# ICC: An Interconnect Controller
## for the Tofu Interconnect Architecture

August 24, 2010

# Takashi Toyoshima

Next Generation Technical Computing Unit

## Fujitsu Limited

**shaping tomorrow with you**

# Background

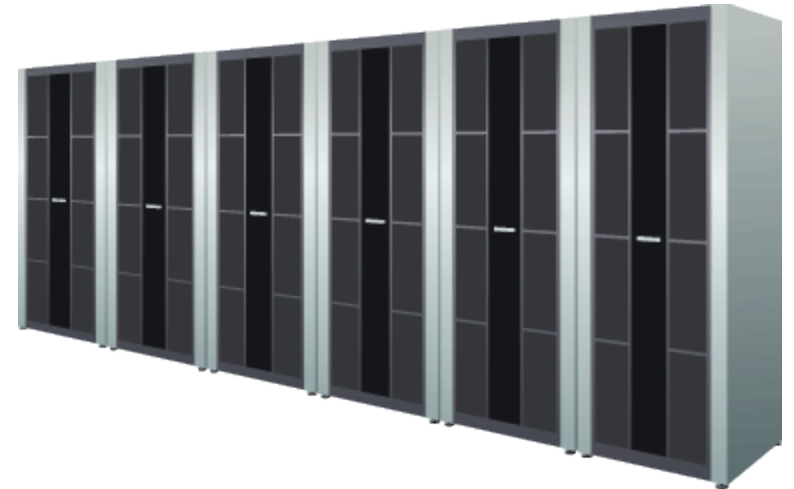- **Requirements for Supercomputing Systems**

  - **Low latency**
    - Communication latency limits the scalability of applications

  - **High bandwidth**
    - Increasing calculation FLOPS requires higher network bandwidth be balanced with FLOPS

  - **RAS – Reliability, Availability and Serviceability**
    - The risk of hardware faults in large systems increases along with the increased number of nodes
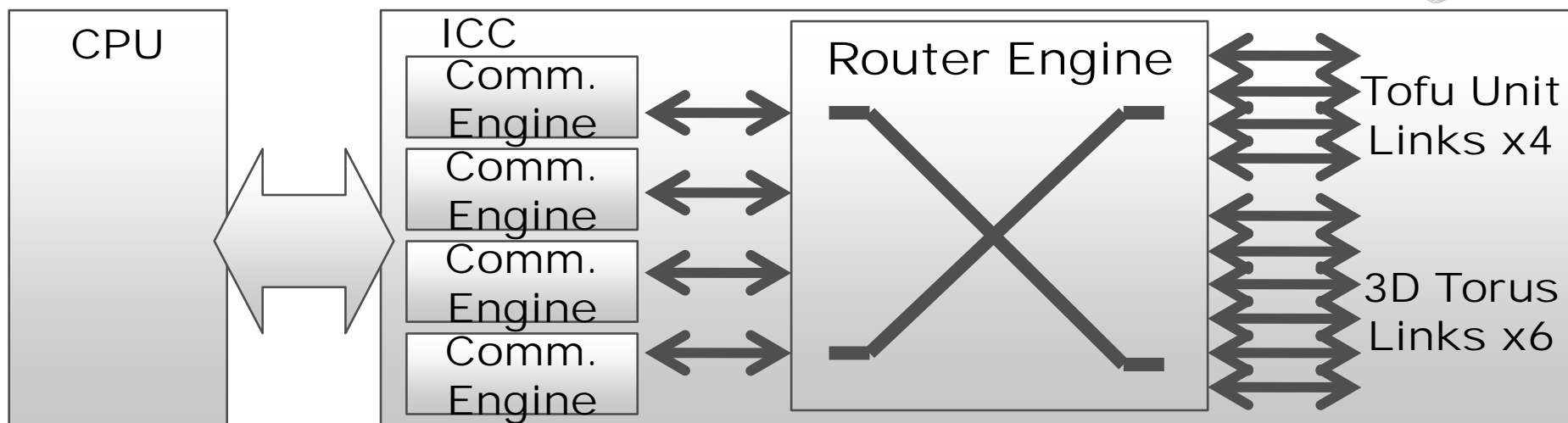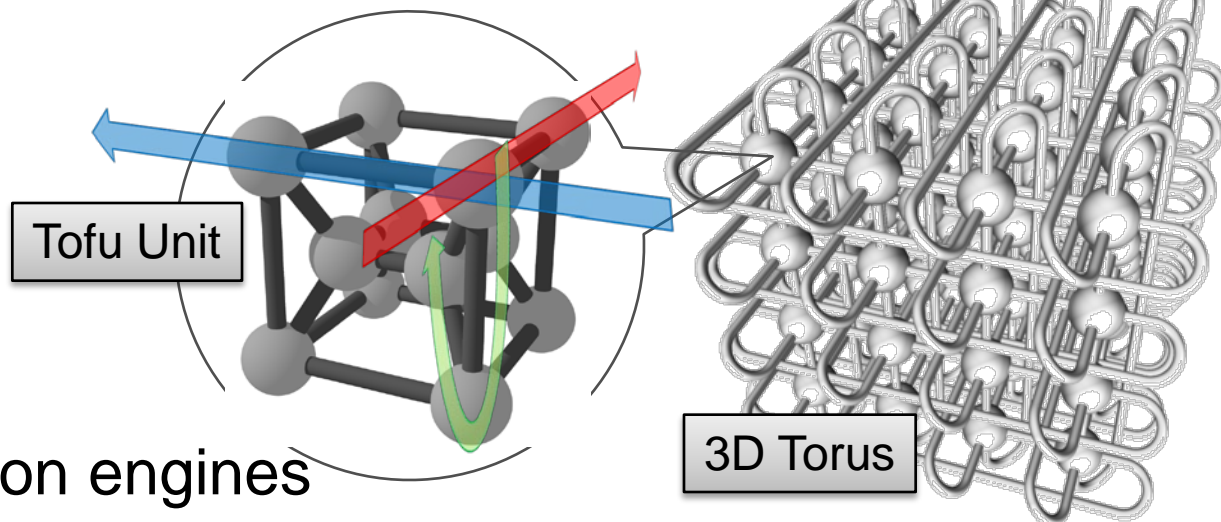
# Fujitsu's New Interconnect Architecture

**FUJITSU**

- ■ 6D Mesh/Torus Interconnect Architecture*
  - ■ Scalability
  - ■ Fault-tolerance
- ■ LSI Features
  - ■ Ten network links
  - ■ Four communication engines

Tofu Unit

3D Torus

CPU

ICC

Comm. Engine

Comm. Engine

Comm. Engine

Comm. Engine

Router Engine

Tofu Unit Links x4

3D Torus Links x6

(*) "Tofu: A 6D Mesh/Torus Interconnect for Exascale Computers", IEEE Computer, vol.42, no.11, Yuichiro Ajima, Shinji Sumimoto, Toshiyuki Shimizu

# **Implementation**

- Implementation

- Features

  - Overview

  - Interface features for latency and throughput

  - Network features for network utilization

- Conclusion

# Specifications

- **Fujitsu's 65nm CMOS Technology**
  - Die size
    - 18.2mm × 18.1mm
  - Transistors
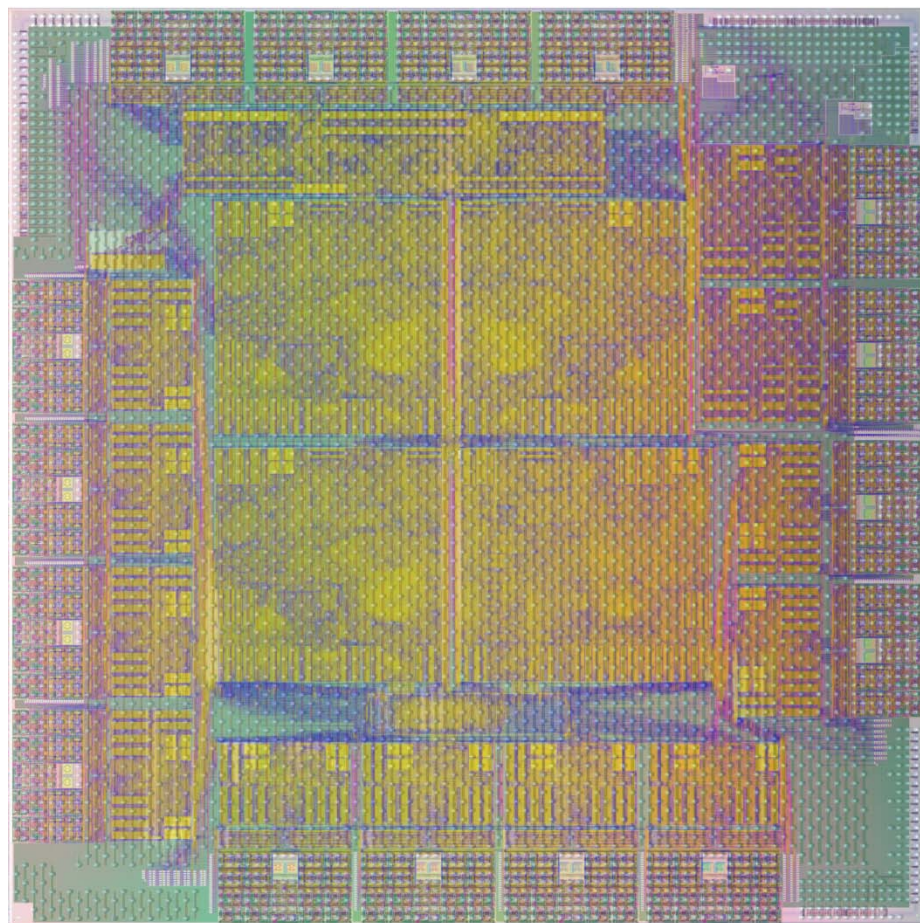    - 48M gates for logic
    - 12M-bit SRAM cells
  - I/O
    - 5GB/s Ports × 16
      - 6.25Gb/s × 8 links / port
  - Misc.
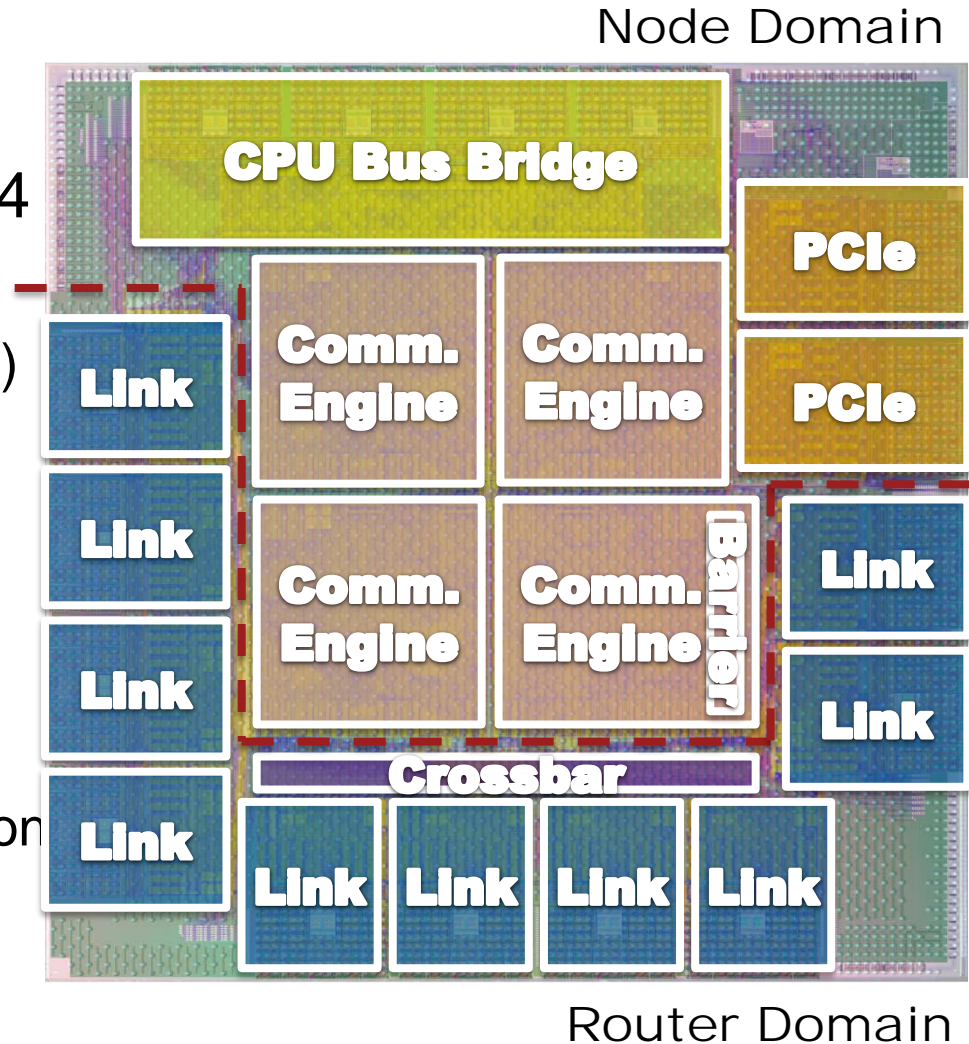    - ASIC design flow
    - 312.5MHz/625.0MHz

# Floor Plan

- ■ **Node Domain**
  - ■ CPU bus bridge
    - • 20GB/s in each direction
  - ■ Communication engines×4
    - • 5GB/s in each direction
    - • Barrier engine (Comm.#0 only)
  - ■ PCIe 2.0 root complex×2
    - • Isolated power domain
- ■ **Router Domain**
  - ■ Crossbar
    - • 14 ports 5GB/s in each direction
  - ■ Link ports×10
    - • 5GB/s in each direction



**Node Domain**

CPU Bus Bridge

PCIe

Link

Comm. Engine

Comm. Engine

PCIe

Link

Link

Comm. Engine

Comm. Engine

Barrier

Link

Link

Link

Link

Crossbar

Link

Link

Link

Link

Link

**Router Domain**
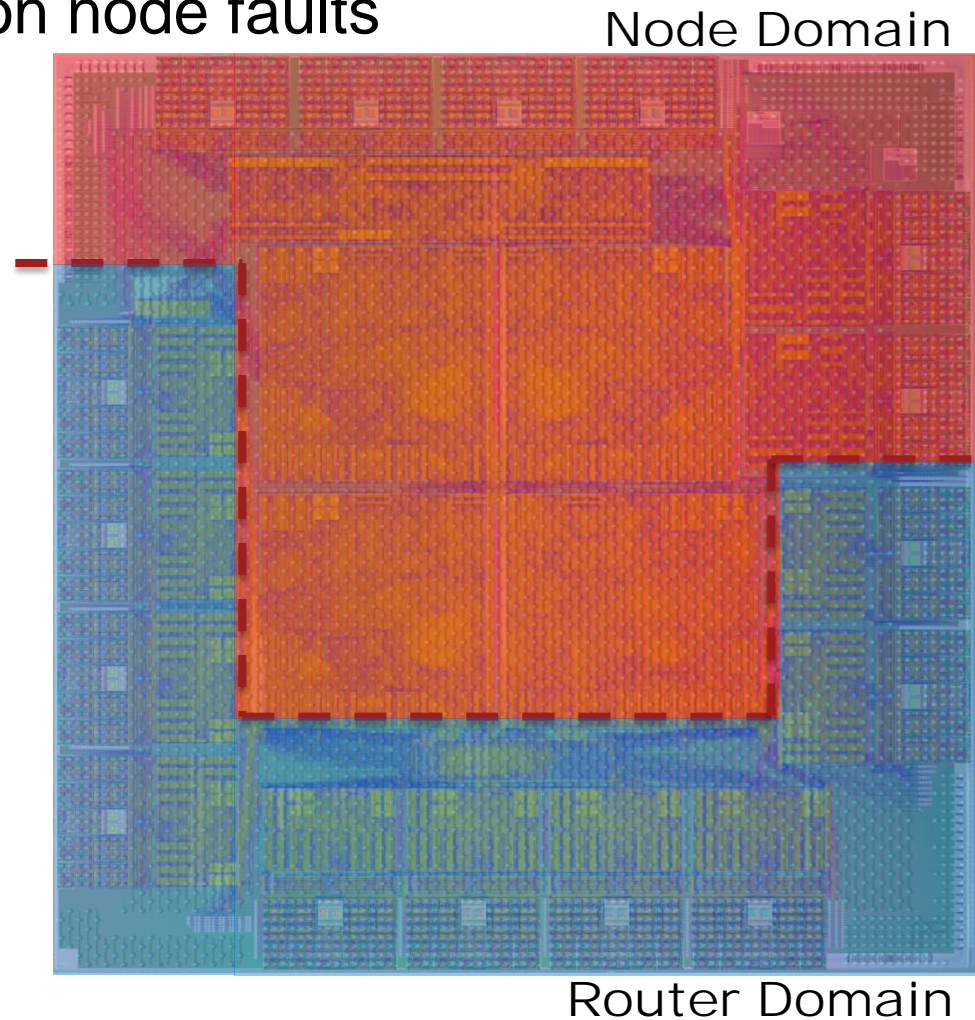
# RAS Features

- **Fault Domain Isolation**
  - Router continues to work on node faults
- **Error Protection**
  - Radiation-hardened FFs
  - ECC protection
    - RAM/Data path
  - Parity error detection
    - Control path
  - CRC protection
    - Data link/Transaction



**Node Domain**

**Router Domain**

# Features
## Overview

- Implementation

- Features

  - Overview

  - Interface features for latency and throughput

  - Network features for network utilization

- Conclusion

# ICC Features

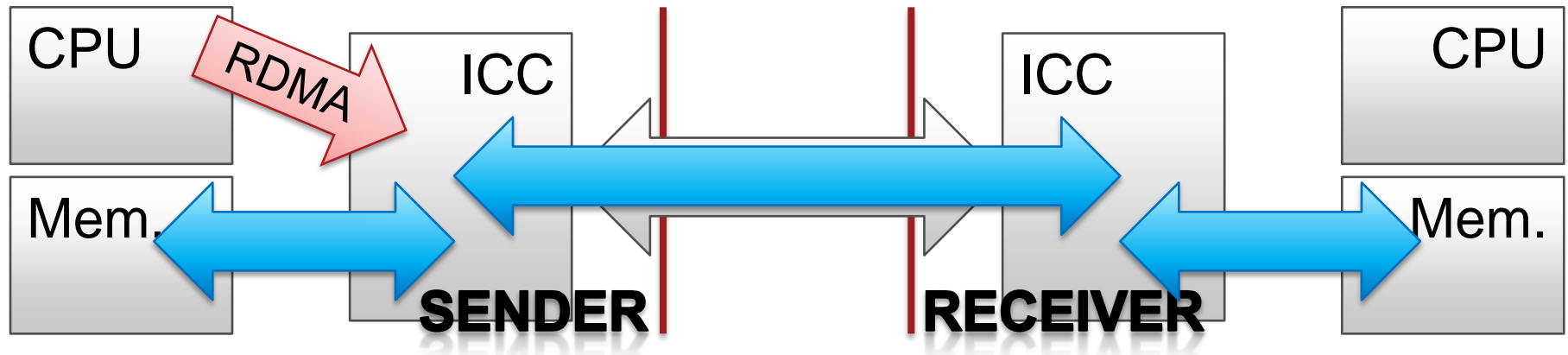| | Latency | Throughput | RAS |
|---|---|---|---|
| **System** | ✓Many Neighbors<br>✓Hop Reduction<br>✓3D Torus View ⭐ | ✓Many Neighbors<br>✓Trunking ⭐<br><br>✓GAP Control ⭐ | ✓Detour Path ⭐<br>✓Subnet Partitioning ⭐ |
| **Network Interface** | ✓RDMA<br> - Quick Start ⭐<br> - Piggyback<br> - Strong Order ⭐<br>✓Stream Offload ⭐<br>✓Barrier Engine ⭐ | ✓Multi-Interfaces ⭐<br> - User Thread x2<br> - Kernel Thread | ✓Radiation-hardened FF<br>✓ECC<br>✓Parity<br>✓CRC |
| **Router Engine** | ✓Cut-through<br>✓Grant Prediction ⭐<br>✓Straight Bypath ⭐ | ✓Straight Bypath ⭐<br>✓New VC Scheduling ⭐ | ✓Node Error Isolation ⭐<br>✓Radiation-hardened FF<br>✓ECC<br>✓Parity<br>✓CRC |

⭐ : Unique Features

Today's topics are highlighted in red

# RDMA: Remote Direct Memory Access

■ Features



| Operation | READ (Get) / WRITE (Put) |
|---|---|
| Length | ~16MB |
| MTU | 256B~1920B |
| Virtual Address | Support (64K set) |

■ Low Latency and High Throughput

■ Command supply throughput and latency

■ Out-of-ordered I/O memory bus

# RDMA Optimizations

**FUJITSU**

- **Sender Techniques**

  - Direct descriptor

    - Quick command supply

  - Piggyback

    - Command embedded communication payload

    - Short message sending without any DMA

|  | Throughput | Latency |
|---|---|---|
| **PIO** | | ✓Good |
| **DMA** | ✓Good | |

**Command Supply Performance**

- **Receiver Techniques**

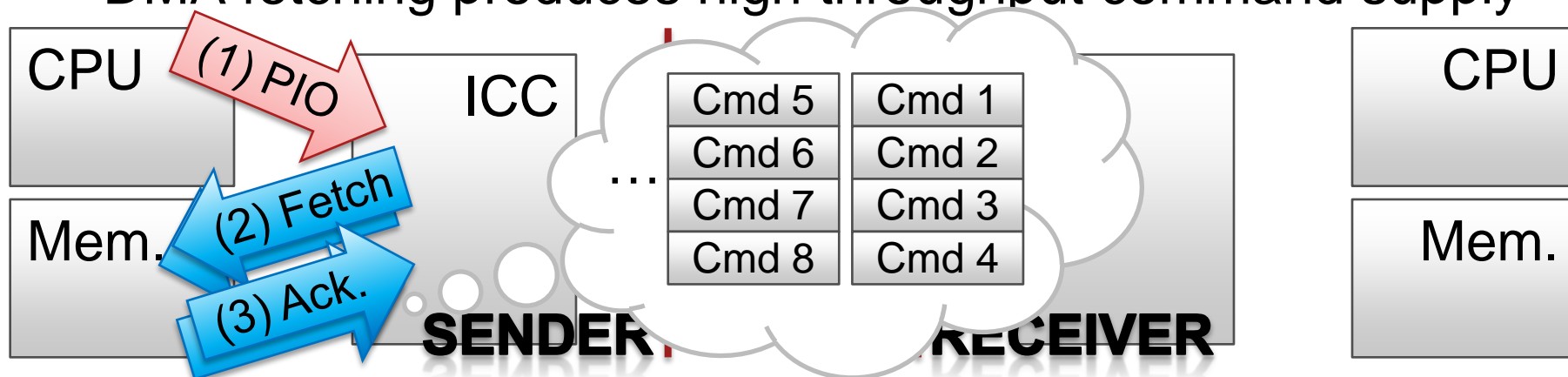- Out-of-ordered I/O memory bus

  - High throughput bus transaction

- Strong ordered store

  - In order completion of DMA transactions for buffer polling
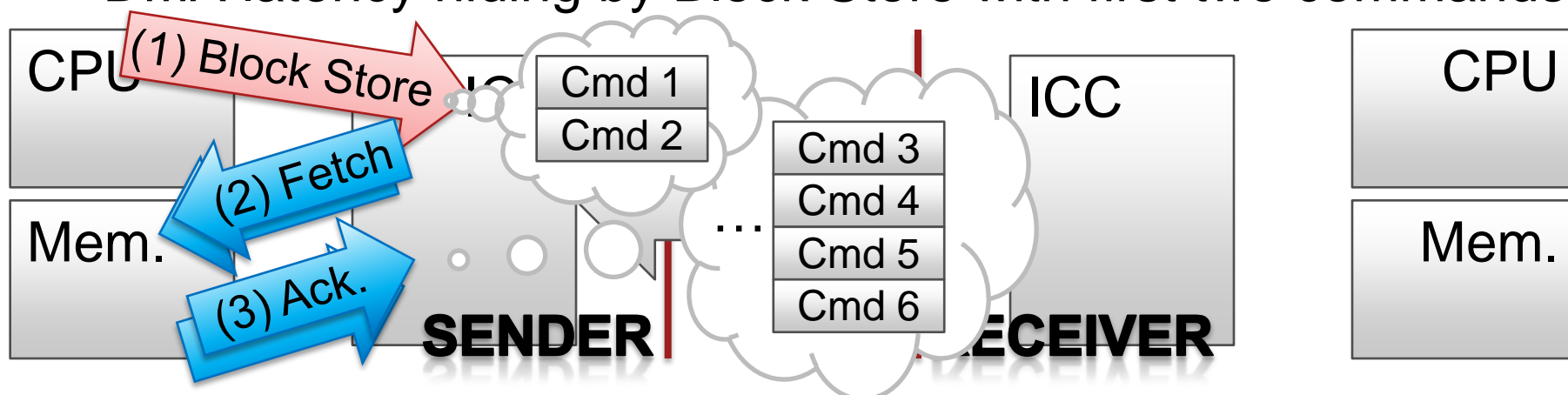
# Direct Descriptor Feature

- ## Normal Command Supply
  - DMA fetching produces high throughput command supply



CPU

(1) PIO

ICC

| Cmd 5 | Cmd 1 |
|-------|-------|
| Cmd 6 | Cmd 2 |
| Cmd 7 | Cmd 3 |
| Cmd 8 | Cmd 4 |

Mem.

(2) Fetch

(3) Ack.

**SENDER**        **RECEIVER**

CPU

Mem.

- ## Direct Descriptor and DMA Command Supply
  - DMA latency hiding by Block Store with first two commands

CPU

(1) Block Store

ICC

| Cmd 1 |
|-------|
| Cmd 2 |

| Cmd 3 |
|-------|
| Cmd 4 |
| Cmd 5 |
| Cmd 6 |

(2) Fetch

Mem.

(3) Ack.

**SENDER**        **RECEIVER**
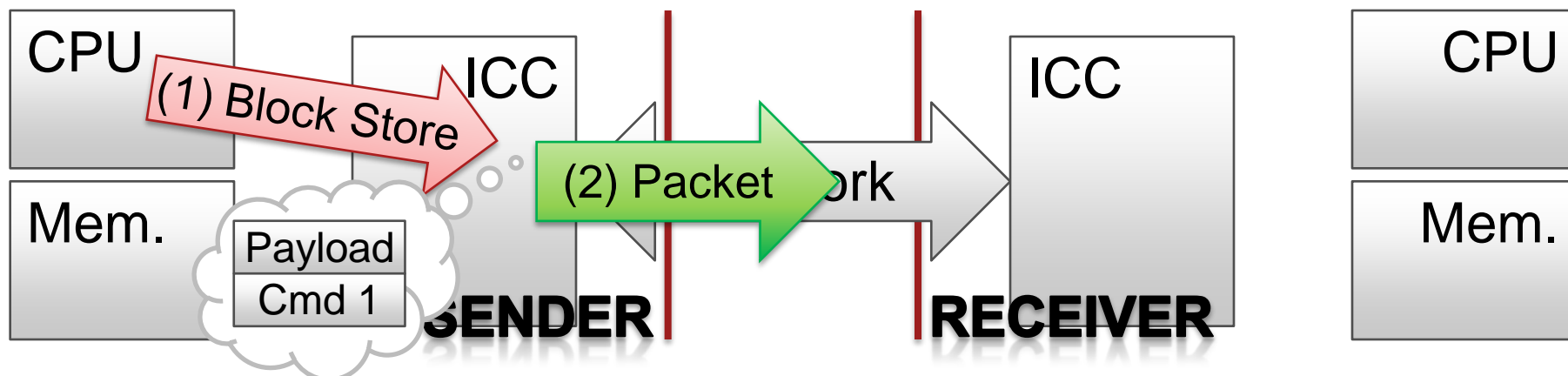
ICC

CPU

Mem.

# Piggyback Feature

■ Normal Payload Supply

- User messages (payload) should be fetched by DMA
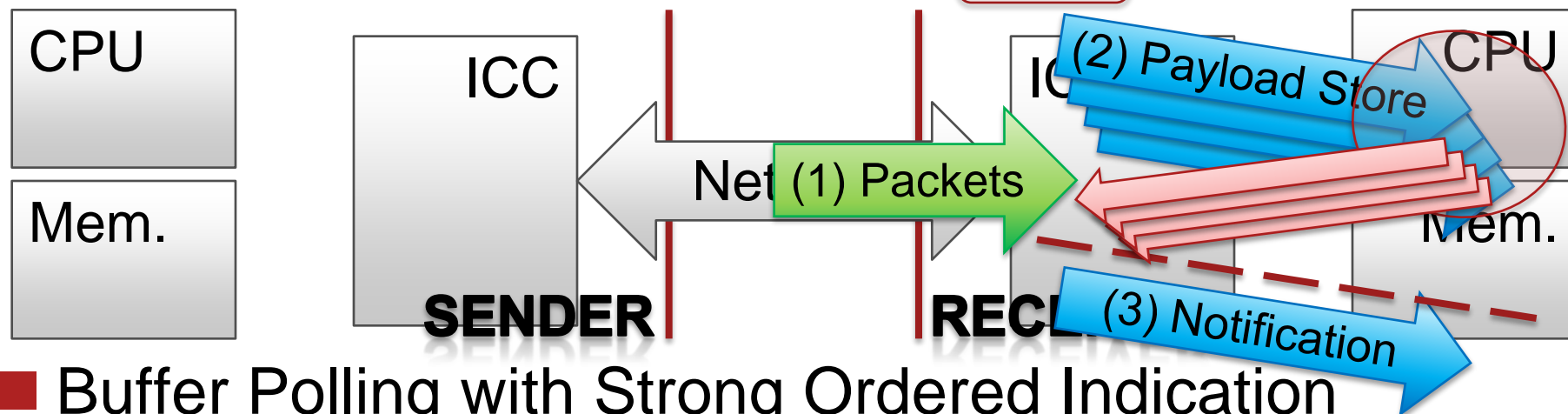


■ Piggyback Payload Supply

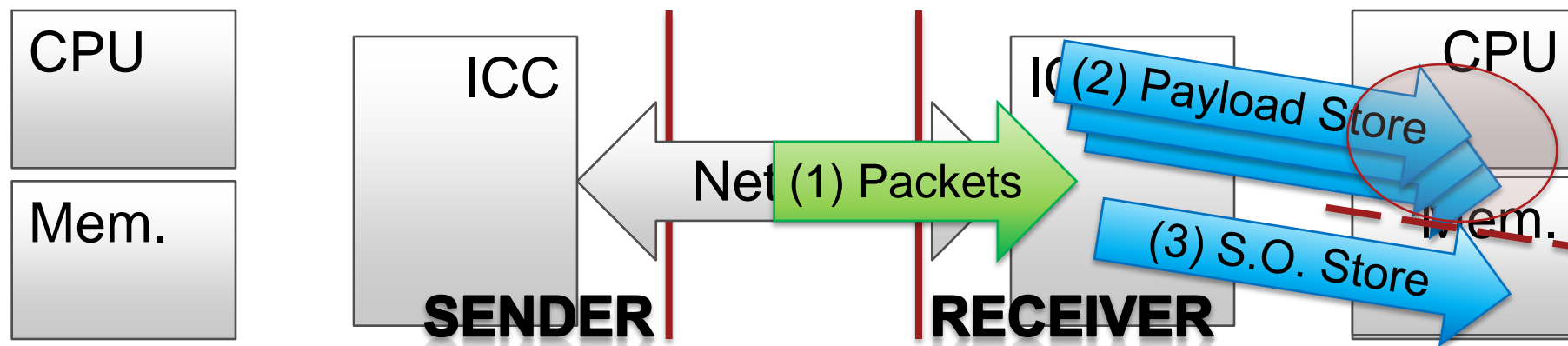- User messages (payload) are embedded in commands

# Out-of-Ordered I/O Memory Bus

■ Completion Notification Polling

  ■ ICC notifies after the completion of O-o-O DMA Stores
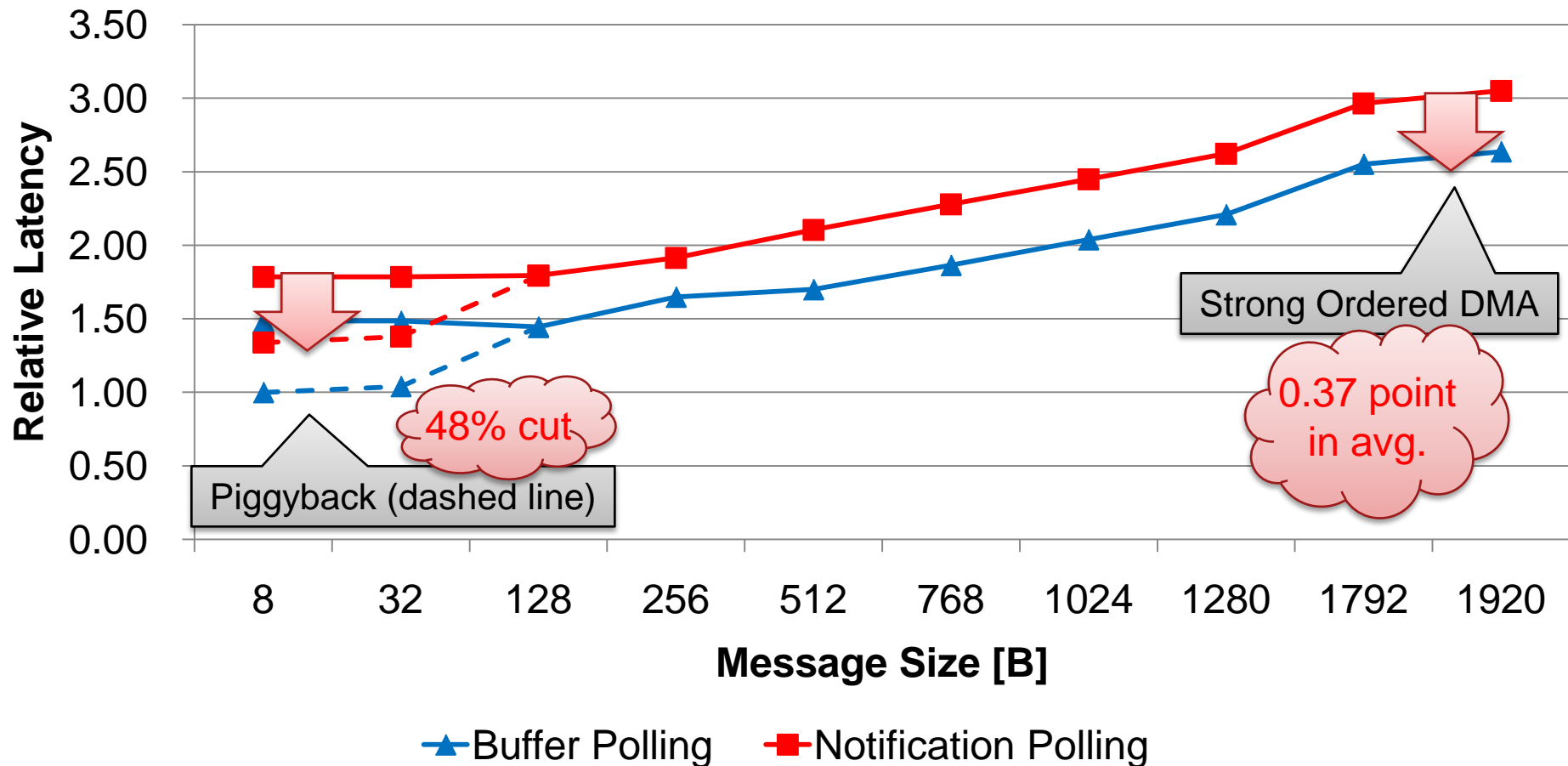
  CPU

  Mem.

  ICC

  SENDER

  Net (1) Packets

  IC...

  RECE...

  (2) Payload Store

  (3) Notification

  CPU

  Mem.

■ Buffer Polling with Strong Ordered Indication

  ■ Memory controller guarantees specified DMA ordering

  CPU

  Mem.

  ICC

  SENDER

  Net (1) Packets

  IC...

  RECEIVER

  (2) Payload Store

  (3) S.O. Store

  CPU

  Mem.

# RDMA Performance

- Hardware Measured Results
  - Piggyback achieves low latency in short message
  - Strong ordered packet makes buffer polling possible

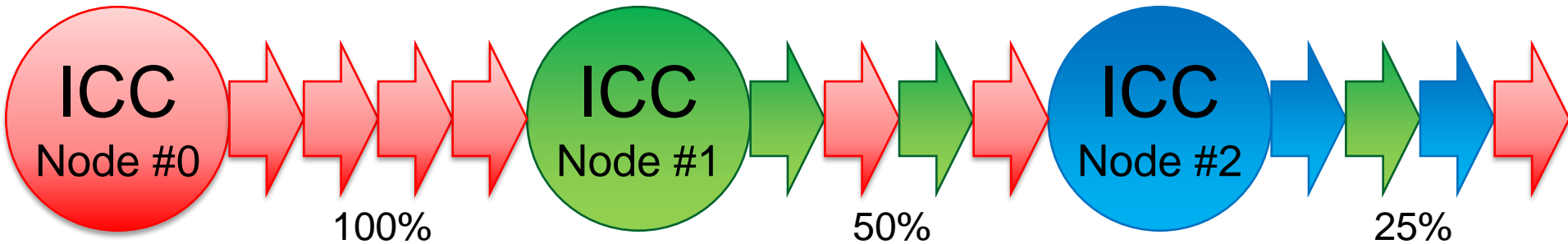Chart: Relative Latency vs Message Size [B]

- 48% cut
- Piggyback (dashed line)
- Strong Ordered DMA
- 0.37 point in avg.

Legend: Buffer Polling, Notification Polling

# Network Utilization Problems

■ Global Unfairness of Throughput

■ Arbitrations with local fairness cause global unfairness



ICC Node #0 → 100% → ICC Node #1 → 50% → ICC Node #2 → 25%

■ Non-uniform Application Traffic in Time and Space

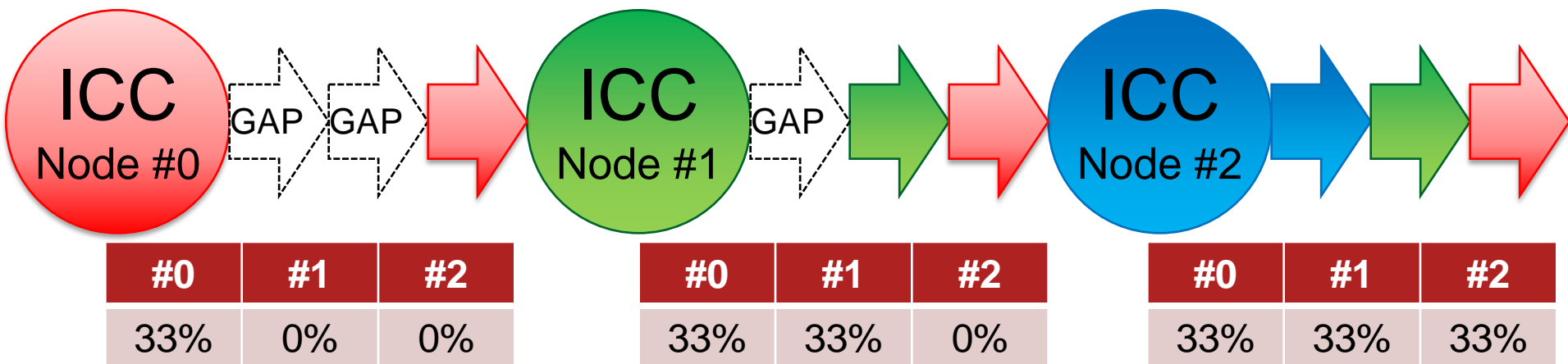■ Bandwidth of idle links needs to be used effectively



Phase A — Idle

Phase B

ICC — Idle — Bottleneck

# Global Unfairness of Throughput

- Local Fairness of arbitration cause global unfairness
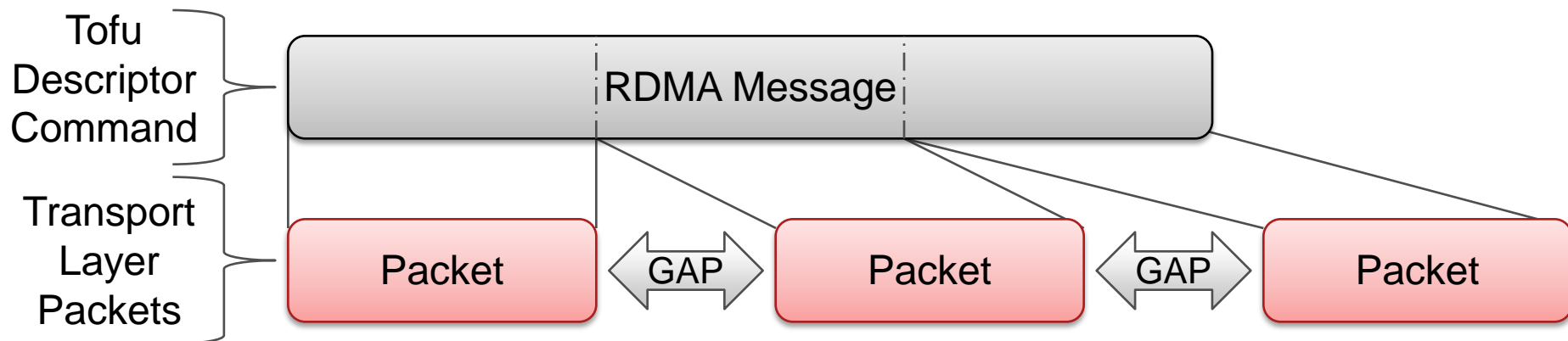  - Arbiters on every junction treat all incoming traffic fairly



| #0 | #1 | #2 |
|------|------|------|
| 100% | 0% | 0% |

| #0 | #1 | #2 |
|------|------|------|
| 50% | 50% | 0% |

| #0 | #1 | #2 |
|------|------|------|
| 25% | 25% | 50% |

halved        halved

ICC Node #2

Ideal Throughput

| #0 | #1 | #2 |
|------|------|------|
| 33% | 33% | 33% |

# Injection Rate Control

**FUJITSU**

■ Software Specify the Inter-Packet GAP Parameter



| #0 | #1 | #2 |
|----|----|----|
| 33% | 0% | 0% |

| #0 | #1 | #2 |
|----|----|----|
| 33% | 33% | 0% |

| #0 | #1 | #2 |
|----|----|----|
| 33% | 33% | 33% |

■ Communication engine works to control injection rate

- Insert temporal gaps between transmitting packets
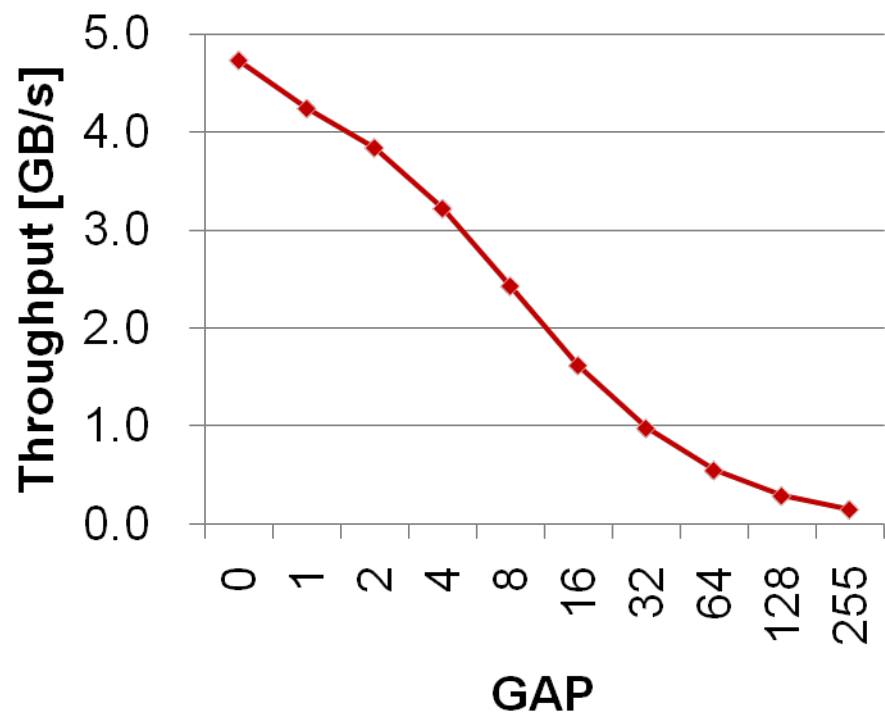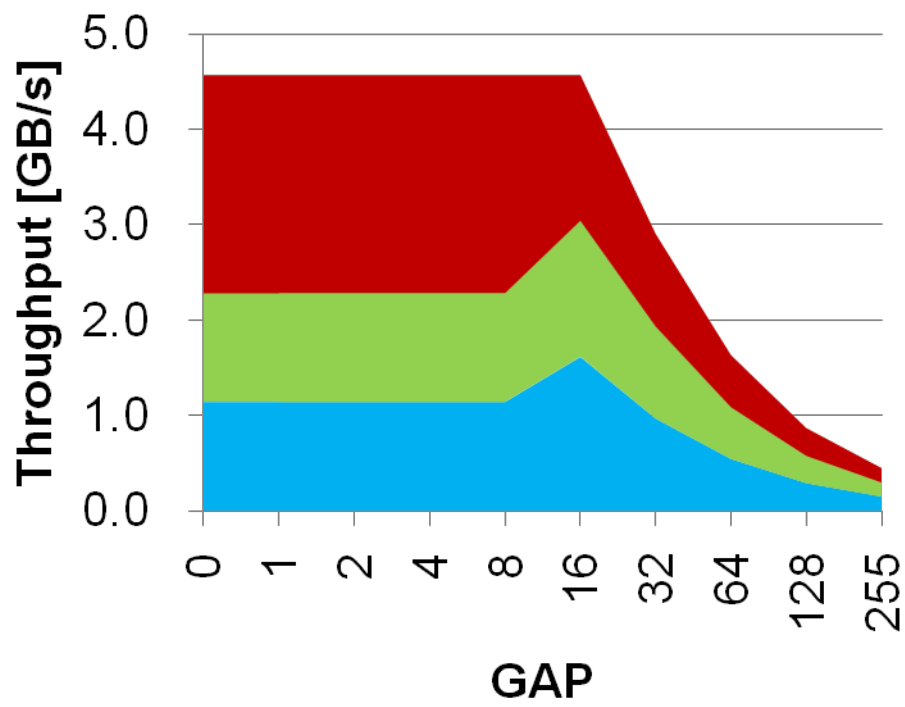- Interval can be specified by the user

Tofu Descriptor Command

RDMA Message

Transport Layer Packets

Packet ⟷ GAP ⟷ Packet ⟷ GAP ⟷ Packet

# Throughput Performance

**FUJITSU**

■ Hardware Measured Results

■ Software can specify fine grained GAP parameters: 0-255

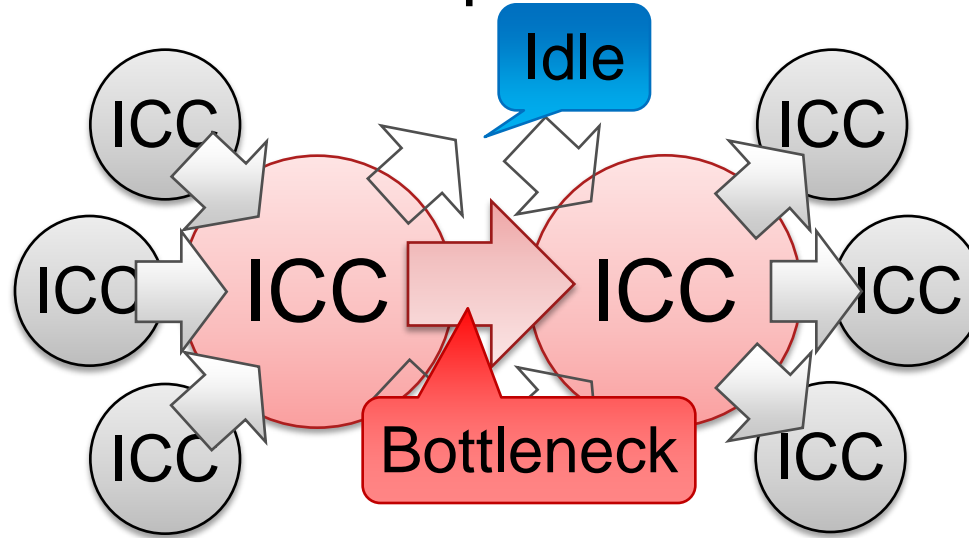■ GAP works to control throughput effectively

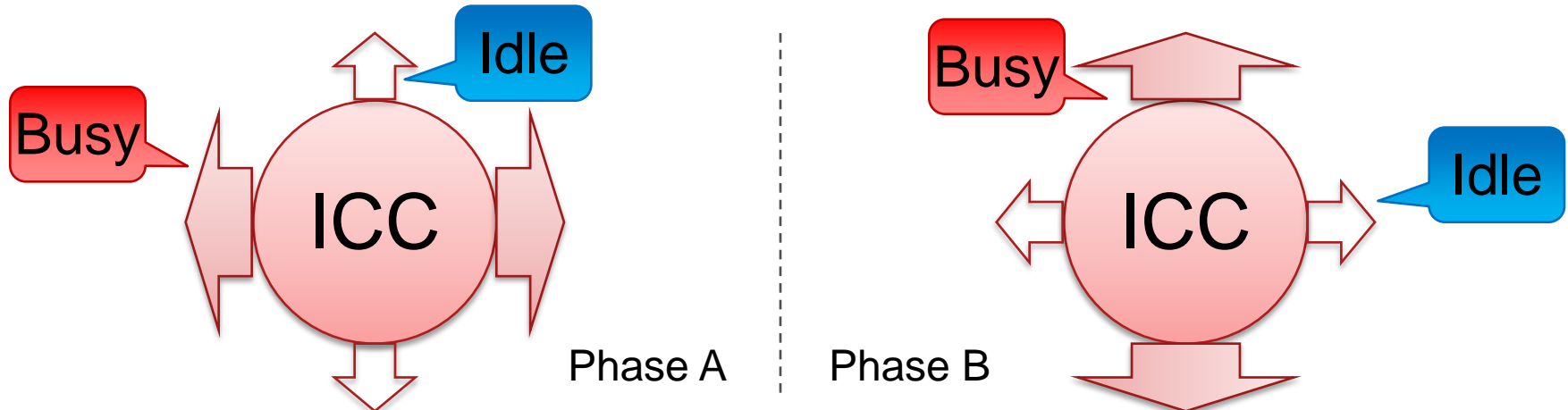## GAP Sensitivity

## Stacked Throughput



■ Node #0  ■ Node #1  ■ Node #2

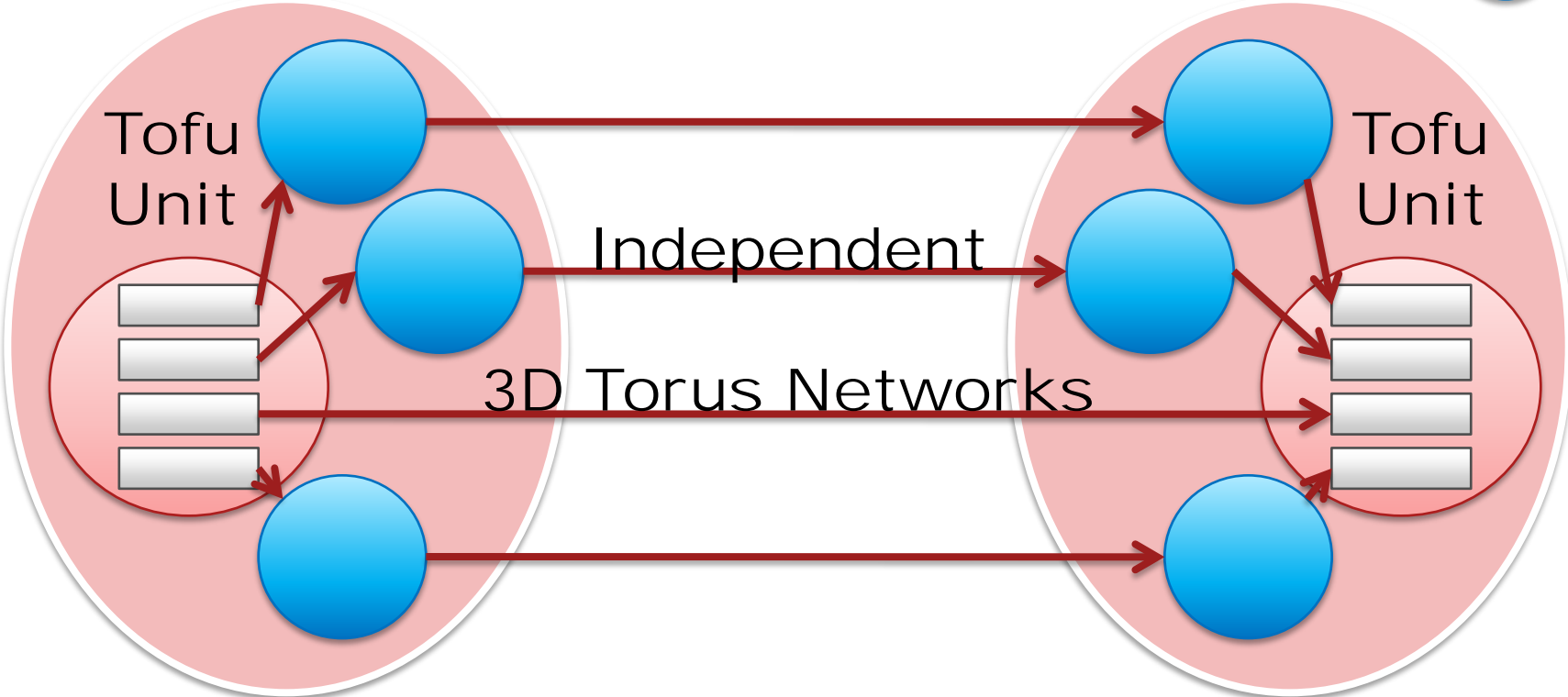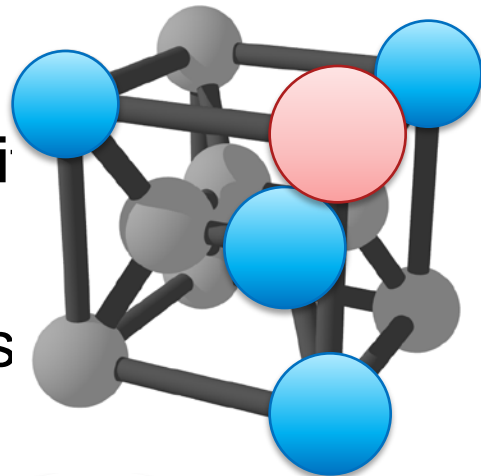# Non-uniform Application Traffic

- Non-uniform Traffic in Space



- Non-uniform Traffic in Time
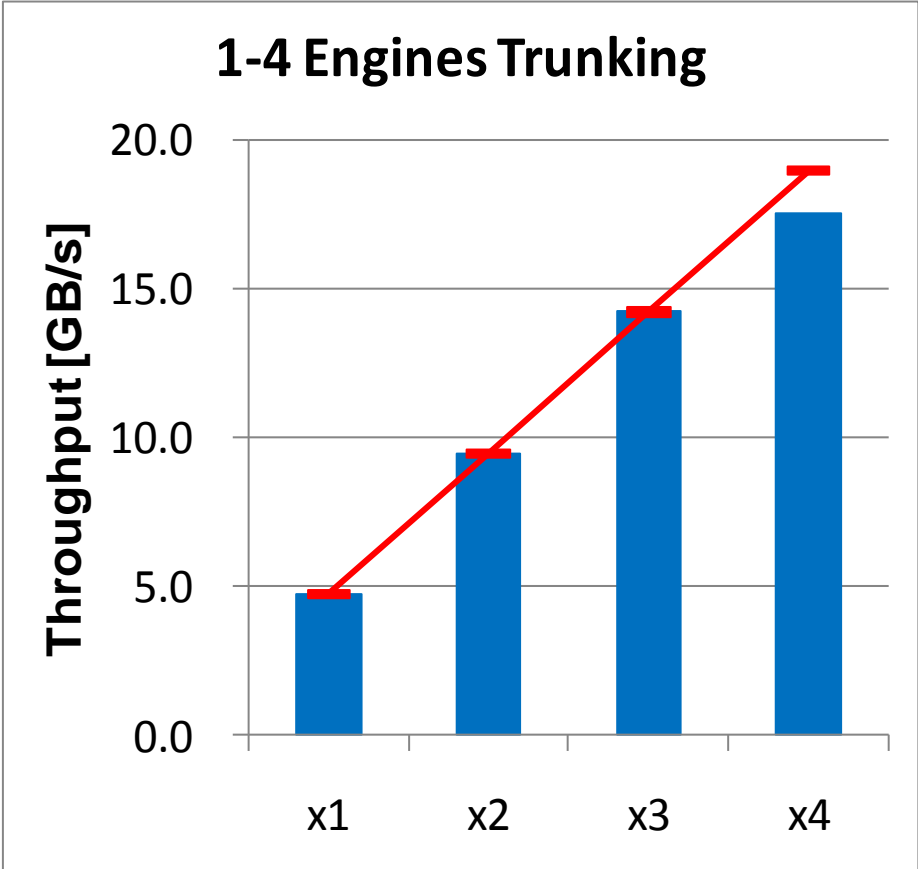
# Trunking Communication
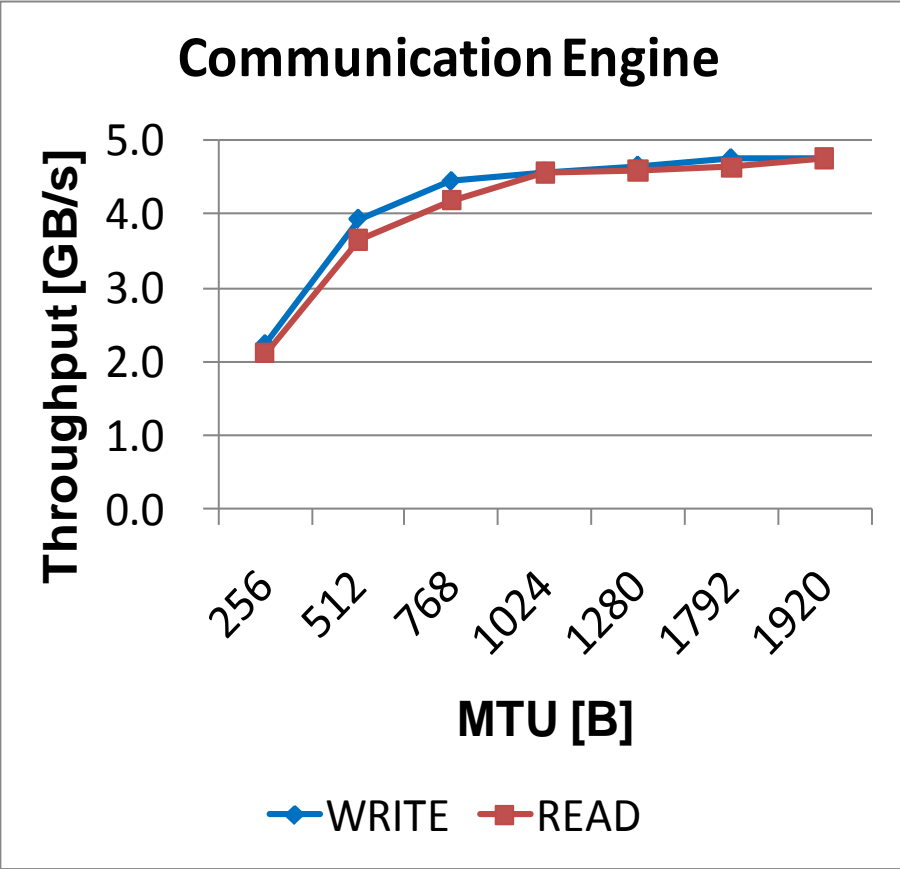
- ■ Trunking Independent Idle Paths
  - ■ Nodes have **four** neighborhoods in Tofu Unit
    - • Independent links and 3D-Torus networks
  - ■ Each node has **four** communication engines
    - • Up to ×4 throughput

**Tofu Unit**

**Independent**

**3D Torus Networks**

**Tofu Unit**

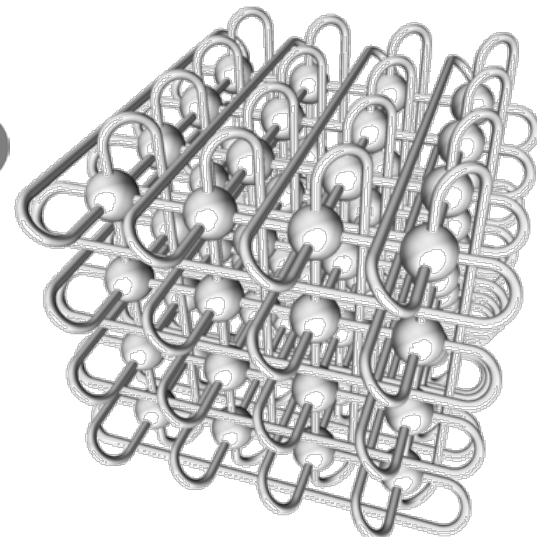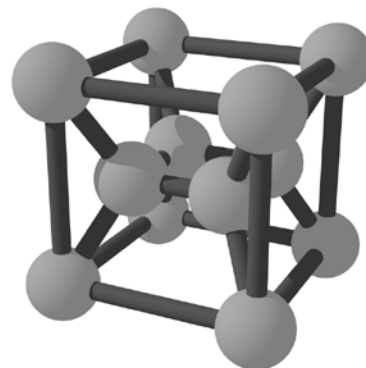# Trunking Performance

- **Results**
  - Communication engines achieve good performance
  - Trunking mechanisms scale up to four engines

# **Conclusion**
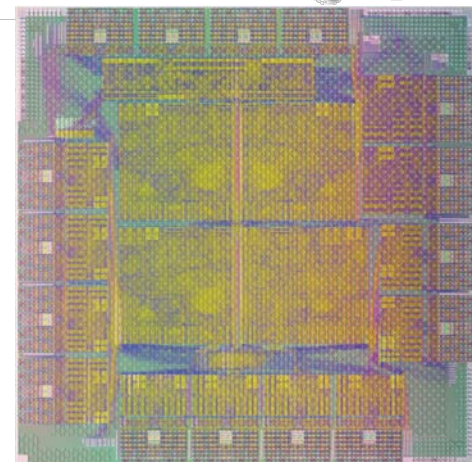
- Implementation

- Features

  - Overview

  - Interface features for latency and throughput

  - Network features for network utilization

- Conclusion

# Concluding Remarks

- **Tofu: A 6D mesh/torus interconnect architecture**
  - Interconnect for Fujitsu's Peta/Exascale computing systems
  - Low latency, High bandwidth and RAS

- **Features**
  - High-throughput and low-latency RDMA
    - Direct Descriptor and Piggyback
    - Out of Ordered I/O Memory Bus
  - Network features for network utilization
    - Network injection rate control
    - Trunking up to four times throughput

**FUJITSU**

# Thanks to…
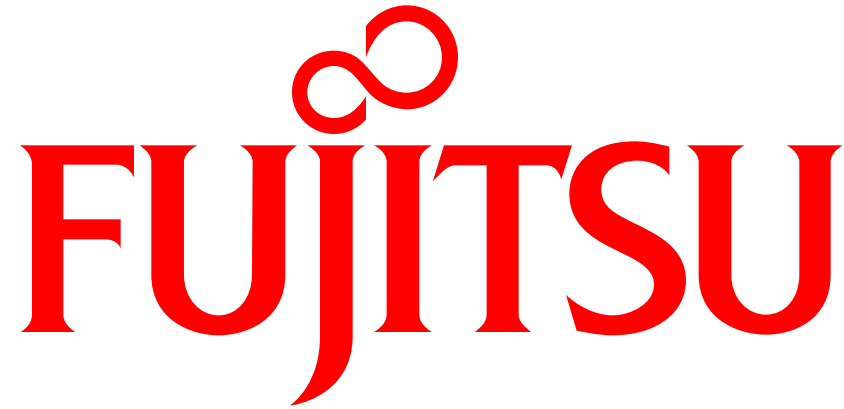
## Next Generation Technical Computing Unit

Aiichiro Inoue, Yuji Oinaga

## Tofu Architecture Team

Toshiyuki Shimizu, Yuichiro Ajima, Tomohiro Inoue, Shinya Hiramoto

## ICC Design Team

Takeo Asakawa, Akira Asato, Takumi Maruyama,

Koichiro Takayama, Koichi Yoshimi, Osamu Moriyama,

Masao Yoshikawa, Shinichi Iwasaki, Takekazu Tabata,

Yoshiro Ikeda, Yuzo Takagi,

Yoshihito Matsushita, Toshihiko Kodama, Satoshi Nakagawa,

Masato Inokai, Shigekatsu Sagi, Ikuto Hosokawa,

Yaroku Sugiyama, Takahide Yoshikawa