# Enabling Efficient GPU Communication over Multiple NICs with FuseLink

**Zhenghang Ren**, Yuxuan Li, Zilong Wang, Xinyang Huang, Wenxue Li, Kaiqiang Xu, Xudong Liao, Yijun Sun, Bowen Liu, Han Tian, Junxue Zhang, Mingfei Wang, Zhizhen Zhong, Guyue Liu, Ying Zhang, Kai Chen

# Advancement of high-bandwidth GPU communication



InfiniBand/Ethernet Switches

NIC  NIC  NIC  NIC

GPU  GPU  GPU  GPU

GPU Switch

GPU  GPU  GPU  GPU

NIC  NIC  NIC  NIC

InfiniBand/Ethernet Switches

Dedicated Fabrics
~TB/s

50-100 GBps

## Dedicated interconnect

- Relatively small groups
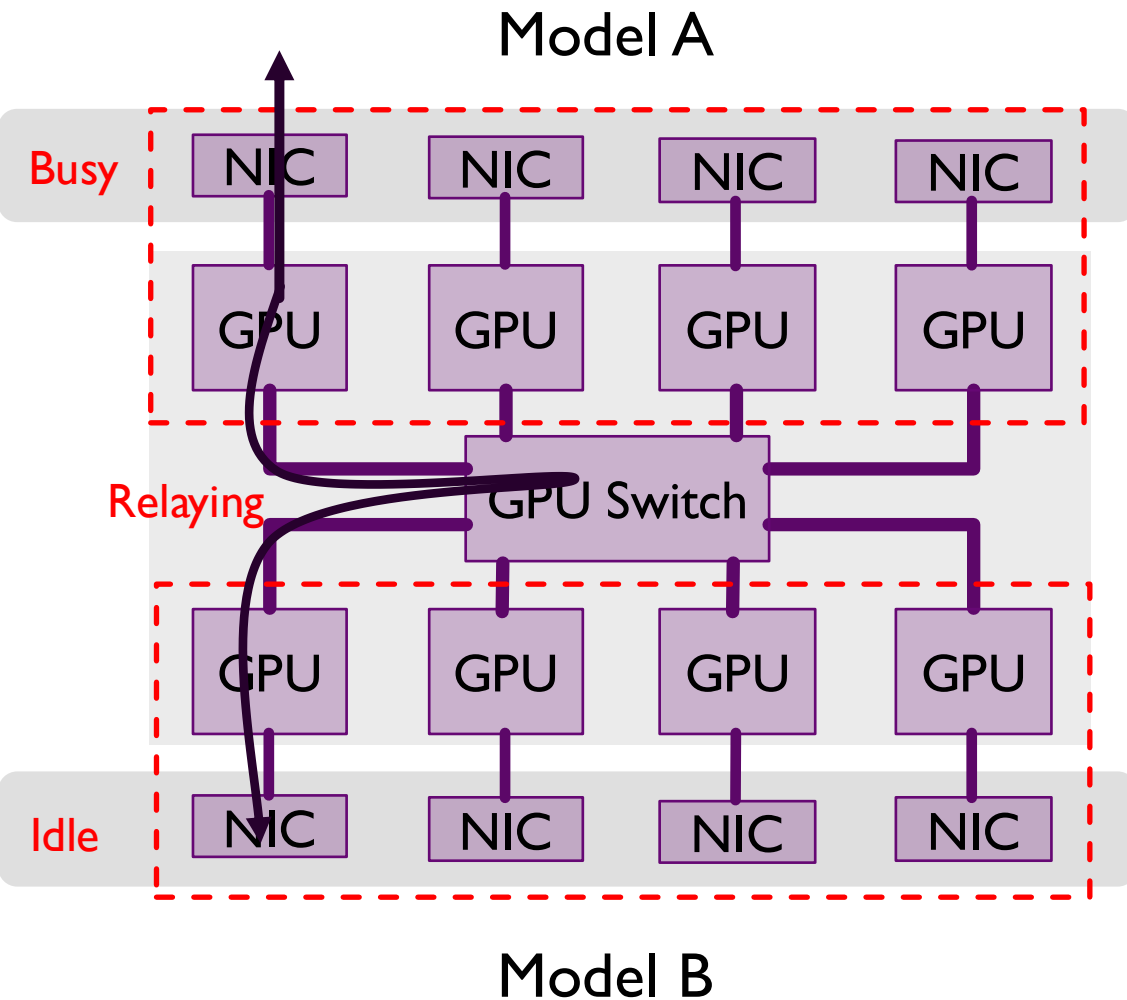- Fast growing bandwidth

## Network over general NICs

- Large scale
- Slow growing bandwidth

Can be tens of GB for every GPU!
KV Cache/Embeddings/Gradients,…

When can we multiplex the NICs and to accelerate inter-server communication?

With PCIe limitation, can we leverage **dedicated fabrics** to improve bandwidth over **general NICs**?

# Observation: traffic imbalance is common

Model A

**Busy**



**Relaying**

GPU Switch

**Idle**

Model B

The GPU server has two independent LLM serving deployments with inter-server context (KVCache) transmission.

- Models serve different requests
- Half NICs busy and half idle.

4

# Opportunity: multi-NIC acceleration under imbalanced traffic

**Biased Traffic Volume**

NICs with less traffic being idle

**Using indirect NICs**

Data traverse via CPU/UPI becomes bottleneck

**Delayed Communication**

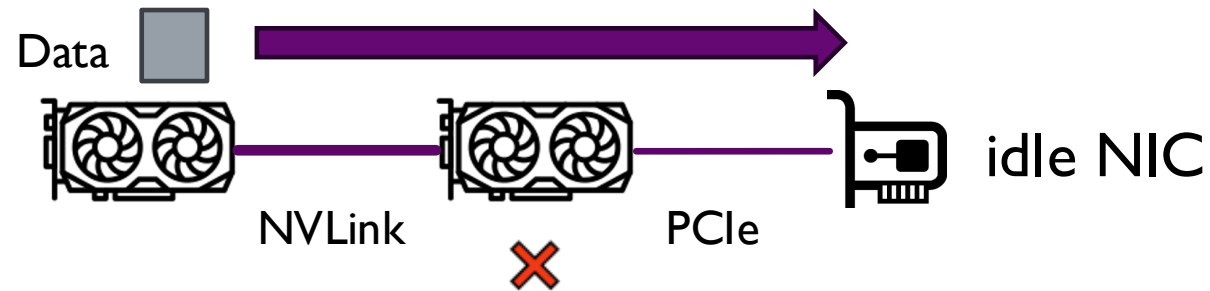Late GPUs increase transmission tail

**Suboptimal NICs throughput**

Accelerate via multi-NIC transmission

Can arbitrary ML workloads effectively leverage all available NICs, aggregating links as a "FuseLink" for inter-server communication ?

# Design goals

- Efficiently use idle NICs through dedicated intra-server fabrics

- Seamless integration into existing systems (NCCL, Gloo,…)

- Avoid contention & interruption among GPUs

# Challenge: incompatibility of dedicated fabrics and NICs

Data ▪

idle NIC

NVLink    PCIe

✖

### TX Direction

Hardware-offloaded NIC network stack cannot transmit data across GPU interconnects

### RX Direction

NIC traffic cannot be modified to be routed across GPUs

# Efficient indirect NIC communication: memory & network combined



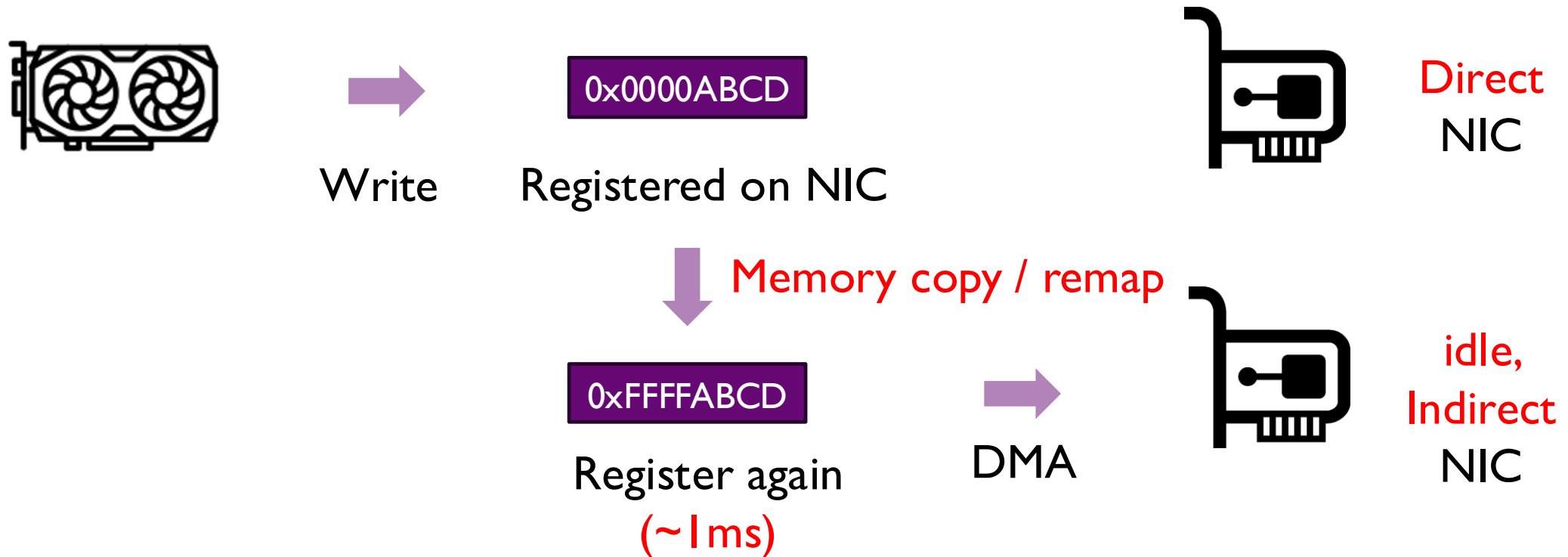Write    Registered on NIC    DMA

`0x0000ABCD`

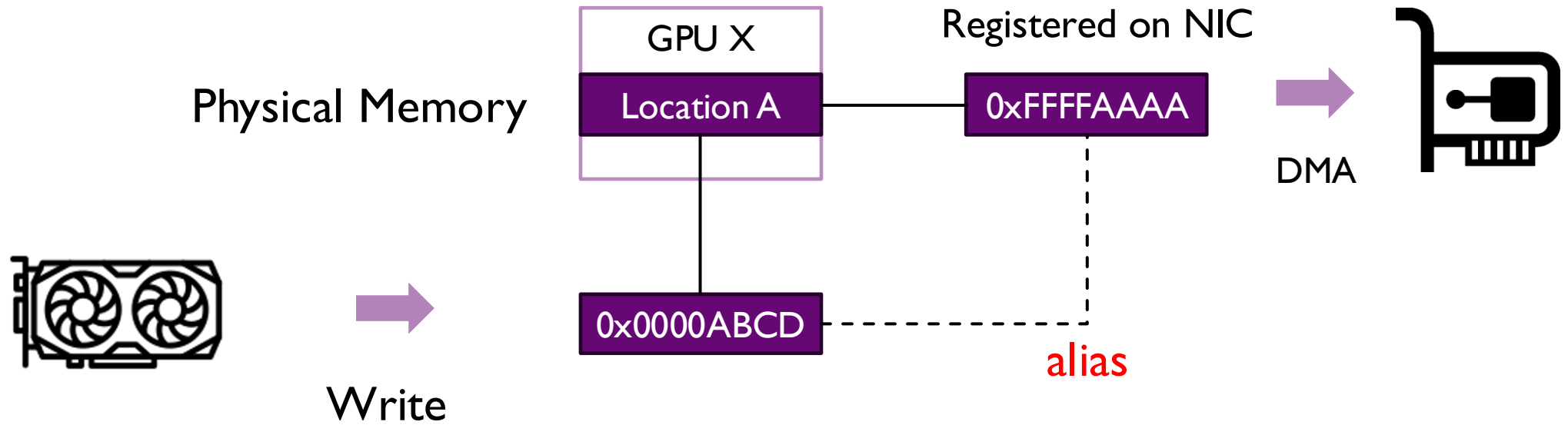GPU program fills network buffer registered on NICs

NIC reads registered (pinned) memory through PCIe

# Efficient indirect NIC communication: memory & network combined

Straight forward solution: **copy data** to intermediate GPU close to the NIC
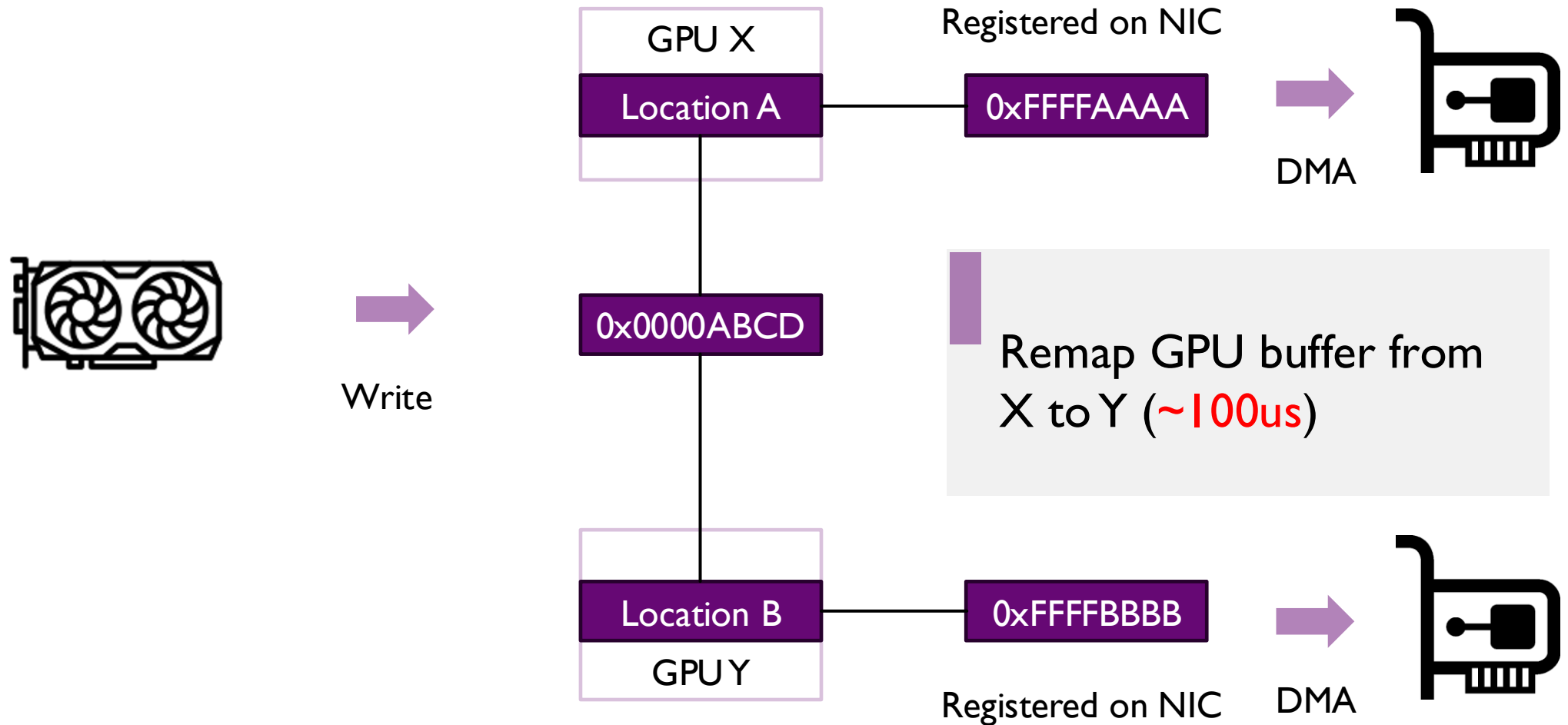Not performant because of extra copy and frequent CPU interruption

Write → 0x0000ABCD    Direct NIC
Registered on NIC

Memory copy / remap

0xFFFFABCD → DMA    idle, Indirect NIC
Register again
(~1ms)

10

# Efficient indirect NIC communication: memory & network combined

Physical Memory

GPU X

Registered on NIC

Location A ———— 0xFFFFAAAA →→ DMA

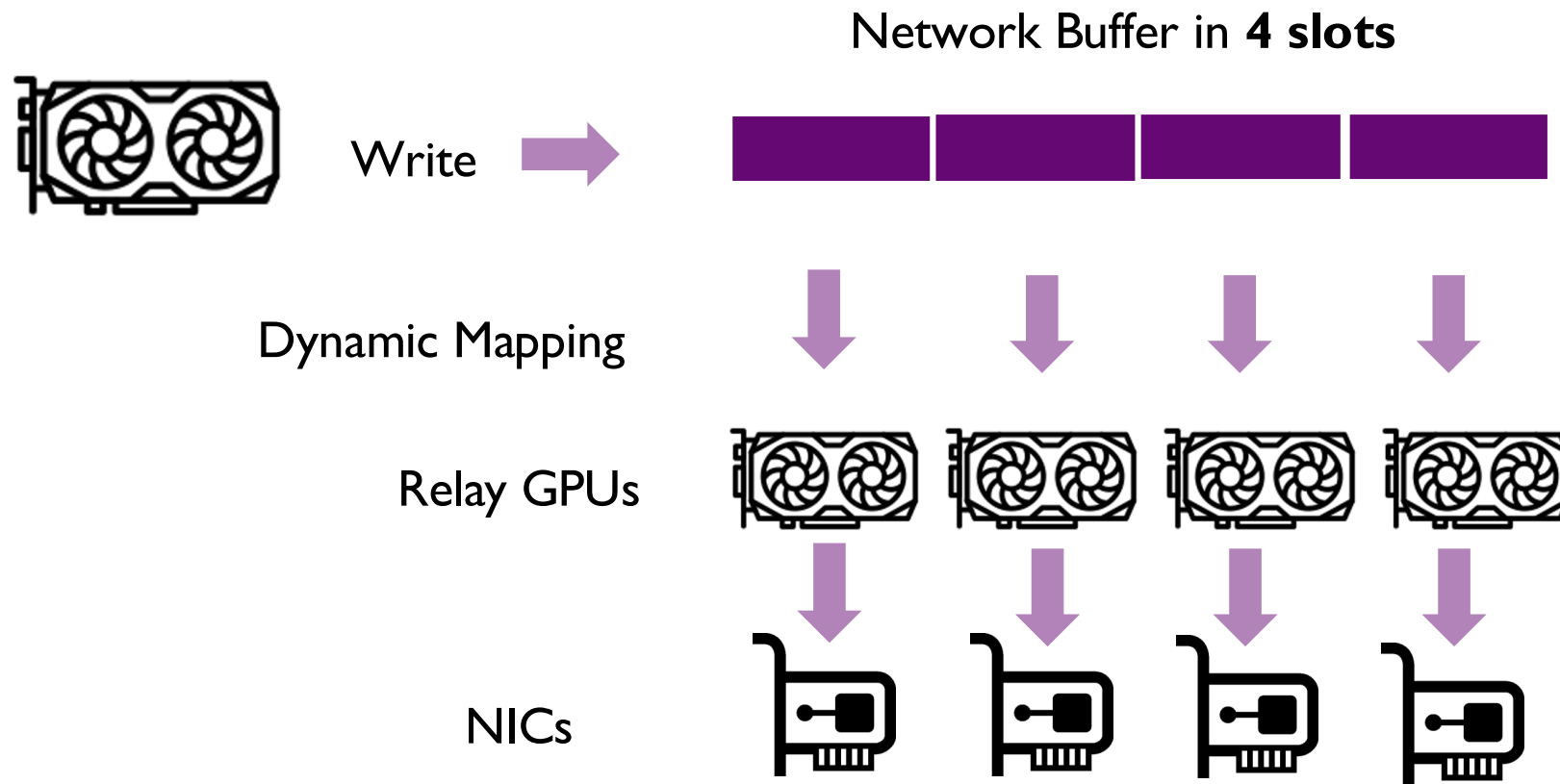0x0000ABCD - - - - - - - - ┘

alias

Write

GPUs and NICs has different memory address with the same backend (alias)
No need to change memory registration (~1ms)

11

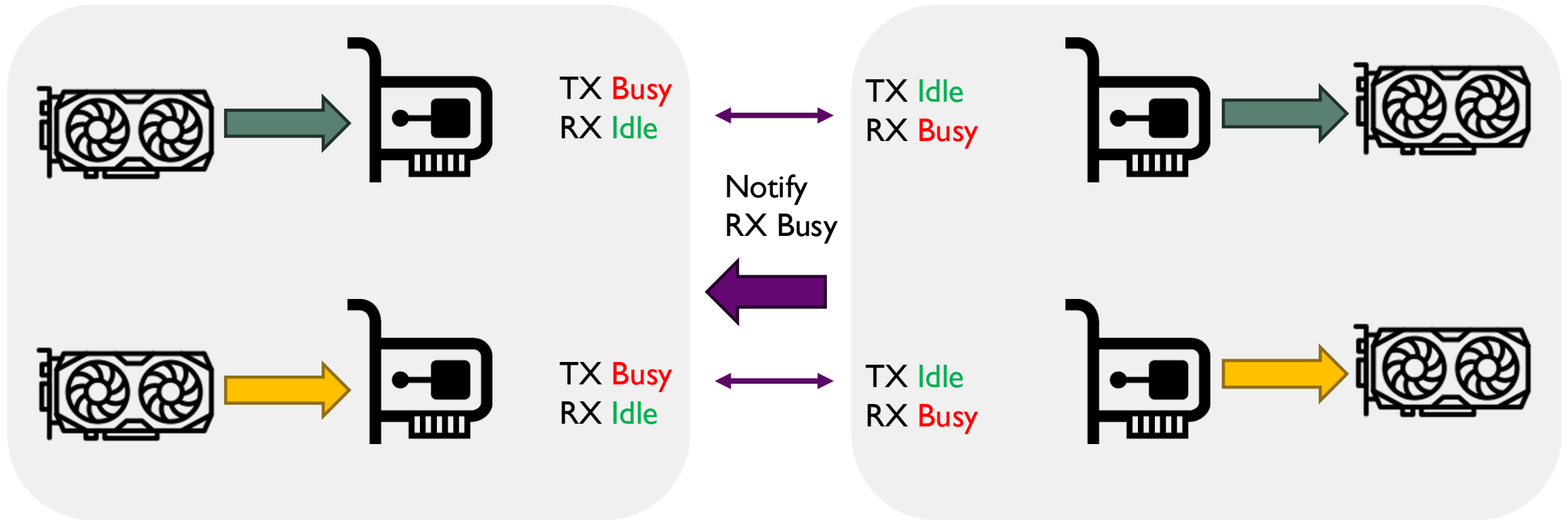# Efficient indirect NIC communication: memory & network combined

GPU X

Location A

Registered on NIC

0xFFFFAAAA

DMA

0x0000ABCD

Write

Remap GPU buffer from X to Y (~100us)

Location B

0xFFFFBBBB

DMA

GPU Y

Registered on NIC

# Efficient indirect NIC communication: memory & network combined

Network Buffer in **4 slots**

Write

Dynamic Mapping

Relay GPUs

NICs

# Dynamic load balance without contention & interruption



TX Busy
RX Idle

TX Idle
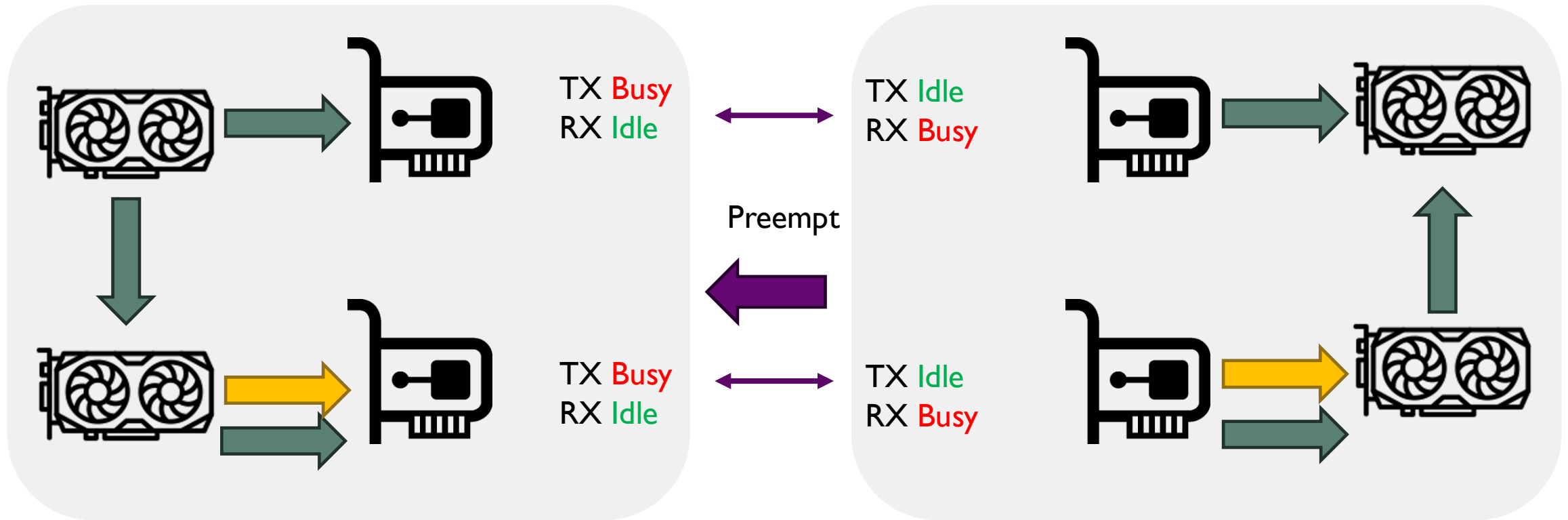RX Busy

Notify
RX Busy

TX Busy
RX Idle

TX Idle
RX Busy

- Record NIC TX/RX workload status within each server
- Exchange TX/RX status, agree on NIC selection

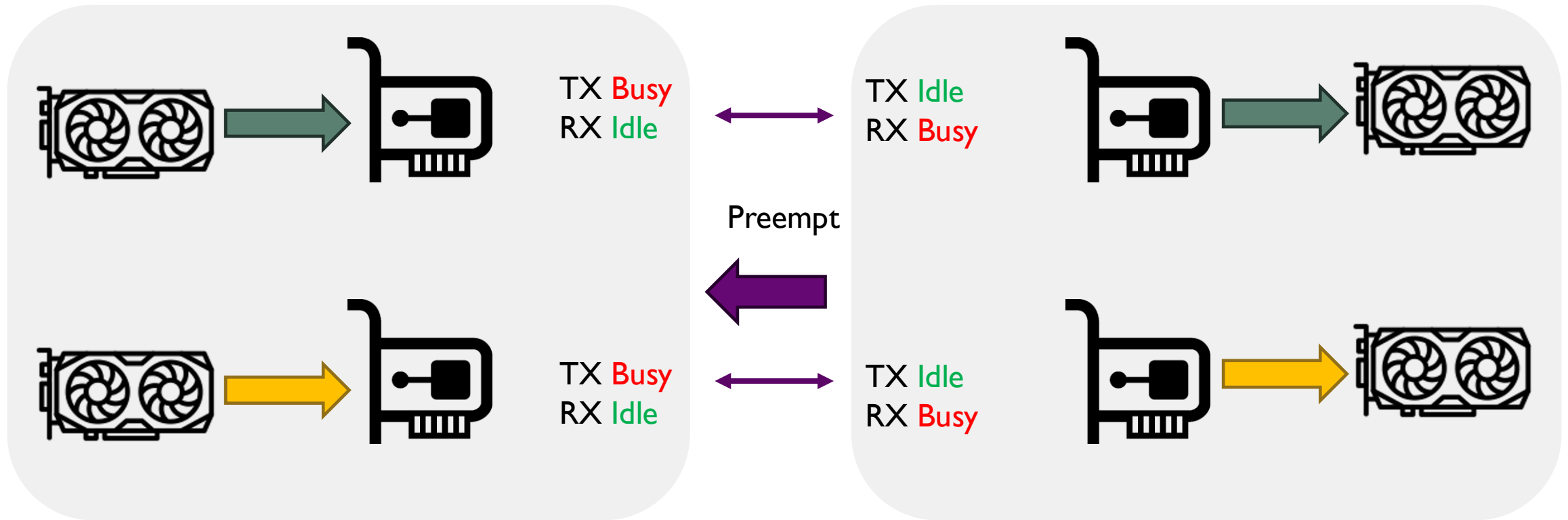14

# Dynamic load balance without contention & interruption

TX **Busy**
RX Idle

Notify
Idleness

TX Idle
RX **Busy**

Dedicated fabrics

Dedicated fabrics

TX Idle
RX Idle

TX Idle
RX Idle

- RX & TX being idle
- Using two NICs through intra-server relaying

15

# Dynamic load balance without contention & interruption



TX **Busy**
RX **Idle**

TX **Idle**
RX **Busy**
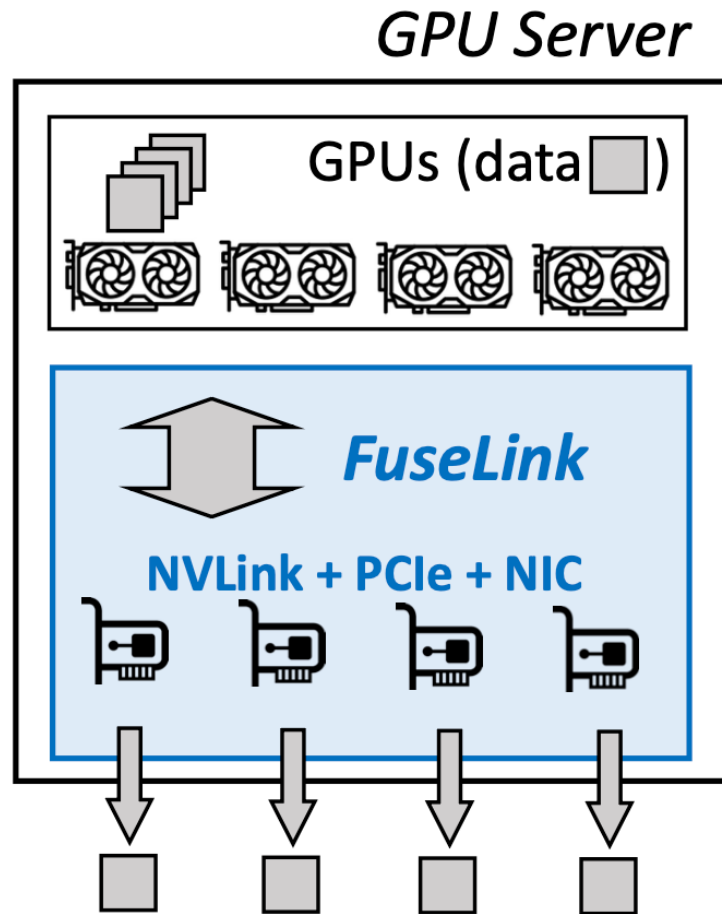
Preempt

TX **Busy**
RX **Idle**

TX **Idle**
RX **Busy**

- NIC contention detected
- Preempt relay traffic and fallback to direct NIC

16

# Dynamic load balance without contention & interruption



TX **Busy**
RX **Idle**

← Preempt →

TX **Idle**
RX **Busy**

Preempt

TX **Busy**
RX **Idle**

TX **Idle**
RX **Busy**

- NIC contention detected
- Preempt relay traffic and fallback to direct NIC
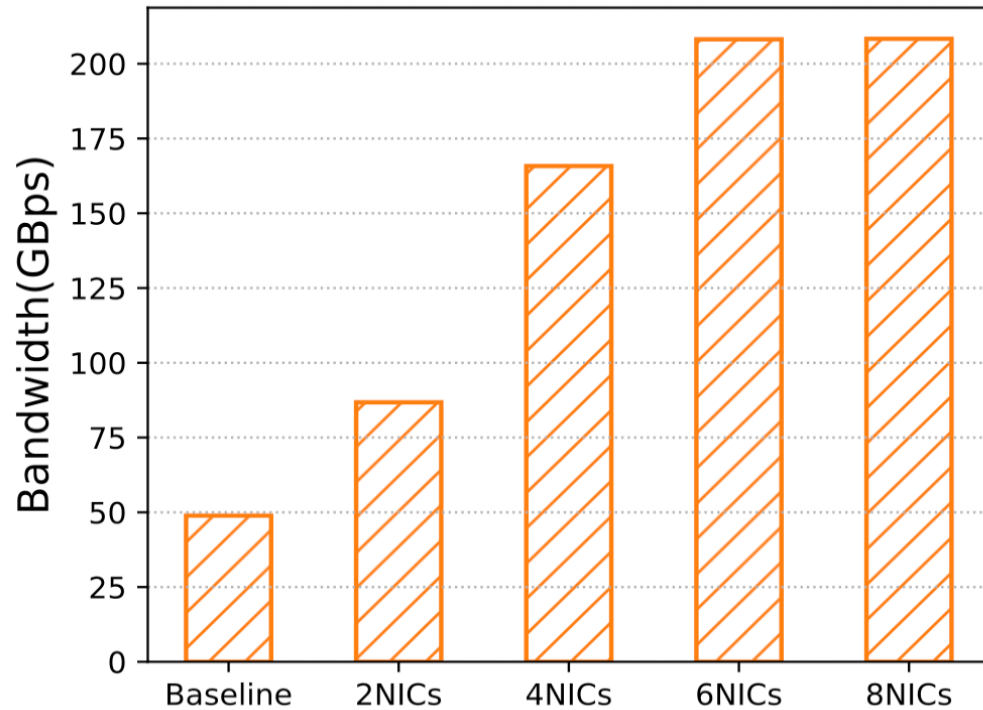
# System Overview



## FuseLink system components

- **NIC idleness detection**: mark NICs as idle when nothing is being sent/received.

- **NVLink + NIC transport**: send traffic to intermediate GPUs and transmit through idle NICs efficiently.

- **Contention elimination**: ensure GPUs fully occupy direct NICs during communication.
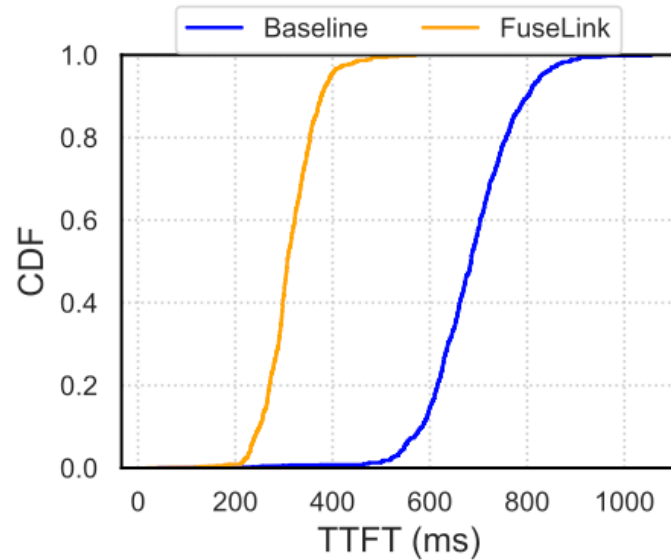
18

# Evaluation: FuseLink bandwidth over NVLink + NIC

Each server: eight GPU, eight-lane NVLink (~160 GBps) and 400 Gbps (50 GBps) NICs.
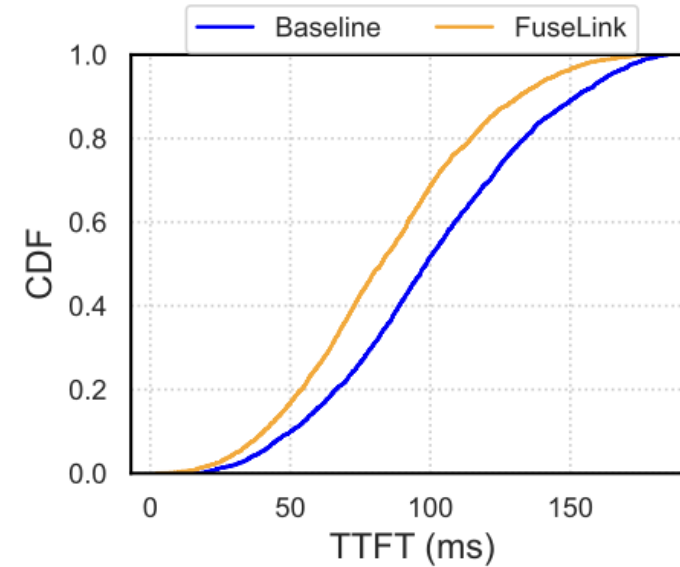


FuseLink bandwidth for two GPUs using different number of NICs

Bandwidth reaches limit when both direct NIC and NVLink are fully utilized.

19

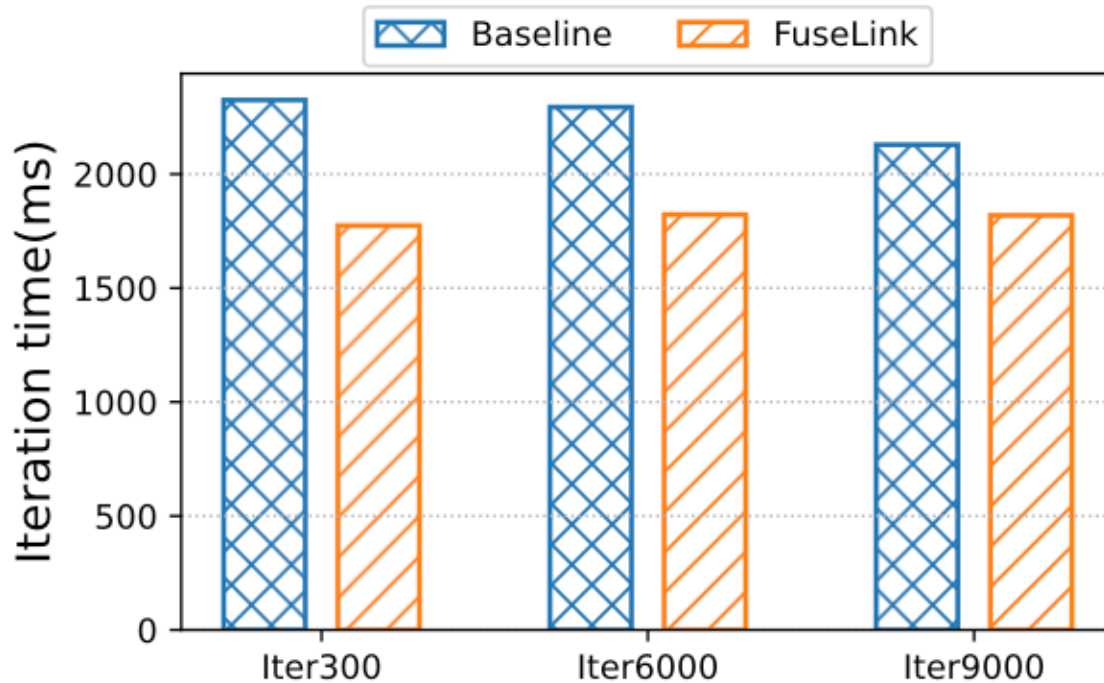# Evaluation: LLM serving of independent instances
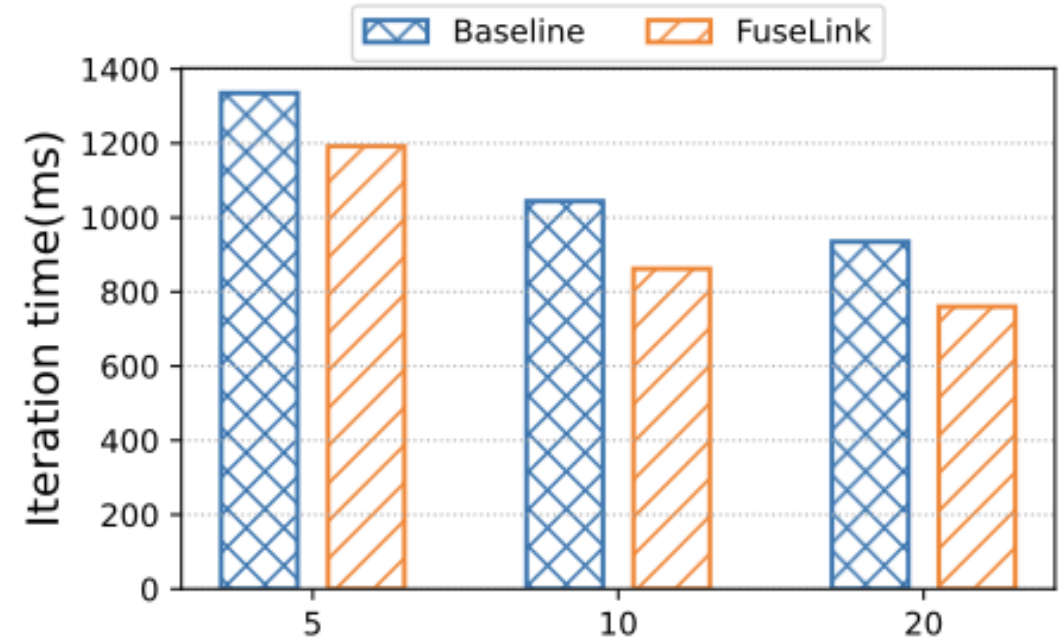


Eight serving instances



Two serving instances

Improvement: accelerated data transfer & reduced waiting time

# Evaluation: expert-parallelism and imbalanced embedding transmission



Each server has two expert shards
Accelerating imbalanced all-to-all in
Mixtral 8x22B expert-parallel training

Accelerating imbalanced embedding
transmission when training DLRM

# Limitations and future works

- Appliable to other GPUs?
  Yes, only need P2P memory access & virtual memory mapping feature.


- Fine-grained load balancing?
  Per-chunk load balancing to per-packet load balancing.

# Thank You!

Code: https://github.com/axio-project/FuseLink
Contact: zrenak@cse.ust.hk