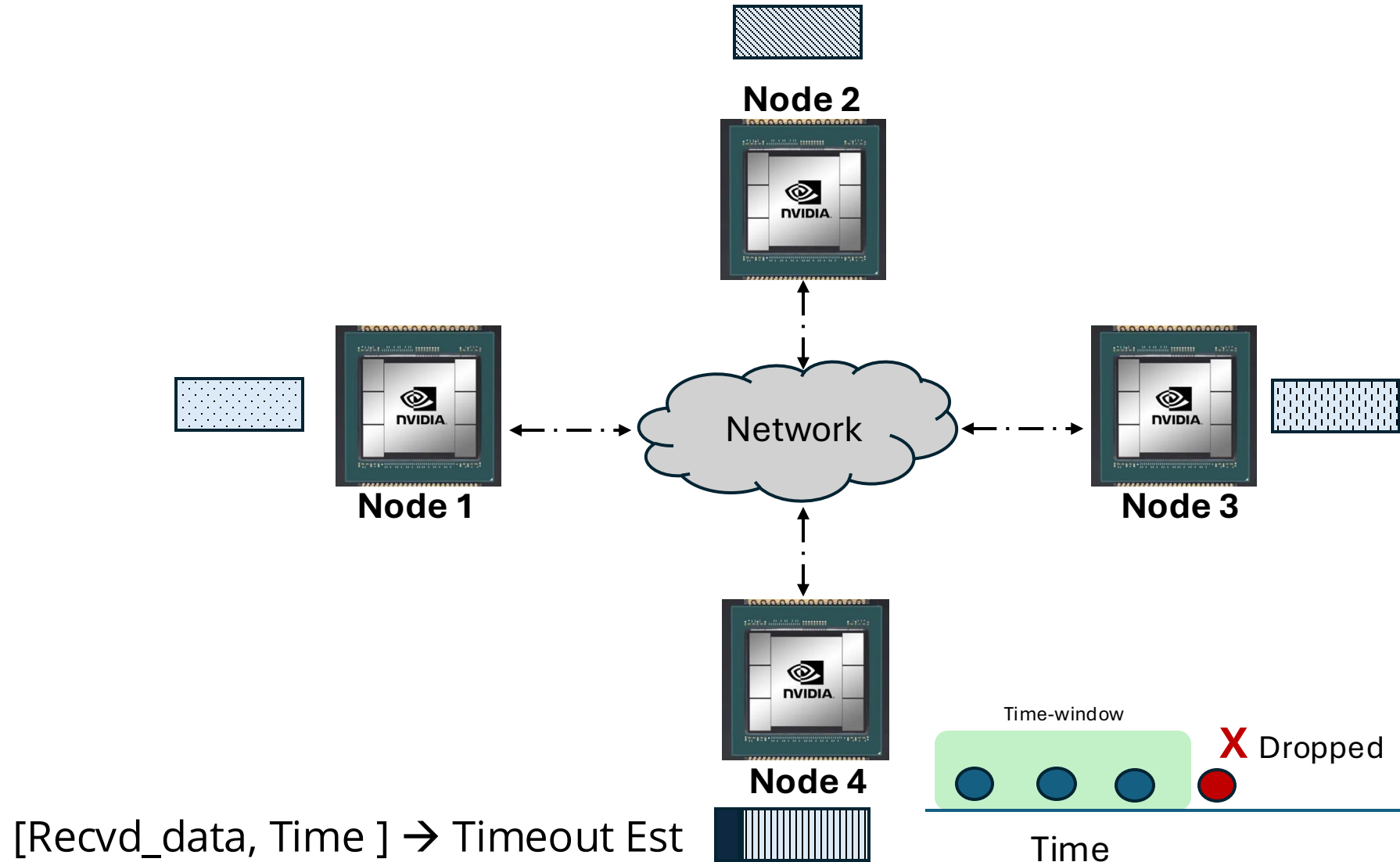To handle dynamic delays

*Observation 3:*
Ring AllRreduce exacerbates the loss

➡

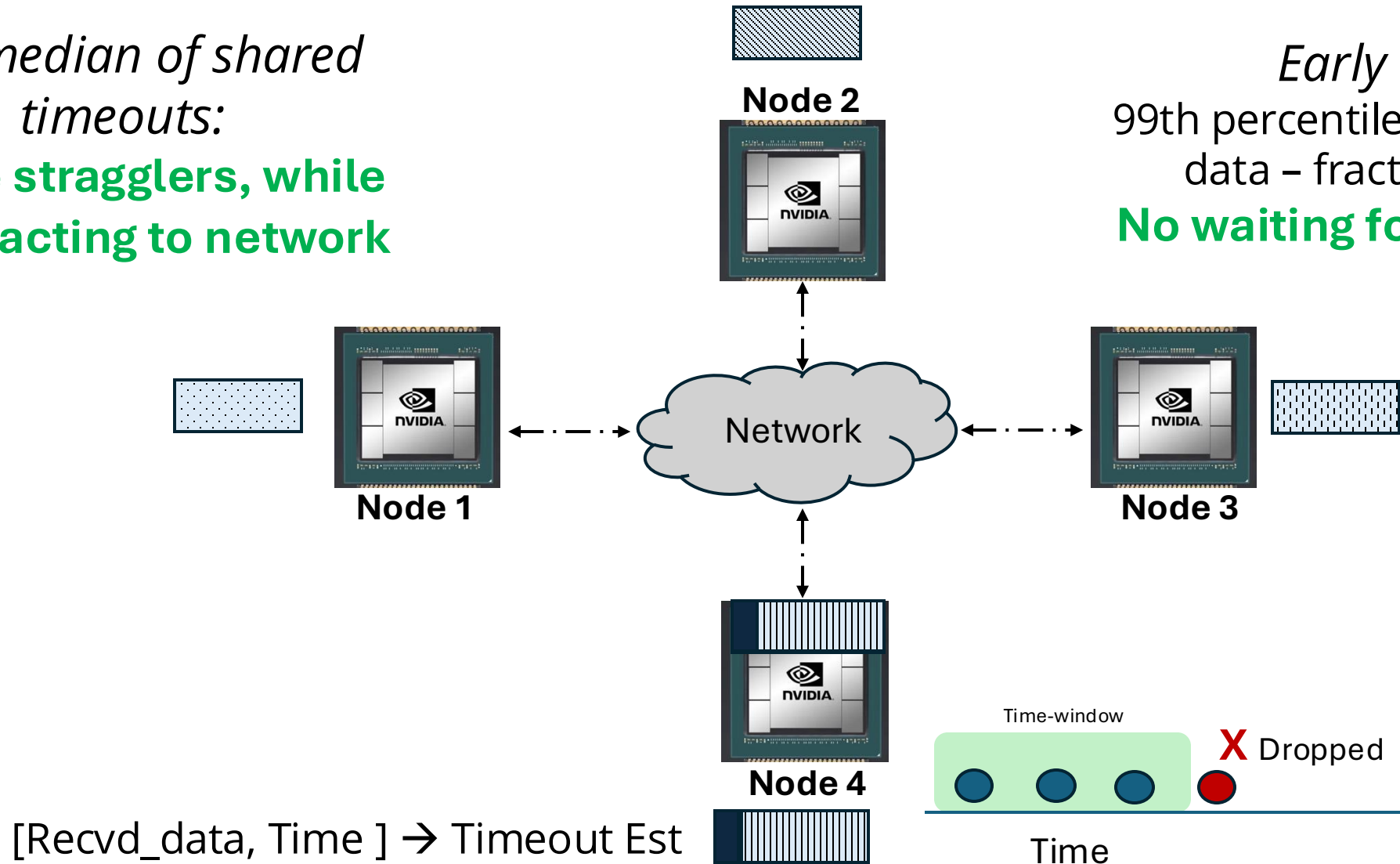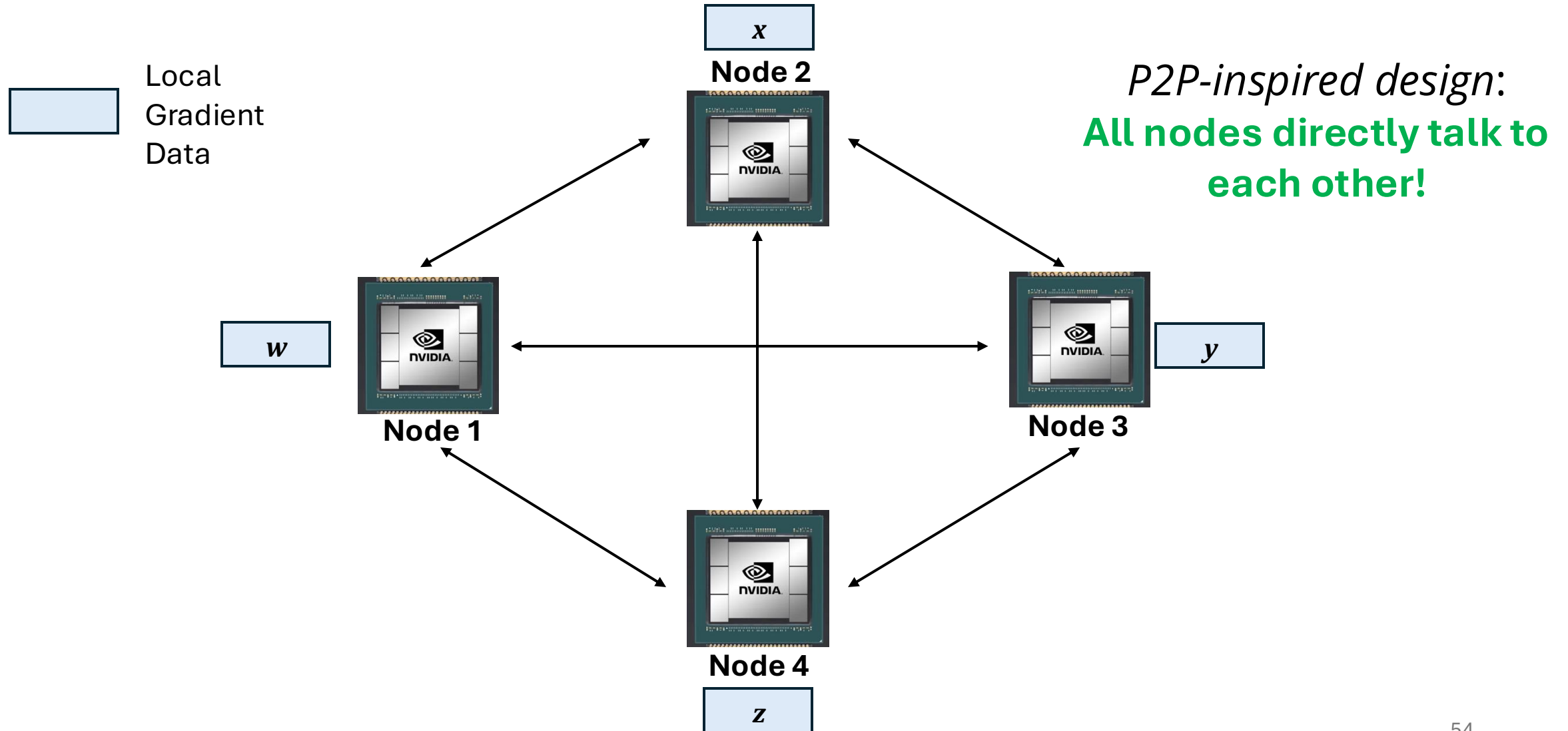*3. Transpose AllReduce*
To minimize loss

# Adaptive Timeouts



**Node 2**

**Node 1**

Network

**Node 3**

**Node 4**

[Recvd_data, Time ] → Timeout Est

Time-window

**X** Dropped

Time

51

# Adaptive Timeouts

*Use median of shared timeouts:*
**Ignore stragglers, while still reacting to network**

**Node 2**

Network

**Node 1**

**Node 3**

*Early timeout:*
99th percentile arrived & no new data – fraction of timeout
**No waiting for dropped data**

**Node 4**

[Recvd_data, Time ] → Timeout Est

Time-window

**X** Dropped

Time

52

# OptiReduce Design

**Observation 1:**
Vanilla UDP fires away and adds undue losses

*1. Unreliable Bounded Transport*
To capitalize & bound the loss

**Observation 2:**
Static timeouts can add delays or increase loss

*2. Adaptive Timeouts*
To handle dynamic delays

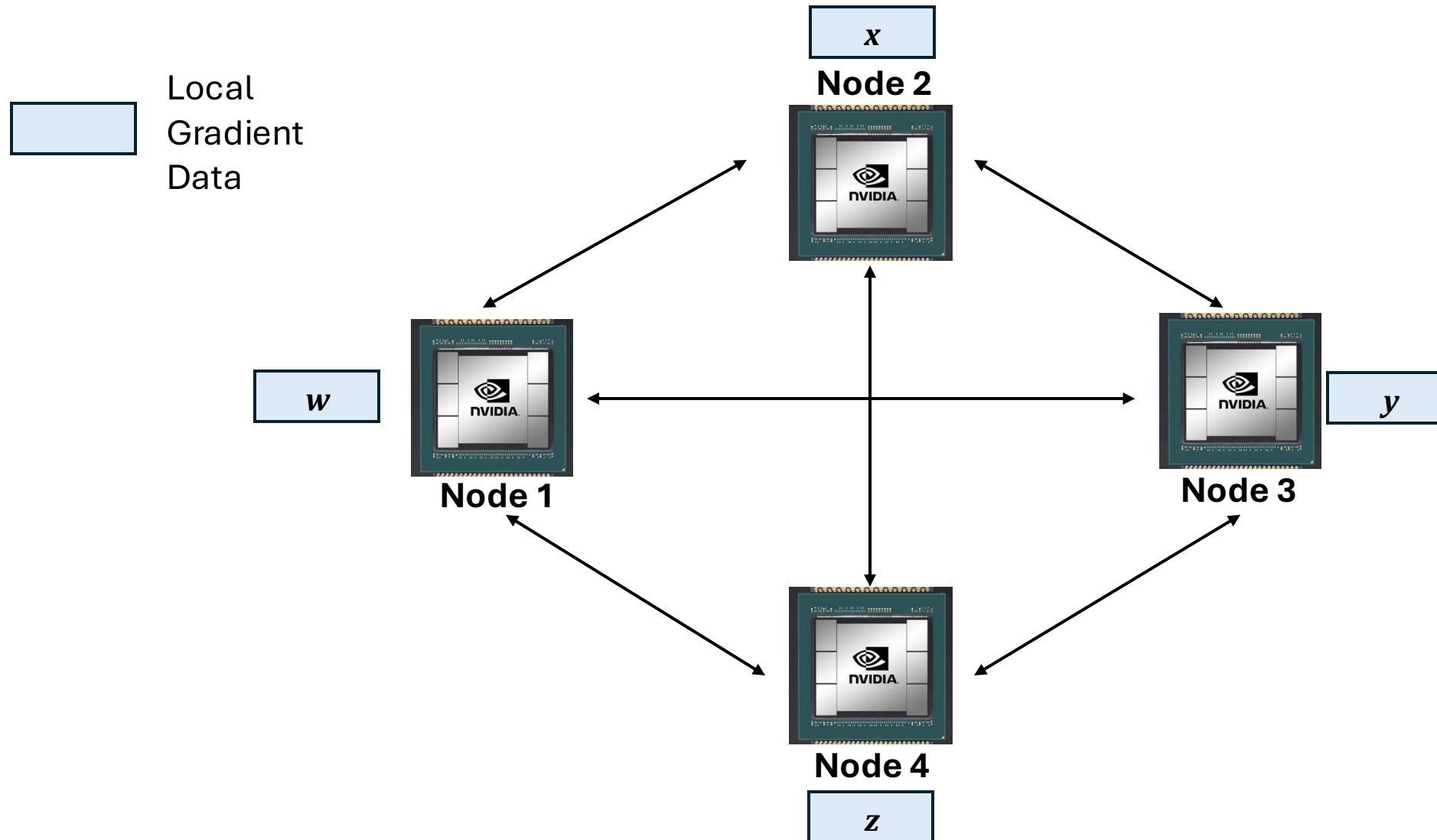**Observation 3:**
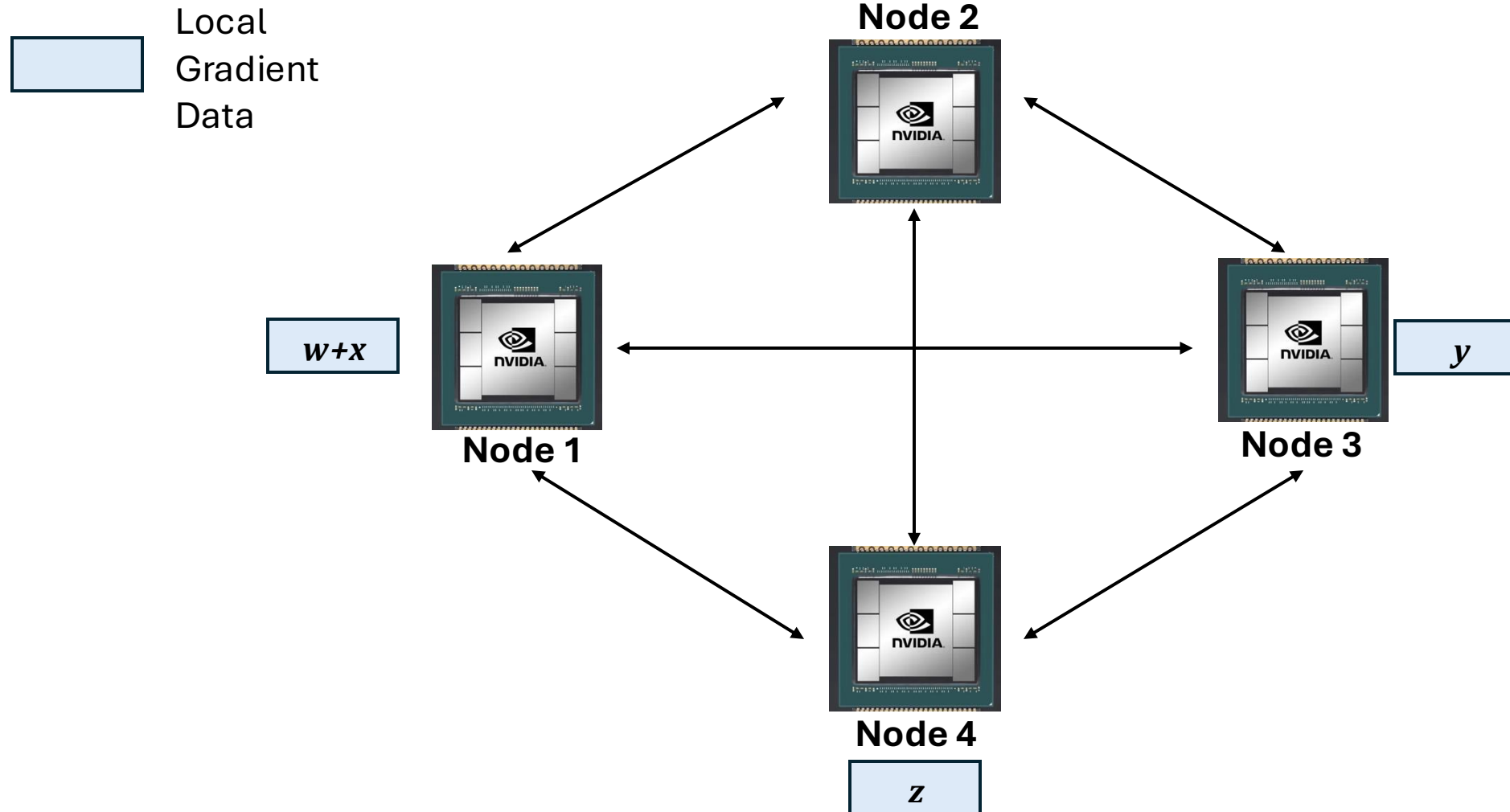Ring AllRreduce exacerbates the loss

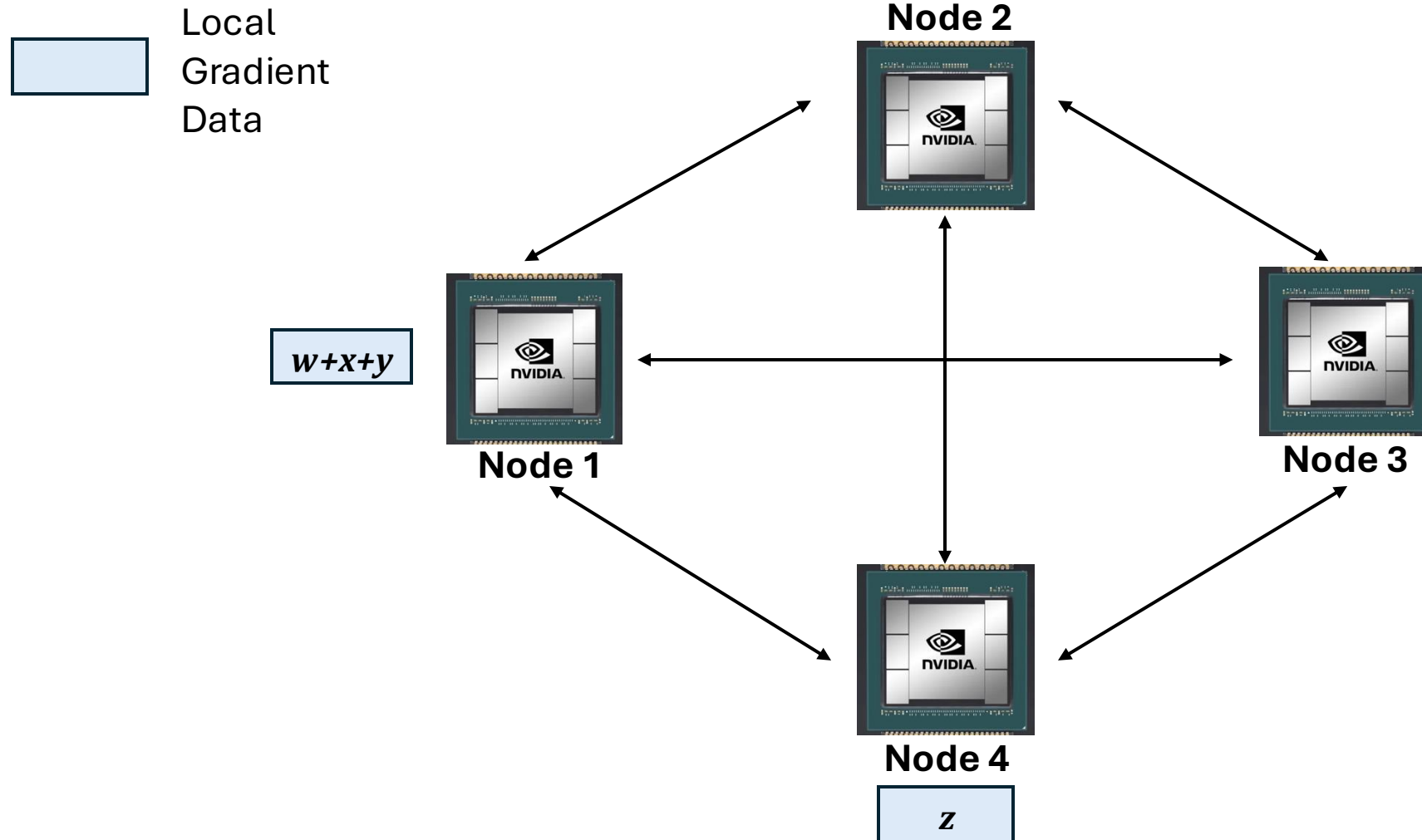*3. Transpose AllReduce*
To minimize loss

# Transpose AllReduce (TAR)



Local Gradient Data

*P2P-inspired design*: **All nodes directly talk to each other!**

# Transpose AllReduce (TAR)

# Transpose AllReduce (TAR)



Local Gradient Data

**Node 2**

**Node 1**

$w+x$

**Node 3**

$y$

**Node 4**

$z$

# Transpose AllReduce (TAR)



Local Gradient Data

**Node 2**

**Node 1**

$w+x+y$

**Node 3**

**Node 4**

$z$

# Transpose AllReduce (TAR)



Local Gradient Data

$w+x+y+z$

Aggregated Gradients

**Node 1**

**Node 2**

**Node 3**

**Node 4**

*Same number of rounds as Ring AllReduce*

# Transpose AllReduce (TAR)



Local Gradient Data

Unreliable Link

$x$

Node 2

$w$

Node 1

$y$

Node 3

$z$

Node 4

# Transpose AllReduce (TAR)

# Transpose AllReduce (TAR)

Local Gradient Data

Unreliable Link

**Node 2**

**Node 1**

$w+x+y$

$z$

**Node 3**

**Node 4**

$z$

# Transpose AllReduce (TAR)



**Node 2**

**Node 1**

**Node 3**

**Node 4**

Local Gradient Data

Unreliable Link

$w+x$

$z$

$z$

*Effect of this approach:*
**Ring's MSE 7x larger than TAR's!**

Only a **single node's data** lost!

# Optimizations for Sustaining Accuracy

- Safeguards against excessive loss

- Randomized Hadamard Transform (RHT)

- Encode data before sending in the network
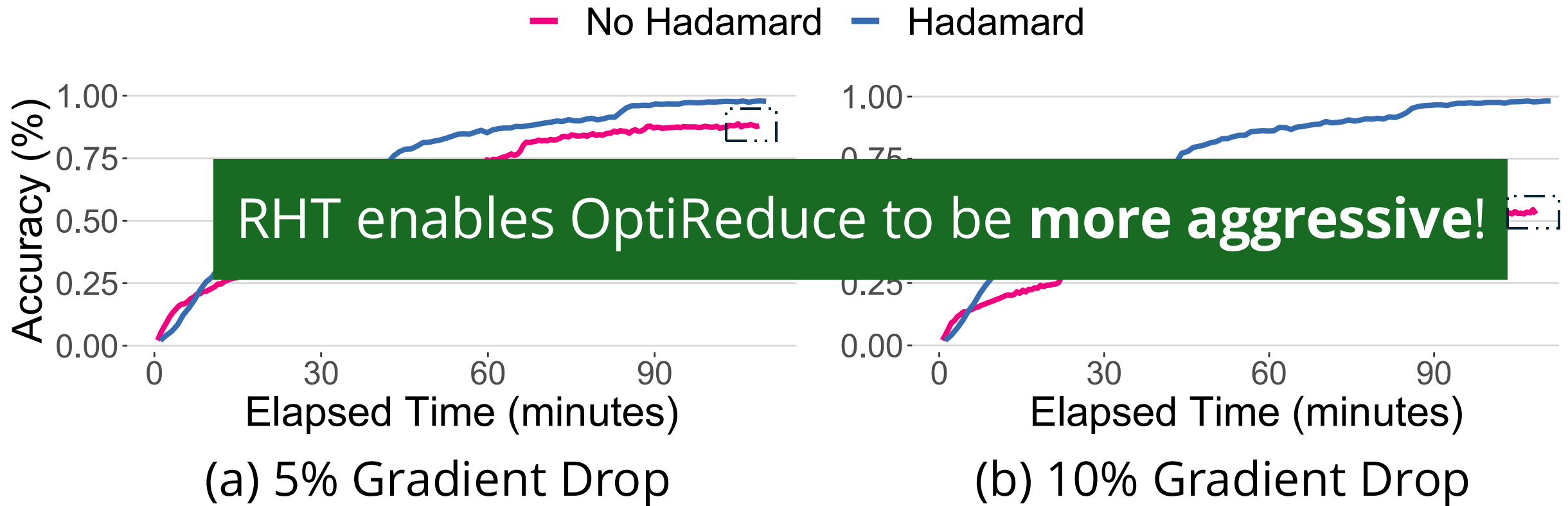
- Approximate lost gradients, essentially recovering data

| Send Data | | | |
|---|---|---|---|
| 1.0 | | 5 | |
| 1.5 | | -13 | |
| 2.0 | HT | -5 | |
| 2.5 | | 5 | |
| 3.0 | | 3 | |
| 3.5 | | 17 | |
| 4.0 | | -5 | |
| 4.5 | | 1 | |

| Tail Drop | | Recv Data | |
|---|---|---|---|
| 5 | | 0.9 | |
| -13 | | 1.4 | |
| -5 | HT | 2.1 | |
| 5 | | 2.4 | |
| 3 | | 2.9 | |
| 17 | | 3.4 | |
| -5 | | 4.1 | |
| 0 | | 4.6 | |

"Hadamard transform applied to a vector of length 8"

# Effect of RHT on Accuracy



(a) 5% Gradient Drop

(b) 10% Gradient Drop

Model: VGG-19

# Effect of RHT on Accuracy



(a) 5% Gradient Drop

(b) 10% Gradient Drop

Model: VGG-19

RHT enables OptiReduce to be **more aggressive**!

# Evaluating OptiReduce

▬ Gloo Ring   ▬ Gloo BCube   ▬ NCCL Ring   ▬ NCCL Tree   ▬ TAR+TCP   ▬ OptiReduce

**Baselines**
- NCCL
- Gloo

**Environments**
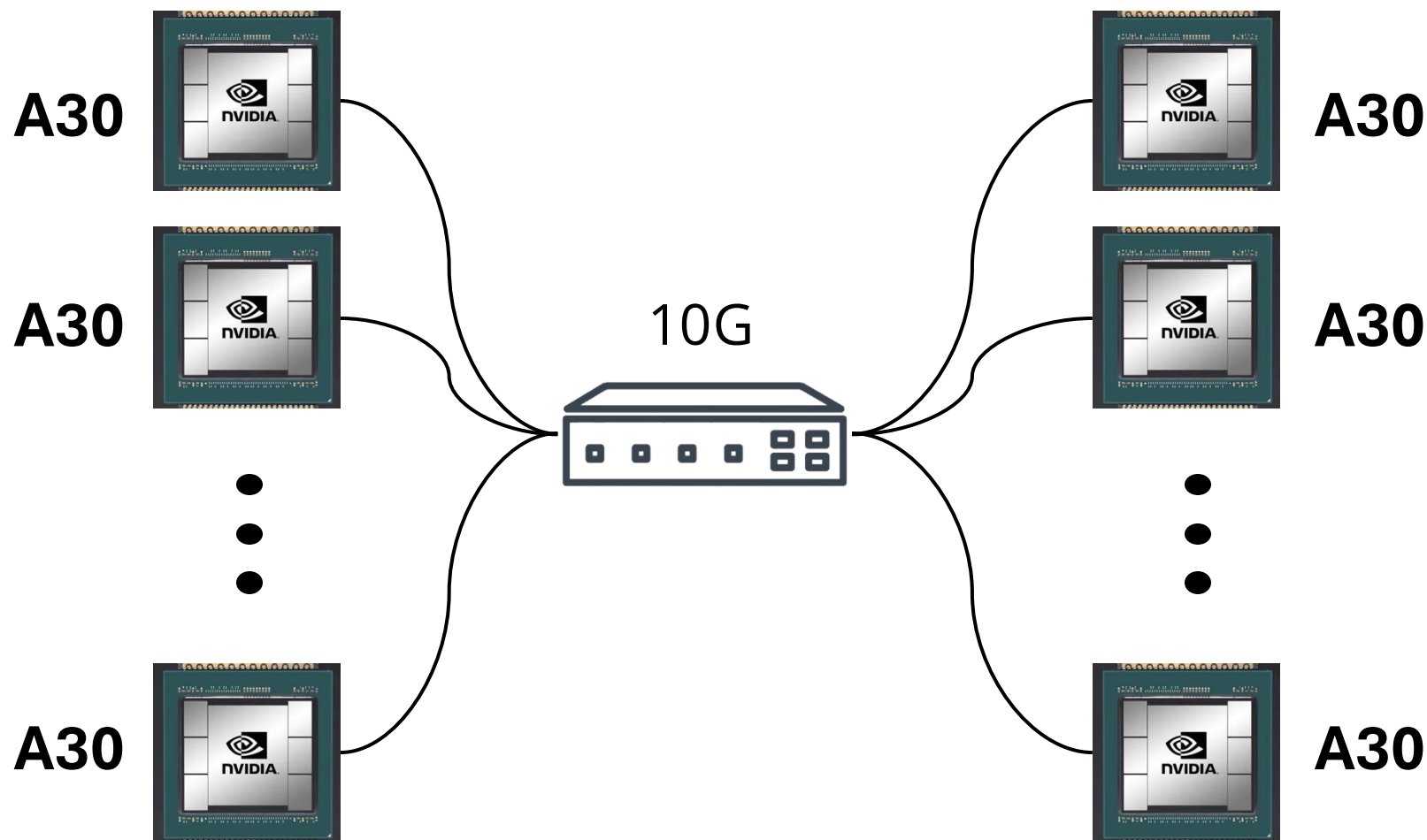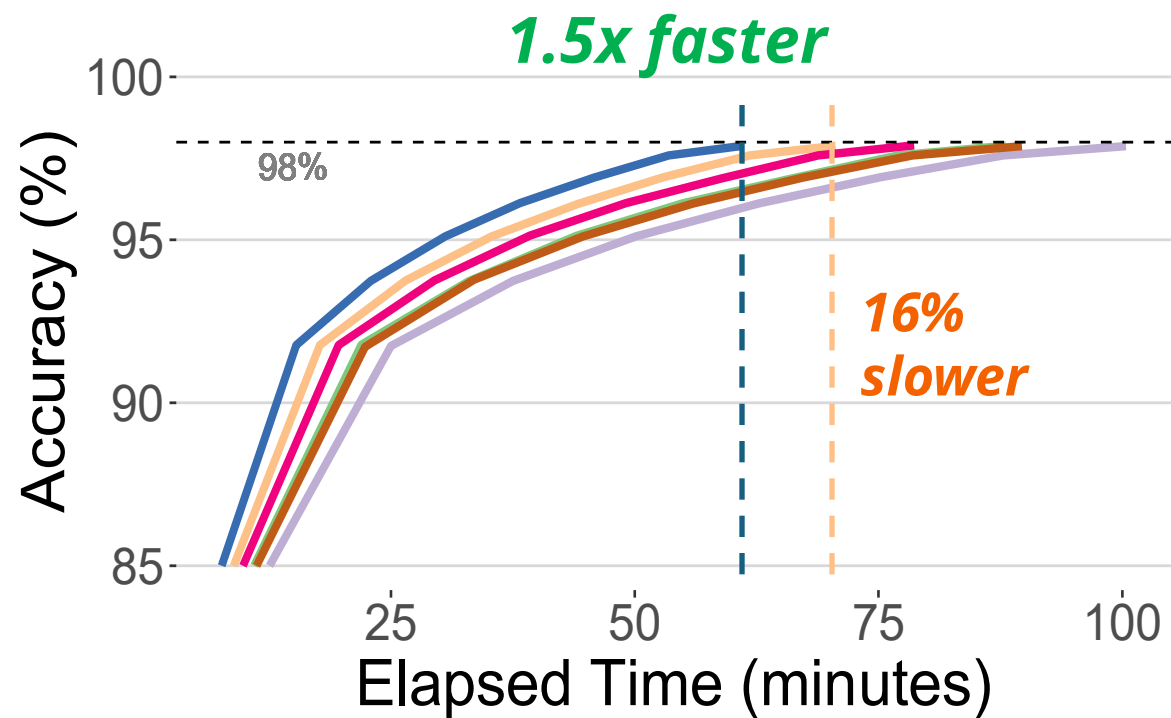- Cloudlab
- Local Setup

**Metrics**
- TTA
- Throughput

# Cloudlab: A Public Cloud Environment

# Cloudlab: A Public Cloud Environment



Legend: Gloo Ring, Gloo BCube, NCCL Ring, NCCL Tree, TAR+TCP, OptiReduce

**1.5x faster**

**16% slower**

98%

Accuracy (%) vs. Elapsed Time (minutes)

**TTA GPT-2**

# Our Local Setup



V100

V100

V100

V100

V100

V100

25G

Background
Traffic

# Our Local Setup



(a) Low Tail Env

(b) High Tail Env

(a) $P_{99/50} = 1.5$

(b) $P_{99/50} = 3$

# Time-to-Accuracy (TTA)



(a) Low Tail Env

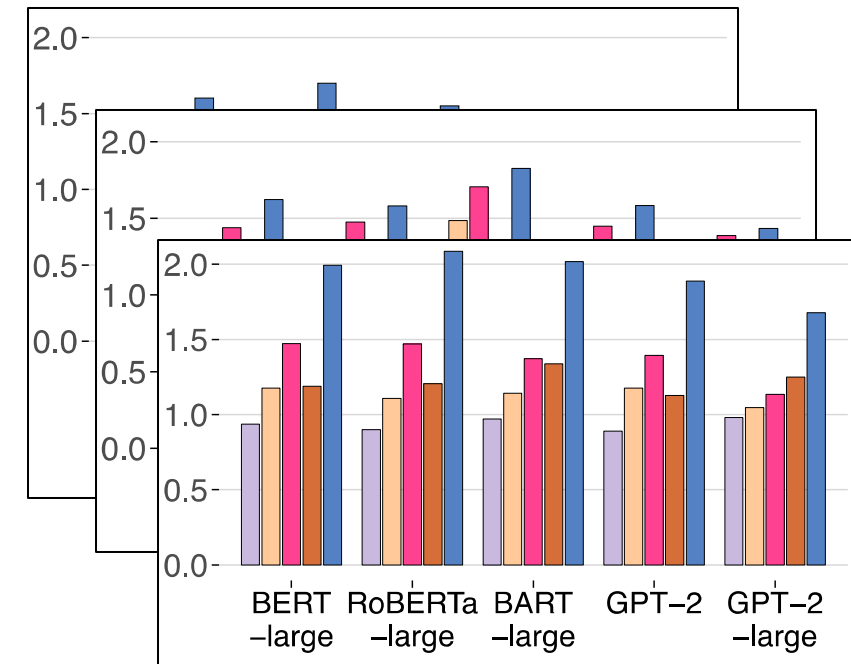(b) High Tail Env

**TTA: GPT-2**

# Training Throughput (Speedup)



**Training Throughput: Llama-3.2**

# More Evaluations in the Paper

- OptiReduce **scalability** results
  - **24 nodes** evaluation in local cluster
  - **144 nodes** in simulations
- Comparison with **other schemes**
  - **Parameter Server** approach
  - **Quantization** and **Sparsification** schemes
  - **In-network aggregation** (INA) approaches
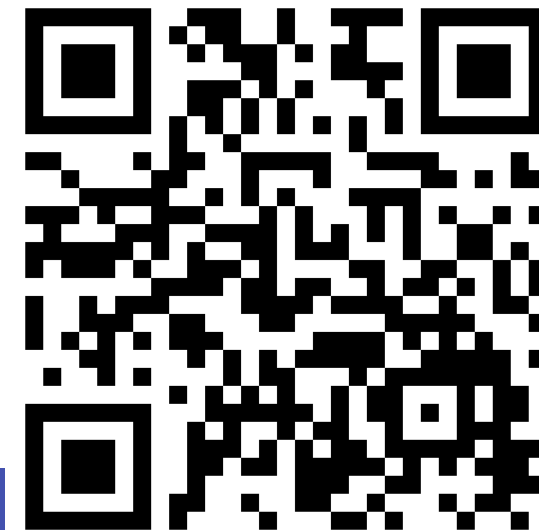- More Models and Datasets
  - More **LLMs** and **vision models**

# Conclusion

- **AllReduce bottleneck** in cloud training
- OptiReduce → **Time-bounded AR** for **Cloud**
- More than **2× Speedup** in **high-tail**
- Try **OptiReduce** – follow the link!

**optireduce.github.io**

# Thank You!

optireduce.github.io

**OptiReduce - Optimizing Large-Scale ML Training**

Home    Getting Started    Installation    Usage Guide    Benchmark    Technical Details    Contributing    References

Home

Why OptiReduce?

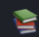🚀 Faster Training

💡 Key Features

🔧 Technical Highlights

Choose your path with OptiReduce

👉 Quick Start

📚 Learn More

Research



**OPTIREDUCE**

OPTIMIZNG LARGE-SCALE ML TRAINING