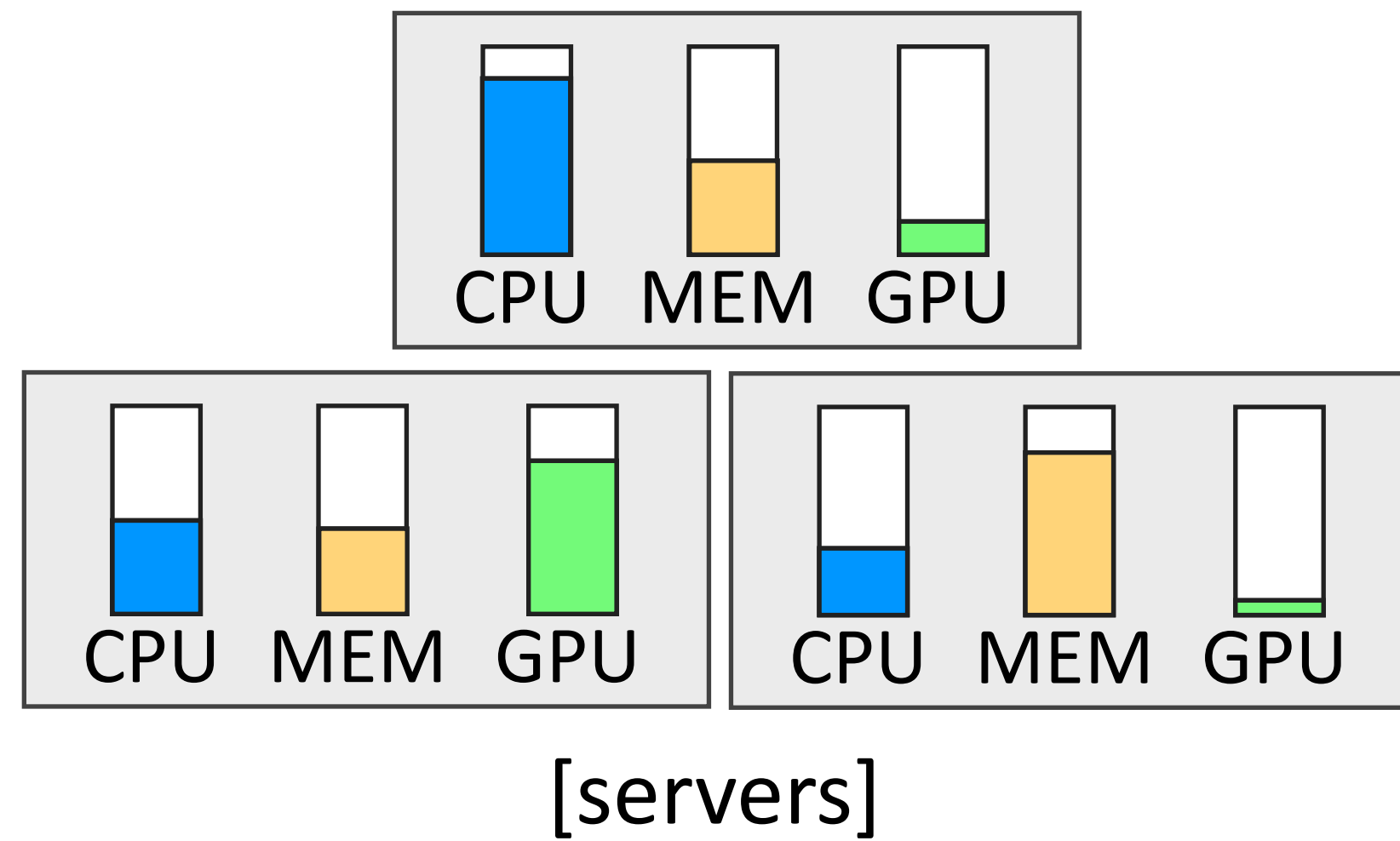# MIND: In-Network Memory Management for Disaggregated Data Centers
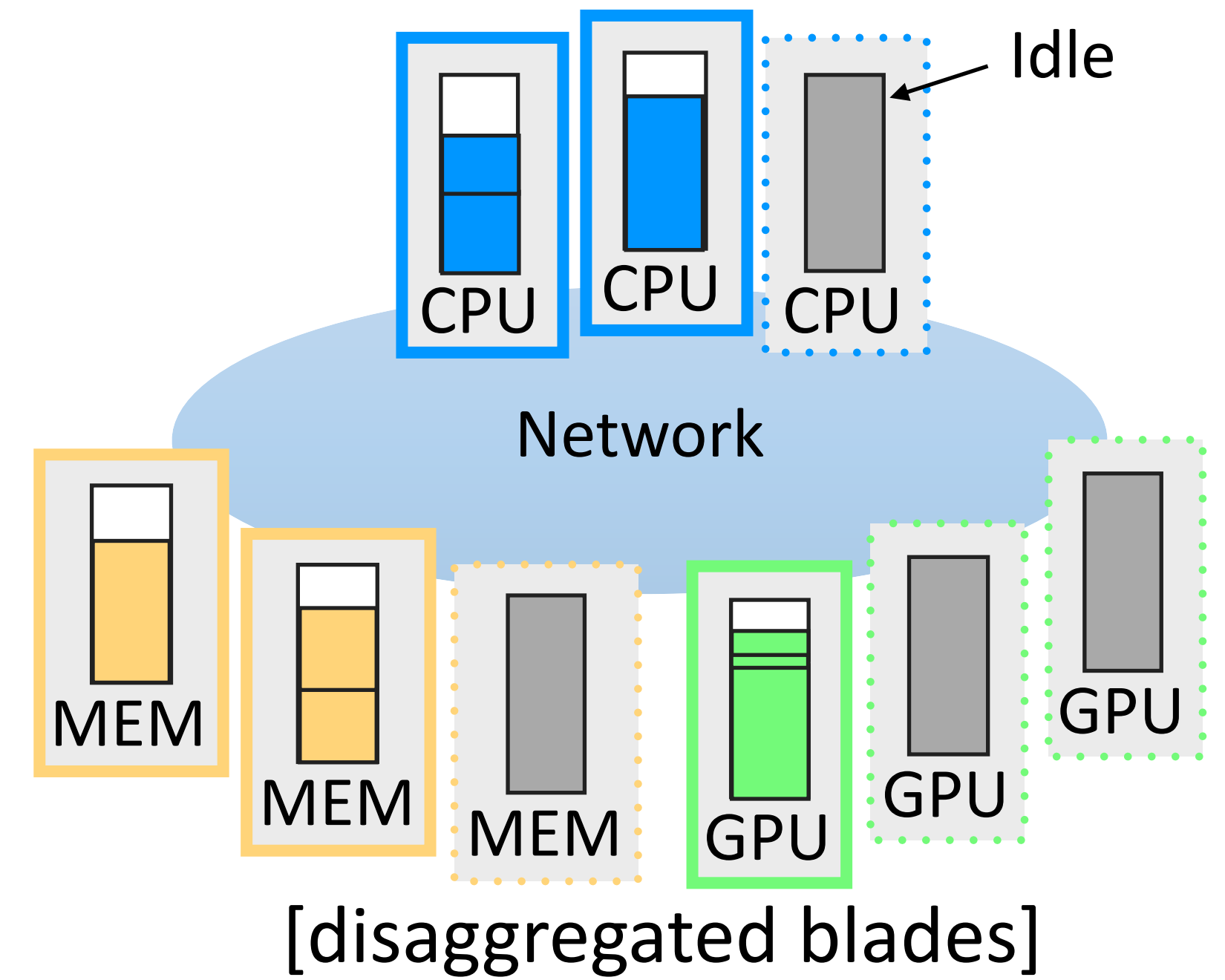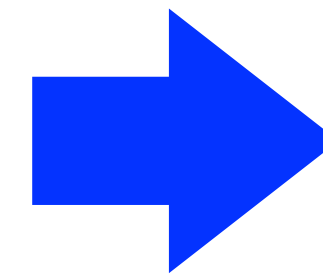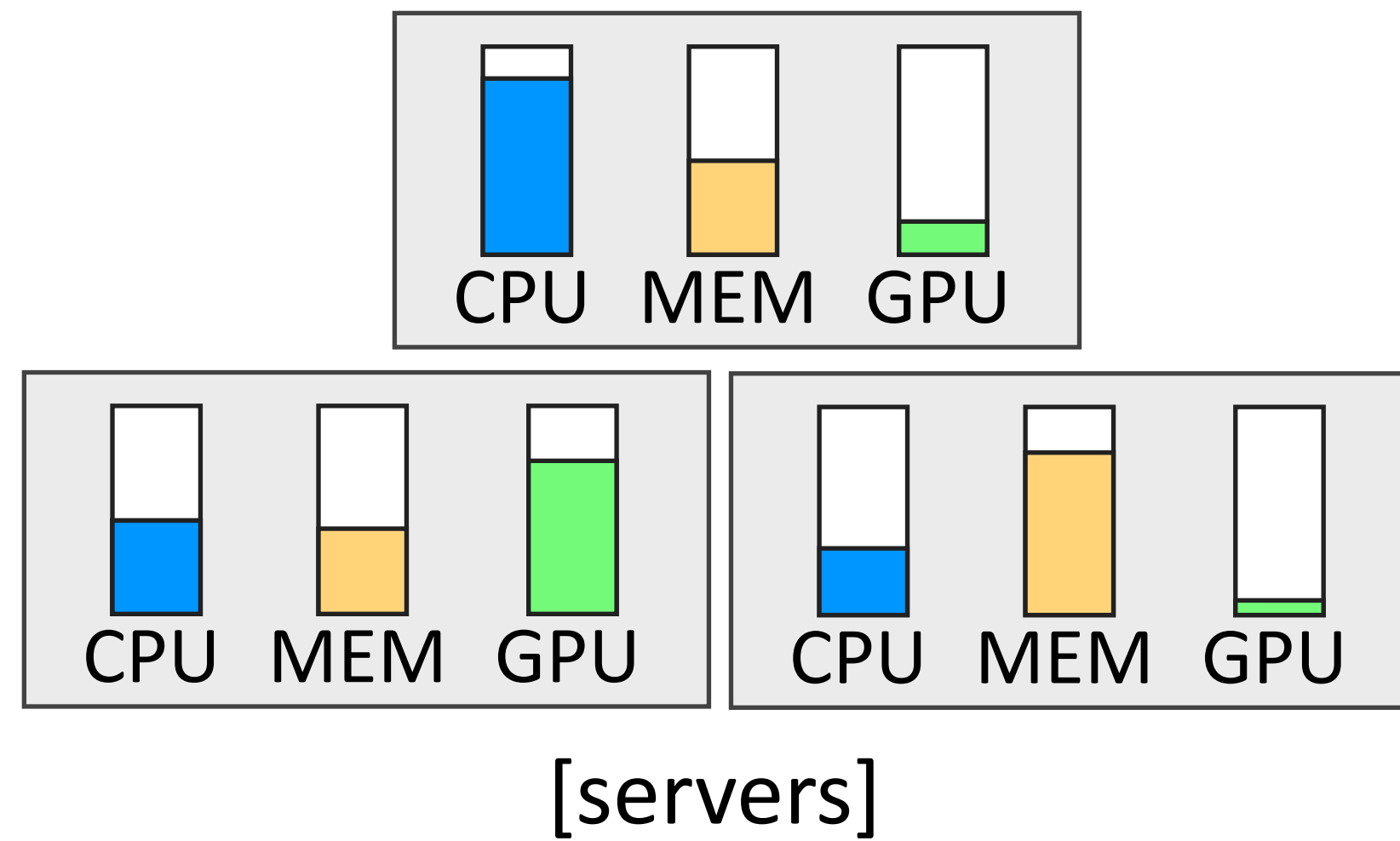
*Seung-seob Lee, Yanpeng Yu, Yupeng Tang,*
*Anurag Khandelwal, Lin Zhong, Abhishek Bhattacharjee*

*Yale University*

# Resource Disaggregation



[servers]

# Resource Disaggregation



[servers]

[disaggregated blades]

# Resource Disaggregation



[servers]

[disaggregated blades]
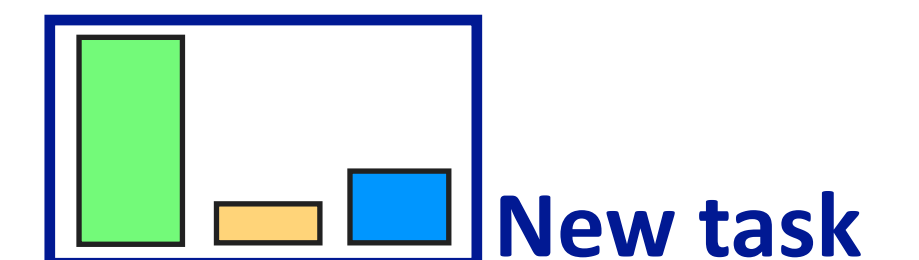
- **Benefits of resource disaggregation**
  - High resource utilization
  - Easy to manage (modularity)
  - Elastic scalability

# Memory Disaggregation

- **Need for memory disaggregation**

**Low utilization**
(as low as 30%)

**Unbalanced usage**
(> 70% of the time in clusters)

**Energy consumption**
(up to 46% of avg. system)

# Memory Disaggregation

- **Need for memory disaggregation**

  **Low utilization**
  (as low as 30%)

  **Unbalanced usage**
  (> 70% of the time in clusters)

  **Energy consumption**
  (up to 46% of avg. system)

- **Timely problem**

  *Faster & harder to disaggregate* →

  | | Storage (NVMe SSD) | Network | Memory (DDR4) | L3 cache |
  |---|---|---|---|---|
  | **Data rate** | 5 GB/s (PCIe 4.0) | 25 GB/s (PCIe 4.0) | 25.6 GB/s (per channel) | 370 GB/s |
  | **Latency** | 45 / 450 us (R/W) | ~5 us (4KB) | <50 ns | 11 ns |

  snowflake
  amazon S3

  *Storage disaggregation is already popular*

# Memory Disaggregation Goals



[disaggregated blades]

# Memory Disaggregation Goals

- **Performance (CPU ↔ memory)**: high throughput & low latency



[disaggregated blades]

# Memory Disaggregation Goals

- **Performance (CPU ↔ memory)**: high throughput & low latency

- **Transparent elasticity**: flexible resource allocation



[disaggregated blades]

# Existing Approaches

# Existing Approaches



*Distributed shared memory (DSM)*

# Existing Approaches



*Distributed shared memory (DSM)*

# Existing Approaches



*Distributed shared memory (DSM)*

*Recent disaggregated memory schemes*

# Existing Approaches



*Distributed shared memory (DSM)*

*Recent disaggregated memory schemes*

# Key Insight — In-network Memory management

- **Central location:** global view & processing directly in the data path



Compute blades
(disaggregated)

Memory blades
(disaggregated)

# Key Insight — In-network Memory management

- **Central location:** global view & processing directly in the data path

- **Programmable switching ASIC:** flexible processing at line rate



Programmable &
running @line-rate

Compute blades
(disaggregated)

Memory blades
(disaggregated)

# Key Insight — In-network Memory management

- **Central location:** global view & processing directly in the data path

- **Programmable switching ASIC:** flexible processing at line rate

- **Similarity between network functions and memory management**

| Networking | Memory management |
|---|---|
| • IP forwarding | • Address translation |
| • IP assignment | • Memory allocation |
| • Access control | • Memory protection |
| • Multicast/broadcast | • Cache invalidations |

# MIND: In-network Memory Management

# MIND: In-network Memory Management

# MIND: In-network Memory Management

# MIND: In-network Memory Management

# MIND: In-network Memory Management

# Challenges in Building In-network MMU

- **Limited amount of in-network resources**
  - *Limited size of in-network memory*
    - Tens of MB → not sufficient to store metadata in a traditional way
      - *E.g., page table: 4 KB for 4GB of memory → 1M entries*

# Challenges in Building In-network MMU

- **Limited amount of in-network resources**

  - *Limited size of in-network memory*
    - Tens of MB → not sufficient to store metadata in a traditional way
      - *E.g., page table: 4 KB for 4GB of memory → 1M entries*

  - *Limited computation capability*
    - Switching ASIC → not sufficient to directly port traditional MMU functions
      - *E.g., limited number of operations for line rate processing*
        *↔ complicated cache coherence protocol*

# 3 Principles for System Design

- **P1 - Decouple memory management functionalities**
  → *Each function has the data structure suitable to its purpose*

# 3 Principles for System Design

■ **P1 - Decouple memory management functionalities**
 → *Each function has the data structure*
 *suitable to its purpose*

■ **P2 - Leverage global view of network**
 → *Make better decisions for memory management*

# 3 Principles for System Design



- **P1 - Decouple memory management functionalities**
  → *Each function has the data structure suitable to its purpose*

- **P2 - Leverage global view of network**
  → *Make better decisions for memory management*

- **P3 - Exploit network-centric hardware primitives**
  → *Reuse network hardware highly optimized for network functions*

# MIND overview: challenges and solutions



Compute blades

Control plane

| Memory allocation | Permission assignment | Directory entry allocation |

Data plane

Cache directory → Memory protection → Address translation

Memory blades

Programmable switch

# MIND overview: challenges and solutions



New access

**CPU** ... **CPU**
Local DRAM as cache
NIC

**CPU** ... **CPU**
Local DRAM as cache
NIC

Compute blades

## Programmable switch

### Control plane
| Memory allocation | Permission assignment | Directory entry allocation |

### Data plane
| Cache directory | → | Memory protection | → | Address translation |

MEM MEM
MEM MEM
NIC

Memory blades

# MIND overview: challenges and solutions



New access

Compute blades

Programmable switch

**Control plane**
- Memory allocation
- Permission assignment
- Directory entry allocation

**Data plane**
- Cache directory
- Memory protection
- Address translation

coherence

Memory blades

CPU ... CPU
Local DRAM as cache
NIC

MEM MEM
NIC

# MIND overview: challenges and solutions

# MIND overview: challenges and solutions



New access

CPU · · · CPU

Local DRAM as cache

NIC

CPU · · · CPU

Local DRAM as cache

NIC

coherence

Compute blades

**Programmable switch**

**Control plane**

| Memory allocation | Permission assignment | Directory entry allocation |

**Data plane**

Cache directory → Memory protection → Address translation

MEM MEM

MEM MEM

NIC

Memory blades

# MIND overview: challenges and solutions



New access

CPU ... CPU

Local DRAM as cache

NIC

CPU ... CPU

Local DRAM as cache

NIC

Compute blades

**Control plane**

Memory allocation

Permission assignment

Directory entry allocation

**Data plane**

Cache directory

Memory protection

Address translation

Programmable switch

MEM MEM

MEM MEM

NIC

Memory blades

# MIND overview: challenges and solutions



New access

CPU ... CPU

Local DRAM as cache

NIC

CPU ... CPU

Local DRAM as cache

NIC

Compute blades

**Control plane**

| Memory allocation | Permission assignment | Directory entry allocation |

**Data plane**

| Cache directory | Memory protection | Address translation |

**4KB page: too many entries**
**Huge page: load balancing**

Programmable switch

MEM MEM

MEM MEM

NIC

Memory blades

# MIND overview: challenges and solutions



New access

CPU … CPU

Local DRAM as cache

NIC

CPU … CPU

Local DRAM as cache

NIC

Compute blades

Control

Memory allocation

Permission assignment

Data p

Cache directory

Memory protection

Address translation

- **Static huge partition**
- **Fine-grained range-based translation**

4KB page: too many entries
Huge page: load balancing

MEM MEM

MEM MEM

NIC

Memory blades

Programmable switch

# MIND overview: challenges and solutions



New access

CPU ··· CPU
Local DRAM as cache
NIC

CPU ··· CPU
Local DRAM as cache
NIC

Compute blades

**Control plane**

| Memory allocation | Permission assignment | Directory entry allocation |

**Data plane**

Cache directory → Memory protection → Address translation

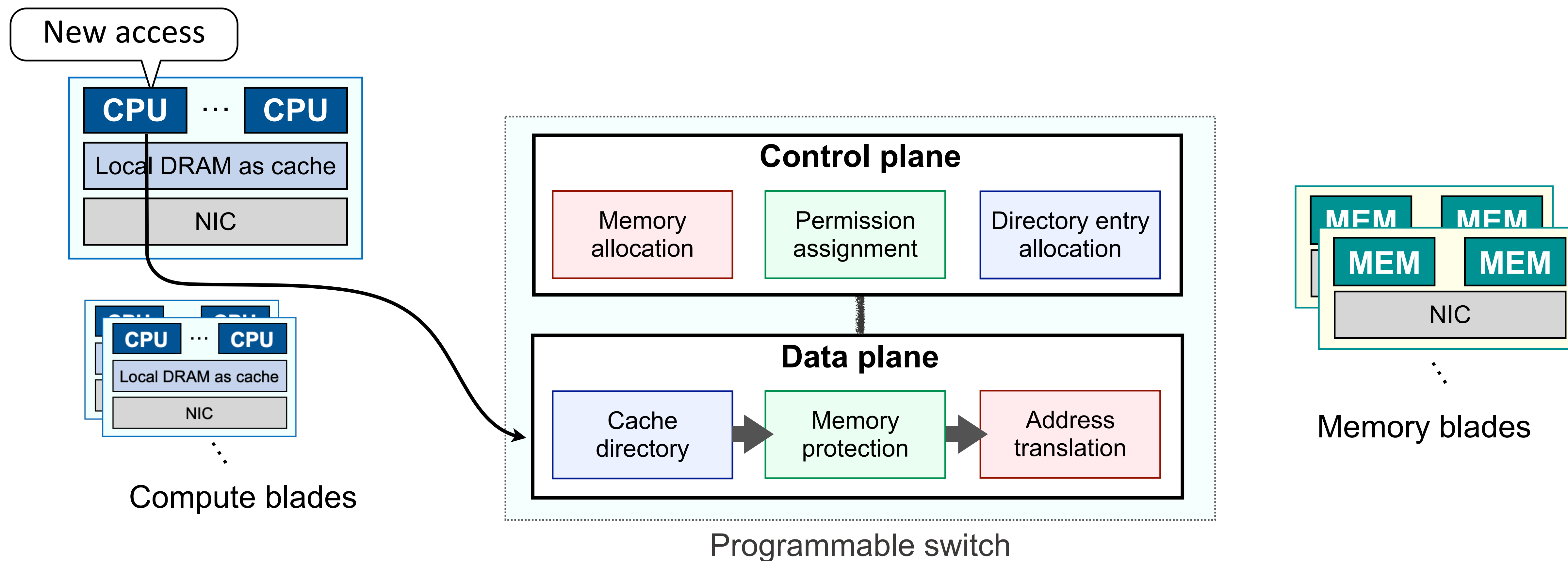Programmable switch

MEM MEM
MEM MEM
NIC

Memory blades

# MIND overview: challenges and solutions



New access

CPU ··· CPU

Local DRAM as cache

NIC

CPU ··· CPU

Local DRAM as cache

NIC

Compute blades

**Control plane**

| Memory allocation | Permission assignment | Directory entry allocation |

**4KB page: too many entries**
**Huge page: inflexible protection**

| Cache directory | Memory protection | Address translation |

Programmable switch

MEM   MEM

MEM   MEM
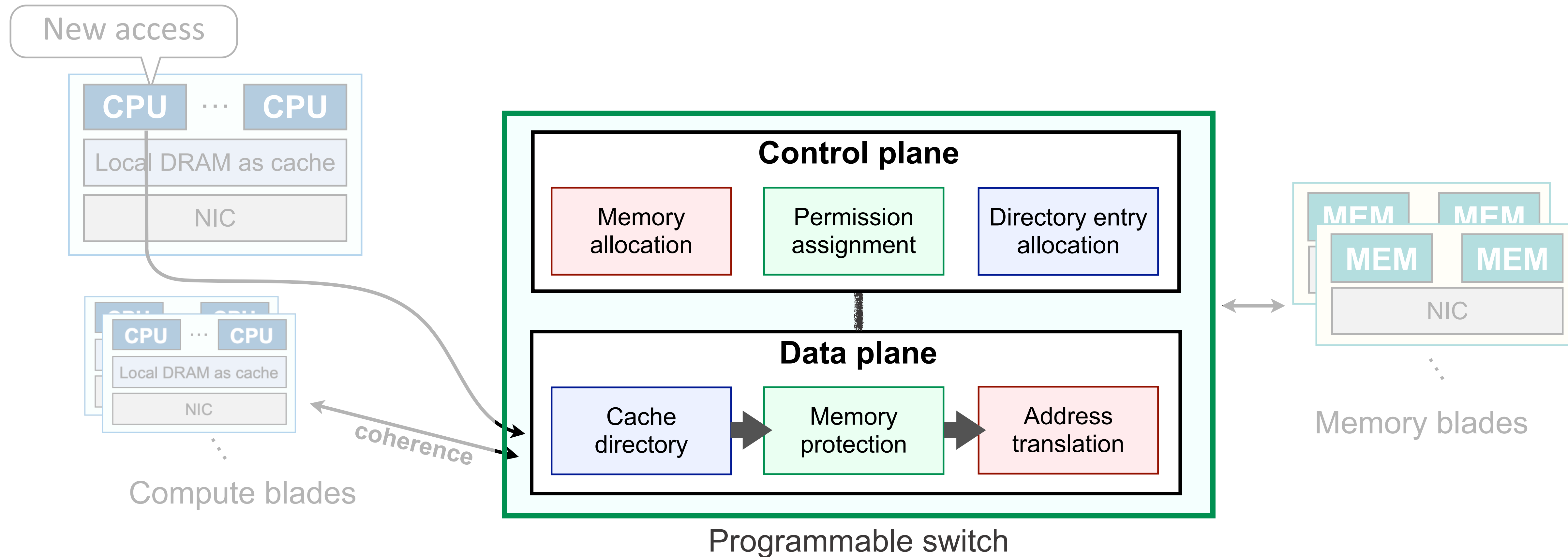
NIC

Memory blades

# MIND overview: challenges and solutions



New access

CPU ... CPU

Local DRAM as cache

NIC

CPU ... CPU

Local DRAM as cache

NIC

Compute blades

**Control plane**

Memory allocation

**VMA-granularity protection**

ry entry allocation

assignment

4KB page: too many entries
Huge page: inflexible protection

Cache directory

Memory protection

Address translation

Programmable switch

MEM MEM

MEM MEM

NIC

Memory blades

# MIND overview: challenges and solutions



New access

CPU ... CPU
Local DRAM as cache
NIC

CPU ... CPU
Local DRAM as cache
NIC

Compute blades

**Control plane**

Memory allocation | Permission assignment | Directory entry allocation

**Data plane**

Cache directory → Memory protection → Address translation

Programmable switch

MEM MEM
MEM MEM
NIC

Memory blades

# MIND overview: challenges and solutions



New access

CPU ⋯ CPU

Local DRAM as cache

NIC

CPU ⋯ CPU

Local DRAM as cache

NIC

Compute blades

**Control plane**

Memory allocation

Permission assignment

Directory entry allocation

**Small cache line: too many entries**
**Large cache line: false sharing**

Cache directory

Memory protection

Address translation

Programmable switch

MEM MEM

MEM MEM

NIC

Memory blades

# MIND overview: challenges and solutions



New access

CPU ... CPU

Local DRAM as cache

NIC

CPU ... CPU

Local DRAM as cache

NIC

Compute blades

- **Access granularity**
- **Coherence granularity**
  ↳**Dynamic resizing**

*with theoretical bound of O(log₂) overhead*

**Small cache line: too many entries**
**Large cache line: false sharing**

allocation

Cache directory

Memory protection

Address translation

Directory entry allocation

Programmable switch

MEM MEM

MEM MEM

NIC

Memory blades

# MIND overview: challenges and solutions



New memory allocation

CPU ⋯ CPU

Local DRAM as cache

NIC

Compute blades

Control plane

Memory allocation

Permission assignment

Directory entry allocation

Data plane

Cache directory → Memory protection → Address translation

Programmable switch

MEM MEM MEM MEM

NIC

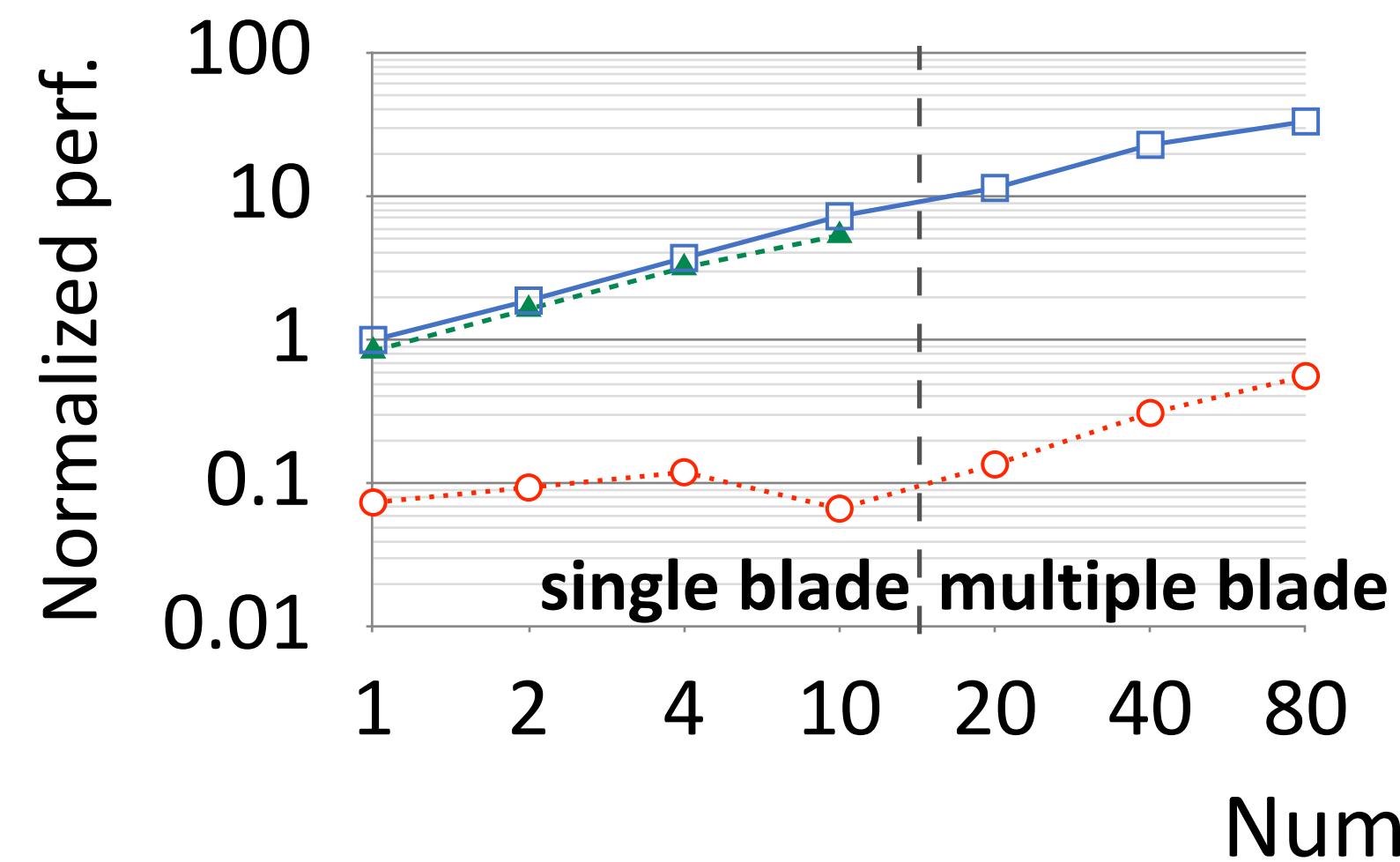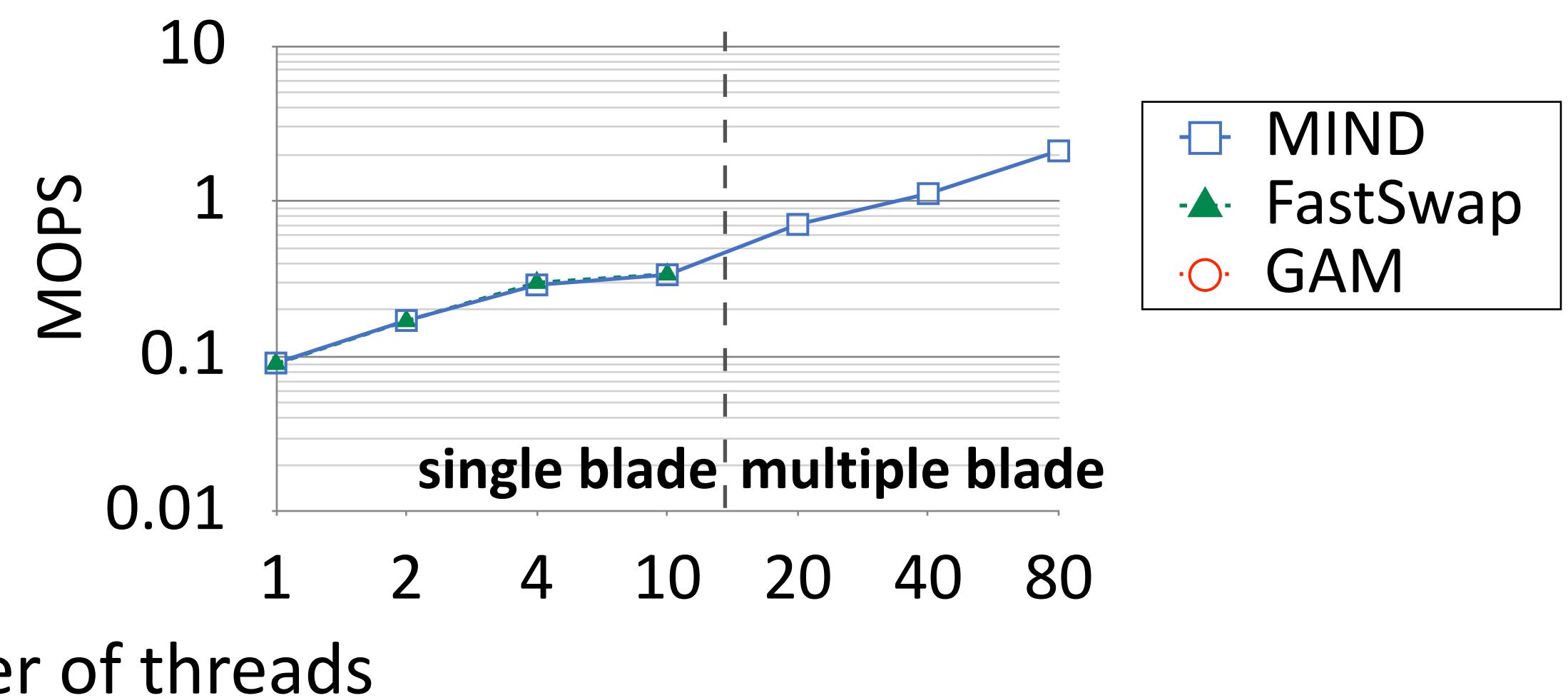Memory blades

# Performance Evaluation

- **Workloads with low contention**
  - TensorFlow (ResNet50), YCSB workload C (read only)
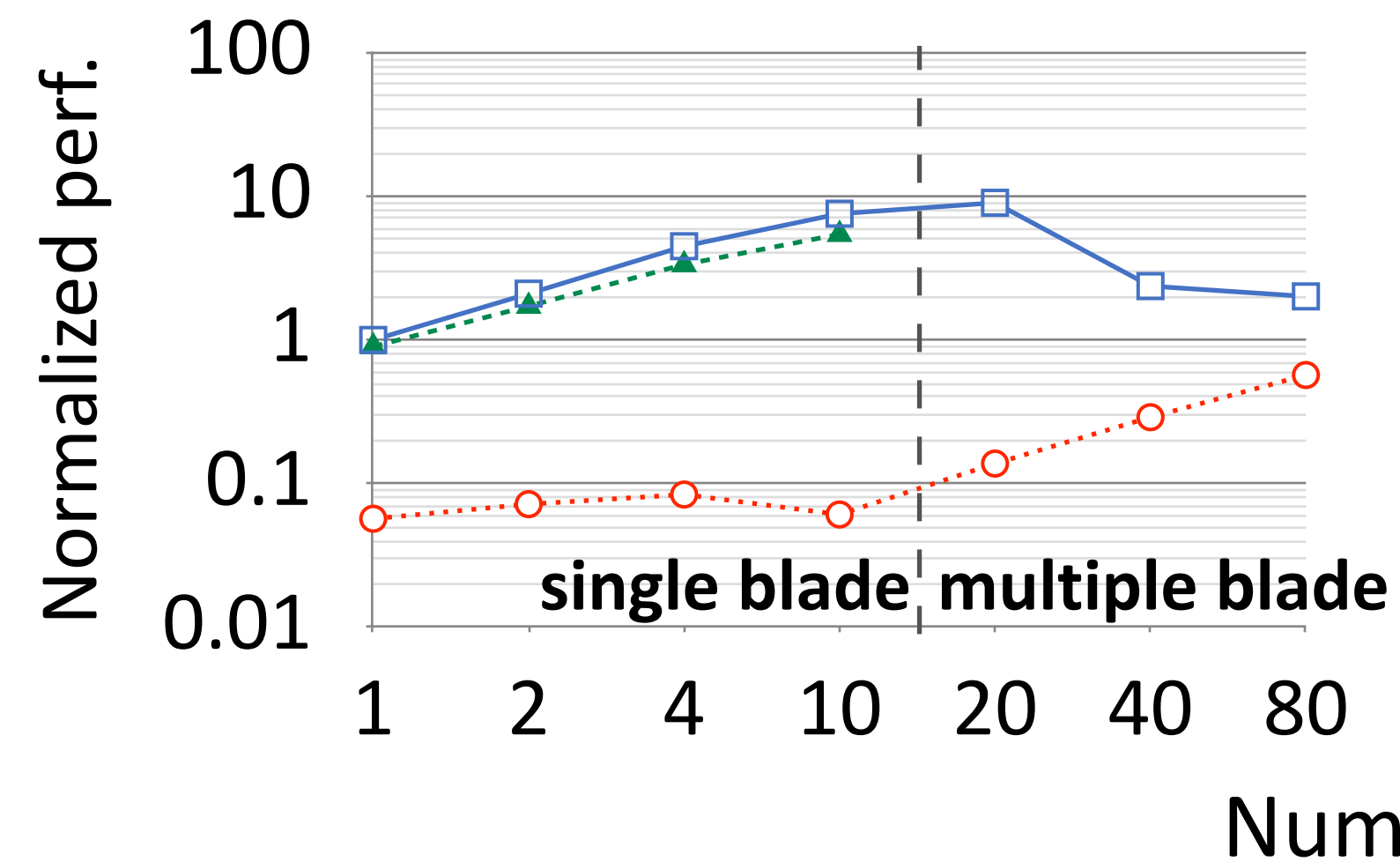
**TensorFlow** (trace-based for GAM)
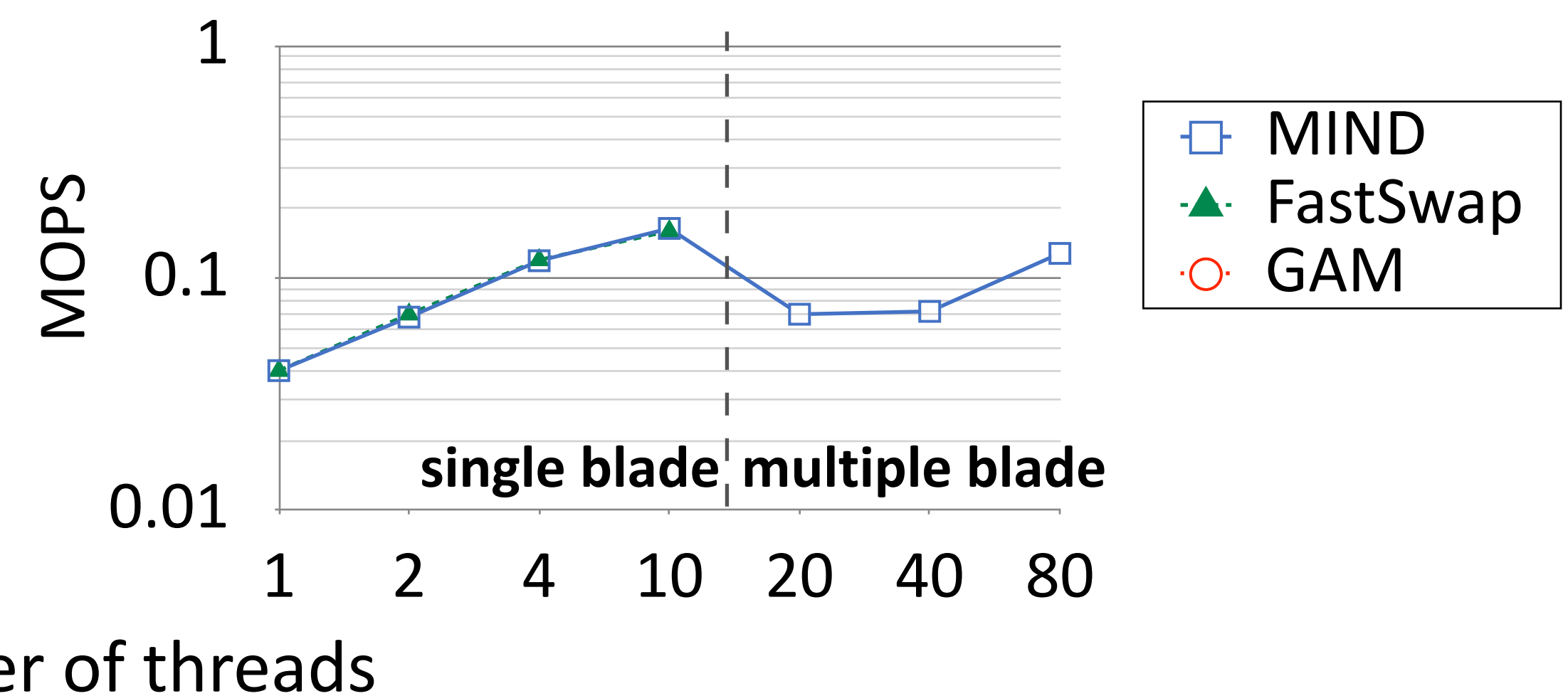
**YCSB workload C** (native)

# Performance Evaluation

- **Workloads with high contention**
  - GraphChi (PageRank), YCSB workload A (50 % update, 50 % read)

**GraphChi** (trace-based for GAM)

**YCSB workload A** (native)

# Conclusion & Summary

- **Trade-off between resource elasticity and performance**
  in memory disaggregation

# Conclusion & Summary

- **Trade-off between resource elasticity and performance** in memory disaggregation

- We designed MIND, an **in-network MMU** by leveraging programmable network

# Conclusion & Summary

- **Trade-off between resource elasticity and performance** in memory disaggregation

- We designed MIND, an **in-network MMU** by leveraging programmable network

- Our prototype of MIND can **match the performance of prior proposals and provide transparent elasticity**