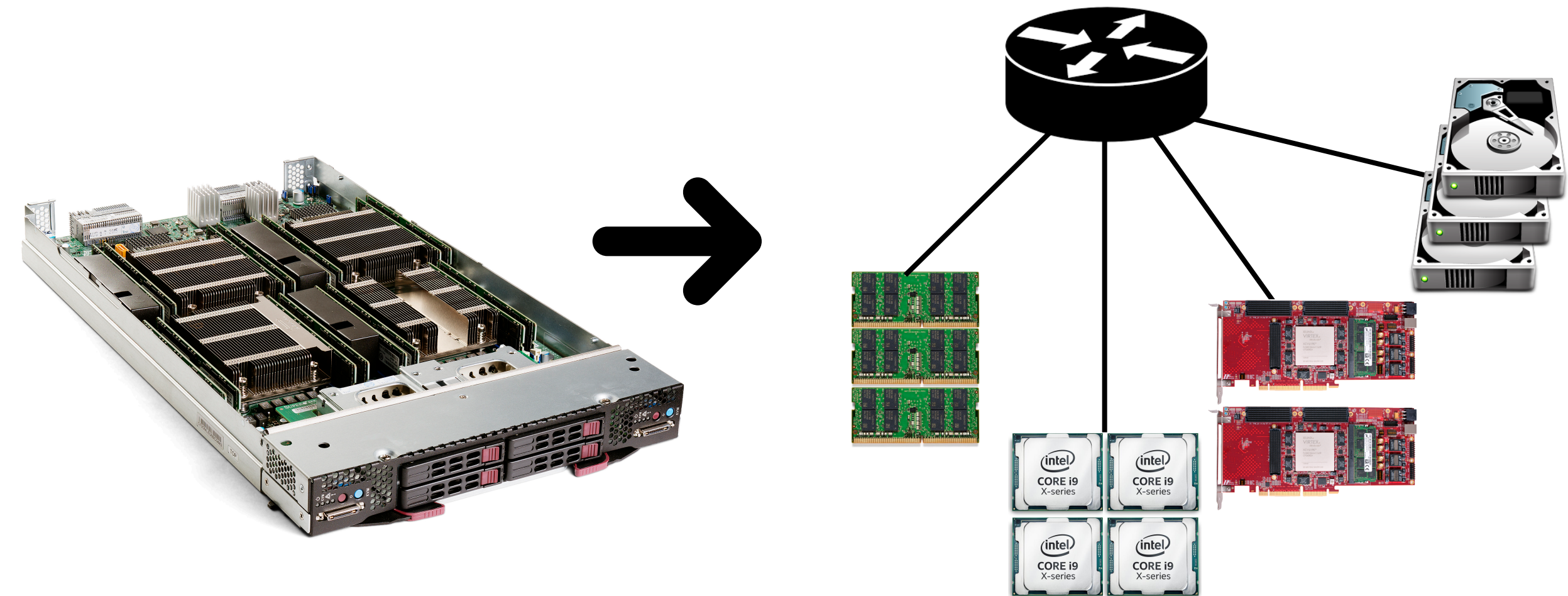# EDM: An Ultra-Low Latency Ethernet Fabric for Memory Disaggregation

Weigao Su, Vishal Shrivastav
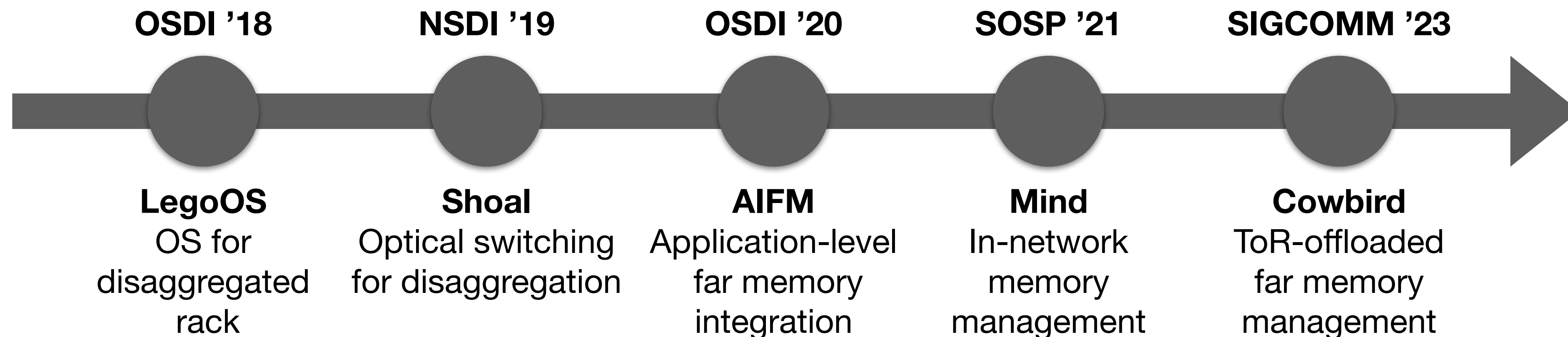
**PURDUE UNIVERSITY**

# Why memory disaggregation?

- The need for memory is surging
- Constraints of individual servers
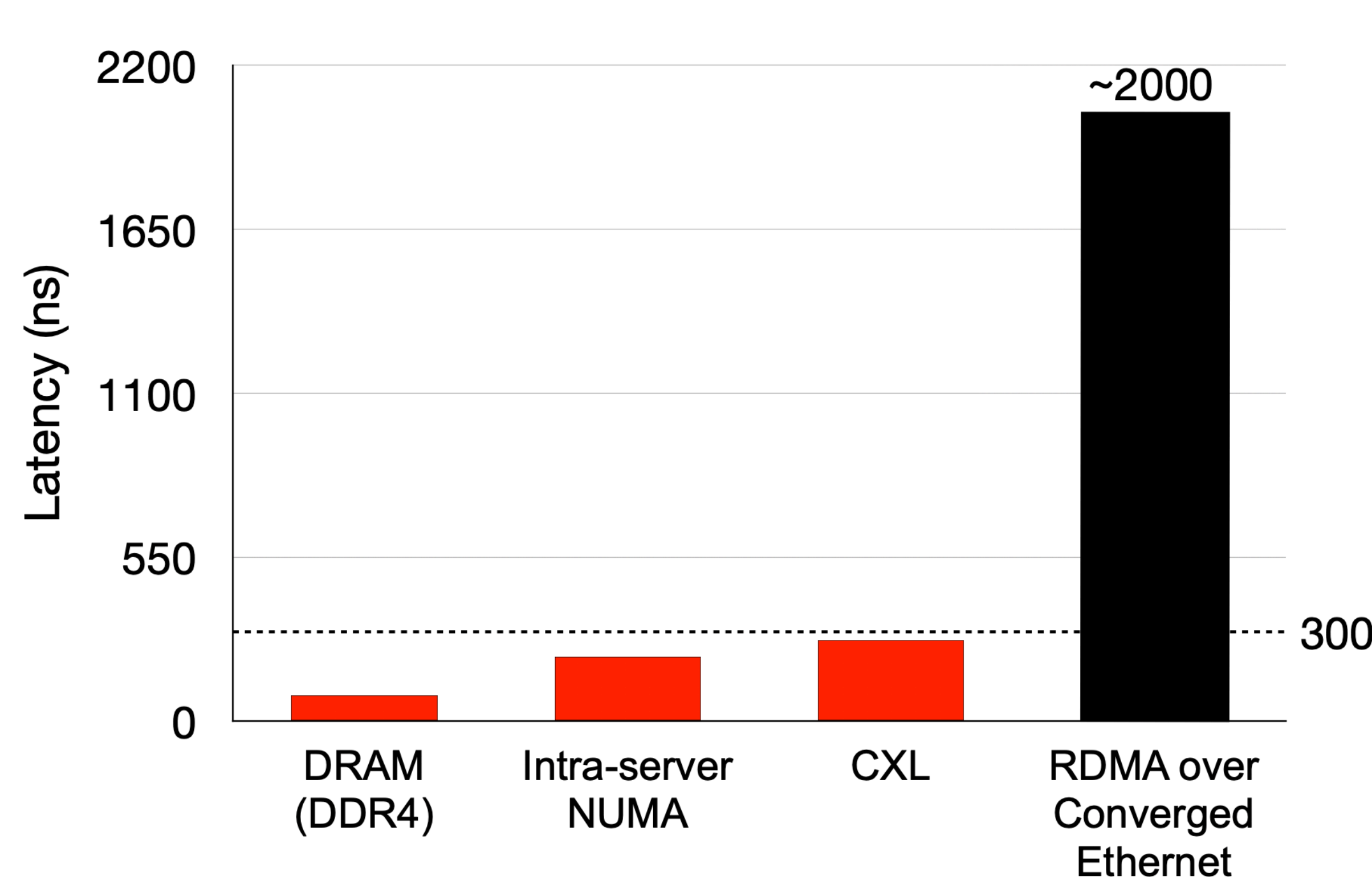- Fine-grained pooling, elastic scaling

# Why is Ethernet promising?

- Dominant datacenter network fabric
- High bandwidth (Terabit Ethernet link)
- Low management cost, distance scaling…

| **OSDI '18** | **NSDI '19** | **OSDI '20** | **SOSP '21** | **SIGCOMM '23** |
|:---:|:---:|:---:|:---:|:---:|
| **LegoOS** | **Shoal** | **AIFM** | **Mind** | **Cowbird** |
| OS for disaggregated rack | Optical switching for disaggregation | Application-level far memory integration | In-network memory management | ToR-offloaded far memory management |

2

# However, the latency in Ethernet is prohibitive, prompting proposals of _separate_ fabric to carry memory traffic

Custom processor interconnect, PCIe, Infiniband, etc.
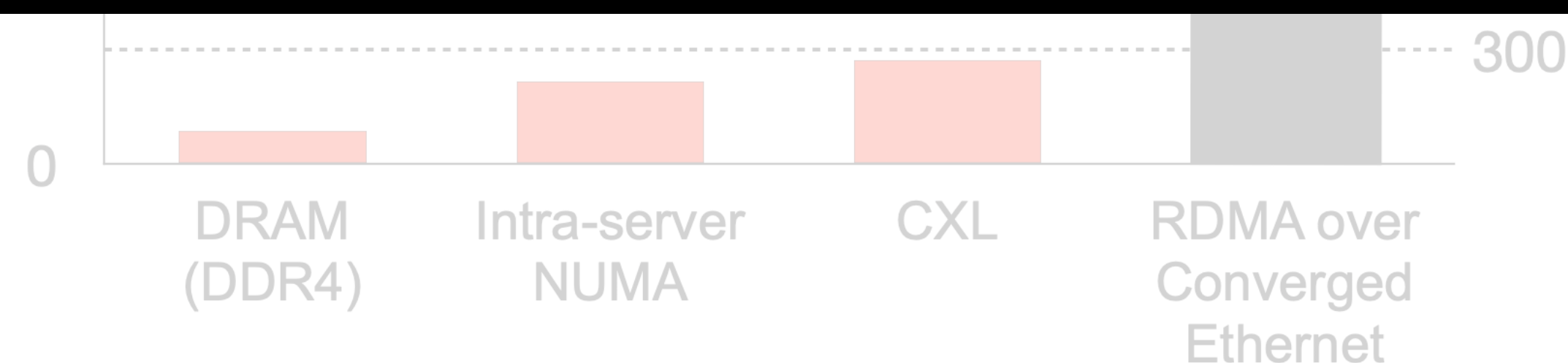


But, separate fabrics for different traffic makes the network **costly** and **harder to manage**

**However, the latency in Ethernet is prohibitive, prompting proposals of _separate_ fabric to carry memory traffic**

Custom processor interconnect, PCIe, Infiniband, etc.

A **low latency** Ethernet fabric would allow us to have a **single unified** network fabric to carry all kinds of traffic (memory, storage, IP, …)

_… easier to manage, lower cost, statistical bandwidth multiplexing_

300

0

DRAM (DDR4)   Intra-server NUMA   CXL   RDMA over Converged Ethernet

The monolithic server model where a server is the unit of deployment, operation, and failure is meeting its limits in the face of several recent hardware and application trends. To improve resource utilization, elasticity, het-

often limits the speed of new hardware adoption.
We believe that datacenters should break mono-lithic servers and organize hardware devices like CPU, DRAM, and disks as independent, failure-isolated

Ricardo Bianchini
Microsoft Azure
USA

But, separate fabrics for different traffic makes the network **costly** and **harder to manage**

# Research goal

Achieving near **intra-server memory access latency** over rack-scale **Ethernet**

(while maintaining high bandwidth utilization)

# Memory Disaggregation over Ethernet

**Compute Node**

Application

**Available Solutions for NIC interconnect latency**

Faster link (e.g., UPI); integrate processor with network controller, etc.

**Memory Node**

Memory Controller

Transport protocol (e.g., TCP/IP or RDMA)

Create/read packet

Reliability

Congestion Control

**Switch**

Layer 2 forwarding

Parser

DST Port

Match-Action

Queue

L2 PKT

L2 PKT

ETH MAC

ETH PHY

Transport protocol (e.g., TCP/IP or RDMA)

Create/read packet

Reliability

Congestion Control

ETH MAC

ETH PHY

ETH MAC

ETH PHY

Ethernet link

# Memory Disaggregation over Ethernet

**Compute Node**

Application

**Available Solutions for NIC interconnect latency**

Faster link (e.g., UPI); integrate processor with network controller, etc.

**Memory Node**

Memory Controller

Transport protocol (e.g., TCP/IP or RDMA)

Create/read packet

Reliability

Congestion Control

**Switch**

Layer 2 forwarding

Parser

DST Port

Match-Action

Queue

L2 PKT

L2 PKT

ETH MAC

ETH MAC

ETH PHY

ETH MAC

ETH PHY

Transport protocol (e.g., TCP/IP or RDMA)

Create/read packet

Reliability

Congestion Control

ETH MAC

ETH PHY
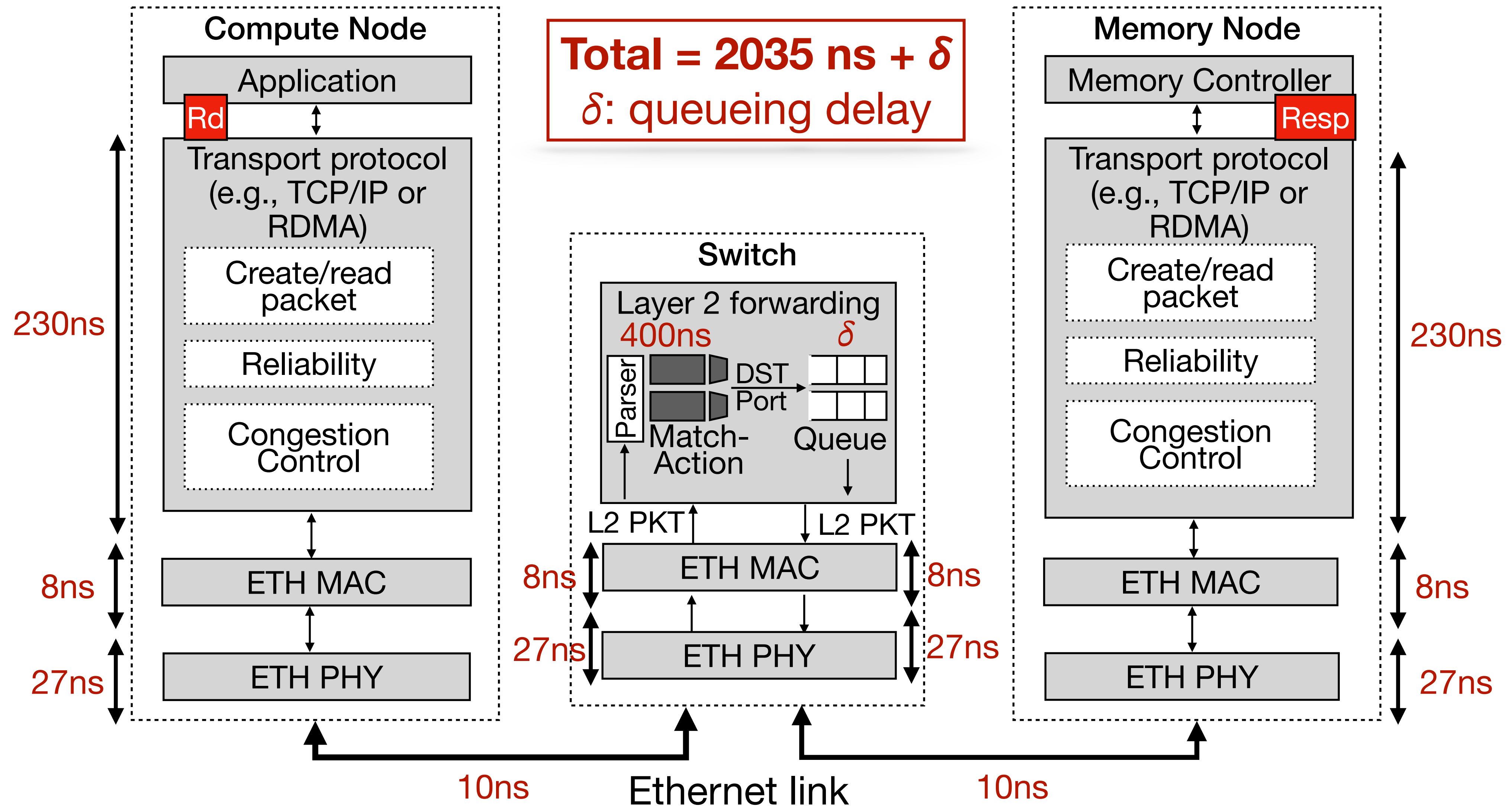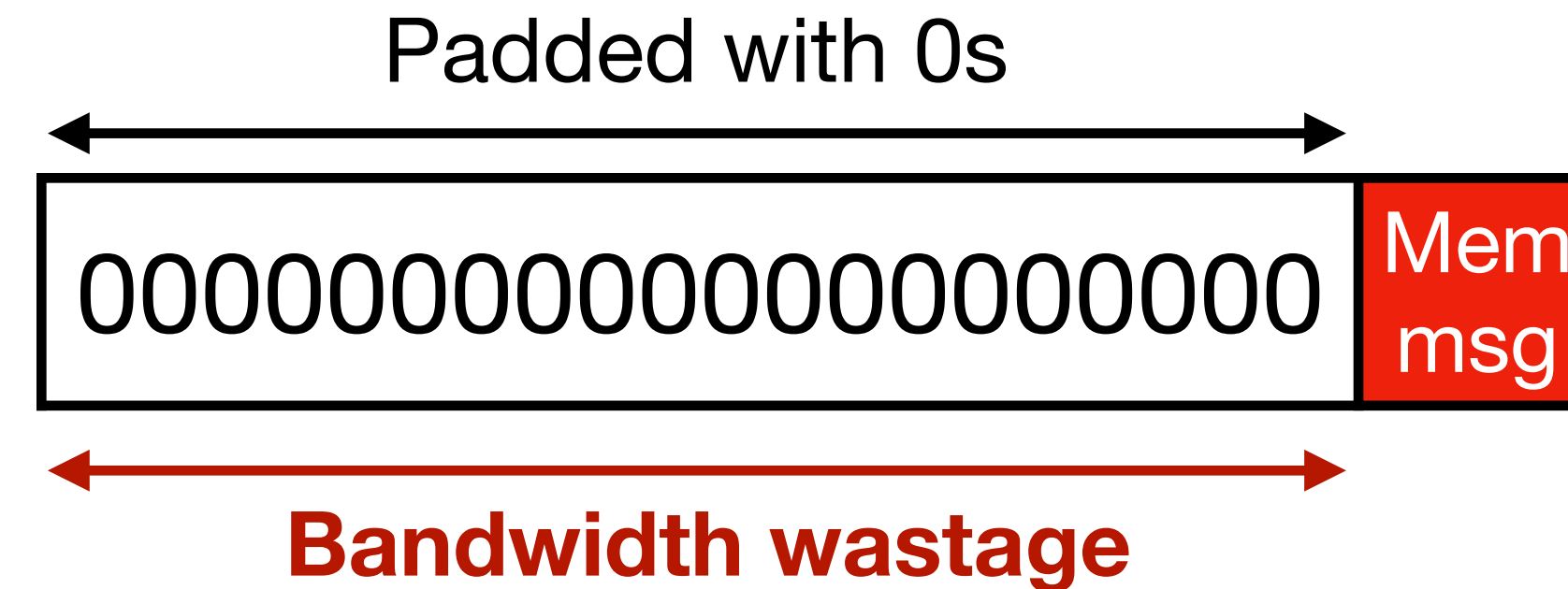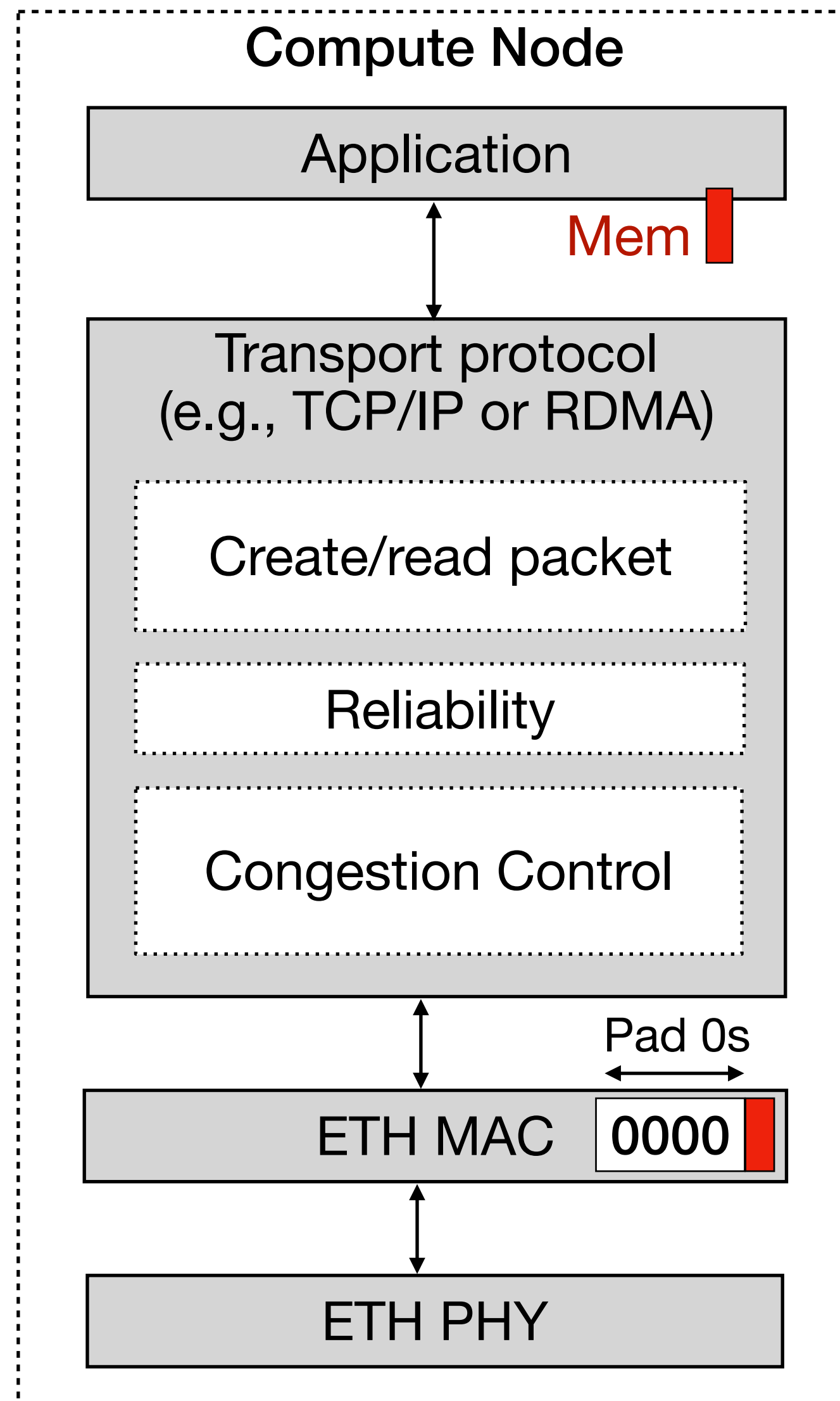
Ethernet link

7

# Latency Overheads of Existing Memory Disaggregation over Ethernet

An example of remote read request over RDMA

# Other Overheads

## Compute Node

**Application**

Mem ▮

**Transport protocol
(e.g., TCP/IP or RDMA)**

Create/read packet

Reliability

Congestion Control

Pad 0s

ETH MAC  | 0000 ▮

ETH PHY

---

1. **Ethernet MAC enforces minimum 64B frame**
… but memory messages can be much smaller
(e.g., read requests are typically 8-16B)

Padded with 0s

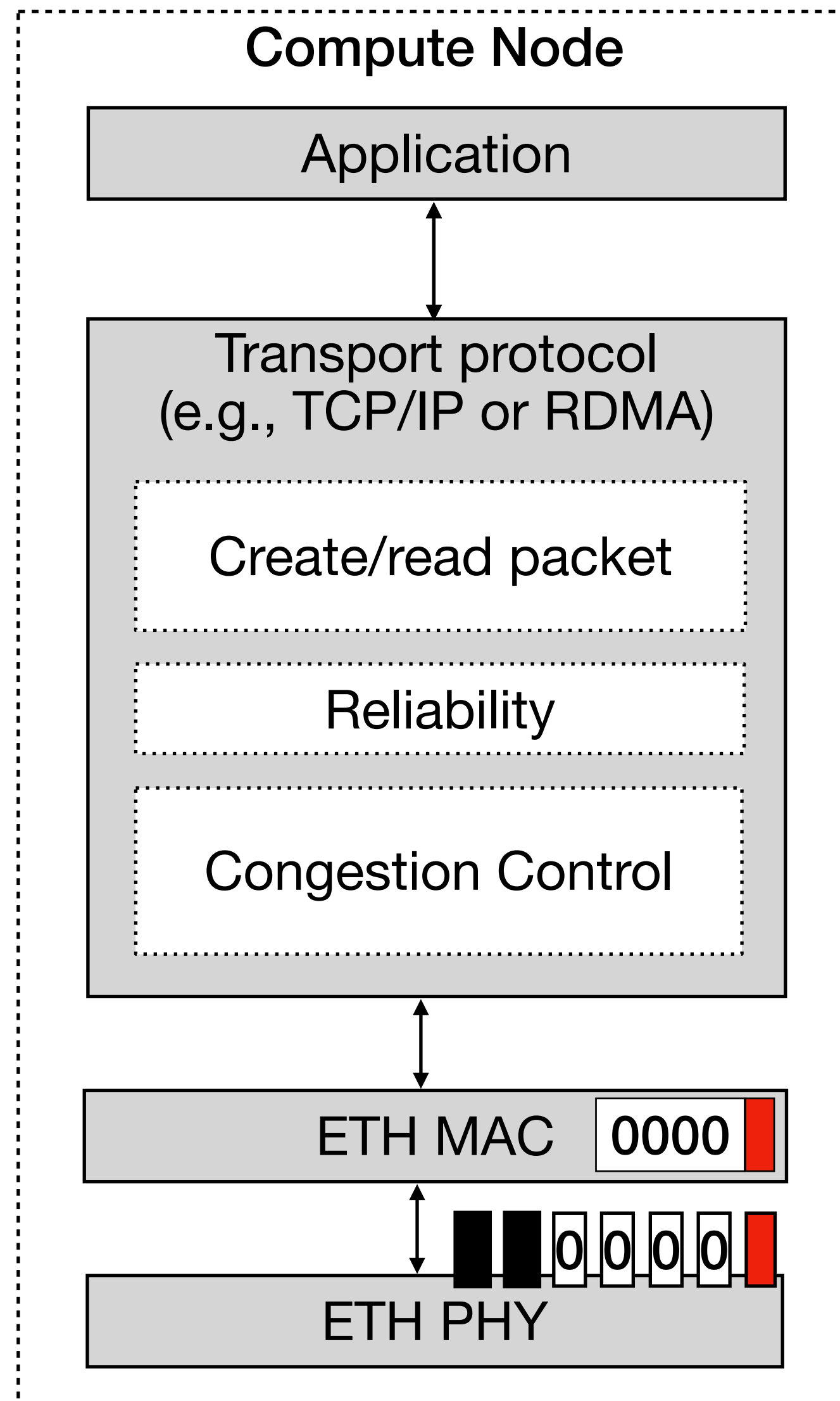00000000000000000000 | Mem msg

**Bandwidth wastage**

# Other Overheads



1. **Ethernet MAC enforces minimum 64B frame**
   … but memory messages can be much smaller (e.g., read requests are typically 8-16B)

2. <span style="color:red">**Ethernet MAC enforces minimum of 12 bytes Inter-frame gap (IFG)**</span>
   … high overhead for small memory messages

# Other Overheads

## Compute Node

Application 1 | Application 2

**Mem** | **IP**

Transport protocol
(e.g., TCP/IP or RDMA)

Create/read packet
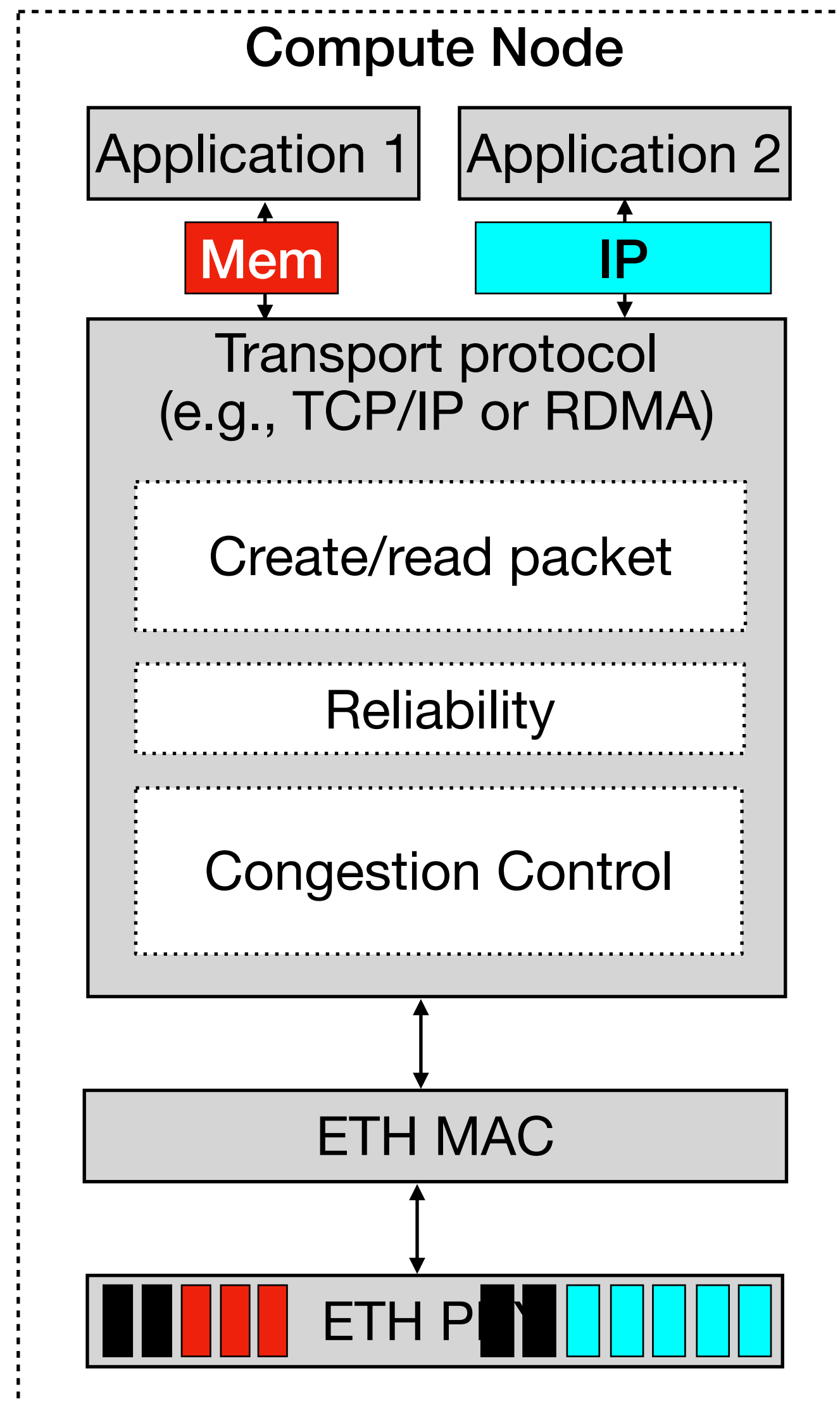
Reliability

Congestion Control

ETH MAC

ETH PHY

1. **Ethernet MAC enforces minimum 64B frame**
   … but memory messages can be much smaller (e.g., read requests are typically 8-16B)

2. **Ethernet MAC enforces minimum of 12 bytes Inter-frame gap (IFG)**
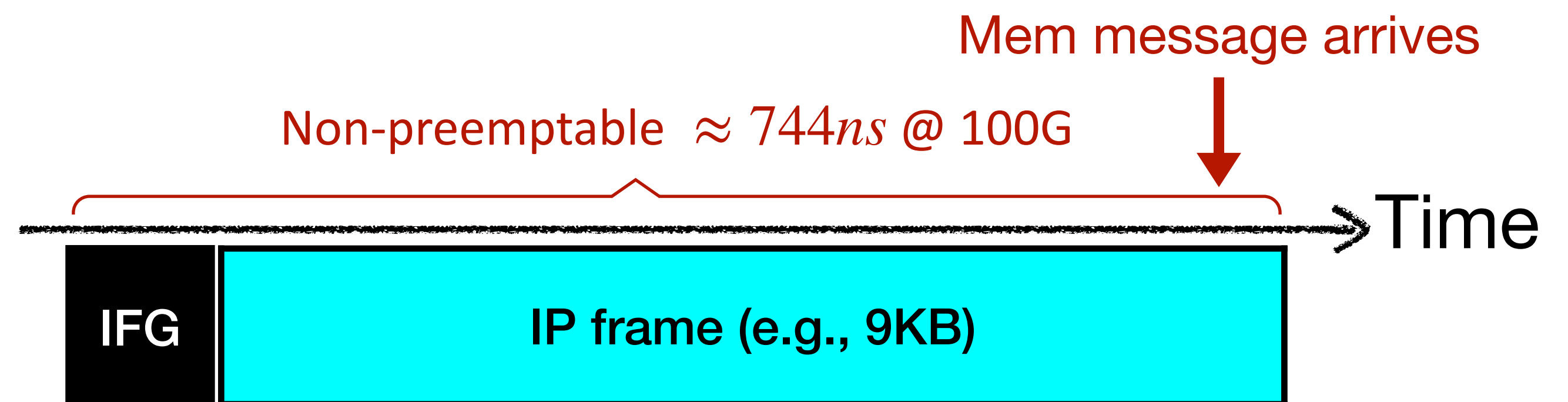   … high overhead for small memory messages

3. **Ethernet MAC does not allow intra-frame preemption**
   … a large non-memory frame may block the transmission of a small memory message

Mem message arrives

Non-preemptable $\approx 744ns$ @ 100G

Time

IFG | IP frame (e.g., 9KB)

## Other Overheads



**Compute Node**

Application 1    Application 2

Transport protocol
(e.g., TCP/IP or RDMA)

Create/read packet

Reliability

Congestion Control

ETH MAC

ETH PHY

1. **Ethernet MAC enforces minimum 64B frame**
   … but memory messages can be much smaller
   (e.g., read requests are typically 8-16B)

2. **Ethernet MAC enforces minimum of 12 bytes Inter-frame gap (IFG)**
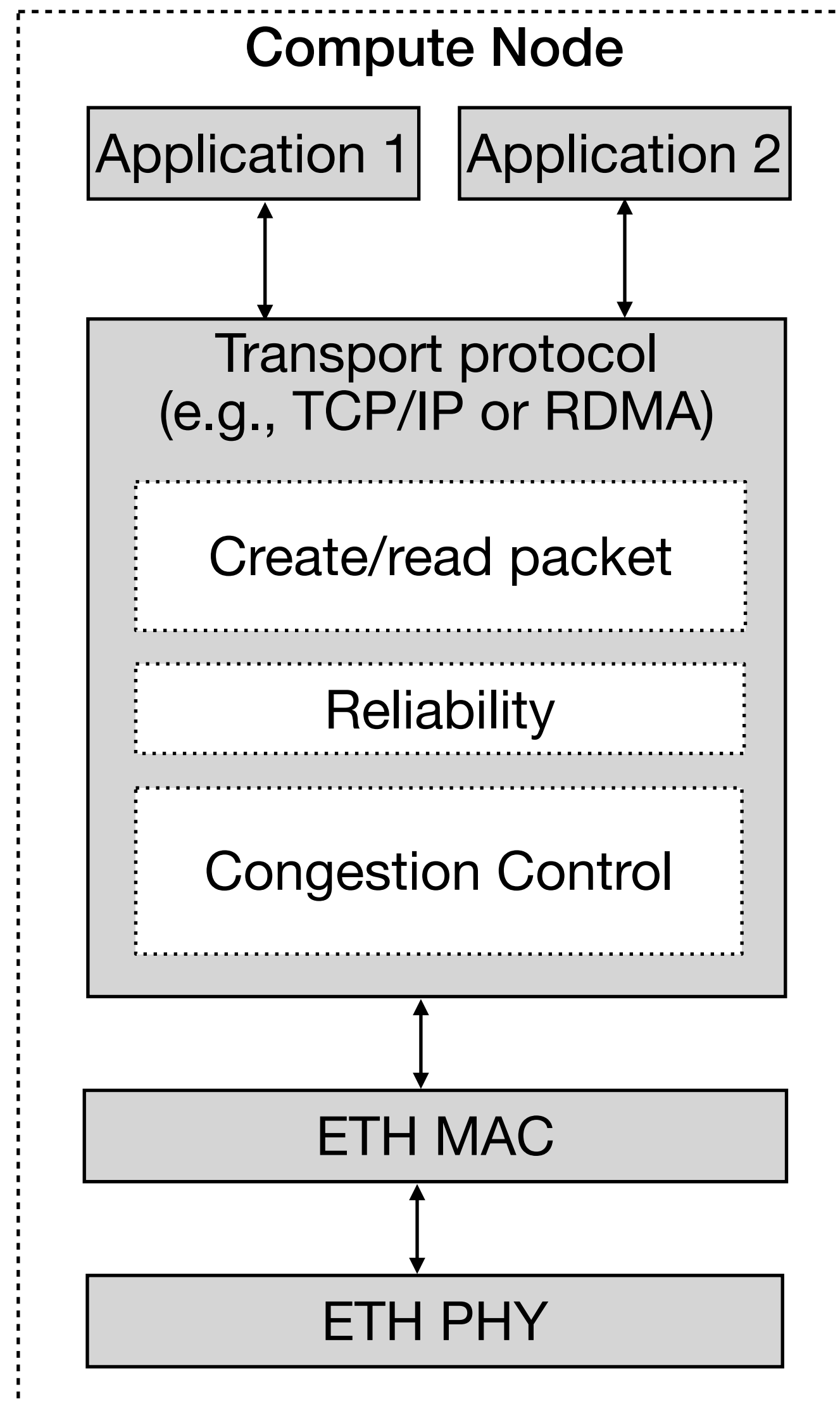   … high overhead for small memory messages

3. **Ethernet MAC does not allow intra-frame preemption**
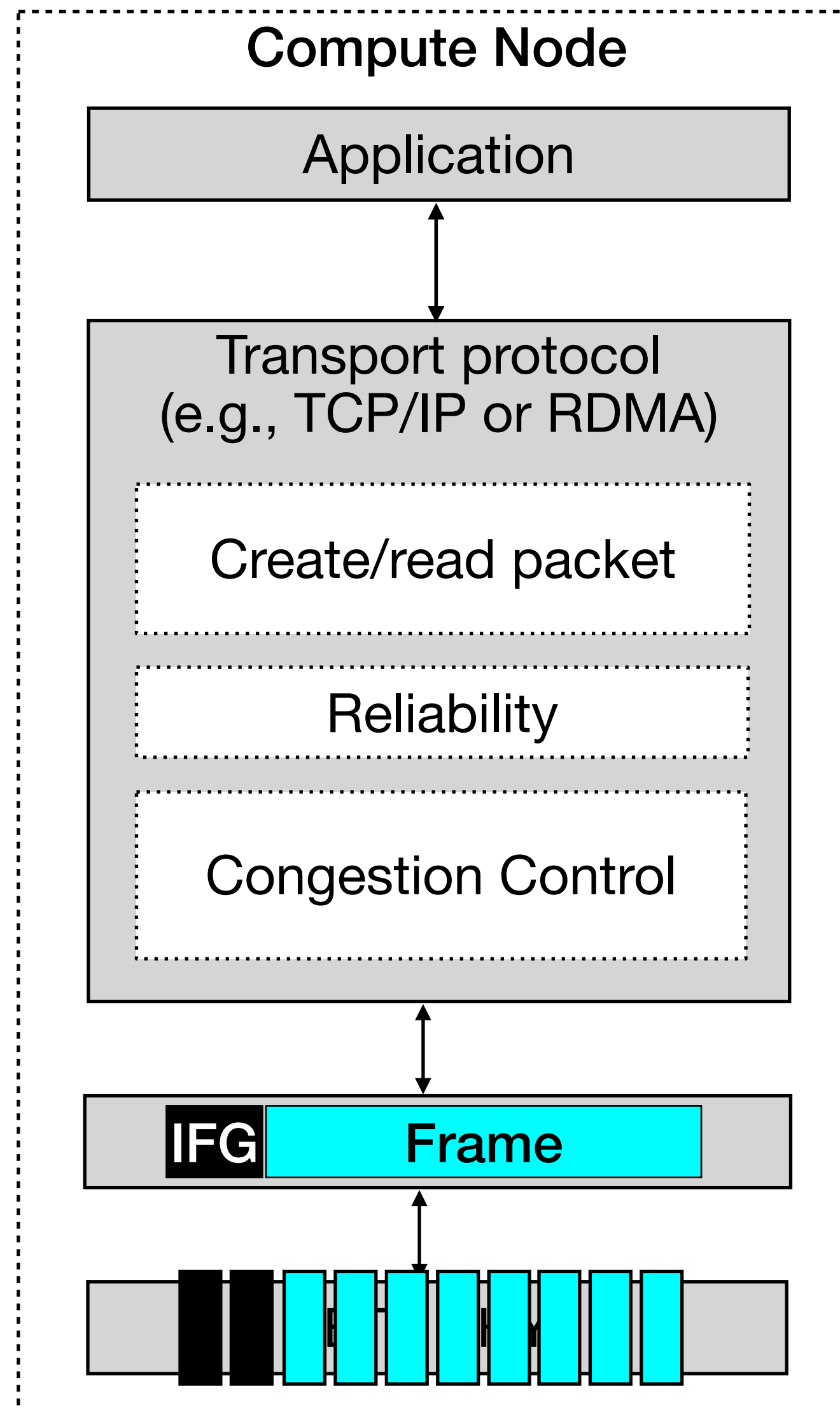   … a large non-memory frame may block the transmission of a small memory message

## Root cause: MAC layer processing

# Design Choice # 1:

Implement the entire
**protocol for remote memory access**
within Ethernet's **Physical layer**
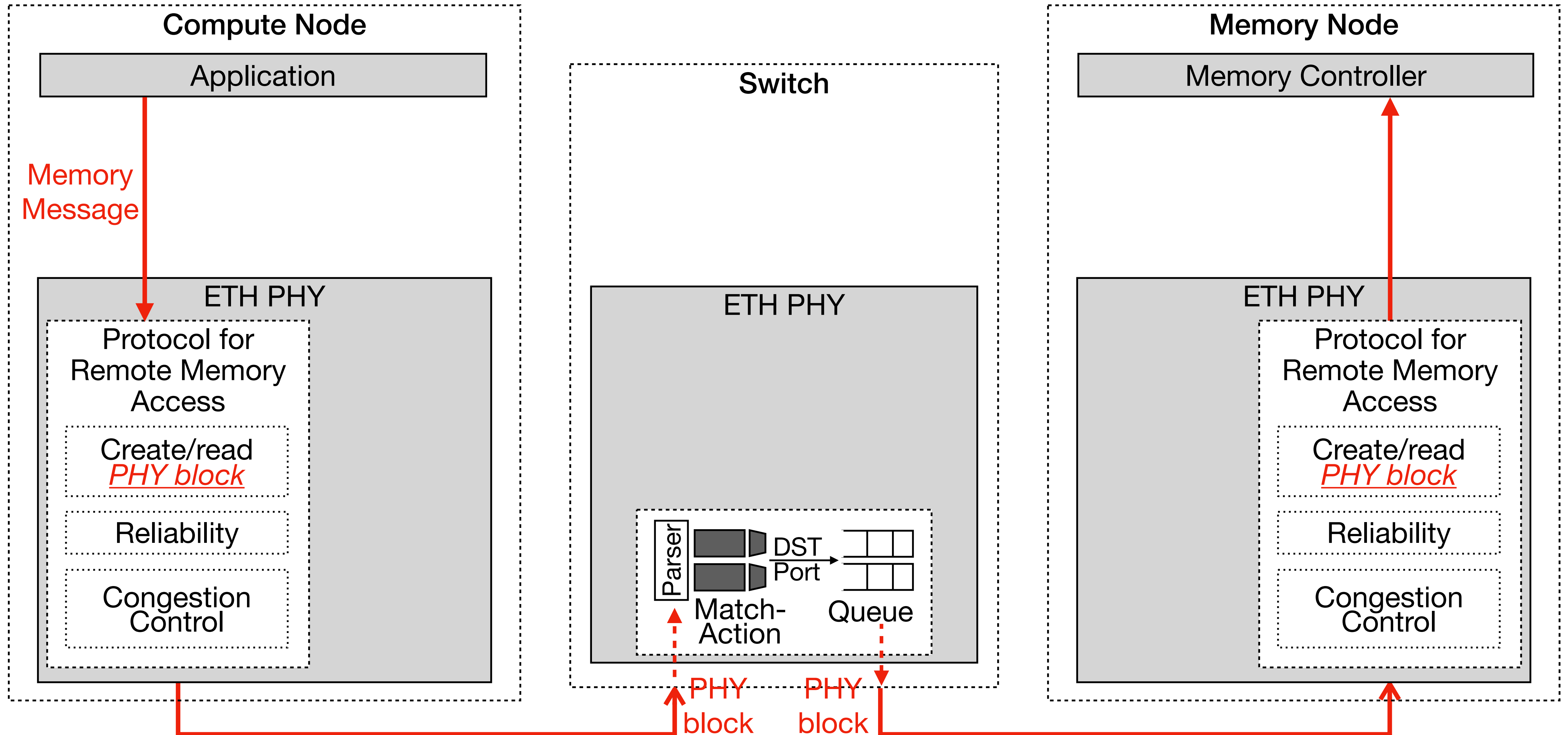
# Rationale for Remote Memory Protocol in PHY



Compute Node

Application

Transport protocol
(e.g., TCP/IP or RDMA)

Create/read packet

Reliability

Congestion Control

IFG Frame

**Ethernet PHY already reformats a MAC layer frame into a series of 66-bit PHY blocks**
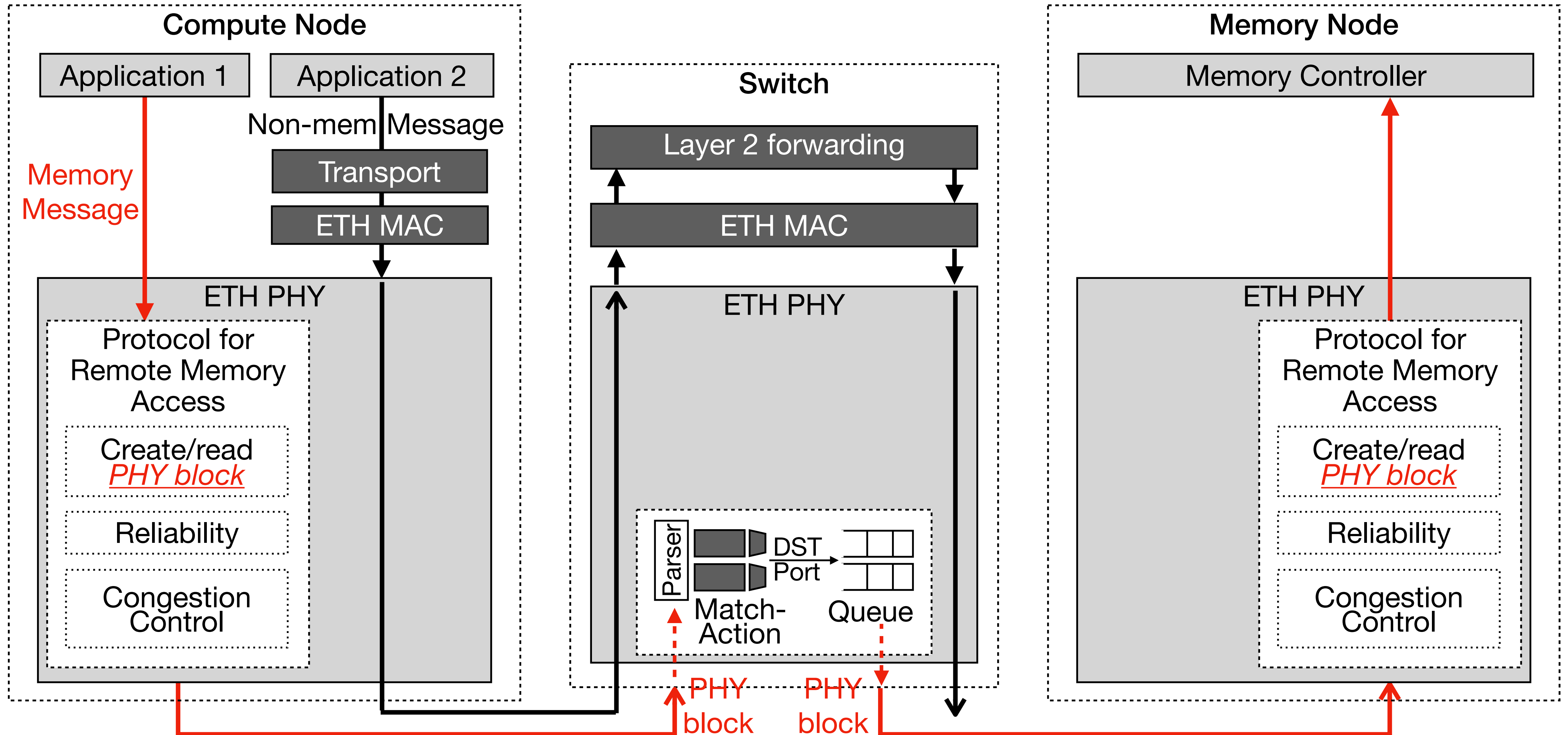
… thus, unlike the MAC layer that works at a **frame** granularity, PHY works at fine-grained **block** granularity

- 66 bit PHY block vs. 64 byte minimum MAC frame size
- PHY also has access to IFG blocks
- Message interleaving can be done at block granularity in PHY rather than at frame granularity in MAC

# Architecture of Remote Memory Protocol in the PHY



**Compute Node**

Application

Memory Message

ETH PHY

Protocol for Remote Memory Access

Create/read *PHY block*

Reliability

Congestion Control

**Switch**

ETH PHY

Parser | Match-Action | DST Port | Queue

PHY block | PHY block

**Memory Node**

Memory Controller

ETH PHY

Protocol for Remote Memory Access

Create/read *PHY block*

Reliability

Congestion Control

# Architecture of Remote Memory Protocol in the PHY

**Compute Node**

| Application 1 | Application 2 |
|---|---|

Non-mem Message

Transport

ETH MAC

**Memory Message**

ETH PHY

Protocol for Remote Memory Access
- Create/read *PHY block*
- Reliability
- Congestion Control

**Switch**

Layer 2 forwarding

ETH MAC

ETH PHY

Parser

DST Port

Match-Action

Queue

PHY block

PHY block

**Memory Node**

Memory Controller

ETH PHY

Protocol for Remote Memory Access
- Create/read *PHY block*
- Reliability
- Congestion Control

# Benefits of Remote Memory Protocol in the PHY



1. **Ethernet PHY operates at a fine data granularity of 66-bit PHY blocks.**
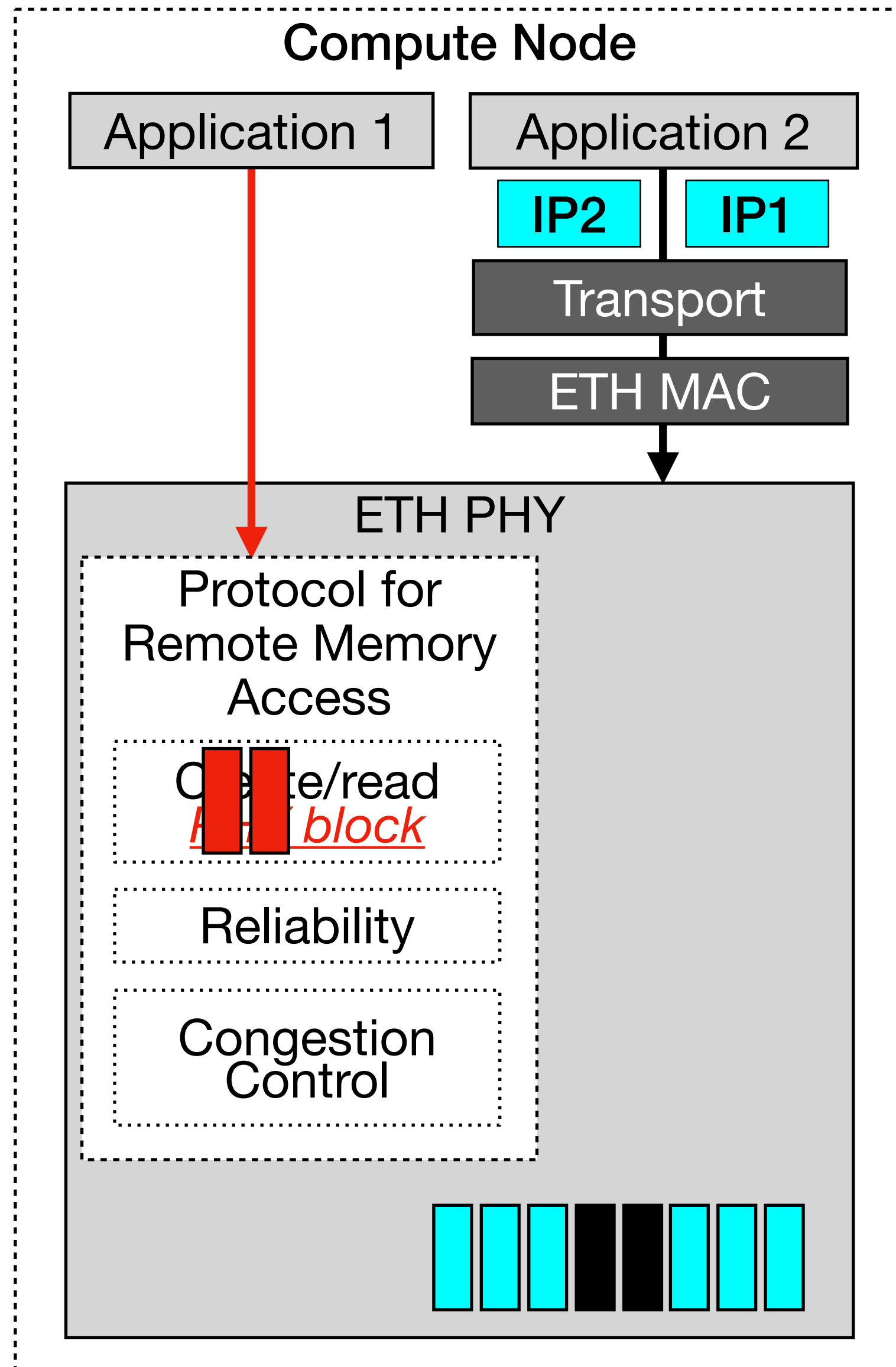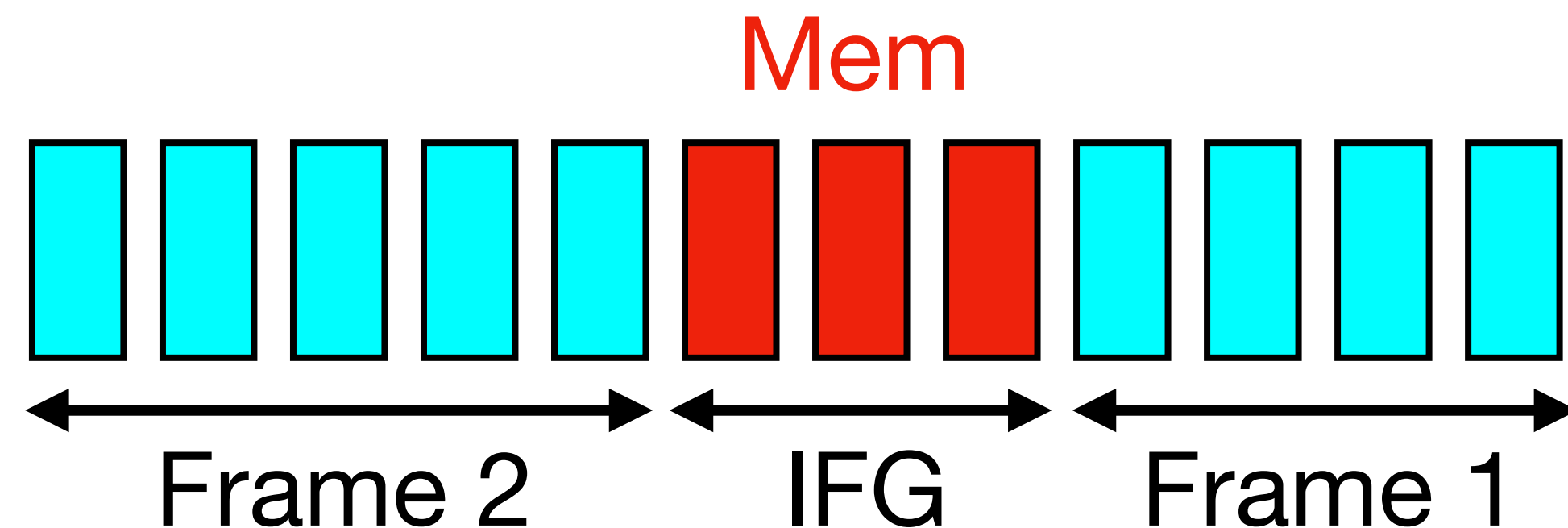   Avoids bandwidth wastage for small memory messages.
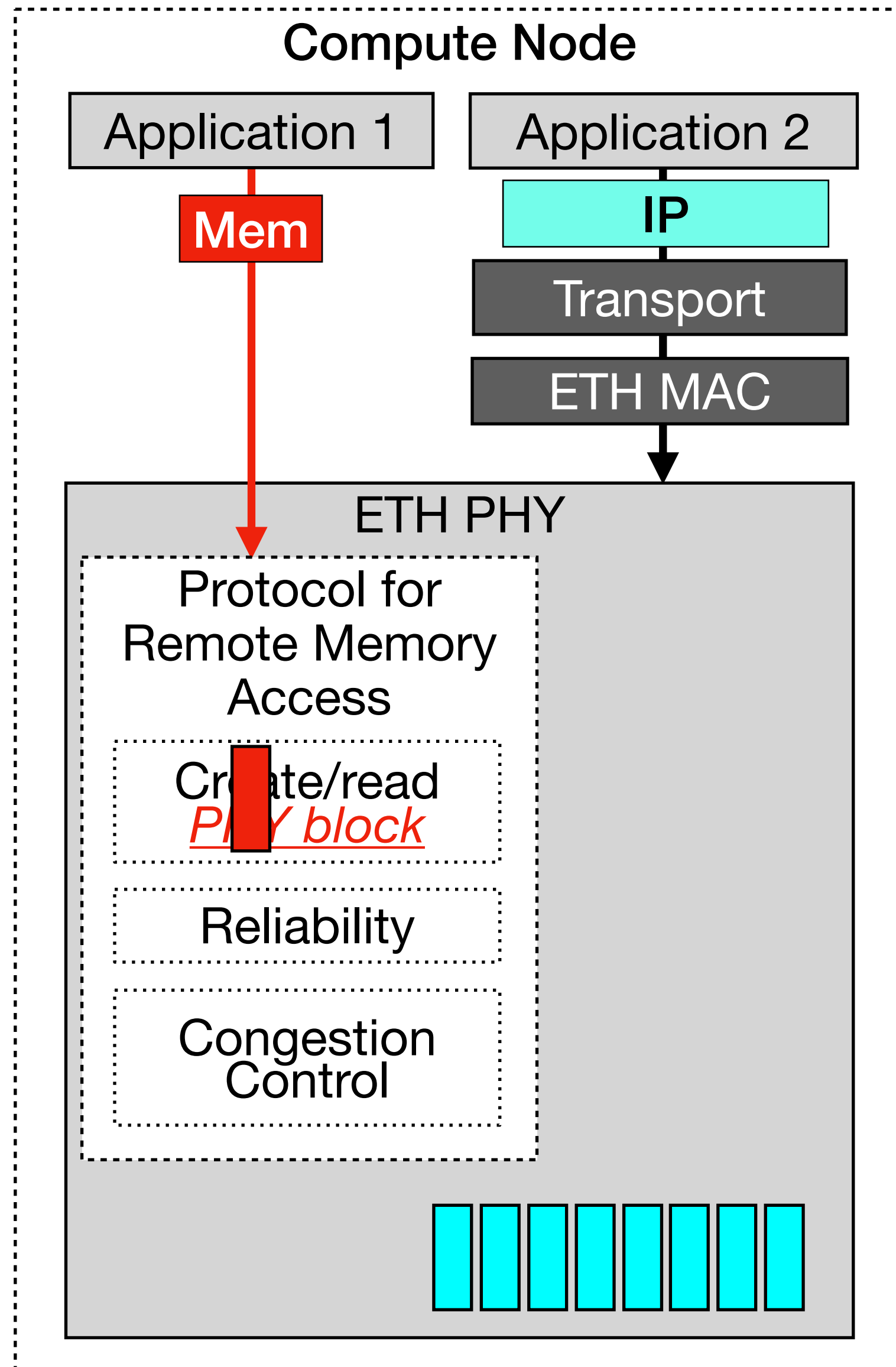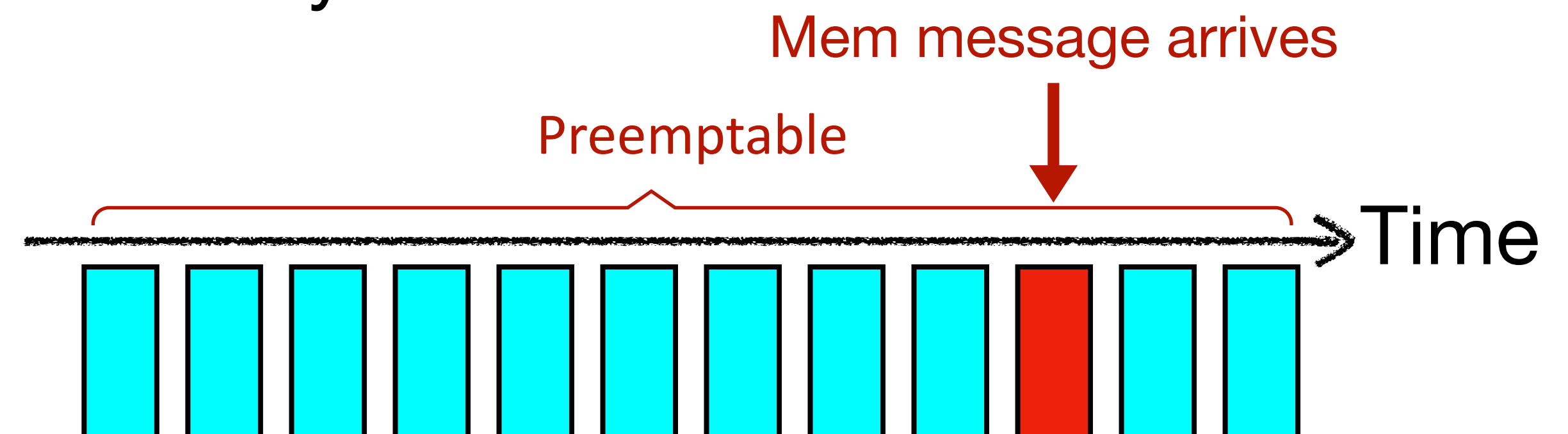
# Benefits of Remote Memory Protocol in the PHY



1. **Ethernet PHY operates at a fine data granularity of 66-bit PHY blocks.**
   Avoids bandwidth wastage for small memory messages.

2. **Ethernet PHY has access to IFG bits.**
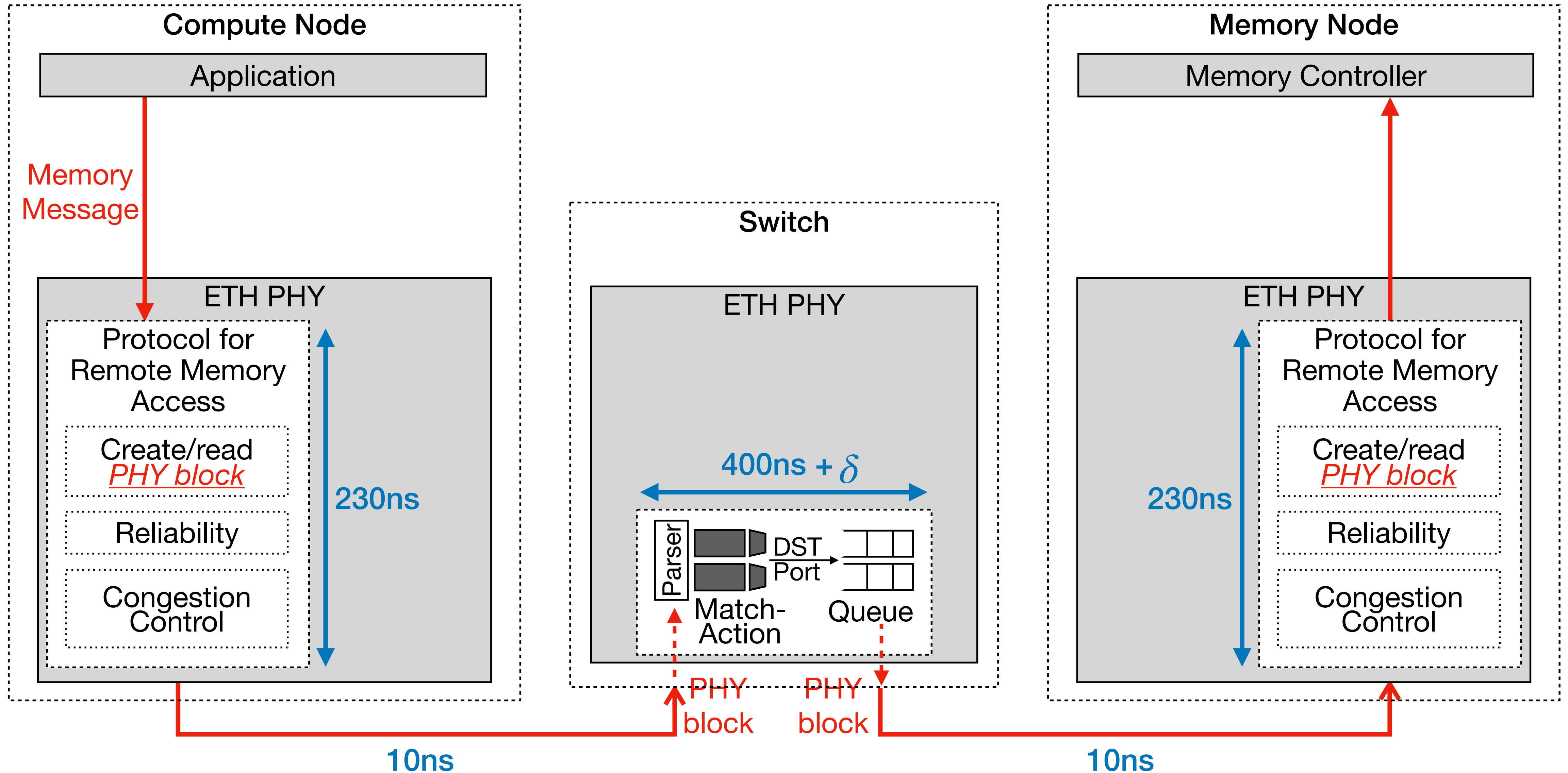   Can repurpose IFG to carry memory messages.

# Benefits of Remote Memory Protocol in the PHY



1. **Ethernet PHY operates at a fine data granularity of 66-bit PHY blocks.**
   Avoids bandwidth wastage for small memory messages.

2. **Ethernet PHY has access to IFG bits.**
   Can repurpose IFG to carry memory messages.

3. **Ethernet PHY enables intra-frame preemption**
   Avoids blocking of small memory message by a large non-memory frame.

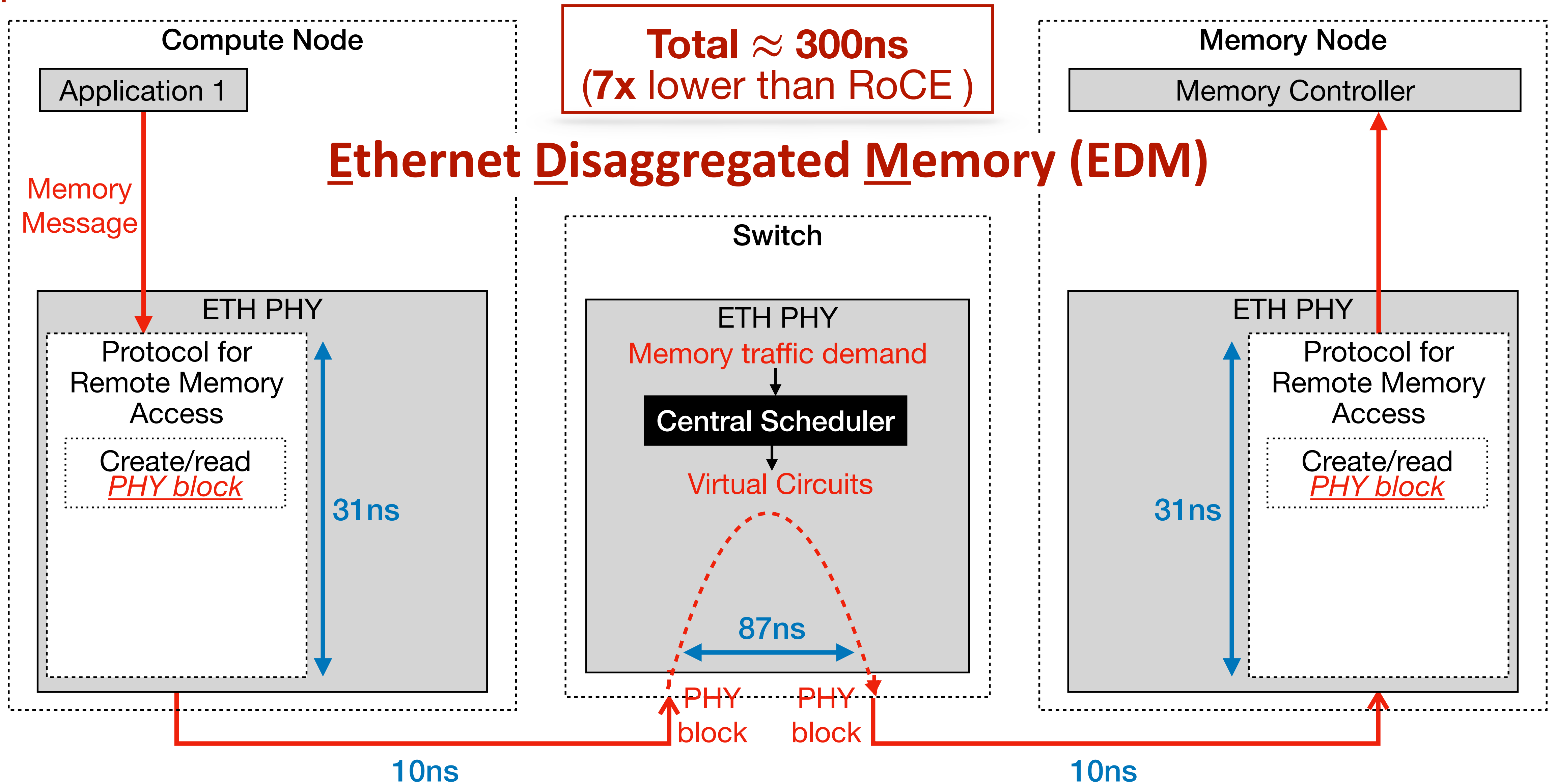# Remote Memory Protocol in the PHY : What about latency?

**Compute Node**

| Application |

Memory
Message

**ETH PHY**

Protocol for
Remote Memory
Access

Create/read
*PHY block*

Reliability

Congestion
Control

230ns

**Switch**

**ETH PHY**

$400ns + \delta$

Parser | DST Port

Match-Action

Queue

PHY
block

PHY
block

**Memory Node**

| Memory Controller |

**ETH PHY**

Protocol for
Remote Memory
Access

230ns

Create/read
*PHY block*

Reliability

Congestion
Control

10ns

10ns

# Design Choice # 2:

# Implement a
# <span style="color:darkred">centralized memory traffic scheduler</span>
# in the PHY of the switch

# Central Scheduler in the Switch PHY



**Compute Node**

Application 1

Memory Message

**Total ≈ 300ns**
(**7x** lower than RoCE )

## Ethernet Disaggregated Memory (EDM)

**Memory Node**

Memory Controller

ETH PHY

Protocol for Remote Memory Access

Create/read *PHY block*

31ns

**Switch**

ETH PHY

Memory traffic demand

**Central Scheduler**

Virtual Circuits

87ns

PHY block    PHY block

ETH PHY

Protocol for Remote Memory Access

Create/read *PHY block*

31ns

10ns

10ns

22

# Overview of Central Scheduler

- **Step 1:** Nodes send their memory message demands {src->dst} to the switch scheduler
- **Step 2:** Switch scheduler creates virtual circuits by forming a Matching based on demand
  - Naive maximal matching ~O(N); EDM uses **Parallel Iterative Matching (PIM)**[1] ~O(log(N))
- **Step 3:** Nodes exchange memory messages over established circuits



[1] Anderson et al. "High Speed switch scheduling for Local Area Networks". TOCS 1993.
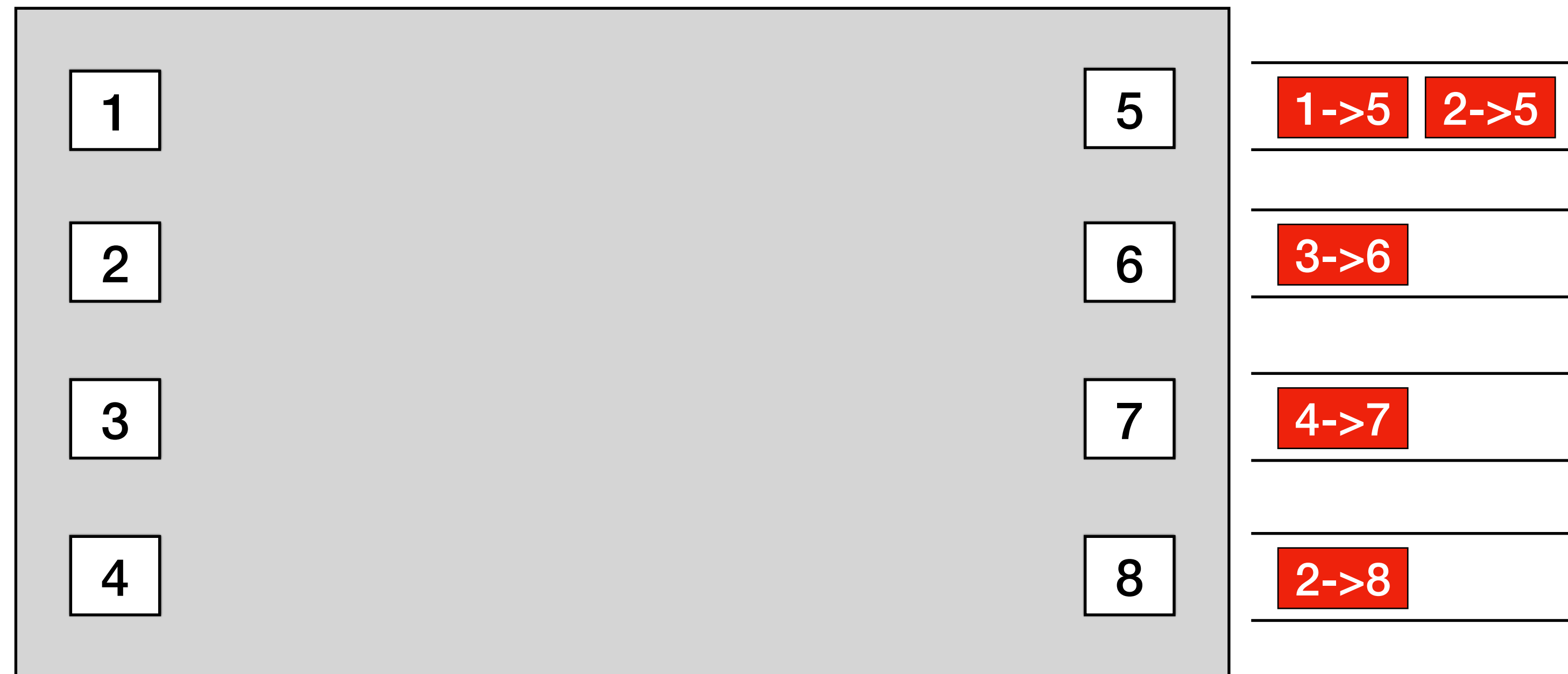
# Practical Central Scheduler

- **Challenge 1:** Low latency for memory messages under <span style="color:darkred">bandwidth contention</span>

- **Challenge 2:** Accurate, low overhead memory <span style="color:darkred">traffic demand estimation</span>

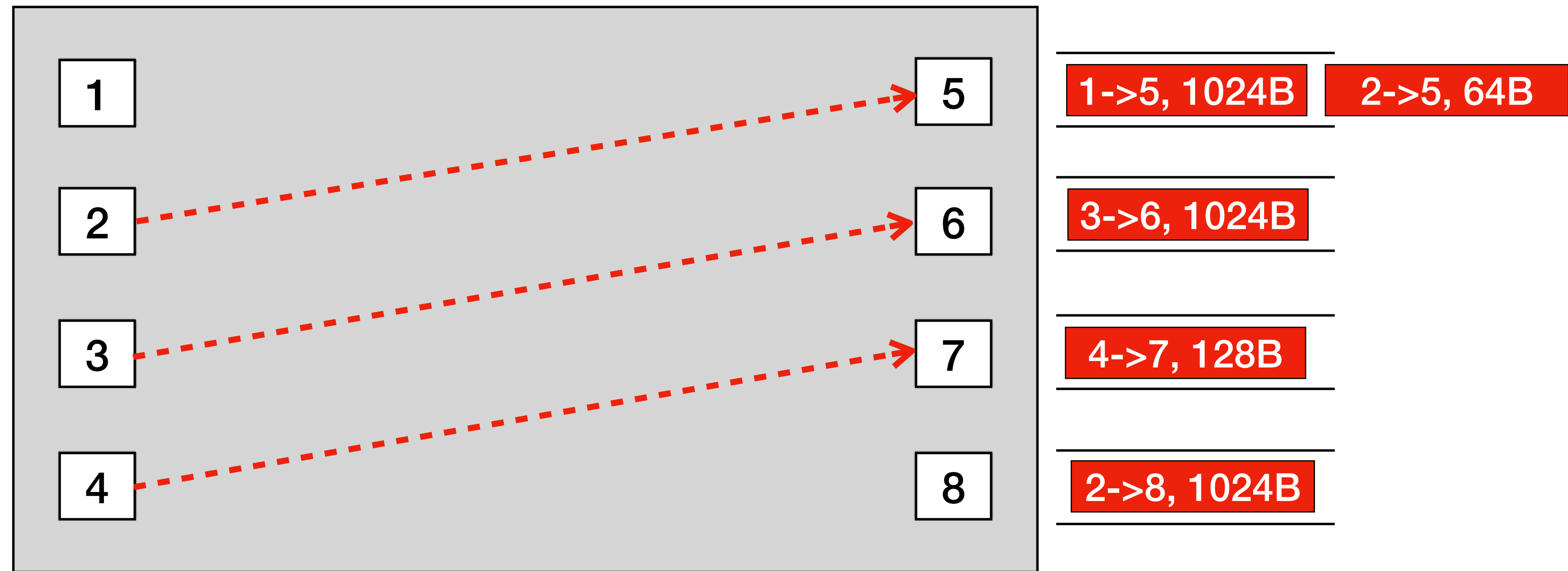- **Challenge 3:** <span style="color:darkred">Line rate,</span> low latency scheduling hardware pipeline

# Challenge # 1: Achieve low latency under bandwidth contention

**Solution:** Augment PIM with priority scheduling

- First Come First Serve (FCFS) for *light-tailed* traffic distribution
- Shortest Remaining Processing First (SRPT) for *heavy-tailed* traffic distribution

# Challenge # 1: Achieve low latency under bandwidth contention

**Solution:** Augment PIM with priority scheduling
- First Come First Serve (FCFS) for *light-tailed* traffic distribution
- Shortest Remaining Processing First (SRPT) for *heavy-tailed* traffic distribution



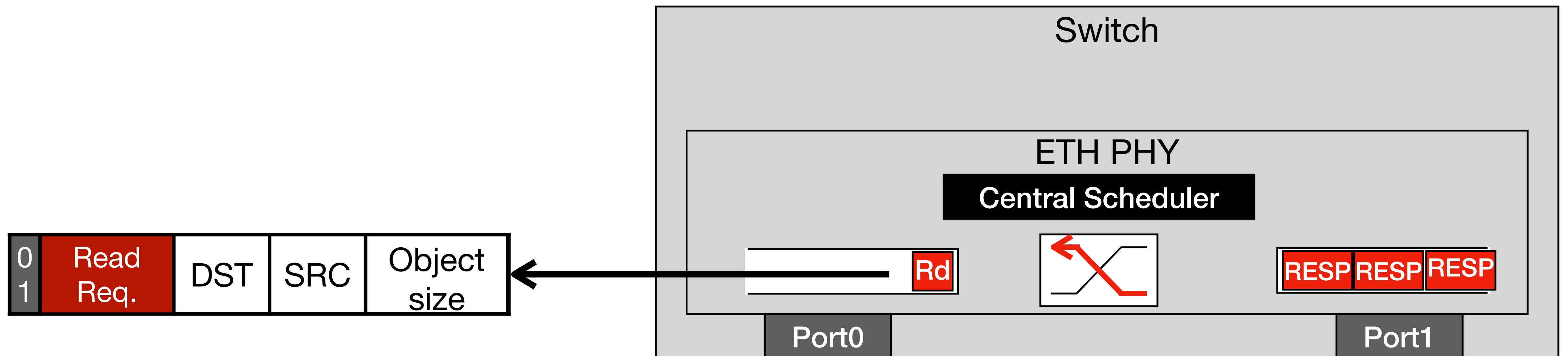| 1 | → | 5 | | 1->5, 1024B | 2->5, 64B |
| 2 | → | 6 | | 3->6, 1024B | |
| 3 | → | 7 | | 4->7, 128B | |
| 4 | → | 8 | | 2->8, 1024B | |

- For SRPT, each demand message from the nodes also contains the **size of message**
- Demand messages per node are processed in the **increasing order of remaining bytes**
- Matching contention -> **prioritize** demand messages with **smaller remaining bytes**

# Challenge # 2: Acquire accurate, low overhead memory traffic demand matrix

**Solution:** Leverage the nature of memory access interface that specifies amount of data to be read or written

- For **reads**, read request implicity contains demand for read reply
  - Zero bandwidth and latency overhead
- For **writes**, send an explicit demand message to switch
  - Small bandwidth overhead (notifications are small)
  - Latency (~RTT/2) is small within a rack

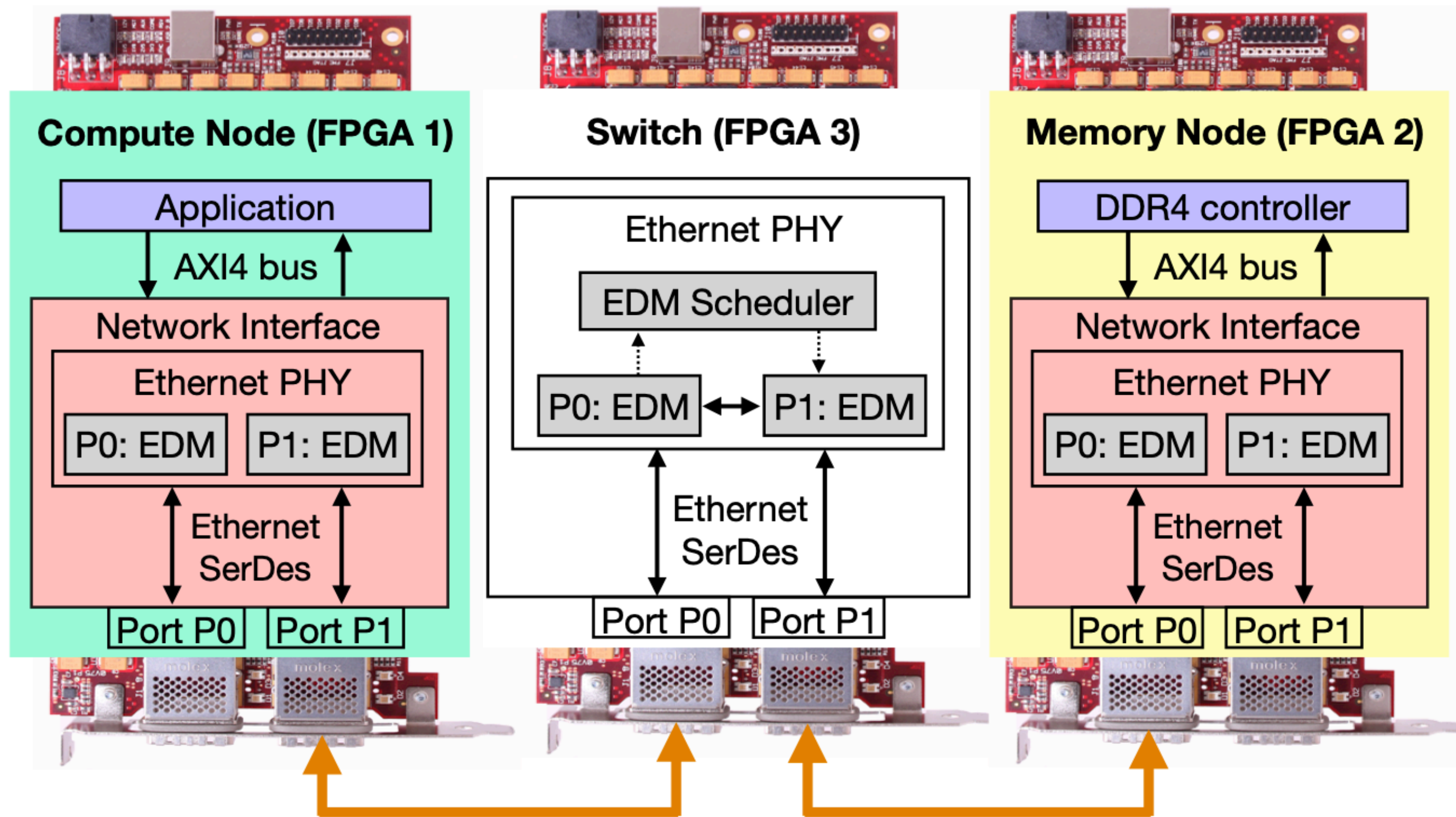**Challenge # 3: Design line-rate, low latency scheduling hardware pipeline**

Naive implementation of **priority-based PIM** would take **O(log(N)) cycles** per PIM iteration (N: number of demand messages)

**Solution:** Leverage hardware parallelism to intelligently trade-off hardware resources for time

- Use combination of constant-time ordered list data structure with a fast priority encoder to implement priority-based PIM

**EDM can implement each iteration of PIM in exactly 3 clock cycles**

# Implementation



**Hardware Testbed**
- Three Xilinx Alveo U200 FPGAs
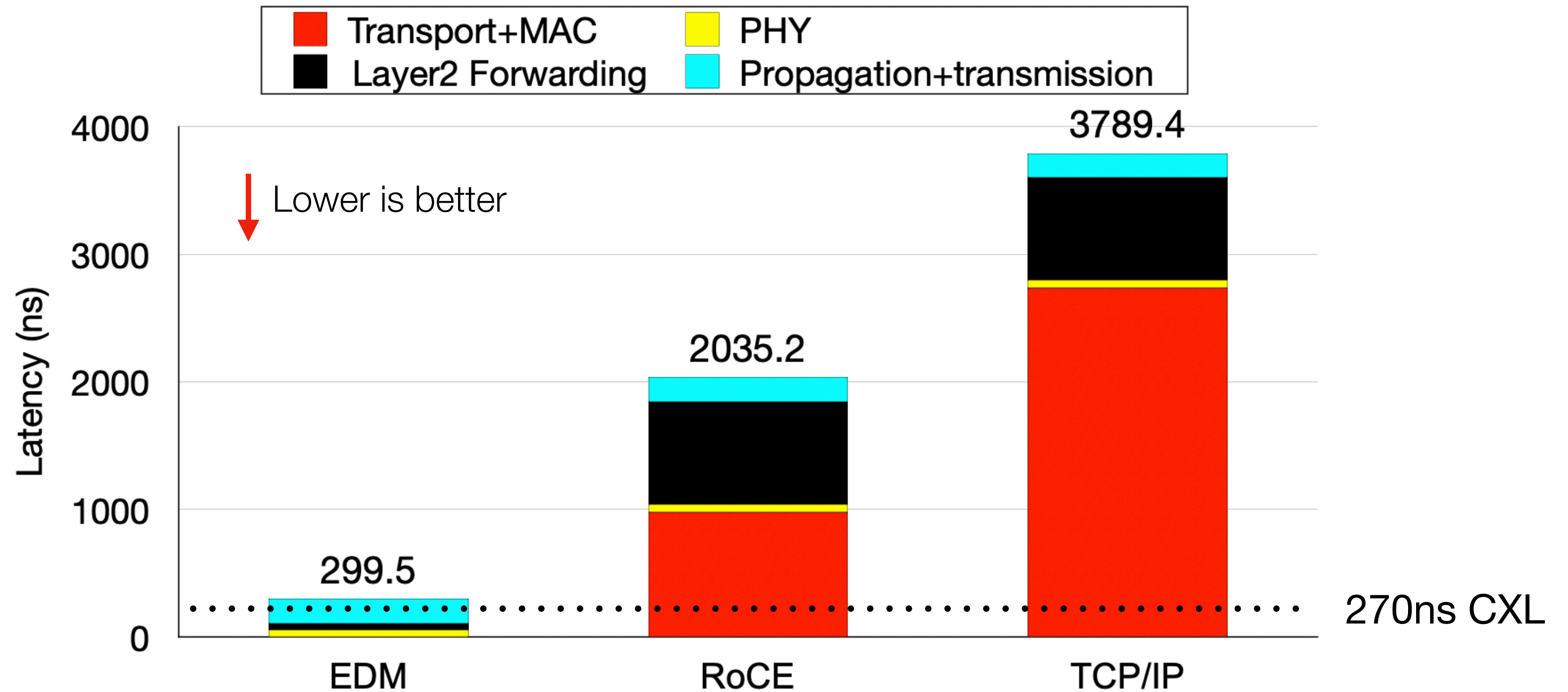- Open-source 25GbE (Corundum)
- Synopsys ASIC RTL compiler

**Network Simulator**
- A single rack with 144 nodes
- Fed with real-world traces
- Compare against 6 classes of scheduling / congestion control

# Evaluation
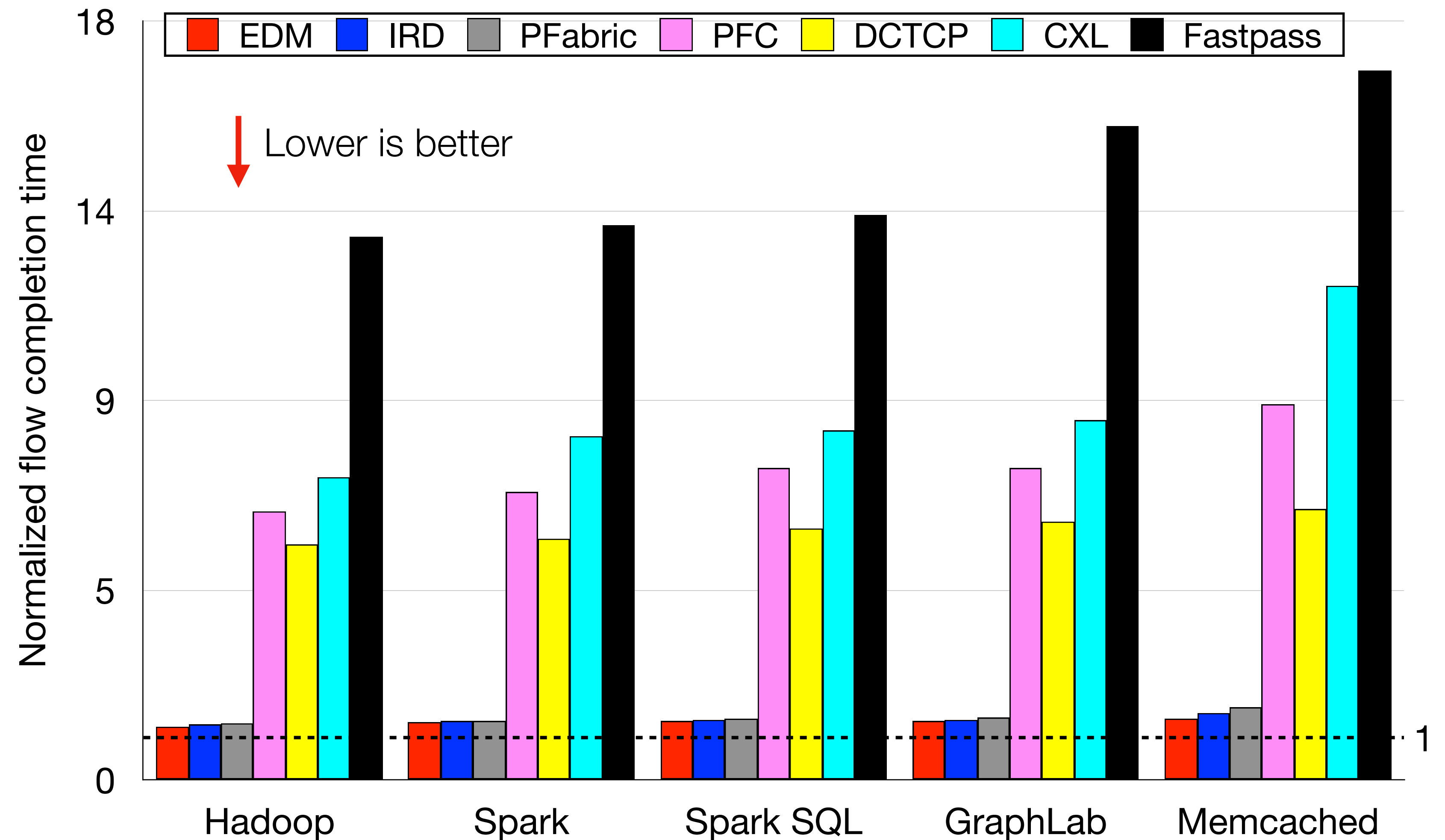
- End-to-end unloaded latency

# Evaluation

- Disaggregated workloads in a loaded network

| Experiment name | Dataset |
|---|---|
| Hadoop, Spark | Generator @sortbench mark.org |
| Spark SQL | Big Data Benchmark@ Berkeley |
| GraphLab | Movie rating data @Netflix |
| Memcached | KV-store@YCSB |

# Summary

- EDM is a low latency Ethernet fabric for memory disaggregation.

- EDM uses two ideas for low latency w/ high bandwidth utilization:

  - EDM implements the **protocol for remote memory access** entirely in the **Ethernet PHY.**

  - EDM implements a **fast, centralized memory traffic scheduler** in the switch's PHY.

- EDM incurs a latency of **~300ns (7x lower than RoCE)** in an unloaded network, and **< 1.3x** its unloaded latency under heavy network loads.

# Thank you !

**Code:** https://github.com/wegul/EDM