# Active Access: A Mechanism for High-Performance Distributed Data-Centric Computations

**MACIEJ BESTA, TORSTEN HOEFLER**

# REMOTE MEMORY ACCESS (RMA) PROGRAMMING

# REMOTE MEMORY ACCESS (RMA) PROGRAMMING

**Process p**

Memory

A

# REMOTE MEMORY ACCESS (RMA) PROGRAMMING

**Process p**

Memory

A

**Process q**

Memory

B

# REMOTE MEMORY ACCESS (RMA) PROGRAMMING
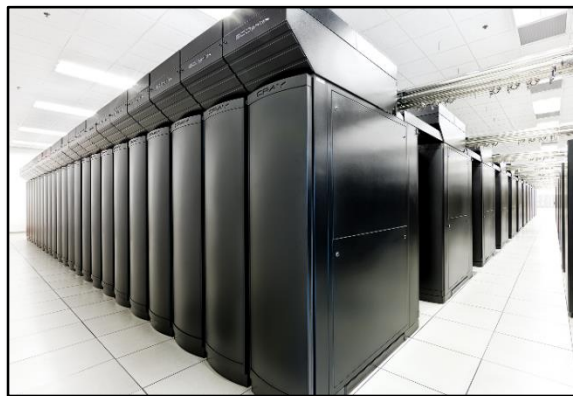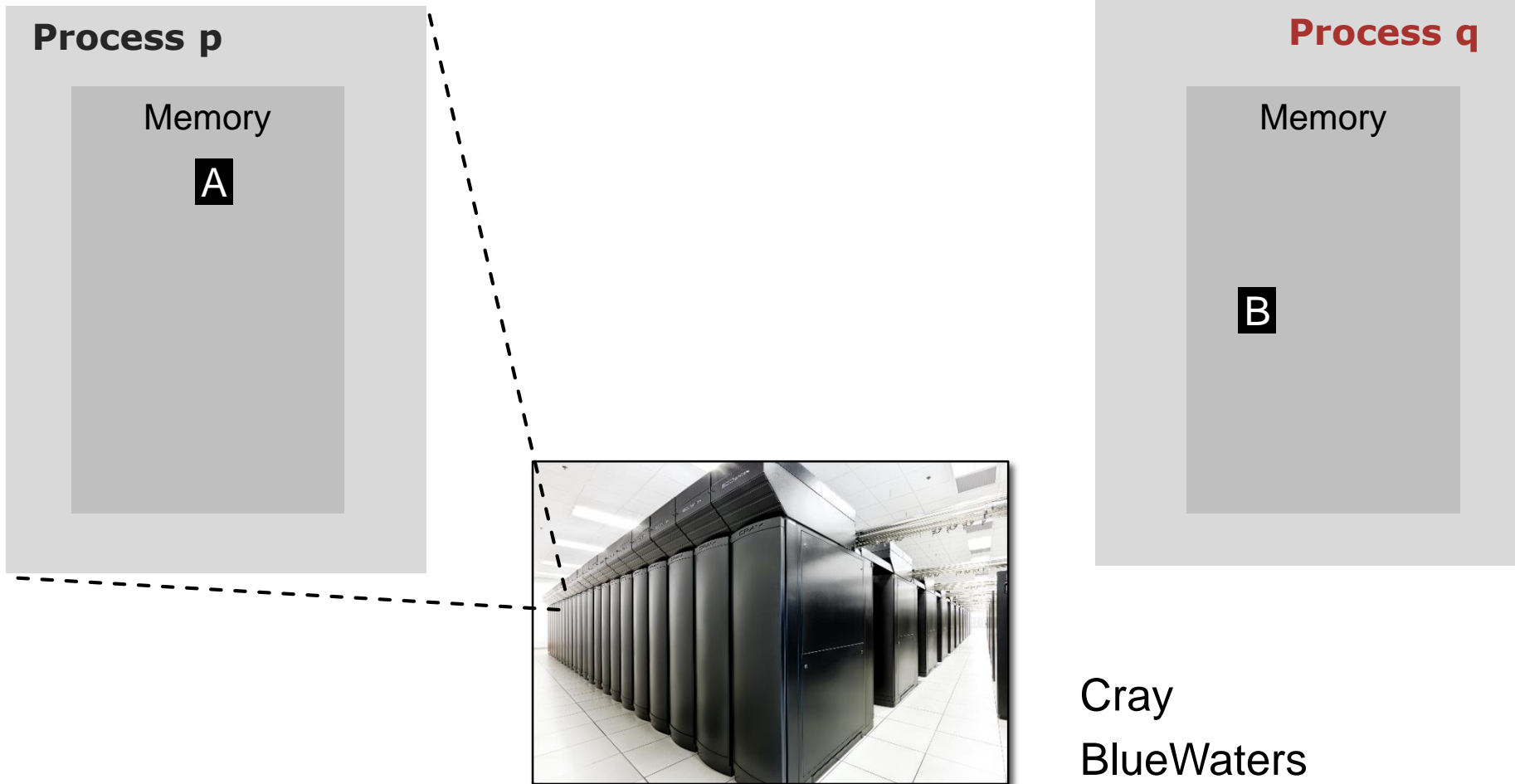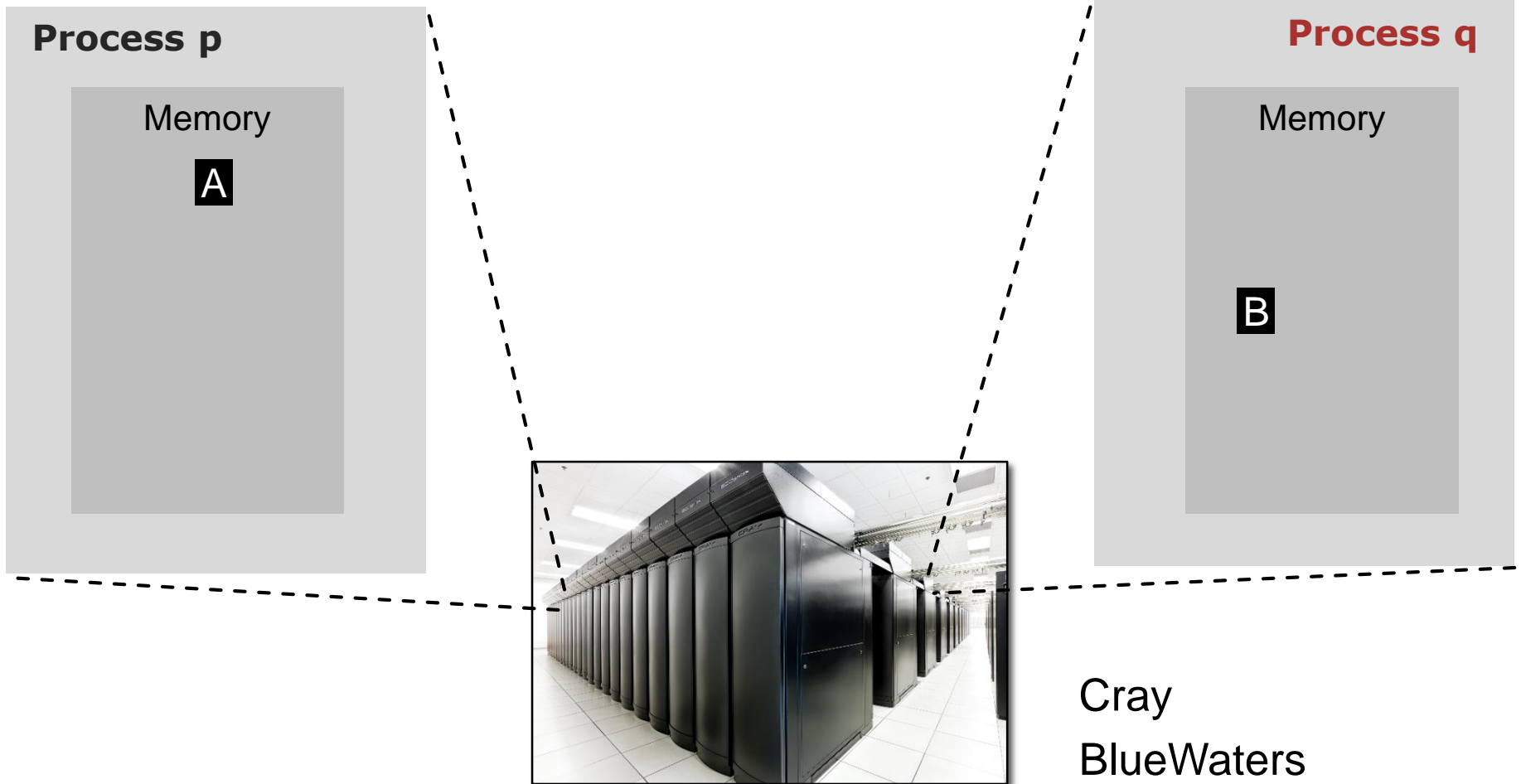
**Process p**

Memory

A

**Process q**

Memory

B



Cray
BlueWaters

# REMOTE MEMORY ACCESS (RMA) PROGRAMMING

**Process p**

Memory

A

**Process q**

Memory

B



Cray
BlueWaters

# REMOTE MEMORY ACCESS (RMA) PROGRAMMING
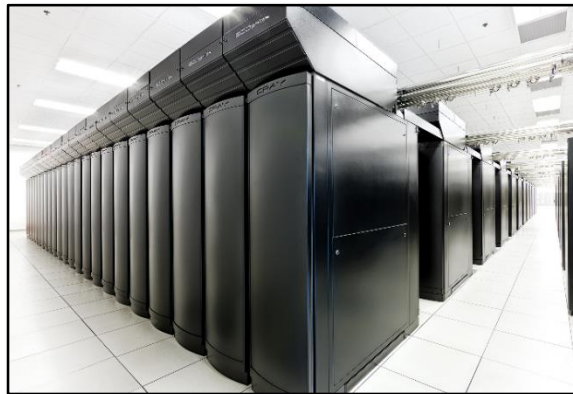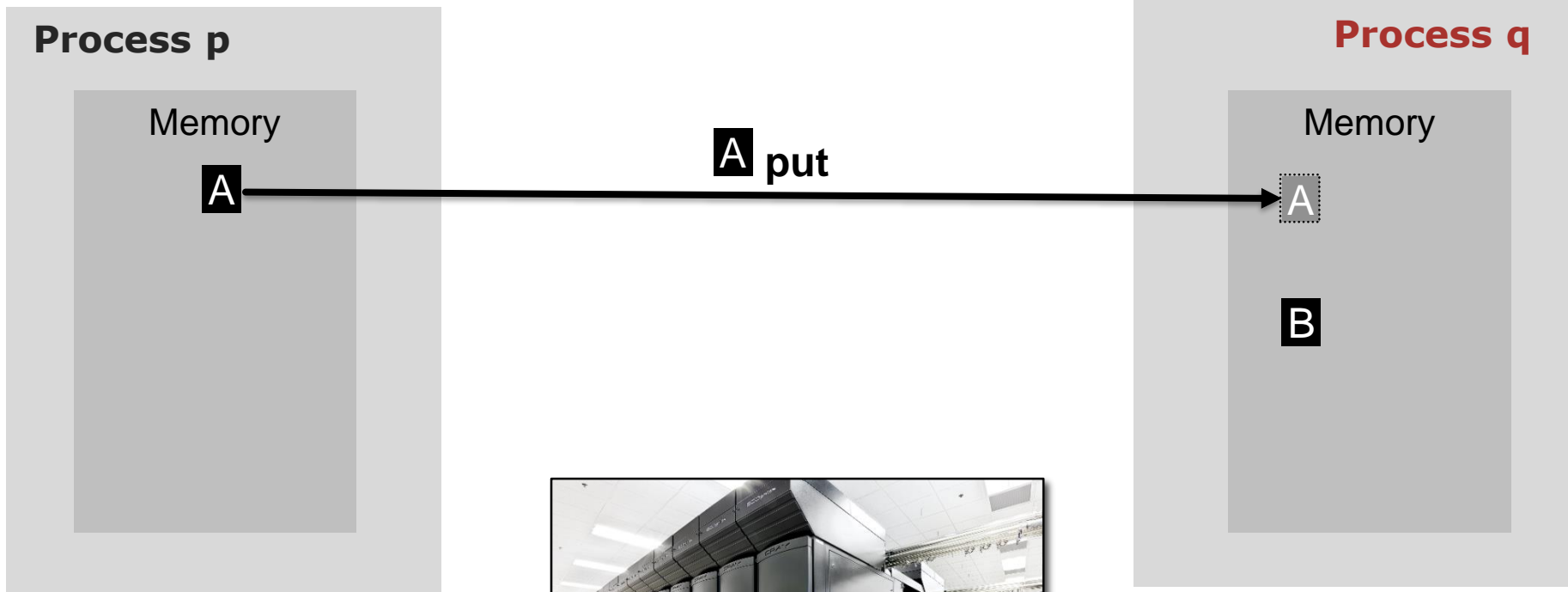
**Process p**

Memory

A

**Process q**

Memory

B

Cray
BlueWaters

# REMOTE MEMORY ACCESS (RMA) PROGRAMMING

**Process p**

Memory

A

**Process q**

Memory

B

Cray
BlueWaters

# REMOTE MEMORY ACCESS (RMA) PROGRAMMING

Process p

Memory

A

A **put**

Process q

Memory

A

B

Cray
BlueWaters

# REMOTE MEMORY ACCESS (RMA) PROGRAMMING



**Process p**

**Process q**

Memory

Memory

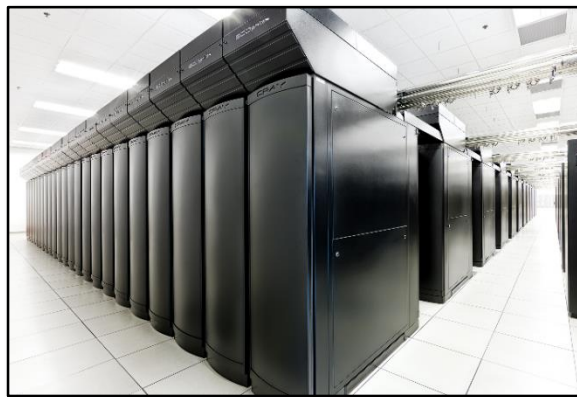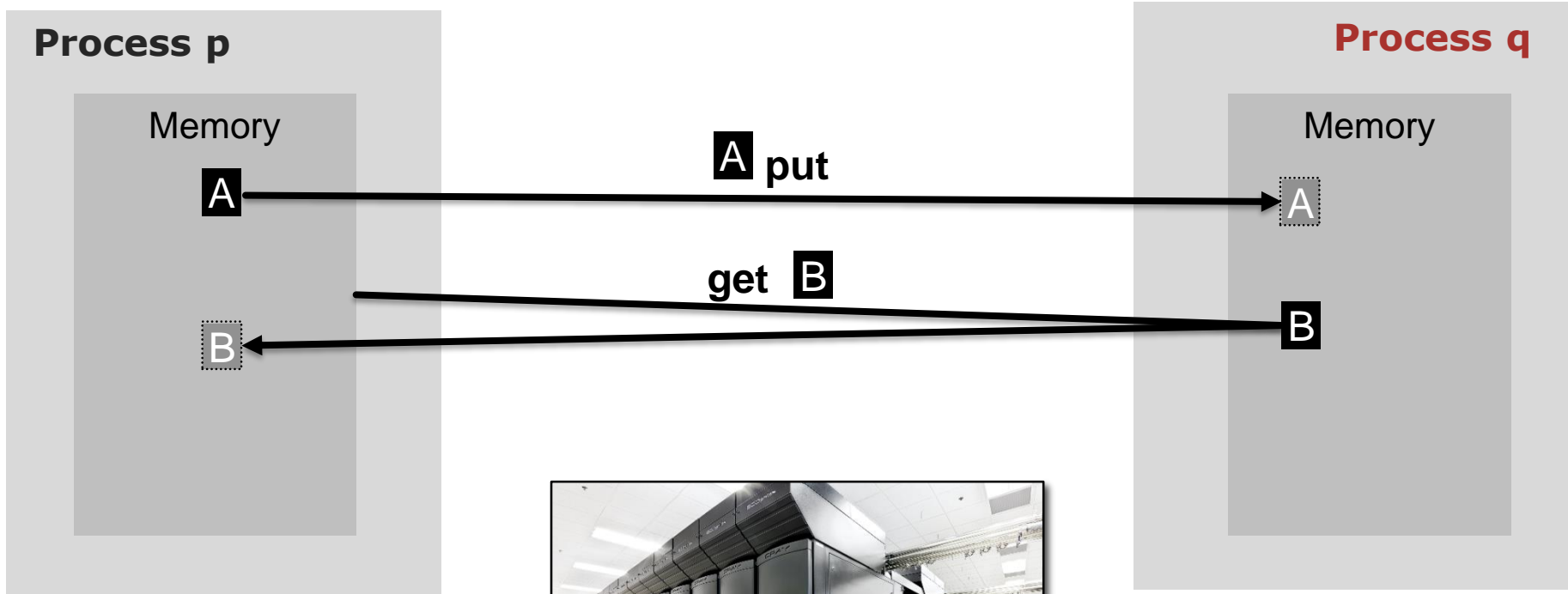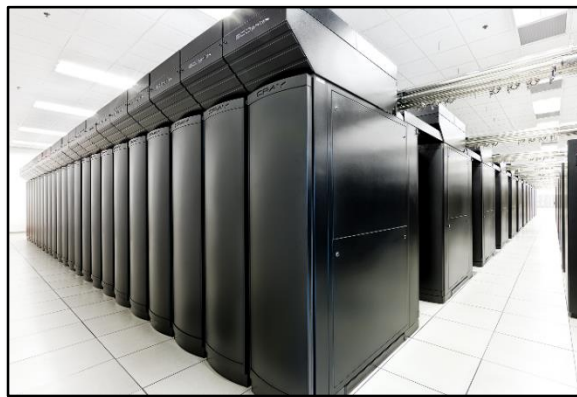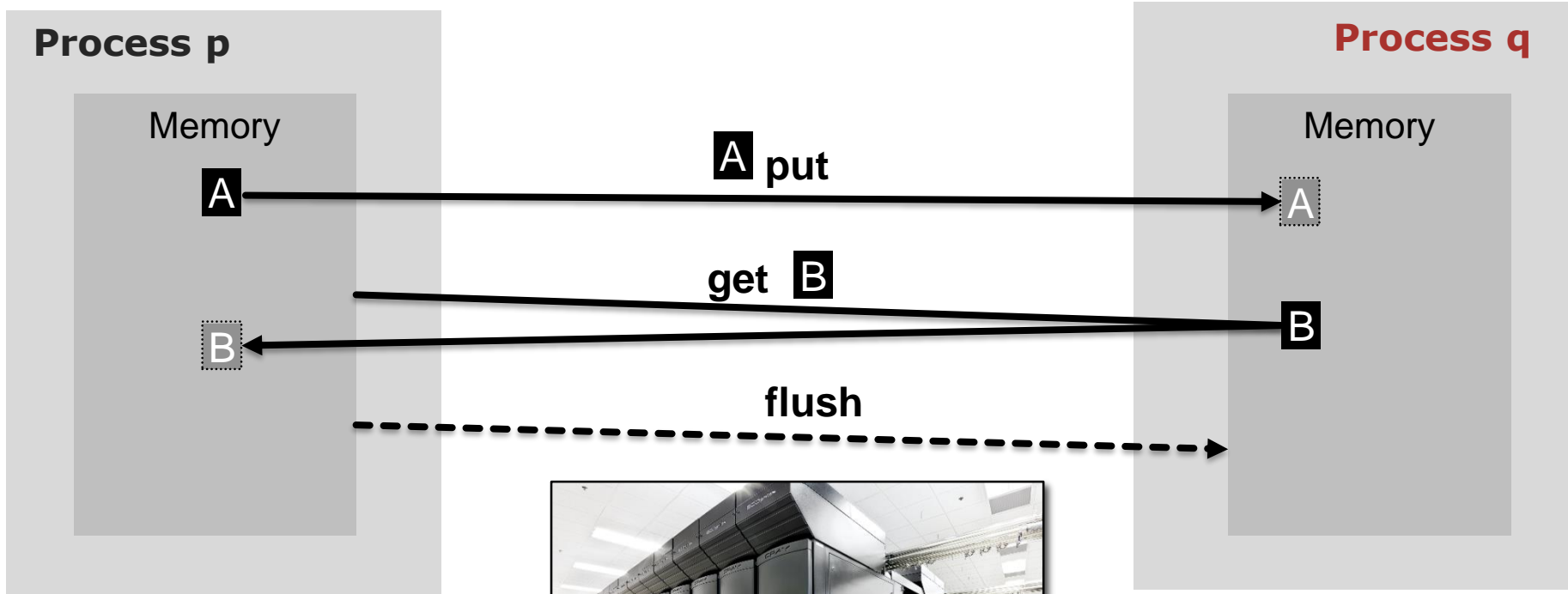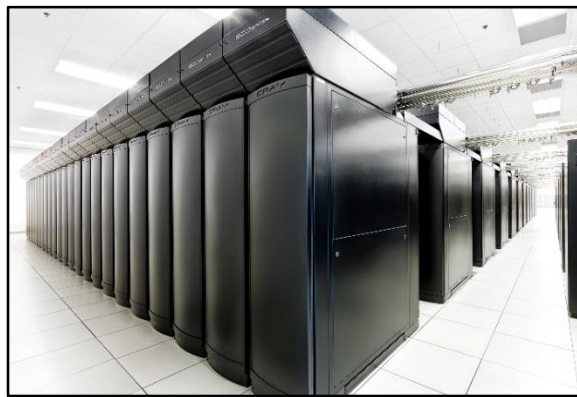A **put** A → A

**get** B

B B → B

Cray
BlueWaters

# REMOTE MEMORY ACCESS (RMA) PROGRAMMING



Cray
BlueWaters
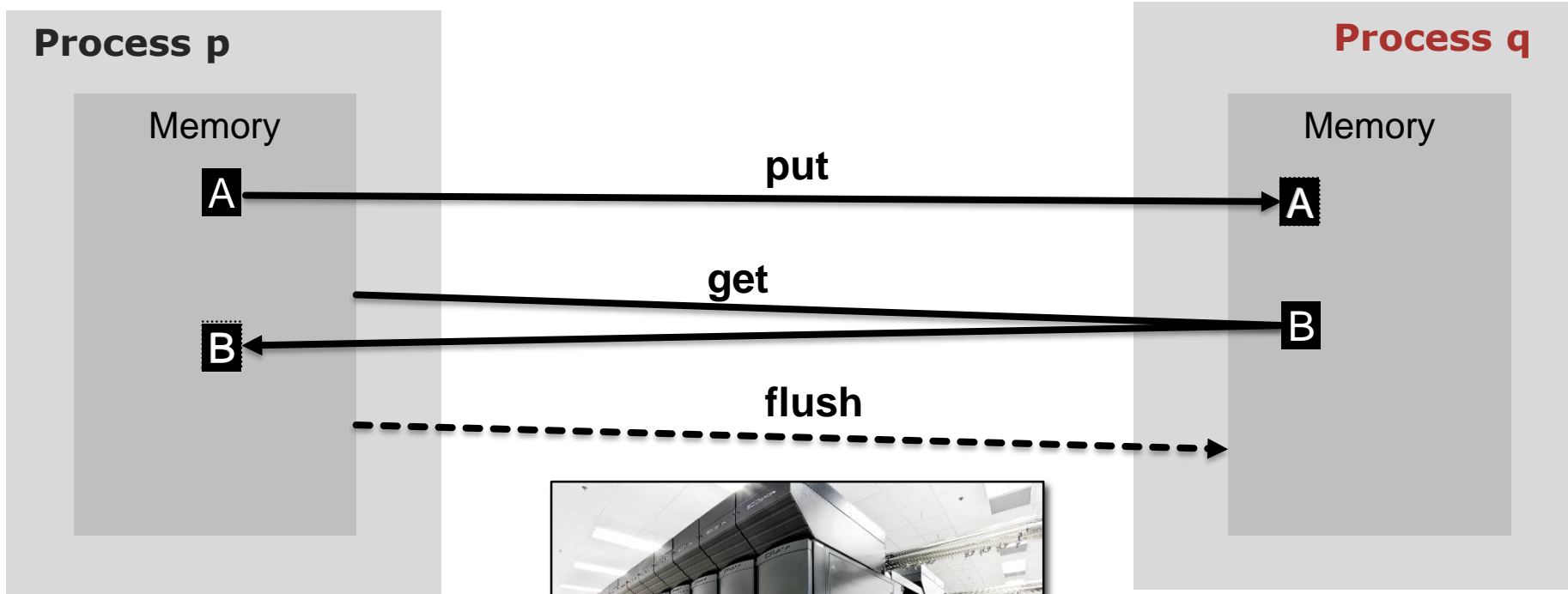
# REMOTE MEMORY ACCESS (RMA) PROGRAMMING



Process p

Process q

Memory

Memory

A → put → A

get

B → B

flush

Cray
BlueWaters

# REMOTE MEMORY ACCESS PROGRAMMING

- Implemented in hardware in NICs in the majority of HPC networks (RDMA)

# REMOTE MEMORY ACCESS PROGRAMMING

- Implemented in hardware in NICs in the majority of HPC networks (RDMA)

# REMOTE MEMORY ACCESS PROGRAMMING

- Implemented in hardware in NICs in the majority of HPC networks (RDMA)

# REMOTE MEMORY ACCESS PROGRAMMING

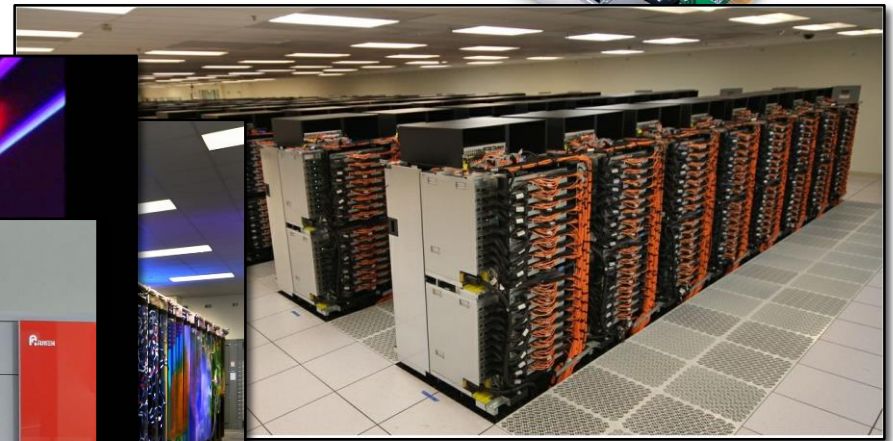- Implemented in hardware in NICs in the majority of HPC networks (RDMA)

# REMOTE MEMORY ACCESS PROGRAMMING

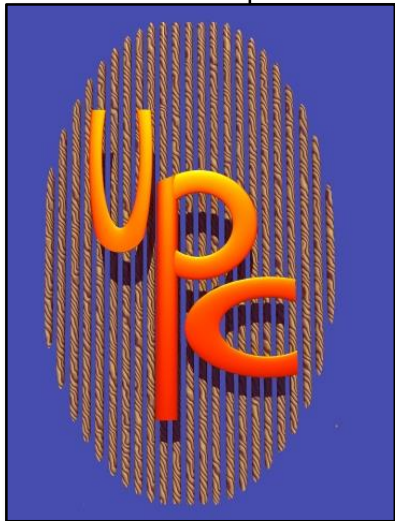- Implemented in hardware in NICs in the majority of HPC networks (RDMA)

# REMOTE MEMORY ACCESS PROGRAMMING

- Supported by many HPC libraries and languages

# REMOTE MEMORY ACCESS PROGRAMMING

- Supported by many HPC libraries and languages

# REMOTE MEMORY ACCESS PROGRAMMING

- Supported by many HPC libraries and languages
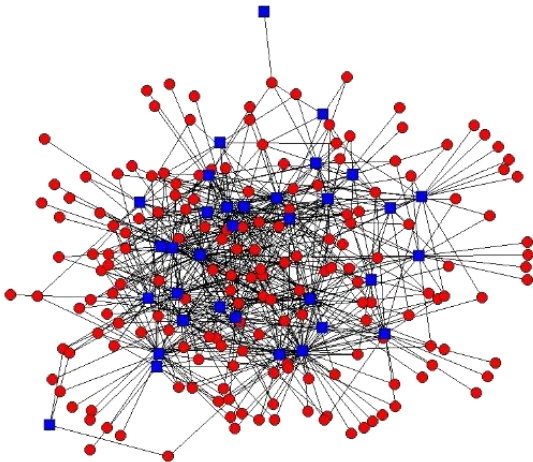
# REMOTE MEMORY ACCESS PROGRAMMING

- Enables significant speedups over message passing in many types of applications, e.g.:

[1] R. Gerstenberger et al. Enabling Highly-Scalable Remote Memory Access Programming with MPI-3 One-Sided. SC13
[2] D. Petrovic et al., High-performance RMA-based broadcast on the Intel SCC. SPAA'12

# REMOTE MEMORY ACCESS PROGRAMMING

- Enables significant speedups over message passing in many types of applications, e.g.:
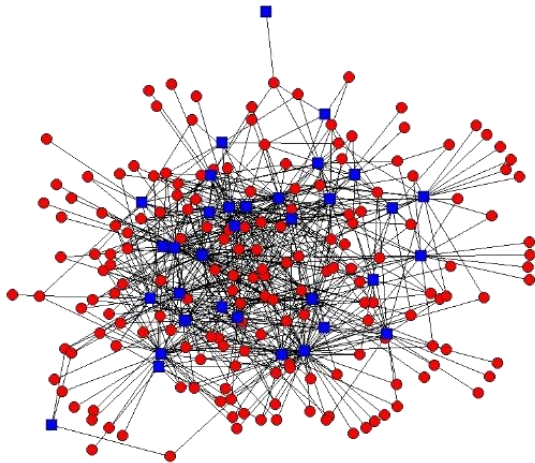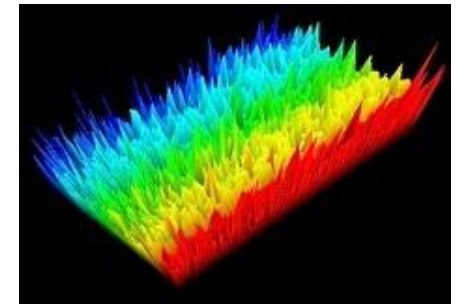  - Speedup of ~1.5 for communication patterns in irregular workloads



[1] R. Gerstenberger et al. Enabling Highly-Scalable Remote Memory Access Programming with MPI-3 One-Sided. SC13
[2] D. Petrovic et al., High-performance RMA-based broadcast on the Intel SCC. SPAA'12

# REMOTE MEMORY ACCESS PROGRAMMING

- Enables significant speedups over message passing in many types of applications, e.g.:
  - Speedup of ~1.5 for communication patterns in irregular workloads
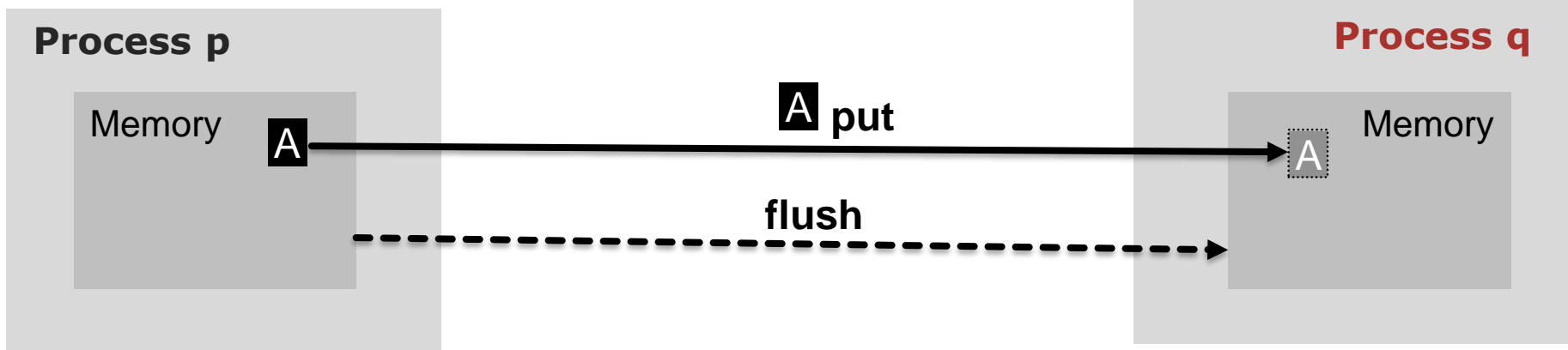  - Speedup of ~1.4-2 in physics computations

$$\frac{1}{\sqrt{2}}\left|\,\raisebox{-2pt}{🐱}\,\right\rangle + \frac{1}{\sqrt{2}}\left|\,\raisebox{-2pt}{🐱}\,\right\rangle$$

[1] R. Gerstenberger et al. Enabling Highly-Scalable Remote Memory Access Programming with MPI-3 One-Sided. SC13
[2] D. Petrovic et al., High-performance RMA-based broadcast on the Intel SCC. SPAA'12
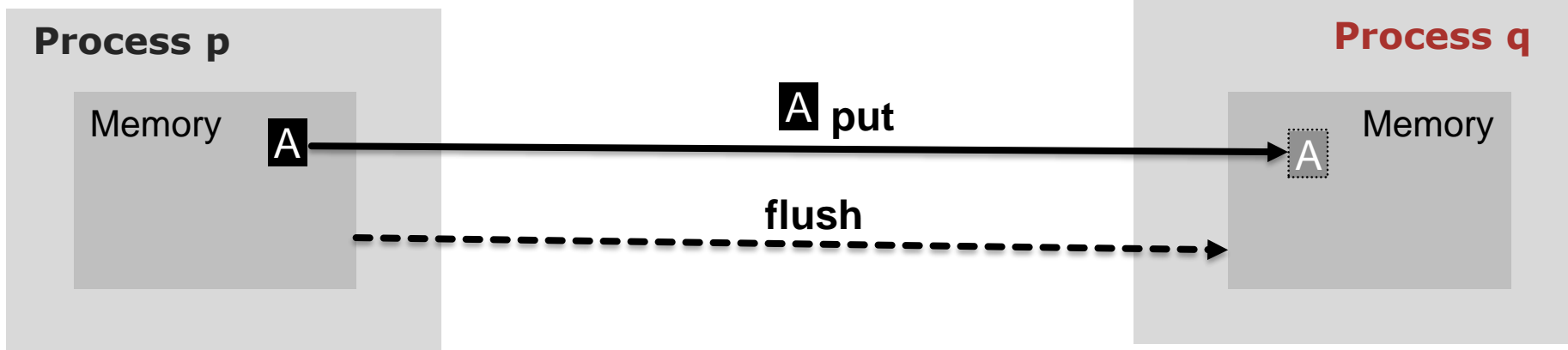
# RMA vs. Message Passing
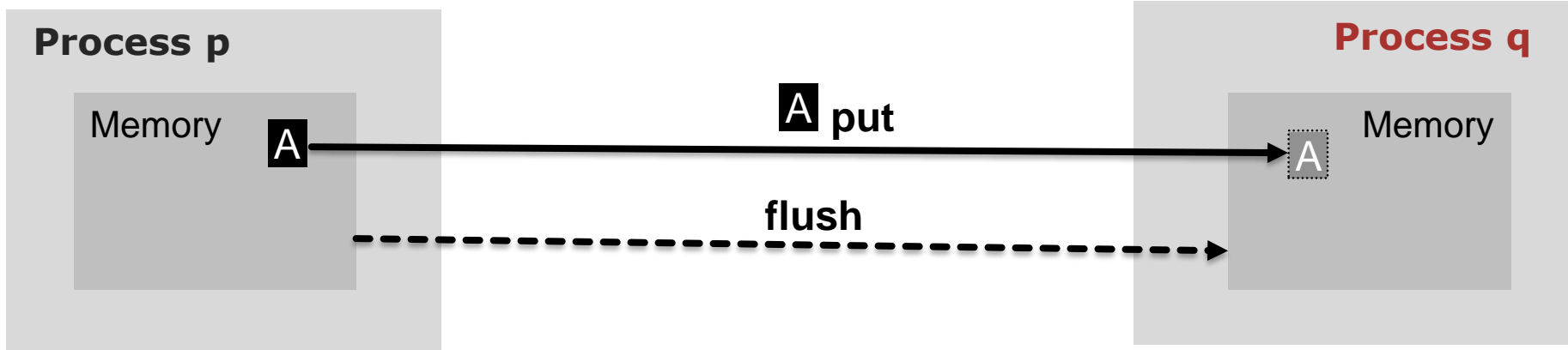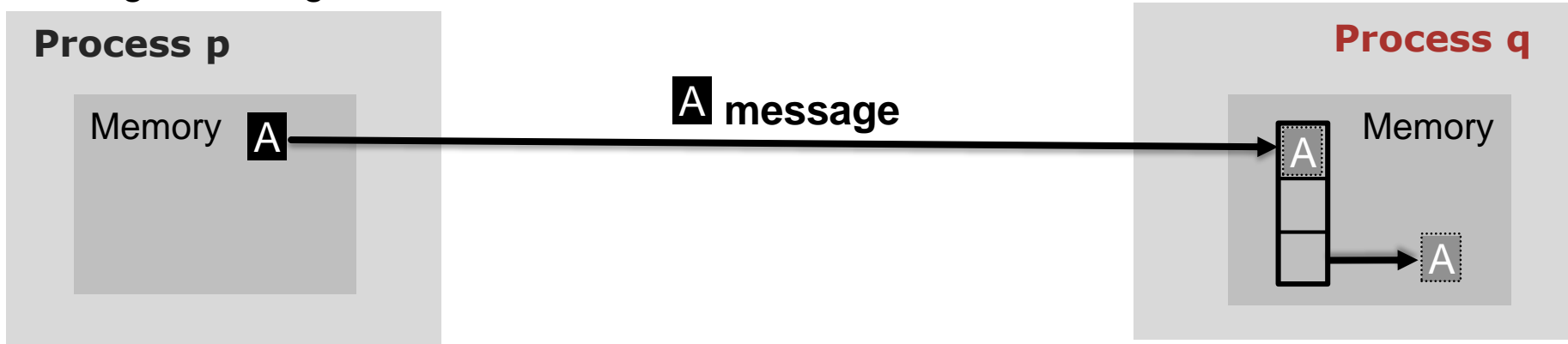
RMA:

# RMA vs. Message Passing
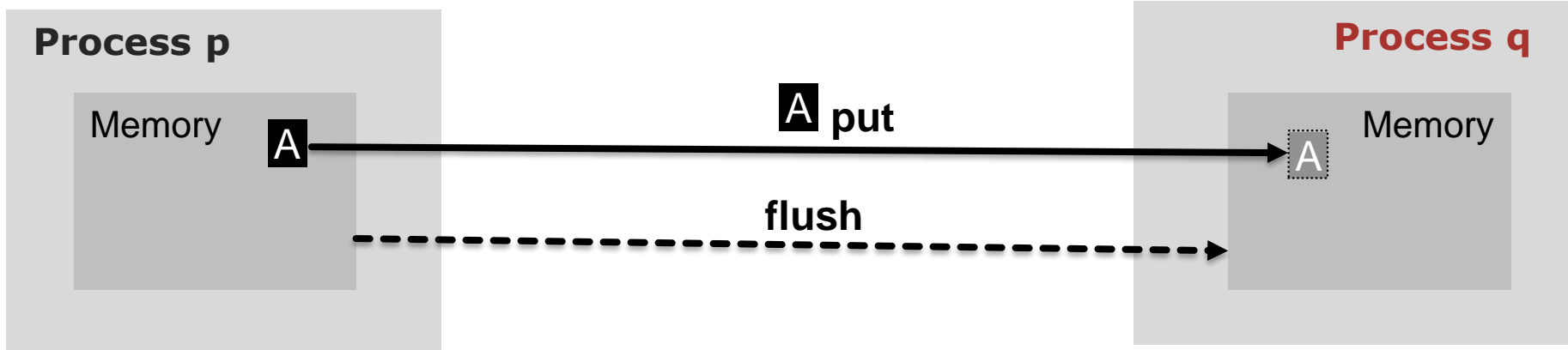
RMA:



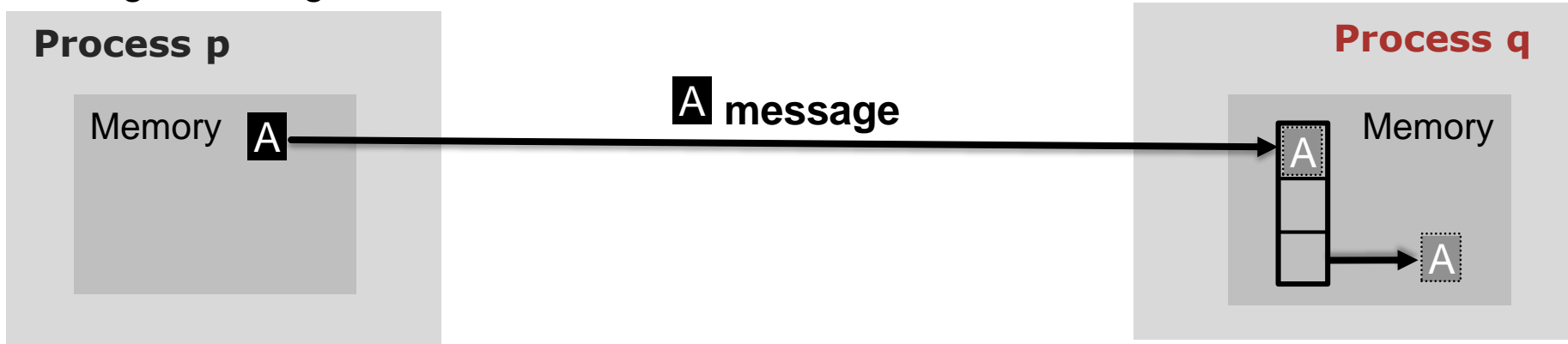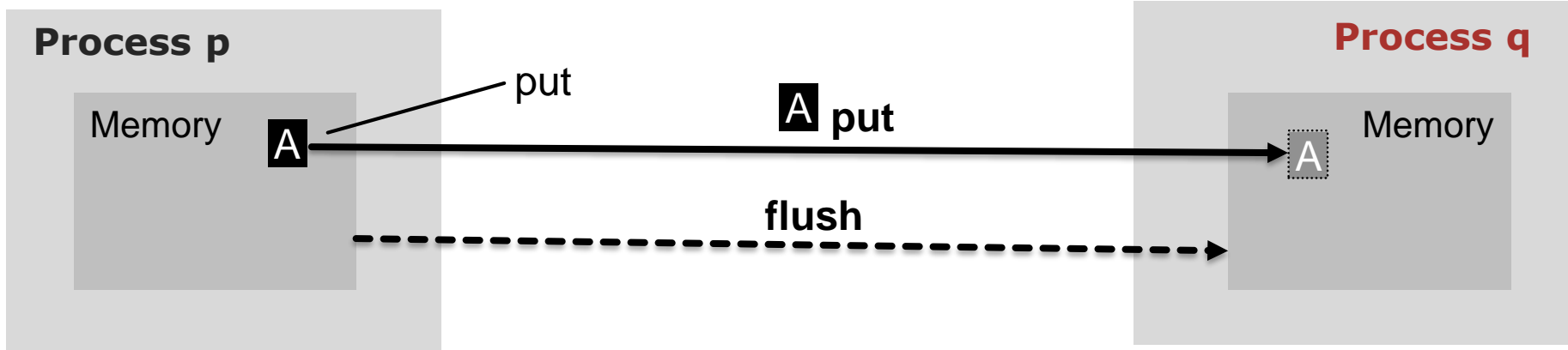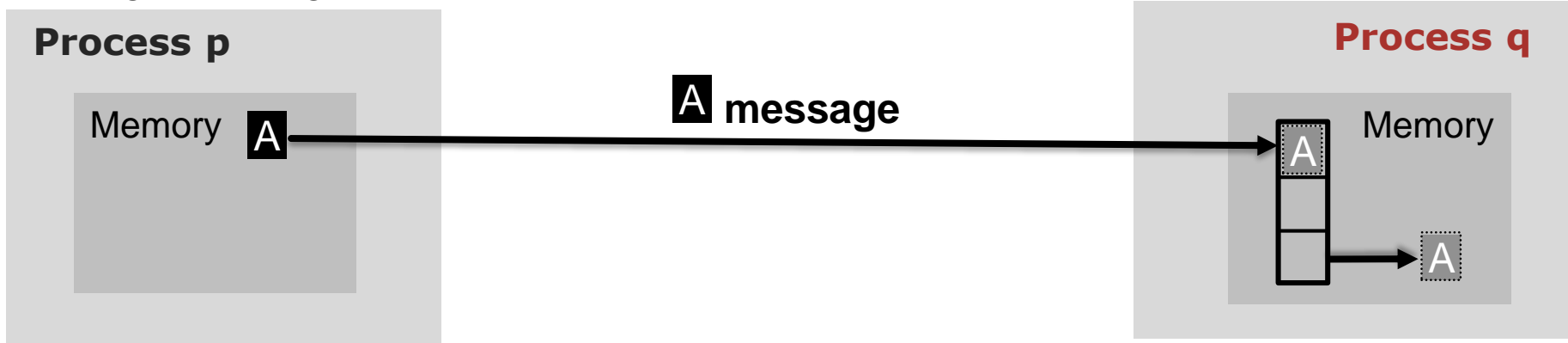Message Passing:

# RMA VS. MESSAGE PASSING

# RMA VS. MESSAGE PASSING

- Communication in RMA is one-sided

RMA:



Message Passing:

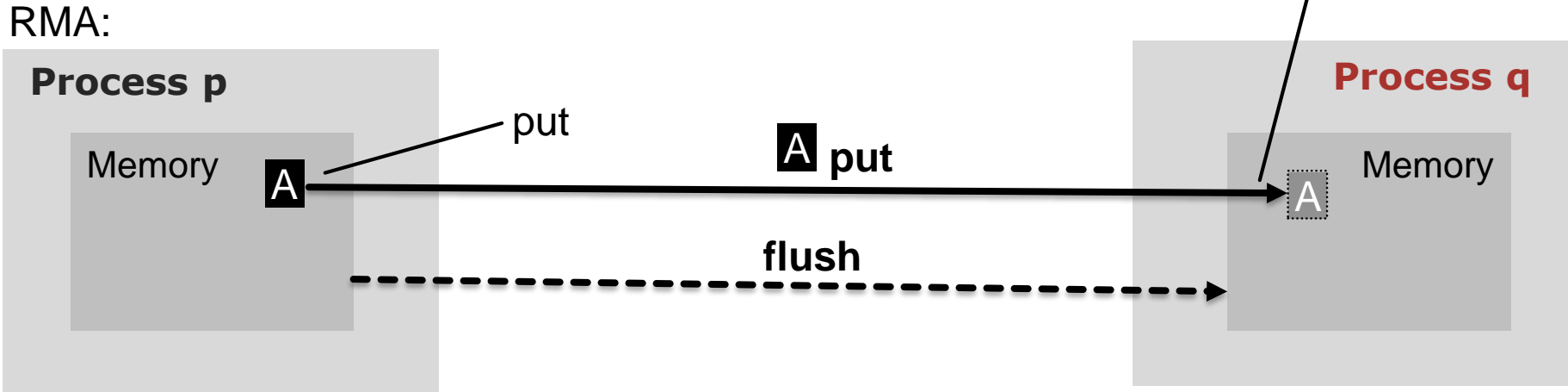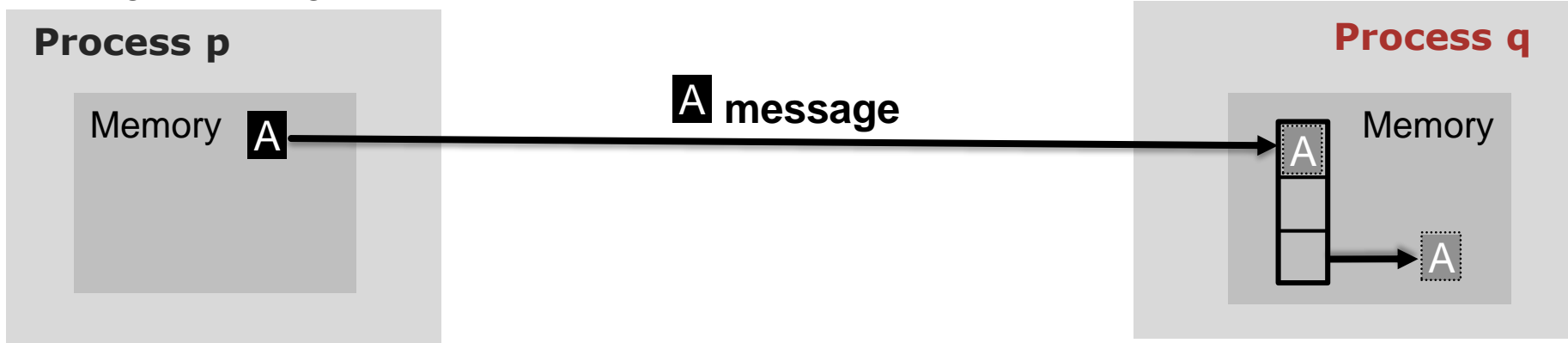# RMA VS. MESSAGE PASSING

▪ Communication in RMA is one-sided

RMA:



Message Passing:

# RMA VS. MESSAGE PASSING

- Communication in RMA is one-sided
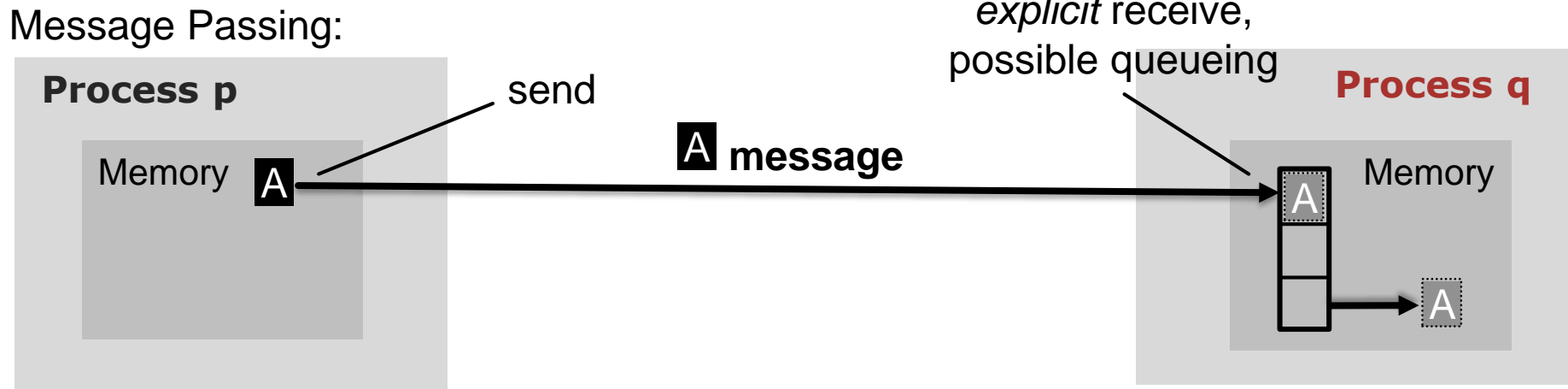
no active participation, direct access to memory

RMA:

**Process p**

Memory  A

put

**A** put

flush

**Process q**

Memory  A

Message Passing:

**Process p**

Memory  A

**A** message

**Process q**

Memory  A

A

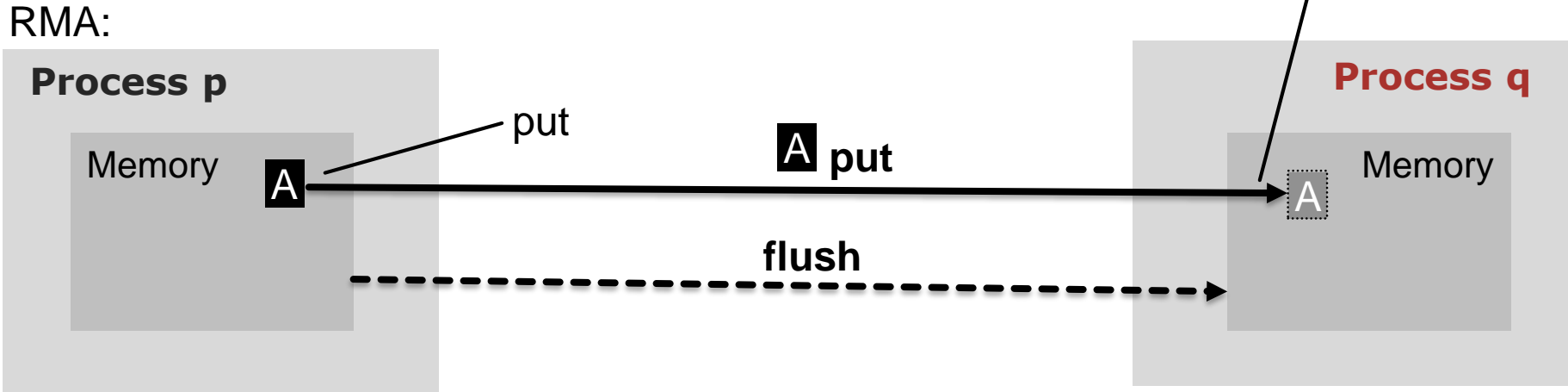# RMA VS. MESSAGE PASSING

- Communication in RMA is one-sided

no active participation, direct access to memory

# REMOTE MEMORY ACCESS PROGRAMMING

[1] R. Gerstenberger et al. Enabling Highly-Scalable Remote Memory Access Programming with MPI-3 One-Sided. SC13

# REMOTE MEMORY ACCESS PROGRAMMING

- Is it ideal?

[1] R. Gerstenberger et al. Enabling Highly-Scalable Remote Memory Access Programming with MPI-3 One-Sided. SC13

# REMOTE MEMORY ACCESS PROGRAMMING

- Is it ideal? **NO!**

[1] R. Gerstenberger et al. Enabling Highly-Scalable Remote Memory Access Programming with MPI-3 One-Sided. SC13

# REMOTE MEMORY ACCESS PROGRAMMING

- Is it ideal? **NO!**
- Consider an insert in a distributed hashtable...

[1] R. Gerstenberger et al. Enabling Highly-Scalable Remote Memory Access Programming with MPI-3 One-Sided. SC13

# REMOTE MEMORY ACCESS PROGRAMMING

- Is it ideal? **NO!**
- Consider an insert in a distributed hashtable...

| Proc p | | Proc q |
|--------|---|--------|

[1] R. Gerstenberger et al. Enabling Highly-Scalable Remote Memory Access Programming with MPI-3 One-Sided. SC13

# REMOTE MEMORY ACCESS PROGRAMMING

- Is it ideal? **NO!**
- Consider an insert in a distributed hashtable...

No hash collision:

Proc p

Proc q

[1] R. Gerstenberger et al. Enabling Highly-Scalable Remote Memory Access Programming with MPI-3 One-Sided. SC13

# REMOTE MEMORY ACCESS PROGRAMMING

- Is it ideal? **NO!**
- Consider an insert in a distributed hashtable...

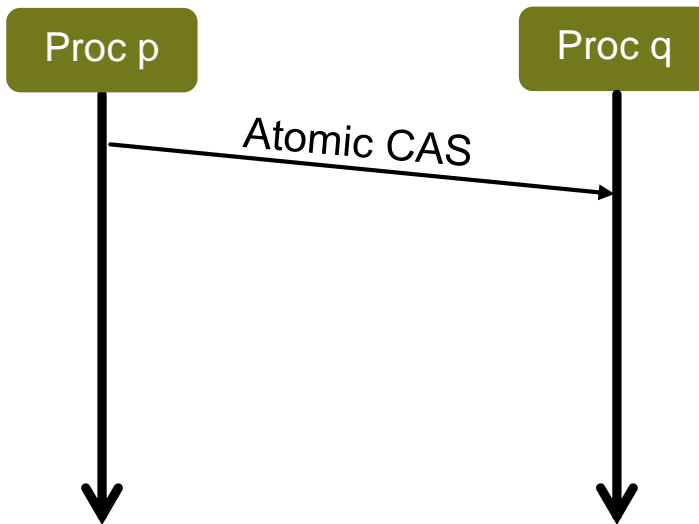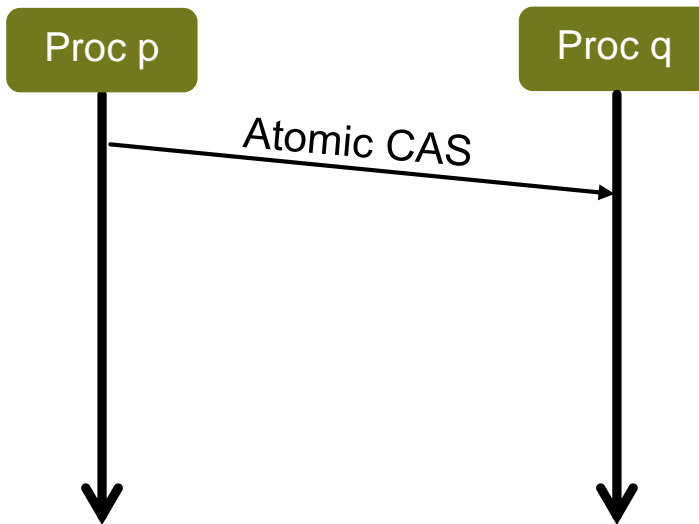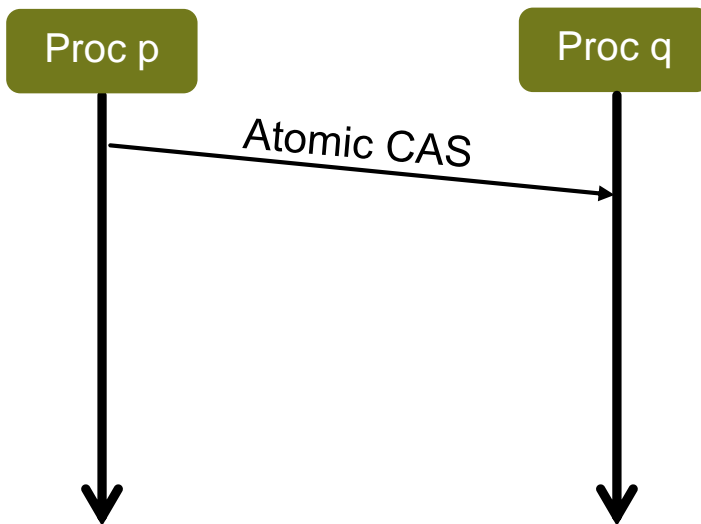No hash collision:

➔ 1 remote atomic
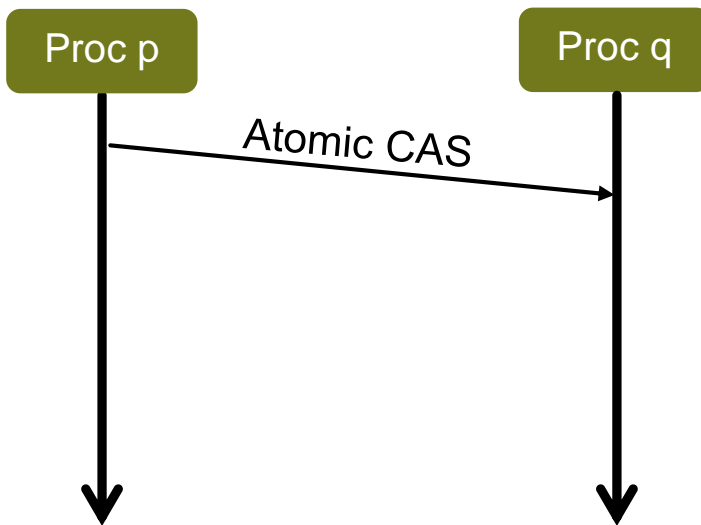
Proc p

Proc q

[1] R. Gerstenberger et al. Enabling Highly-Scalable Remote Memory Access Programming with MPI-3 One-Sided. SC13

# REMOTE MEMORY ACCESS PROGRAMMING

- Is it ideal? **NO!**
- Consider an insert in a distributed hashtable...

No hash collision:

→ 1 remote atomic



Proc p

Proc q

Atomic CAS

[1] R. Gerstenberger et al. Enabling Highly-Scalable Remote Memory Access Programming with MPI-3 One-Sided. SC13
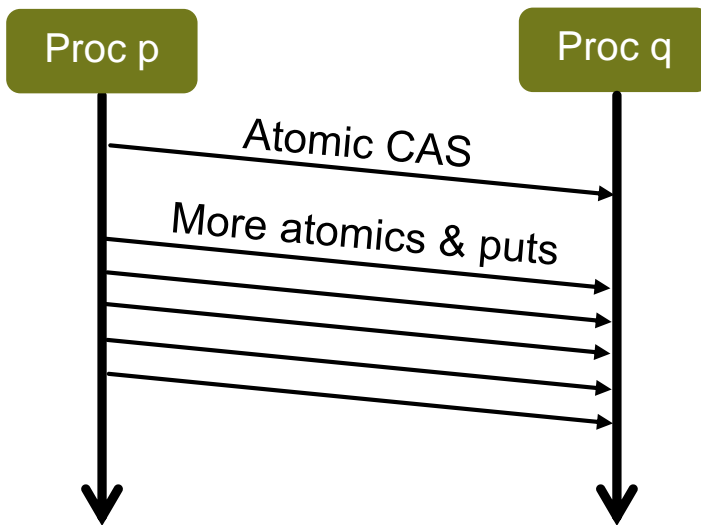
# REMOTE MEMORY ACCESS PROGRAMMING

- Is it ideal? **NO!**
- Consider an insert in a distributed hashtable...

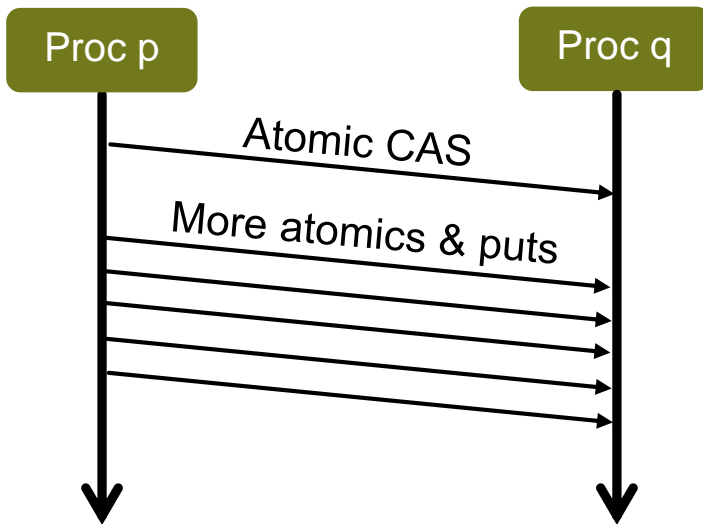No hash collision:

→ 1 remote atomic
→ Up to 5x speedup over MP [1]

Proc p          Proc q

Atomic CAS

[1] R. Gerstenberger et al. Enabling Highly-Scalable Remote Memory Access Programming with MPI-3 One-Sided. SC13

# REMOTE MEMORY ACCESS PROGRAMMING

- Is it ideal?  **NO!**
- Consider an insert in a distributed hashtable...

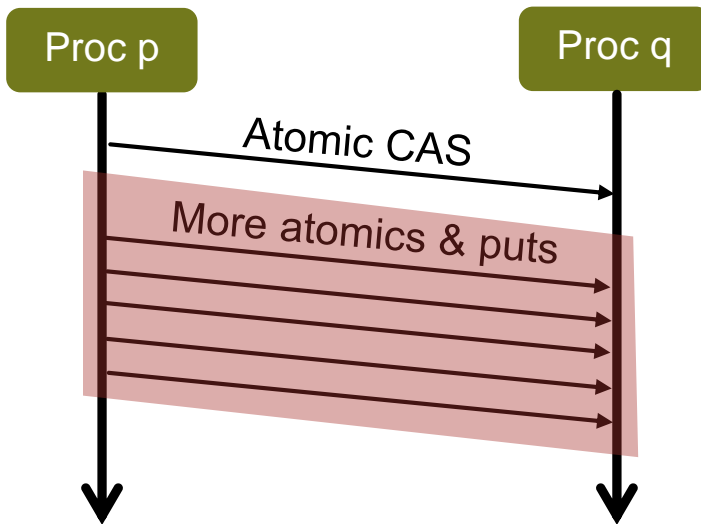No hash collision:

→ 1 remote atomic
→ Up to 5x speedup over MP [1]

A hash collision:

| Proc p | Proc q |
|--------|--------|

Atomic CAS

[1] R. Gerstenberger et al. Enabling Highly-Scalable Remote Memory Access Programming with MPI-3 One-Sided. SC13

# REMOTE MEMORY ACCESS PROGRAMMING
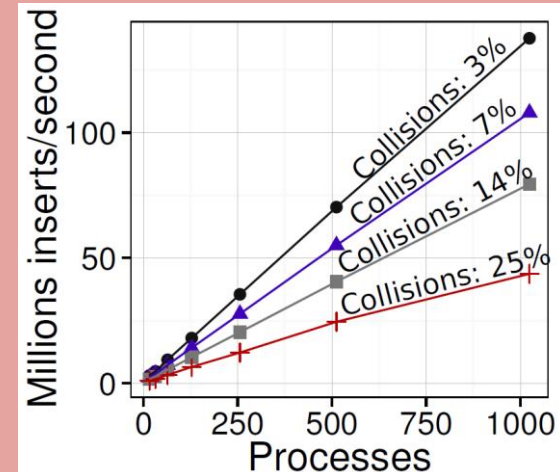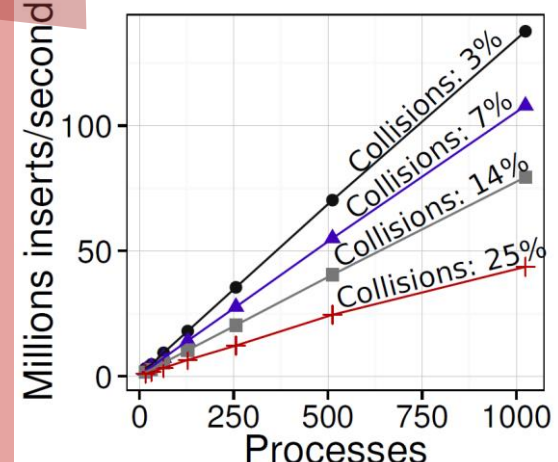
- Is it ideal? **NO!**
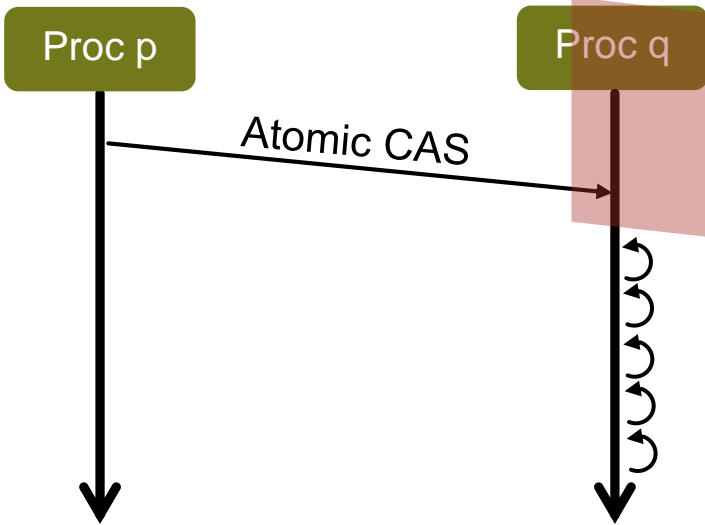- Consider an insert in a distributed hashtable...

No hash collision:

→ 1 remote atomic
→ Up to 5x speedup over MP [1]

A hash collision:

→ 4 remote atomics + 2 remote puts

Proc p        Proc q

*Atomic CAS*

[1] R. Gerstenberger et al. Enabling Highly-Scalable Remote Memory Access Programming with MPI-3 One-Sided. SC13
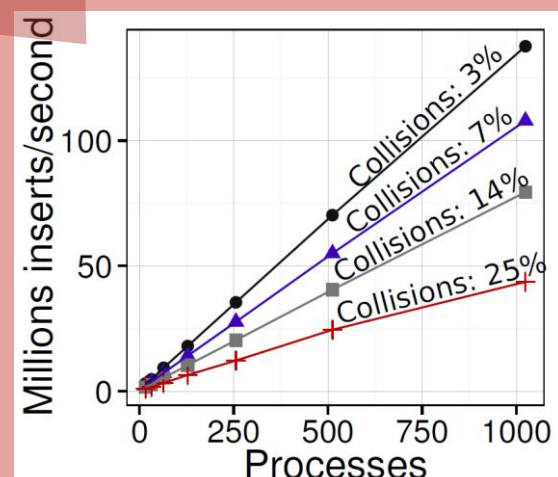
# REMOTE MEMORY ACCESS PROGRAMMING

- Is it ideal? **NO!**
- Consider an insert in a distributed hashtable...

**No hash collision:**

→ 1 remote atomic
→ Up to 5x speedup over MP [1]



Proc p    Proc q

Atomic CAS

More atomics & puts

**A hash collision:**

→ 4 remote atomics + 2 remote puts

[1] R. Gerstenberger et al. Enabling Highly-Scalable Remote Memory Access Programming with MPI-3 One-Sided. SC13

# REMOTE MEMORY ACCESS PROGRAMMING

- Is it ideal? **NO!**
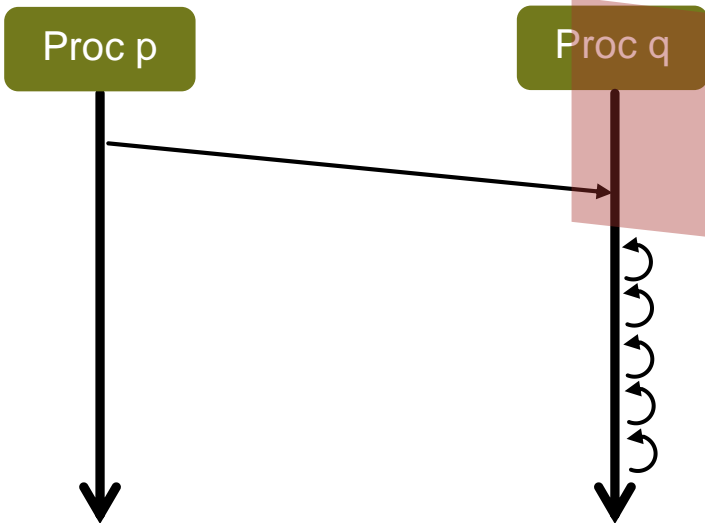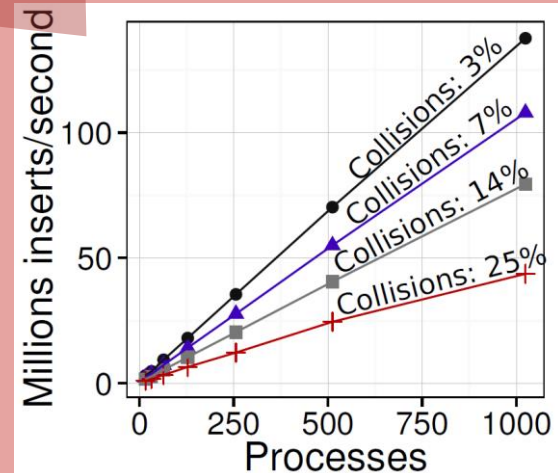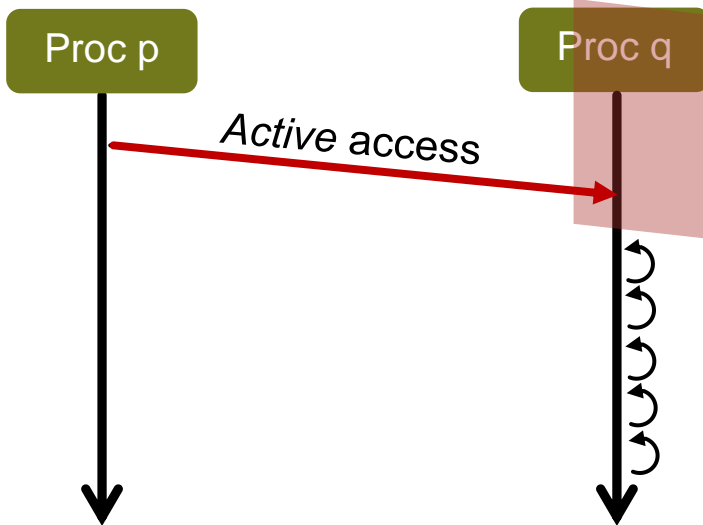- Consider an insert in a distributed hashtable...

No hash collision:

→ 1 remote atomic
→ Up to 5x speedup over MP [1]

A hash collision:

→ 4 remote atomics + 2 remote puts
→ Significant performance drops

[1] R. Gerstenberger et al. Enabling Highly-Scalable Remote Memory Access Programming with MPI-3 One-Sided. SC13
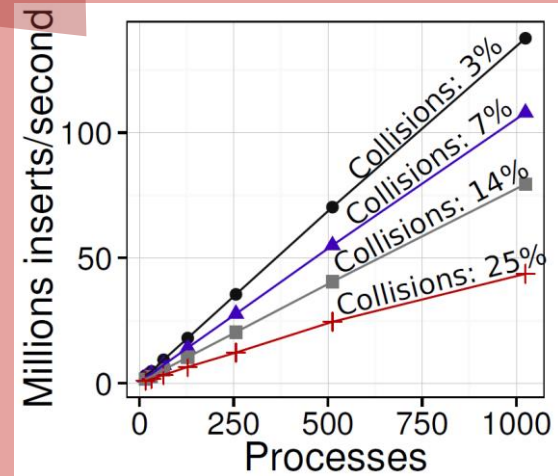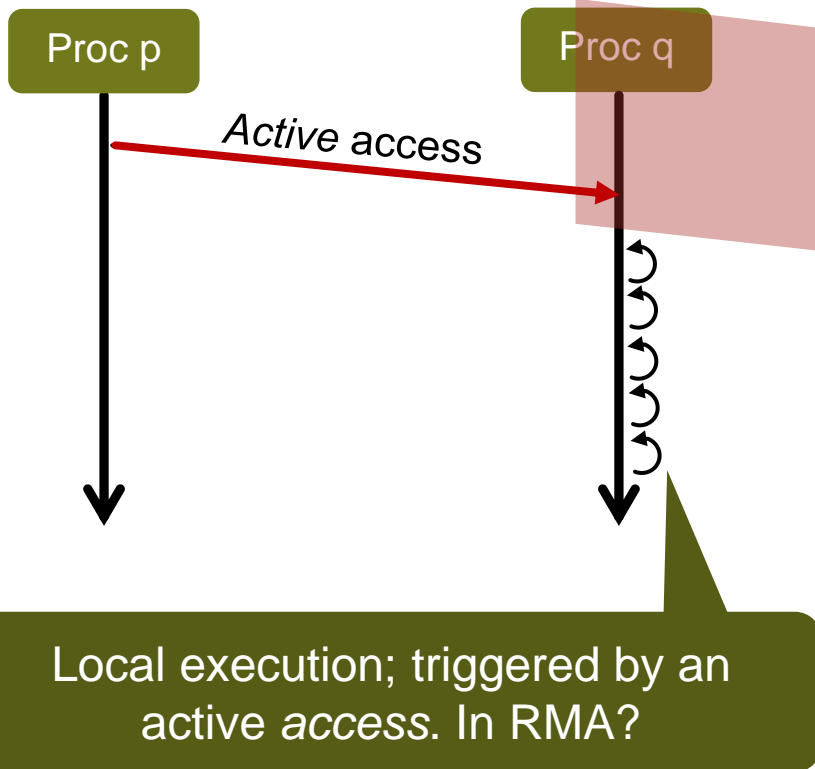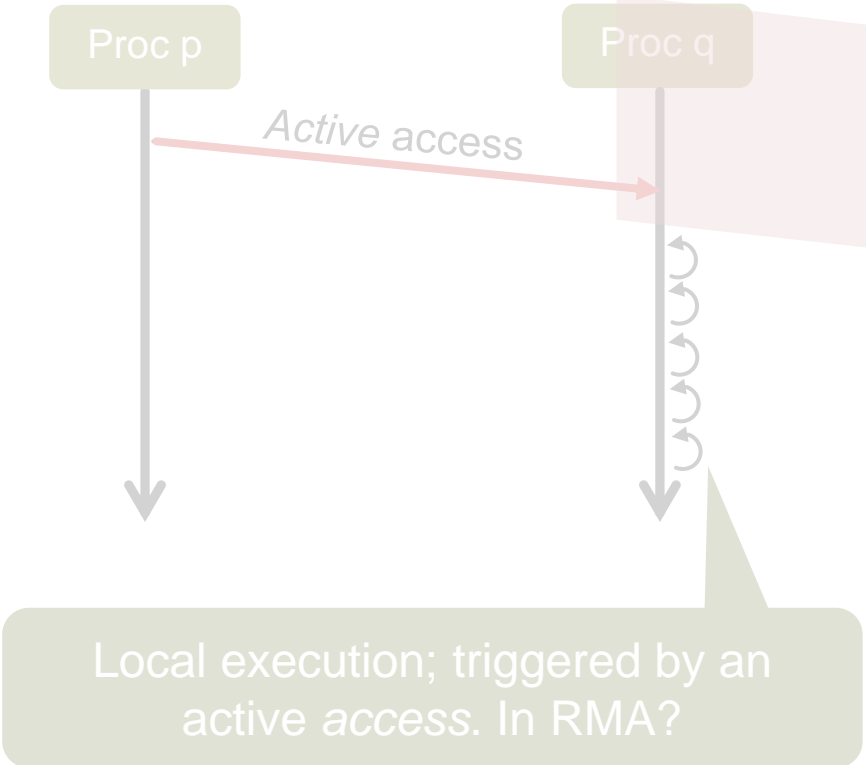
# REMOTE MEMORY ACCESS PROGRAMMING

- Is it ideal? **NO!**
- Consider an insert in a distributed hashtable...

No hash collision:

➔ 1 remote atomic
➔ Up to 5x speedup over MP [1]

| Proc p | | Proc q |

Atomic CAS

More atomics & puts

A hash collision:

➔ 4 remote atomics + 2 remote puts
➔ Significant performance drops

[1] R. Gerstenberger et al. Enabling Highly-Scalable Remote Memory Access Programming with MPI-3 One-Sided. SC13

# REMOTE MEMORY ACCESS PROGRAMMING

- Is it ideal? **NO!**
- Consider an insert in a distributed hashtable...

No hash collision:

→ 1 remote atomic
→ Up to 5x speedup over MP [1]

Proc p

Proc q

Atomic CAS

A hash collision:

→ 4 remote atomics + 2 remote puts
→ Significant performance drops

[1] R. Gerstenberger et al. Enabling Highly-Scalable Remote Memory Access Programming with MPI-3 One-Sided. SC13

# REMOTE MEMORY ACCESS PROGRAMMING

- Is it ideal? **NO!**
- Consider an insert in a distributed hashtable...

No hash collision:

→ 1 remote atomic
→ Up to 5x speedup over MP [1]

Proc p

Proc q

A hash collision:

→ 4 remote atomics + 2 remote puts
→ Significant performance drops



[1] R. Gerstenberger et al. Enabling Highly-Scalable Remote Memory Access Programming with MPI-3 One-Sided. SC13

# REMOTE MEMORY ACCESS PROGRAMMING

- Is it ideal? **NO!**
- Consider an insert in a distributed hashtable...

No hash collision:

➔ 1 remote atomic
➔ Up to 5x speedup over MP [1]

Proc p

Proc q

*Active* access

A hash collision:

➔ 4 remote atomics + 2 remote puts
➔ Significant performance drops

[1] R. Gerstenberger et al. Enabling Highly-Scalable Remote Memory Access Programming with MPI-3 One-Sided. SC13

# REMOTE MEMORY ACCESS PROGRAMMING

- Is it ideal? **NO!**
- Consider an insert in a distributed hashtable...

No hash collision:

➔ 1 remote atomic
➔ Up to 5x speedup over MP [1]

Proc p | Proc q

*Active* access

A hash collision:

➔ 4 remote atomics + 2 remote puts
➔ Significant performance drops

Local execution; triggered by an active *access*. In RMA?



[1] R. Gerstenberger et al. Enabling Highly-Scalable Remote Memory Access Programming with MPI-3 One-Sided. SC13

# REMOTE MEMORY ACCESS PROGRAMMING

- Is ...
- Co...
distributed hashtable...

## How to enable it?

No hash collision:

➔ 1 remote atomic
➔ Up to 5x speedup over MP [1]

| Proc p | | Proc q |

*Active* access

A hash collision:

➔ 4 remote atomics + 2 remote puts
➔ Significant performance drops

Local execution; triggered by an active *access*. In RMA?

[1] R. Gerstenberger et al. Enabling Highly-Scalable Remote Memory Access Programming with MPI-3 One-Sided. SC13

# REMOTE MEMORY ACCESS PROGRAMMING

- Is t...
- Co...
  distributed hashtable...

**How to enable it?**

No hash collision:

→ 1 remote atomic
→ Up to 5x speedup over MP [1]

| Proc p | Proc q |

*Active* access

A hash collision:

→ 4 remote atomics + 2 remote puts
→ Significant performance drops

**Use "active" semantics**

Local executi...
active *access.* in RMA?

[1] R. Gerstenberger et al. Enabling Highly-Scalable Remote Memory Access Programming with MPI-3 One-Sided. SC13

# REMOTE MEMORY ACCESS PROGRAMMING

- Is ...
- Co...
distributed hashtable...

**How to enable it?**

No hash collision:

➜ 1 remote atomic
➜ Up to 5x speedup over MP [1]

Proc p

*Active* access

Proc ...

**Use and extend I/O MMUs and their paging capabilities**

A hash collision:
remote puts
drops

**Use "active" semantics**

Local execution
active *access* in RMA?



[1] R. Gerstenberger et al. Enabling Highly-Scalable Remote Memory Access Programming with MPI-3 One-Sided. SC13

# "ACTIVE" REMOTE MEMORY ACCESS PROGRAMMING

# "ACTIVE" REMOTE MEMORY ACCESS PROGRAMMING

Local execution; triggered by an active *access*. In RMA?

NO!

# "ACTIVE" REMOTE MEMORY ACCESS PROGRAMMING

Local execution; triggered by an active *access*. In RMA?

How to enable it?

# "ACTIVE" REMOTE MEMORY ACCESS PROGRAMMING

Local execution; triggered by an active *access*. In RMA?

How to enable it?

Use "active" semantics

# "Active" Remote Memory Access Programming

Local execution; triggered by an active *access*. In RMA?

How to enable it?

Use "active" semantics

Use and extend I/O MMUs and their paging capabilities

# USE SEMANTICS FROM ACTIVE MESSAGES (AM) [1]

[1] T. von Eicken et al. Active messages: a mechanism for integrated communication and computation. ISCA'92.

# USE SEMANTICS FROM ACTIVE MESSAGES (AM) [1]

**Process p**

**Process q**

[1] T. von Eicken et al. Active messages: a mechanism for integrated communication and computation. ISCA'92.

# USE SEMANTICS FROM ACTIVE MESSAGES (AM) [1]

Process p

Process q

Memory

[1] T. von Eicken et al. Active messages: a mechanism for integrated communication and computation. ISCA'92.

# USE SEMANTICS FROM ACTIVE MESSAGES (AM) [1]

**Process p**

**Process q**

Memory

Handler A

. . .

Handler Z

[1] T. von Eicken et al. Active messages: a mechanism for integrated communication and computation. ISCA'92.

# USE SEMANTICS FROM ACTIVE MESSAGES (AM) [1]

**Process p**

**Process q**

Memory

A's addr: Handler A

. . .

Z's addr: Handler Z

[1] T. von Eicken et al. Active messages: a mechanism for integrated communication and computation. ISCA'92.

# USE SEMANTICS FROM ACTIVE MESSAGES (AM) [1]

**Process p**

*Active message*

**Process q**

Memory

A's addr: Handler A

. . .

Z's addr: Handler Z

[1] T. von Eicken et al. Active messages: a mechanism for integrated communication and computation. ISCA'92.

# USE SEMANTICS FROM ACTIVE MESSAGES (AM) [1]

**Process p**

Z's addr

Active message

**Process q**

Memory

A's addr:    Handler A

. . .

Z's addr:    Handler Z

[1] T. von Eicken et al. Active messages: a mechanism for integrated communication and computation. ISCA'92.

# USE SEMANTICS FROM ACTIVE MESSAGES (AM) [1]



[1] T. von Eicken et al. Active messages: a mechanism for integrated communication and computation. ISCA'92.

# USE SEMANTICS FROM ACTIVE MESSAGES (AM) [1]

**Process p**

Active message
Z's addr | Payload

**Process q**

Memory

A's addr:　Handler A

...

Z's addr:　Handler Z

[1] T. von Eicken et al. Active messages: a mechanism for integrated communication and computation. ISCA'92.

# USE SEMANTICS FROM ACTIVE MESSAGES (AM) [1]

**Process p**

Active message
Z's addr | Payload

**Process q**

Memory

A's addr: | Handler A

. . .

Z's addr: | Handler Z

[1] T. von Eicken et al. Active messages: a mechanism for integrated communication and computation. ISCA'92.

# USE SEMANTICS FROM ACTIVE MESSAGES (AM) [1]



Process p

Active message

Z's addr | Payload

Process q

Memory

A's addr: Handler A

...

Z's addr: Handler Z

[1] T. von Eicken et al. Active messages: a mechanism for integrated communication and computation. ISCA'92.

# USE SEMANTICS FROM ACTIVE MESSAGES (AM) [1]



IBM

Myricom

AM++[2]
GASNet [3]

**Process p**

Active message
Z's addr | Payload

**Process q**

Memory

A's addr: Handler A

...

Z's addr: Handler Z

[1] T. von Eicken et al. Active messages: a mechanism for integrated communication and computation. ISCA'92.
[2] J. J. Willcock et al. AM++: A generalized active message framework. PACT'10.
[3] D. Bonachea, GASNet Specification, v1.1. Berkeley Technical Report. 2002.

# USE SEMANTICS FROM ACTIVE MESSAGES (AM) [1]

**IBM**

**Myricom**

AM++[2]
GASNet [3]

**Process p**

Active message
Z's addr | Payload

**Process q**

Memory

A's addr: Handler A

. . .

Z's addr: Handler Z

We need *active* puts/gets:
- Invoke a handler upon accessing a given page
- Preserve one-sided RMA behavior

[1] T. von Eicken et al. Active messages: a mechanism for integrated communication and computation. ISCA'92.
[2] J. J. Willcock et al. AM++: A generalized active message framework. PACT'10.
[3] D. Bonachea, GASNet Specification, v1.1. Berkeley Technical Report. 2002.

# USE SEMANTICS FROM ACTIVE MESSAGES (AM) [1]

Process p

Z's addr | Payload

Active message

Process q

A's addr: Handler A

...

Z's addr: Handler Z

We use it in syntax & semantics to enable the "active" behavior

We need *active* puts/gets.

- Invoke a handler upon accessing a given page
- Preserve one-sided RMA behavior

[1] T. von Eicken et al. Active messages: a mechanism for integrated communication and computation. ISCA'92.
[2] J. J. Willcock et al. AM++: A generalized active message framework. PACT'10.
[3] D. Bonachea, GASNet Specification, v1.1. Berkeley Technical Report. 2002.

# USE INPUT/OUTPUT MEMORY MANAGEMENT UNITS

# USE INPUT/OUTPUT MEMORY MANAGEMENT UNITS

# USE INPUT/OUTPUT MEMORY MANAGEMENT UNITS

# USE INPUT/OUTPUT MEMORY MANAGEMENT UNITS

# USE INPUT/OUTPUT MEMORY MANAGEMENT UNITS

# USE INPUT/OUTPUT MEMORY MANAGEMENT UNITS

# USE INPUT/OUTPUT MEMORY MANAGEMENT UNITS

# USE INPUT/OUTPUT MEMORY MANAGEMENT UNITS

# USE INPUT/OUTPUT MEMORY MANAGEMENT UNITS

# USE INPUT/OUTPUT MEMORY MANAGEMENT UNITS

# USE INPUT/OUTPUT MEMORY MANAGEMENT UNITS

# USE INPUT/OUTPUT MEMORY MANAGEMENT UNITS

# USE INPUT/OUTPUT MEMORY MANAGEMENT UNITS

# USE INPUT/OUTPUT MEMORY MANAGEMENT UNITS

Main memory

Physical addresses

Physical addresses

IOMM

MMU

! We propose it as a way to implement the "active" behavior

Device addresses

I/O devices

CPU

AMD

IBM

ARM®

Sun microsystems

SOLARFLARE®

(intel)

+ PCI EXPRESS®

# IOMMUs AND RMA

# IOMMUs AND RMA

NIC

# IOMMUs and RMA

IOMMU

NIC

# IOMMUs AND RMA

MSI

IOMMU

CPU

NIC

# IOMMUS AND RMA

# IOMMUS AND RMA

MSI

IOMMU

CPU

NIC

SMT cores

Main memory

# IOMMUs AND RMA

# IOMMUS AND RMA

# IOMMUs AND RMA

# IOMMUs AND RMA

MSI

IOMMU

CPU

An RDMA
packet

**1**

NIC

SMT cores

**2**

PCIe packets

Main memory

Remapping structures

Dev-to-PT

PT

# IOMMUs AND RMA

# IOMMUs and RMA

# IOMMUs AND RMA

# IOMMUs and RMA

# IOMMUS AND RMA

# IOMMUs and RMA

# IOMMUs AND RMA

# IOMMUs AND RMA

# IOMMUs AND RMA

MSI

An RDMA packet

**1**

**2** PCIe packets

NIC

**3** IOMMU

Dev-to-PT cache **4**

IOTLB **6**

CPU

SMT cores

Main memory

Remapping structures

Dev-to-PT **5**

PT **7**

W
R **8**

System-wide fault log

Fault entry → ... → Fault entry

# IOMMUS AND RMA

# IOMMUs AND RMA

# IOMMUS AND RMA

# IOMMUs AND RMA

# IOMMUS AND RMA

We could use it somehow. But…

# IOMMUS AND RMA



**MSI** (10)

**An RDMA packet** (1)

**IOMMU** (3)

**CPU**

**NIC**

**Dev-to-PT cache** (4)

(11)

**SMT cores**

**IOTLB** (6)

**PCIe packets** (2)

**Main memory**

**Remapping structures**

**System-wide fault log** (9)

**User handlers** (12)

**Dev-to-PT** (5)

W / R (8)

Fault entry → ... → Fault entry

**Handler A**

...

(7)

**PT**

# IOMMUs and RMA

# IOMMUs AND RMA

We could use it somehow. But…

MSI **10**

An RDMA packet **1**

IOMMU **3**

**4** Dev-to-PT cache

**6** IOTLB

CPU

**11** SMT cores

**NIC**

PCIe packets **2**

No parallelism (single log)... BAD

No multiplexing (single log)... BAD

Main memory

Remapping structures

**5** Dev-to-PT

W
R **8**

**7** PT

System-wide fault log **9**

Fault entry → … → Fault entry

User handlers **12**

Handler A

…

# ACTIVE PUTS

MSI

IOMMU

CPU

An RDMA packet

NIC

PCIe packets

Dev-to-PT cache

IOTLB

SMT cores

Main memory

Remapping structures

Dev-to-PT

PT

W
R

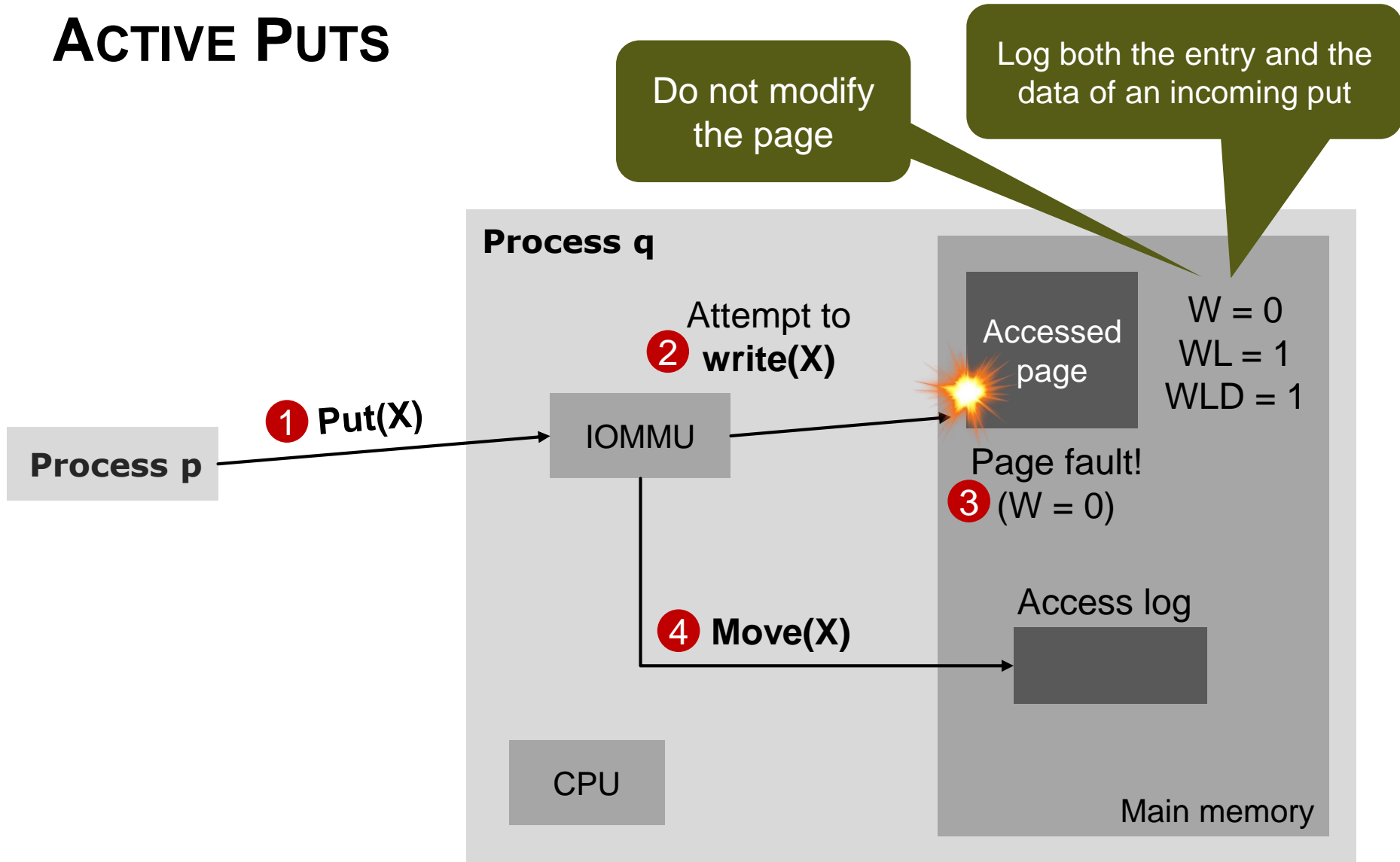System-wide fault log
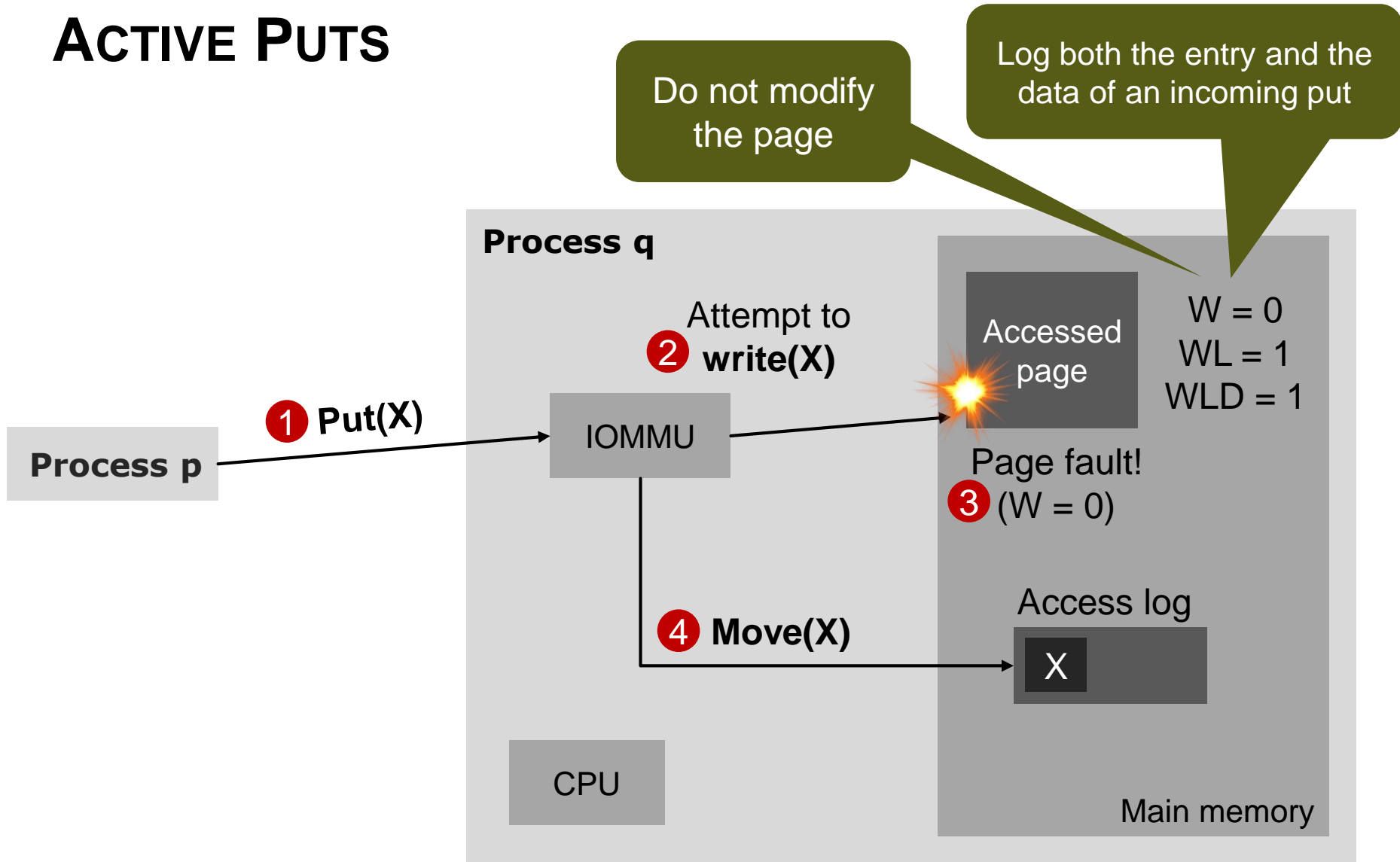
Fault entry → ... → Fault entry

User handlers

Handler A

...

# ACTIVE PUTS

# ACTIVE PUTS

# ACTIVE PUTS

# ACTIVE PUTS



MSI

An RDMA packet

NIC

PCIe packets

IOMMU

Dev-to-PT cache

IOTLB

Access log table +

CPU

SMT cores

Main memory

Remapping structures

Dev-to-PT

PT

W
R

System-wide fault log

Fault entry → ... → Fault entry

Access log (private for each process) +

Fault entry → ... → Fault entry

Request data

Request data

User handlers

Handler A

...

Data can be reused

# ACTIVE PUTS



Stores addresses of each access log

MSI

IOMMU

Dev-to-PT cache

Access log table +

IOTLB

CPU

SMT cores

An RDMA packet

NIC

PCIe packets

Main memory

Remapping structures

W
R

Dev-to-PT

PT

System-wide fault log

Fault entry → … → Fault entry

Access log (private for each process) +

Fault entry → … → Fault entry

Request data

Request data

User handlers

Handler A

…

Data can be reused

# ACTIVE PUTS

An RDMA packet

NIC

PCIe packets

IOMMU

Dev-to-PT cache

IOTLB

Stores addresses of each access log

Access log table +

MSI

CPU

SMT cores

Main memory

Remapping structures

Dev-to-PT

PT

| W |
| R |
| *WL* + |

System-wide fault log

Fault entry → ... → Fault entry

*Access log (private for each process)* +

Fault entry → ... → Fault entry

Request data

Request data

User handlers

Handler A

...

Data can be reused

# ACTIVE PUTS

An RDMA packet

NIC

PCIe packets

IOMMU

Dev-to-PT cache

IOTLB

**Access log table** +

Stores addresses of each access log

MSI

CPU

SMT cores

Decide on keeping/discarding the entry/data

Main memory

Remapping structures

Dev-to-PT

PT

| W |
| R |
| *WL* + |
| *WLD* + |

System-wide fault log

Fault entry → ... → Fault entry

*Access log (private for each process)* +

Fault entry → ... → Fault entry

Request data

Request data

User handlers

Handler A

...

Data can be reused

# ACTIVE PUTS

An RDMA packet

NIC

PCIe packets

**IOMMU**

Dev-to-PT cache

IOTLB

Stores addresses of each access log

MSI

**CPU**

SMT cores

Access log table ⊕

Main memory

Remapping structures

Dev-to-PT

Decide on keeping/discarding the entry/data

| W |
| R |
| *WL* ⊕ |
| *WLD* ⊕ |

System-wide fault log

Fault entry → ... → Fault entry

User handlers

Handler A

...

*Access log (private for each process)* ⊕

Fault entry → ... → Fault entry

Request data

Request data

*IUID* ⊕

Maps each page to an access log

Data can be reused

# ACTIVE PUTS

# ACTIVE PUTS

**Process q**

IOMMU

CPU

Main memory

**Process p**

# ACTIVE PUTS

**Process q**

Accessed page

IOMMU

**Process p**

Access log

CPU

Main memory

# ACTIVE PUTS

**Process p**

**Process q**

IOMMU

Accessed page

W = 0
WL = 1
WLD = 1

Access log

CPU

Main memory

# ACTIVE PUTS

Do not modify
the page

**Process q**

Accessed
page

W = 0
WL = 1
WLD = 1

IOMMU

**Process p**

Access log

CPU

Main memory

# ACTIVE PUTS

Do not modify the page

Log both the entry and the data of an incoming put

**Process q**

Accessed page

W = 0
WL = 1
WLD = 1

IOMMU

**Process p**

Access log

CPU

Main memory

# ACTIVE PUTS

Do not modify the page

Log both the entry and the data of an incoming put

**Process q**

Accessed page

W = 0
WL = 1
WLD = 1

① **Put(X)**

**Process p**

IOMMU

Access log

CPU

Main memory

# ACTIVE PUTS

Do not modify the page

Log both the entry and the data of an incoming put

**Process q**

**① Put(X)**

Attempt to
**② write(X)**

**Process p**

IOMMU

Accessed
page

W = 0
WL = 1
WLD = 1

Access log

CPU

Main memory

# ACTIVE PUTS

Do not modify the page

Log both the entry and the data of an incoming put

**Process q**

Attempt to
2 **write(X)**

Accessed page

W = 0
WL = 1
WLD = 1

1 **Put(X)**

**Process p**

IOMMU

Page fault!
3 (W = 0)

Access log

CPU

Main memory

# ACTIVE PUTS

Do not modify the page

Log both the entry and the data of an incoming put

**Process q**

① **Put(X)**

**Process p**

② Attempt to **write(X)**

IOMMU

Accessed page

W = 0
WL = 1
WLD = 1

Page fault!
③ (W = 0)

④ **Move(X)**

Access log

CPU

Main memory

# ACTIVE PUTS

# ACTIVE PUTS

Do not modify the page

Log both the entry and the data of an incoming put

**Process q**

**1** **Put(X)**

**Process p**

IOMMU

Attempt to
**2** **write(X)**

Accessed page

$W = 0$
$WL = 1$
$WLD = 1$

Page fault!
**3** $(W = 0)$

Access log

**4** **Move(X)**

X

**5** **Process(X)**

CPU

Main memory

# ACTIVE GETS

# ACTIVE GETS

# ACTIVE GETS

**Process q**

IOMMU

CPU

Main memory

**Process p**

# ACTIVE GETS



Process p

Process q

Accessed page

IOMMU

Access log

CPU

Main memory

# ACTIVE GETS

**Process q**

**Process p**

IOMMU

Accessed page

R = 1
RL = 1
RLD = 1

Access log

CPU

Main memory

# ACTIVE GETS

Enable reading from the page

**Process q**

Accessed page

R = 1
RL = 1
RLD = 1

IOMMU

**Process p**

Access log

CPU

Main memory

# ACTIVE GETS

Enable reading from the page

Log both the entry and the data accessed by a get

**Process q**

Accessed page

R = 1
RL = 1
RLD = 1

IOMMU

Access log

**Process p**

CPU

Main memory

# ACTIVE GETS

Enable reading from the page

Log both the entry and the data accessed by a get

**Process q**

Accessed page

R = 1
RL = 1
RLD = 1

1 **Get(X)**

**Process p**

IOMMU

Access log

CPU

Main memory

# ACTIVE GETS

Enable reading from the page

Log both the entry and the data accessed by a get

**Process q**

Accessed page

R = 1
RL = 1
RLD = 1

② **Read(X)**

① **Get(X)**

IOMMU

**Process p**

Access log

CPU

Main memory

# ACTIVE GETS

Enable reading from the page

Log both the entry and the data accessed by a get

**Process q**

Accessed page

R = 1
RL = 1
RLD = 1

② **Read(X)**

① **Get(X)**

**Process p**

IOMMU

③ **Copy(X)**

Access log

X

CPU

Main memory

# ACTIVE GETS

Enable reading from the page

Log both the entry and the data accessed by a get

**Process q**

R = 1
RL = 1
RLD = 1

Accessed page

**2** **Read(X)**

**1** **Get(X)**

**Process p**

IOMMU

**3** **Copy(X)**

Access log

X

CPU

**4** **Process(X)**

Main memory

# ACTIVE GETS

Enable reading from the page

Log both the entry and the data accessed by a get

**Process q**

Accessed page

R = 1
RL = 1
RLD = 1

**2** **Read(X)**

**1** **Get(X)**

IOMMU

**Process p**

Access log

**3** **Copy(X)**

X

**4** **Process(X)**

CPU

Main memory

Sounds like we can reuse most of the existing stuff!

(intel) + PCI EXPRESS

# INTERACTIONS WITH THE CPU

An RDMA packet

NIC

PCIe packets

MSI

IOMMU

Dev-to-PT cache

IOTLB

Access log table

CPU

SMT cores

Main memory

Remapping structures

Dev-to-PT

PT

W
R
*WL* +
*WLD* +
*RL* +
*RLD* +
*IUID* +

System-wide fault log

Fault entry → … → Fault entry

*Access log (private for each process)* +

Fault entry → … → Fault entry

Request data

Request data

User handlers

Handler A

…

# INTERACTIONS WITH THE CPU

MSI

IOMMU

CPU

Dev-to-PT cache

IOTLB

*Access log table*

SMT cores

# INTERACTIONS WITH THE CPU

# INTERACTIONS WITH THE CPU

MSI

IOMMU

CPU

Dev-to-PT
cache

IOTLB

*Access log table*

SMT cores

- Interrupts

# INTERACTIONS WITH THE CPU

MSI

IOMMU                                                                CPU

SMT cores

Dev-to-PT
cache

Access log table

IOTLB

- Interrupts
- Polling

# INTERACTIONS WITH THE CPU

MSI

IOMMU                                                    CPU

Dev-to-PT cache

Access log table

IOTLB

SMT cores

- Interrupts
- Polling
- Direct notifications via scratchpads

# INTERACTIONS WITH THE CPU

MSI

IOMMU

CPU

Dev-to-PT cache

IOTLB

*Access log table*

SMT cores

*Scratchpad memory*

+

- Interrupts
- Polling
- Direct notifications via scratchpads

# INTERACTIONS WITH THE CPU

MSI

IOMMU

CPU

Dev-to-PT cache

Access log table

IOTLB

SMT cores

Scratchpad memory

Var

- Interrupts
- Polling
- Direct notifications via scratchpads

# INTERACTIONS WITH THE CPU

MSI

IOMMU

CPU

SMT cores

Dev-to-PT cache

IOTLB

*Access log table*

*Scratchpad memory*

Var | Handler A

- Interrupts
- Polling
- Direct notifications via scratchpads

# INTERACTIONS WITH THE CPU



- Interrupts
- Polling
- Direct notifications via scratchpads

# INTERACTIONS WITH THE CPU

MSI

IOMMU

CPU

SMT cores

Dev-to-PT cache

Access log table

IOTLB

Scratchpad memory

Var    Handler A

Hyper thread

- Interrupts
- Polling
- Direct notifications via scratchpads

# INTERACTIONS WITH THE CPU

MSI

IOMMU

CPU

SMT cores

log table

IOTLB

## Are we done?

Scratchpad memory

Var    Handler A

Hyper thread

- Interrupts
- Polling
- Direct notifications via scratchpads

# CONSISTENCY

# CONSISTENCY

- A weak consistency model [1]

[1] K. Gharachorloo et al. Memory consistency and event ordering in scalable shared-memory multiprocessors. ISCA '90

# CONSISTENCY

- A weak consistency model [1]
  - Consistency on-demand

REMOTE MEMORY ACCESS (RMA) PROGRAMMING

[1] K. Gharachorloo et al. Memory consistency and event ordering in scalable shared-memory multiprocessors. ISCA '90

# CONSISTENCY

- A weak consistency model [1]
  - Consistency on-demand
- active_flush(int target_id)

[1] K. Gharachorloo et al. Memory consistency and event ordering in scalable shared-memory multiprocessors. ISCA '90

# CONSISTENCY



REMOTE MEMORY ACCESS (RMA) PROGRAMMING

- A weak consistency model [1]
  - Consistency on-demand
- active_flush(int target_id)
  - Enforces the completion of active accesses issued by the calling process and targeted at target_id

[1] K. Gharachorloo et al. Memory consistency and event ordering in scalable shared-memory multiprocessors. ISCA '90

# CONSISTENCY



REMOTE MEMORY ACCESS (RMA) PROGRAMMING

- A weak consistency model [1]
    - Consistency on-demand
- active_flush(int target_id)
    - Enforces the completion of active accesses issued by the calling process and targeted at target_id
    - Implemented with an active get issued at a special *flushing page*

[1] K. Gharachorloo et al. Memory consistency and event ordering in scalable shared-memory multiprocessors. ISCA '90

# CONSISTENCY



- A weak consistency model [1]
  - Consistency on-demand
- active_flush(int target_id)
  - Enforces the completion of active accesses issued by the calling process and targeted at target_id
  - Implemented with an active get issued at a special *flushing page*

**Process p**

Memory

[1] K. Gharachorloo et al. Memory consistency and event ordering in scalable shared-memory multiprocessors. ISCA '90

# CONSISTENCY



- A weak consistency model [1]
  - Consistency on-demand
- active_flush(int target_id)
  - Enforces the completion of active accesses issued by the calling process and targeted at target_id
  - Implemented with an active get issued at a special *flushing page*

| Process p | | Process q | |
|---|---|---|---|
| Memory | | IOMMU | Memory |

[1] K. Gharachorloo et al. Memory consistency and event ordering in scalable shared-memory multiprocessors. ISCA '90

# CONSISTENCY

- A weak consistency model [1]
    - Consistency on-demand
- active_flush(int target_id)
    - Enforces the completion of active accesses issued by the calling process and targeted at target_id
    - Implemented with an active get issued at a special *flushing page*

| Process p | | Process q | |
|---|---|---|---|
| Memory | | IOMMU | Memory |
| | | | Access log |

[1] K. Gharachorloo et al. Memory consistency and event ordering in scalable shared-memory multiprocessors. ISCA '90

# CONSISTENCY

- A weak consistency model [1]
  - Consistency on-demand
- active_flush(int target_id)
  - Enforces the completion of active accesses issued by the calling process and targeted at target_id
  - Implemented with an active get issued at a special *flushing page*

| Process p | Process q |
|---|---|
| Memory | IOMMU | Memory |
| | | Access log |
| | | Flushing page |

[1] K. Gharachorloo et al. Memory consistency and event ordering in scalable shared-memory multiprocessors. ISCA '90

# CONSISTENCY

- A weak consistency model [1]
  - Consistency on-demand
- active_flush(int target_id)
  - Enforces the completion of active accesses issued by the calling process and targeted at target_id
  - Implemented with an active get issued at a special *flushing page*



**Process p**

Memory

X Y Z Some accesses

IOMMU

**Process q**

Memory

Access log

Flushing page

[1] K. Gharachorloo et al. Memory consistency and event ordering in scalable shared-memory multiprocessors. ISCA '90

# CONSISTENCY

- A weak consistency model [1]
  - Consistency on-demand
- active_flush(int target_id)
  - Enforces the completion of active accesses issued by the calling process and targeted at target_id
  - Implemented with an active get issued at a special *flushing page*



**Process p**

Memory

Some accesses

**Process q**

IOMMU

Memory

Access log

X Y Z

Flushing page

[1] K. Gharachorloo et al. Memory consistency and event ordering in scalable shared-memory multiprocessors. ISCA '90

# CONSISTENCY

- A weak consistency model [1]
  - Consistency on-demand
- active_flush(int target_id)
  - Enforces the completion of active accesses issued by the calling process and targeted at target_id
  - Implemented with an active get issued at a special *flushing page*



**Process p**

Memory

Some accesses

IOMMU

**Process q**

Memory

Access log

X Y Z

Active get

Flushing page

[1] K. Gharachorloo et al. Memory consistency and event ordering in scalable shared-memory multiprocessors. ISCA '90

# CONSISTENCY

- A weak consistency model [1]
  - Consistency on-demand
- active_flush(int target_id)
  - Enforces the completion of active accesses issued by the calling process and targeted at target_id
  - Implemented with an active get issued at a special *flushing page*



[1] K. Gharachorloo et al. Memory consistency and event ordering in scalable shared-memory multiprocessors. ISCA '90

# CONSISTENCY

- A weak consistency model [1]
  - Consistency on-demand
- active_flush(int target_id)
  - Enforces the completion of active accesses issued by the calling process and targeted at target_id
  - Implemented with an active get issued at a special *flushing page*

[1] K. Gharachorloo et al. Memory consistency and event ordering in scalable shared-memory multiprocessors. ISCA '90

# CONSISTENCY

- A weak consistency model [1]
  - Consistency on-demand
- active_flush(int target_id)
  - Enforces the completion of active accesses issued by the calling process and targeted at target_id
  - Implemented with an active get issued at a special *flushing page*



**Process p**

Memory

Some accesses

Active get

**Process q**

IOMMU

Memory

Access log

Z

Flushing page

[1] K. Gharachorloo et al. Memory consistency and event ordering in scalable shared-memory multiprocessors. ISCA '90

# CONSISTENCY

- A weak consistency model [1]
  - Consistency on-demand
- active_flush(int target_id)
  - Enforces the completion of active accesses issued by the calling process and targeted at target_id
  - Implemented with an active get issued at a special *flushing page*



**Process p**

Memory

Some accesses

Active get

**Process q**

IOMMU

Memory

Access log

Flushing page

[1] K. Gharachorloo et al. Memory consistency and event ordering in scalable shared-memory multiprocessors. ISCA '90

# CONSISTENCY

- A weak consistency model [1]
  - Consistency on-demand
- active_flush(int target_id)
  - Enforces the completion of active accesses issued by the calling process and targeted at target_id
  - Implemented with an active get issued at a special *flushing page*

[1] K. Gharachorloo et al. Memory consistency and event ordering in scalable shared-memory multiprocessors. ISCA '90

# CONSISTENCY

# CONSISTENCY

# CONSISTENCY



MSI

IOMMU

CPU

Dev-to-PT cache

IOTLB

Access log table +

Flushing buffer +

+

SMT cores

Scratchpad memory

Handler A • • • +

+

Hyper thread +

Contains the addresses of flushing pages

# CONSISTENCY

MSI

IOMMU

CPU

Dev-to-PT cache

IOTLB

Access log table +

Flushing buffer +

+

SMT cores

Scratchpad memory

Handler A · · ·

+

Hyper thread +

+

Contains the addresses of flushing pages

Maps flushing pages to IUIDs and access logs

Let's summarize…

# Let's summarize…



Active Messages

# Let's summarize…

IOMMUs



Active Messages

# Let's summarize…



Active Messages

IOMMUs

Active Puts/Gets

Let's summarize…

Active Messages

IOMMUs

Consistency

Active Puts/Gets

Let's summarize…

Active Messages

IOMMUs

Consistency

Active Puts/Gets

How can we use it?

# ACTIVE ACCESS USE-CASES

## DISTRIBUTED HASHTABLE

# ACTIVE ACCESS USE-CASES
## DISTRIBUTED HASHTABLE

- Used to construct key-value stores (e.g., Memcached [1])

[1] B. Fitzpatrick. Distributed caching with memcached. Linux journal, 2004.

# ACTIVE ACCESS USE-CASES
## DISTRIBUTED HASHTABLE

- Used to construct key-value stores (e.g., Memcached [1])

| Local volume 0 (at process 0) | Local volume 1 (at process 1) | ... | Local volume N-1 (at process N-1) |

[1] B. Fitzpatrick. Distributed caching with memcached. Linux journal, 2004.

# ACTIVE ACCESS USE-CASES

## DISTRIBUTED HASHTABLE

- Used to construct key-value stores (e.g., Memcached [1])

| Local volume 0 (at process 0) | Local volume 1 (at process 1) | | Local volume N-1 (at process N-1) |
|---|---|---|---|
| Table of elements | Table of elements | . . . | Table of elements |

[1] B. Fitzpatrick. Distributed caching with memcached. Linux journal, 2004.

# ACTIVE ACCESS USE-CASES

## DISTRIBUTED HASHTABLE

- Used to construct key-value stores (e.g., Memcached [1])

| Local volume 0 (at process 0) | Local volume 1 (at process 1) | | Local volume N-1 (at process N-1) |
|---|---|---|---|
| Table of elements | Table of elements | ... | Table of elements |
| ↓ | ↓ | | ↓ |
| Overflow heap | Overflow heap | | Overflow heap |

[1] B. Fitzpatrick. Distributed caching with memcached. Linux journal, 2004.

# ACTIVE ACCESS USE-CASES

## DISTRIBUTED HASHTABLE: INSERTS (RMA)

# ACTIVE ACCESS USE-CASES
## DISTRIBUTED HASHTABLE: INSERTS (RMA)

Proc p

Proc q

Table of elements

Overflow heap

# ACTIVE ACCESS USE-CASES
## DISTRIBUTED HASHTABLE: INSERTS (RMA)

# ACTIVE ACCESS USE-CASES
## DISTRIBUTED HASHTABLE: INSERTS (RMA)

# ACTIVE ACCESS USE-CASES
## DISTRIBUTED HASHTABLE: INSERTS (RMA)

# ACTIVE ACCESS USE-CASES
## DISTRIBUTED HASHTABLE: INSERTS (RMA)

# ACTIVE ACCESS USE-CASES
## DISTRIBUTED HASHTABLE: INSERTS (AA)

# ACTIVE ACCESS USE-CASES
## DISTRIBUTED HASHTABLE: INSERTS (AA)

**Proc p**

**Proc q**

Table of elements

FAD (get and increment ptr to the next free cell)

PUT (insert element)

FAD & CAS & PUT (update ptrs)

Overflow heap

# ACTIVE ACCESS USE-CASES
## DISTRIBUTED HASHTABLE: INSERTS (AA)

Proc p

Proc q

PUT (intercepted by the IOMMU)

Table of elements

FAD (get and increment ptr to the next free cell)

PUT (insert element)

FAD & CAS & PUT (update ptrs)

Overflow heap

# ACTIVE ACCESS USE-CASES

## DISTRIBUTED HASHTABLE: INSERTS (AA)

# ACTIVE ACCESS USE-CASES
## DISTRIBUTED HASHTABLE: INSERTS (AA)

Proc p

Proc q

PUT (intercepted by the IOMMU)

Table of elements

Overflow heap

# ACTIVE ACCESS USE-CASES
## VIRTUAL GLOBAL ADDRESS SPACE (V-GAS)

# ACTIVE ACCESS USE-CASES

## VIRTUAL GLOBAL ADDRESS SPACE (V-GAS)

ENABLED

# ACTIVE ACCESS USE-CASES
## VIRTUAL GLOBAL ADDRESS SPACE (V-GAS)

ENABLED

Machine 0

Machine 1

Machine N-1

# ACTIVE ACCESS USE-CASES
## VIRTUAL GLOBAL ADDRESS SPACE (V-GAS)

ENABLED

Machine 0

Machine 1

Machine N-1

NIC

NIC

NIC

# ACTIVE ACCESS USE-CASES
## VIRTUAL GLOBAL ADDRESS SPACE (V-GAS)

ENABLED

| Machine 0 | Machine 1 | | Machine N-1 |

Proc 0

Proc 1

Proc N-1

NIC

NIC

NIC

# ACTIVE ACCESS USE-CASES
## VIRTUAL GLOBAL ADDRESS SPACE (V-GAS)

ENABLED

Machine 0

Proc 0

Memory

NIC

Machine 1

Proc 1

Memory

NIC

Machine N-1

Proc N-1

Memory

NIC

# ACTIVE ACCESS USE-CASES
## VIRTUAL GLOBAL ADDRESS SPACE (V-GAS)

ENABLED

Machine 0

Proc 0

MMU

Memory

NIC

Machine 1

Proc 1

MMU

Memory

NIC

Machine N-1

Proc N-1

MMU

Memory

NIC

# ACTIVE ACCESS USE-CASES
## VIRTUAL GLOBAL ADDRESS SPACE (V-GAS)

ENABLED

Machine 0

Proc 0

MMU

Memory

IOMMU

NIC

Machine 1

Proc 1

MMU

Memory

IOMMU

NIC

Machine N-1

Proc N-1

MMU

Memory

IOMMU

NIC

# ACTIVE ACCESS USE-CASES
## VIRTUAL GLOBAL ADDRESS SPACE (V-GAS)

ENABLED

# ACTIVE ACCESS USE-CASES
## VIRTUAL GLOBAL ADDRESS SPACE (V-GAS)



ENABLED

| Machine 0 | Machine 1 | Machine N-1 |
|-----------|-----------|-------------|
| Proc 0 | Proc 1 | Proc N-1 |
| MMU | MMU | MMU |
| Memory | Memory    V-GAS | Memory |
| IOMMU | IOMMU | IOMMU |
| NIC | NIC | NIC |

# ACTIVE ACCESS USE-CASES
## VIRTUAL GLOBAL ADDRESS SPACE (V-GAS)

ENABLED

| Machine 0 | Machine 1 | Machine N-1 |
|-----------|-----------|-------------|
| Proc 0 | Proc 1 | Proc N-1 |
| MMU | MMU | MMU |
| Memory | Memory    V-GAS | Memory |
| IOMMU | IOMMU | IOMMU |
| NIC | NIC | NIC |

# ACTIVE ACCESS USE-CASES
## VIRTUAL GLOBAL ADDRESS SPACE (V-GAS)

ENABLED

| Machine 0 | Machine 1 | | Machine N-1 |
|---|---|---|---|
| Proc 0 | Proc 1 | | Proc N-1 |
| MMU | MMU | | MMU |
| Memory | Memory | V-GAS | Memory |
| IOMMU | IOMMU | | IOMMU |
| NIC | NIC | | NIC |

# ACTIVE ACCESS USE-CASES
## VIRTUAL GLOBAL ADDRESS SPACE (V-GAS)

ENABLED

Machine 0          Machine 1                          Machine N-1

Proc 0             Proc 1                Local memory   Proc N-1
                                         protection

MMU                MMU                                  MMU

Memory             Memory          V-GAS               Memory

IOMMU              IOMMU                                IOMMU

NIC                NIC                                  NIC

# ACTIVE ACCESS USE-CASES
## VIRTUAL GLOBAL ADDRESS SPACE (V-GAS)

ENABLED

Remote memory
protection

Local memory
protection

Machine 0

Machine 1

Machine N-1

Proc 0

Proc 1

Proc N-1

MMU

MMU

MMU

Memory

Memory

V-GAS

Memory

IOMMU

IOMMU

IOMMU

NIC

NIC

NIC

# ACTIVE ACCESS USE-CASES
## VIRTUAL GLOBAL ADDRESS SPACE (V-GAS)

ENABLED

Remote memory
protection

Machine 0

Machine 1

Machine N-1

Proc 0

Proc 1

Local memory
protection

Proc N-1

MMU

MMU

MMU

Memory

Memory

V-GAS

Memory

IOMMU

IOMMU

IOMMU

Fetch data (used
for logging, fault-
tolerance, etc…)

NIC

NIC

NIC

# PERFORMANCE

- Evaluation on CSCS Monte Rosa
  - 1,496 computing Cray XE6 nodes
  - 47,872 schedulable cores
  - 46TB memory
- 3 microbenchmarks
- 4 use-cases

# PERFORMANCE: MICROBENCHMARKS
## RAW DATA TRANSFER

# PERFORMANCE: MICROBENCHMARKS
## RAW DATA TRANSFER

- Workload simulated with [1]:

# PERFORMANCE: MICROBENCHMARKS
## RAW DATA TRANSFER

- Workload simulated with [1]:

[1] N. Binkert et al. The gem5 simulator. SIGARCH Comput. Archit. News. 2011

# PERFORMANCE: MICROBENCHMARKS
## RAW DATA TRANSFER

- Workload simulated with [1]:



- Data generated with:

[1] N. Binkert et al. The gem5 simulator. SIGARCH Comput. Archit. News. 2011

# PERFORMANCE: MICROBENCHMARKS
## RAW DATA TRANSFER

- Workload simulated with [1]:



- Data generated with:
  - PktGen [2]

[1] N. Binkert et al. The gem5 simulator. SIGARCH Comput. Archit. News. 2011

[2] R. Olsson. PktGen the linux packet generator. Linux Symposium. 2005

# PERFORMANCE: MICROBENCHMARKS
## RAW DATA TRANSFER

- Workload simulated with [1]:



- Data generated with:
  - PktGen [2]
  - Netmap [3]

[1] N. Binkert et al. The gem5 simulator. SIGARCH Comput. Archit. News. 2011

[2] R. Olsson. PktGen the linux packet generator. Linux Symposium. 2005

[3] L. Rizzo. netmap: A novel framework for fast packet i/o. USENIX Annual Technical Conference. 2012

# PERFORMANCE: MICROBENCHMARKS
## RAW DATA TRANSFER

- Workload simulated with [1]:

- Data generated with:
  - PktGen [2]
  - Netmap [3]

[1] N. Binkert et al. The gem5 simulator. SIGARCH Comput. Archit. News. 2011

[2] R. Olsson. PktGen the linux packet generator. Linux Symposium. 2005

[3] L. Rizzo. netmap: A novel framework for fast packet i/o. USENIX Annual Technical Conference. 2012

# PERFORMANCE: LARGE-SCALE CODES
## COMPARISON TARGETS

# PERFORMANCE: LARGE-SCALE CODES
## COMPARISON TARGETS

Active Access

# PERFORMANCE: LARGE-SCALE CODES
## COMPARISON TARGETS

Active Access

AA-Int

# PERFORMANCE: LARGE-SCALE CODES
## COMPARISON TARGETS

Active Access

AA-Poll

AA-Int

# PERFORMANCE: LARGE-SCALE CODES
## COMPARISON TARGETS

Active Access

AA-Poll

AA-Int

AA-SP

# PERFORMANCE: LARGE-SCALE CODES
## COMPARISON TARGETS

Active Access

AA-Poll

AA-Int

AA-SP

RMA

# PERFORMANCE: LARGE-SCALE CODES
## COMPARISON TARGETS

Active Access

AA-Poll

AA-Int

AA-SP

RMA

CRAY

DMAPP

# Performance: Large-Scale Codes
## Comparison Targets

Active Access

AA-Poll

AA-Int

AA-SP

RMA

CRAY

DMAPP

INFINIBAND

# PERFORMANCE: LARGE-SCALE CODES
## COMPARISON TARGETS

Active Access

AA-Poll

AA-Int

AA-SP

RMA

Cell

DMAPP

INFINIBAND

# PERFORMANCE: LARGE-SCALE CODES
## COMPARISON TARGETS

Active Access

AA-Poll

AA-Int

AA-SP

RMA

IBM

Cell

CRAY

DMAPP

INFINIBAND

Mellanox
TECHNOLOGIES

RoCE

# PERFORMANCE: LARGE-SCALE CODES
## COMPARISON TARGETS

**Active Messages**

**Active Access**

AA-Poll

AA-Int

AA-SP

RMA

IBM

Cell

CRAY

DMAPP

INFINIBAND

Mellanox
TECHNOLOGIES

RoCE

# PERFORMANCE: LARGE-SCALE CODES
## COMPARISON TARGETS

Active Access

AA-Poll

AA-Int

AA-SP

Active Messages

AM

RMA

Cell

DMAPP

RoCE

# PERFORMANCE: LARGE-SCALE CODES
## COMPARISON TARGETS

**Active Access**

AA-Poll

AA-Int

AA-SP

**Active Messages**

AM

AM-Exp

RMA



Cell

DMAPP

RoCE

# PERFORMANCE: LARGE-SCALE CODES
## COMPARISON TARGETS

**Active Access**

AA-Poll

AA-Int

AA-SP

**Active Messages**

AM  AM-Onload

AM-Exp

RMA



Cell

DMAPP



RoCE

# PERFORMANCE: LARGE-SCALE CODES
## COMPARISON TARGETS

Active Access

AA-Poll

AA-Int

AA-SP

Active Messages

AM  AM-Onload

AM-Exp  AM-Ints

RMA


Cell

DMAPP



RoCE

# PERFORMANCE: LARGE-SCALE CODES
## COMPARISON TARGETS

Active Access

AA-Poll

AA-Int

AA-SP

Active Messages

AM  AM-Onload

AM-Exp  AM-Ints

**IBM**

DCMF  LAPI

PAMI

RMA

**IBM**

Cell

**CRAY**

DMAPP

INFINIBAND

**Mellanox**
TECHNOLOGIES

RoCE

# PERFORMANCE: LARGE-SCALE CODES
## COMPARISON TARGETS

# PERFORMANCE: LARGE-SCALE CODES
## COMPARISON TARGETS

Active Access

AA-Poll
AA-Int
AA-SP

Active Messages

AM  AM-Onload
AM-Exp  AM-Ints

IBM

DCMF  LAPI
PAMI

RMA

IBM
Cell

CRAY
DMAPP

Mellanox
TECHNOLOGIES

INFINIBAND

RoCE

Myricom MX

AM++GASNet

# PERFORMANCE: LARGE-SCALE CODES
## DISTRIBUTED HASHTABLE

Collisions: 5%

Collisions: 25%

# CONCLUSIONS

# CONCLUSIONS

**Active Access**

# CONCLUSIONS

## Active Access



Alleviates RMA's problems with AMs
while preserving one-sided semantics

# CONCLUSIONS

## Active Access



Uses commodity
& common IOMMUs

Alleviates RMA's problems with AMs
while preserving one-sided semantics

# CONCLUSIONS

## Active Access



Uses commodity
& common IOMMUs

Alleviates RMA's problems with AMs while preserving one-sided semantics

Extends paging capabilities in a distributed environment

# CONCLUSIONS

## Active Access



Uses commodity
& common IOMMUs

Alleviates RMA's problems with AMs while preserving one-sided semantics

Extends paging capabilities in a distributed environment

## Data-centric programming

# CONCLUSIONS

**Active Access**



Uses commodity
& common IOMMUs

Alleviates RMA's problems with AMs
while preserving one-sided semantics

Extends paging capabilities
in a distributed environment

**Data-centric programming**



Addresses of pages guide
the execution of handlers

# CONCLUSIONS

## Active Access



Uses commodity
& common IOMMUs

Alleviates RMA's problems with AMs
while preserving one-sided semantics

Extends paging capabilities
in a distributed environment

## Data-centric programming



Hashtables, logging
schemes, counters, V-GAS,
checkpointing...

Addresses of pages guide
the execution of handlers

# CONCLUSIONS

## Active Access



Alleviates RMA's problems with AMs while preserving one-sided semantics

Uses commodity & common IOMMUs



Extends paging capabilities in a distributed environment

## Data-centric programming



Addresses of pages guide the execution of handlers

Hashtables, logging schemes, counters, V-GAS, checkpointing...

## Performance

# CONCLUSIONS

## Active Access



Alleviates RMA's problems with AMs while preserving one-sided semantics

Uses commodity & common IOMMUs



Extends paging capabilities in a distributed environment

## Data-centric programming



Addresses of pages guide the execution of handlers

Hashtables, logging schemes, counters, V-GAS, checkpointing...

## Performance



Accelerates various distributed codes

# Thank you
# for your attention

# ACTIVE ACCESS USE-CASES
## ACCELERATING LOGGING FOR RMA

# ACTIVE ACCESS USE-CASES

## ACCELERATING LOGGING FOR RMA

- Logging – a popular mechanism for fault-tolerance.

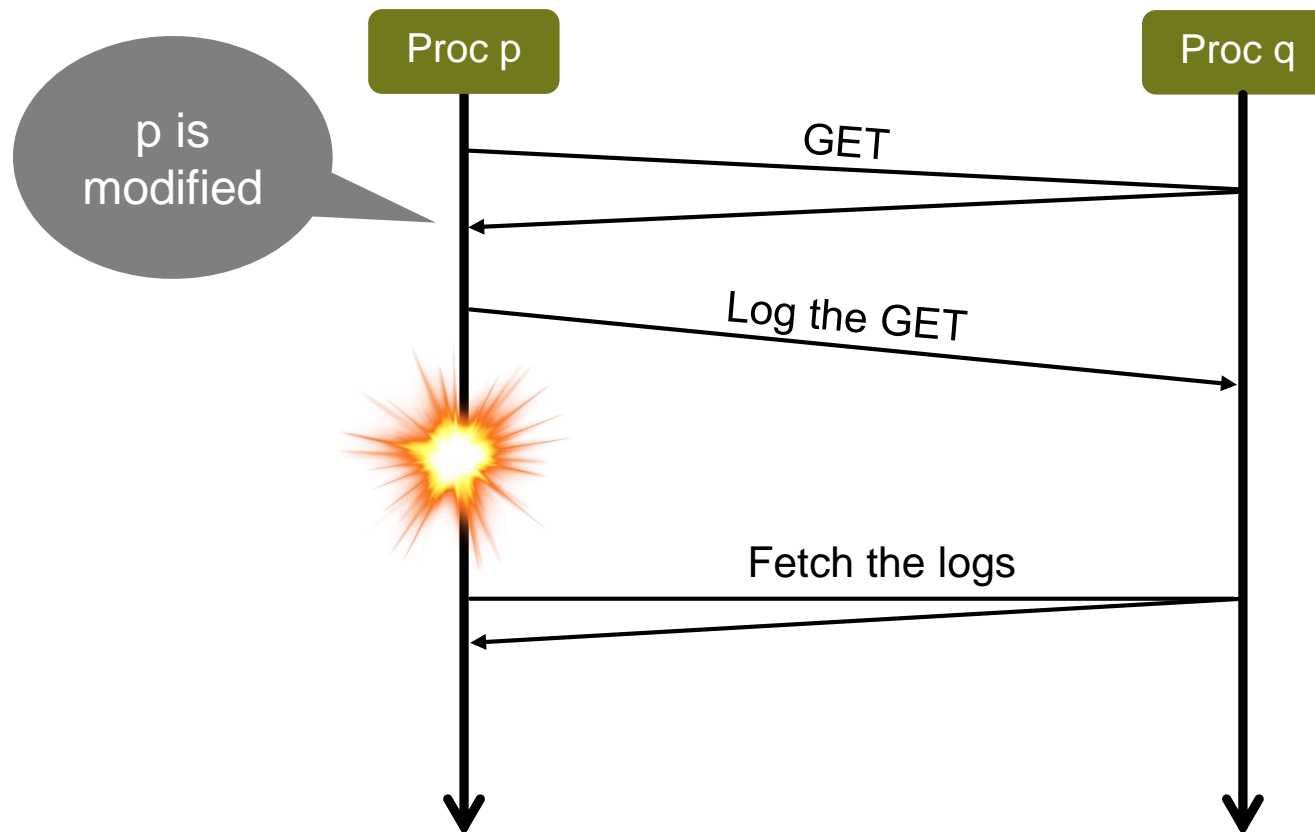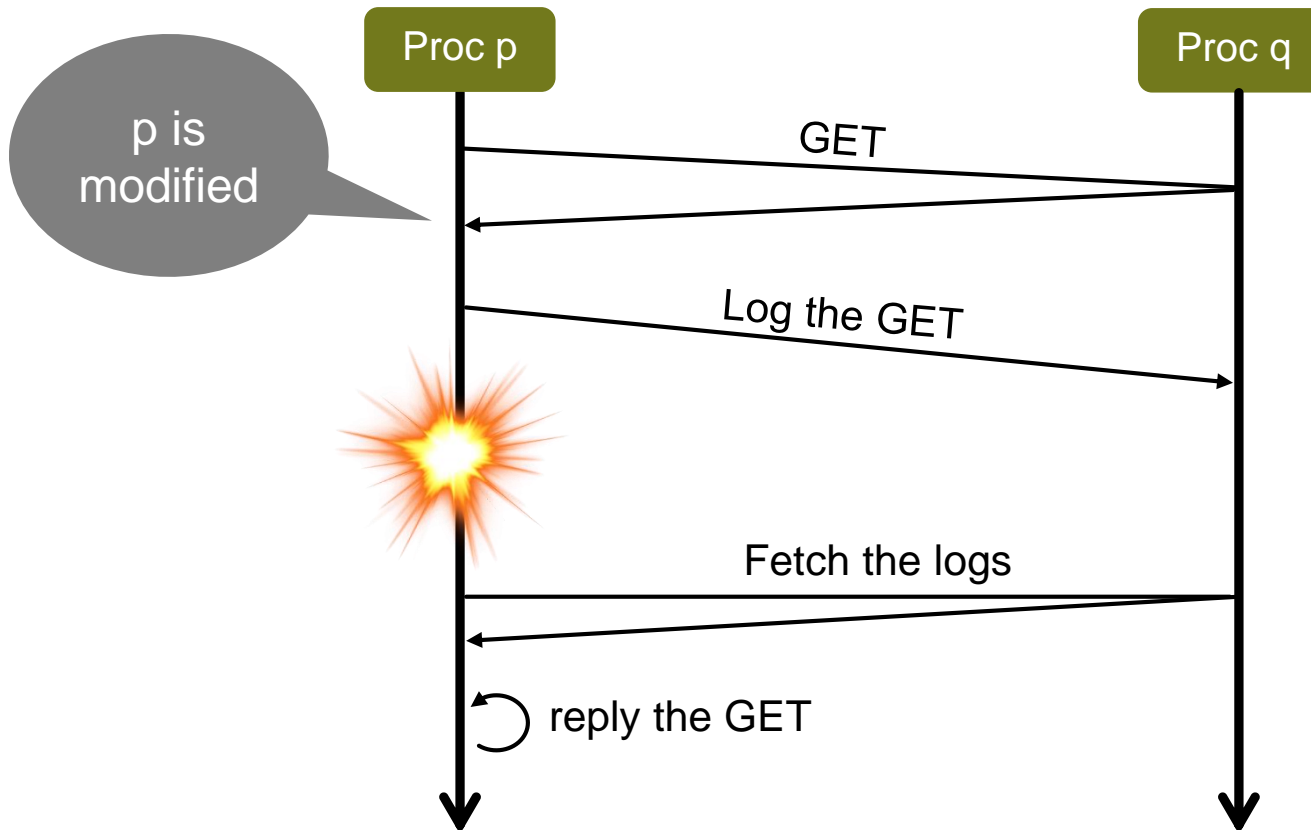# ACTIVE ACCESS USE-CASES

## ACCELERATING LOGGING FOR RMA

- Logging – a popular mechanism for fault-tolerance.
- Remote communication (puts/gets) is logged.

# ACTIVE ACCESS USE-CASES

## ACCELERATING LOGGING FOR RMA

- Logging – a popular mechanism for fault-tolerance.

- Remote communication (puts/gets) is logged.

- Upon a process crash, it is restored and uses the logs to replay its previous actions.

# ACTIVE ACCESS USE-CASES

## ACCELERATING LOGGING FOR RMA

- Logging – a popular mechanism for fault-tolerance.

- Remote communication (puts/gets) is logged.

- Upon a process crash, it is restored and uses the logs to replay its previous actions.
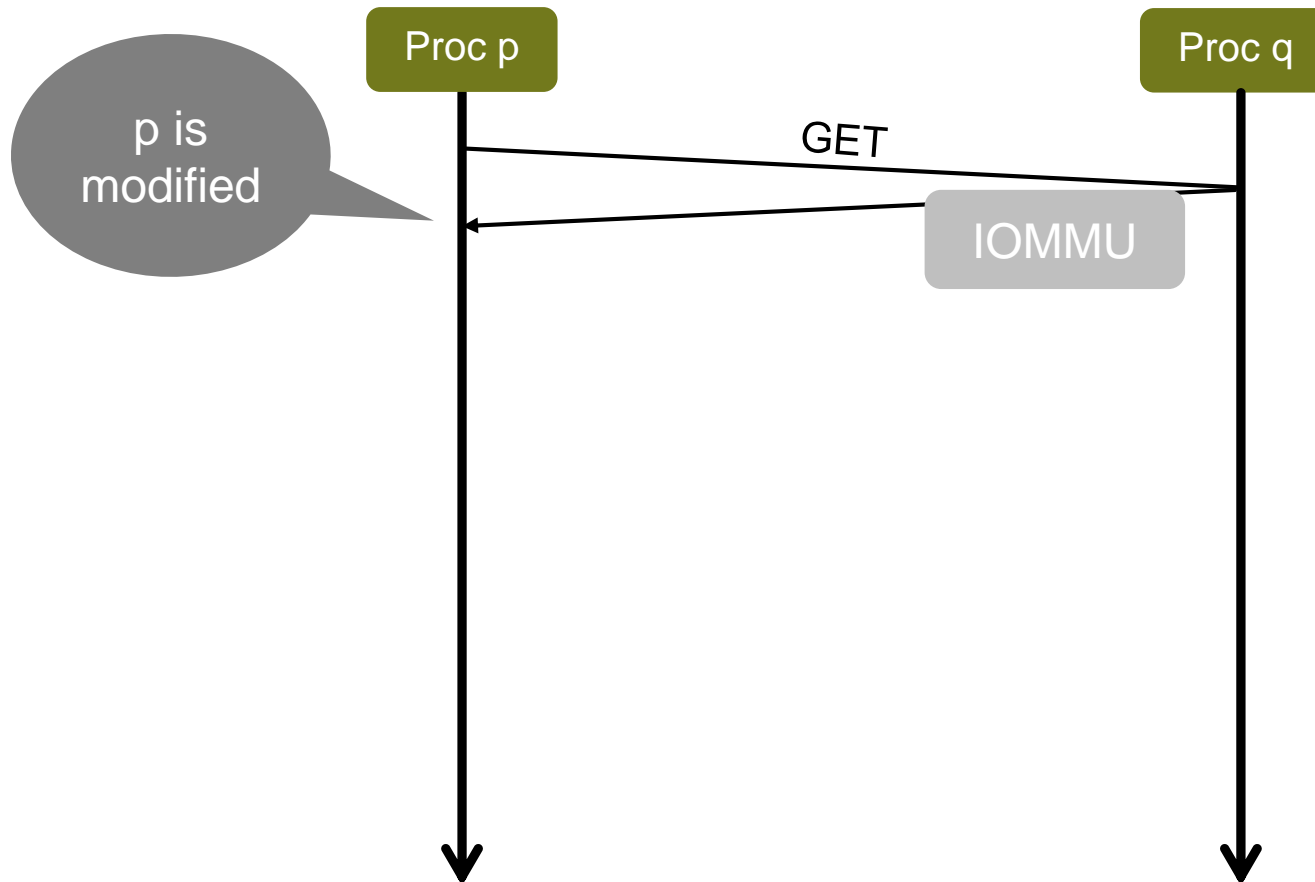
- Logs are stored in volatile memories.

# ACTIVE ACCESS USE-CASES

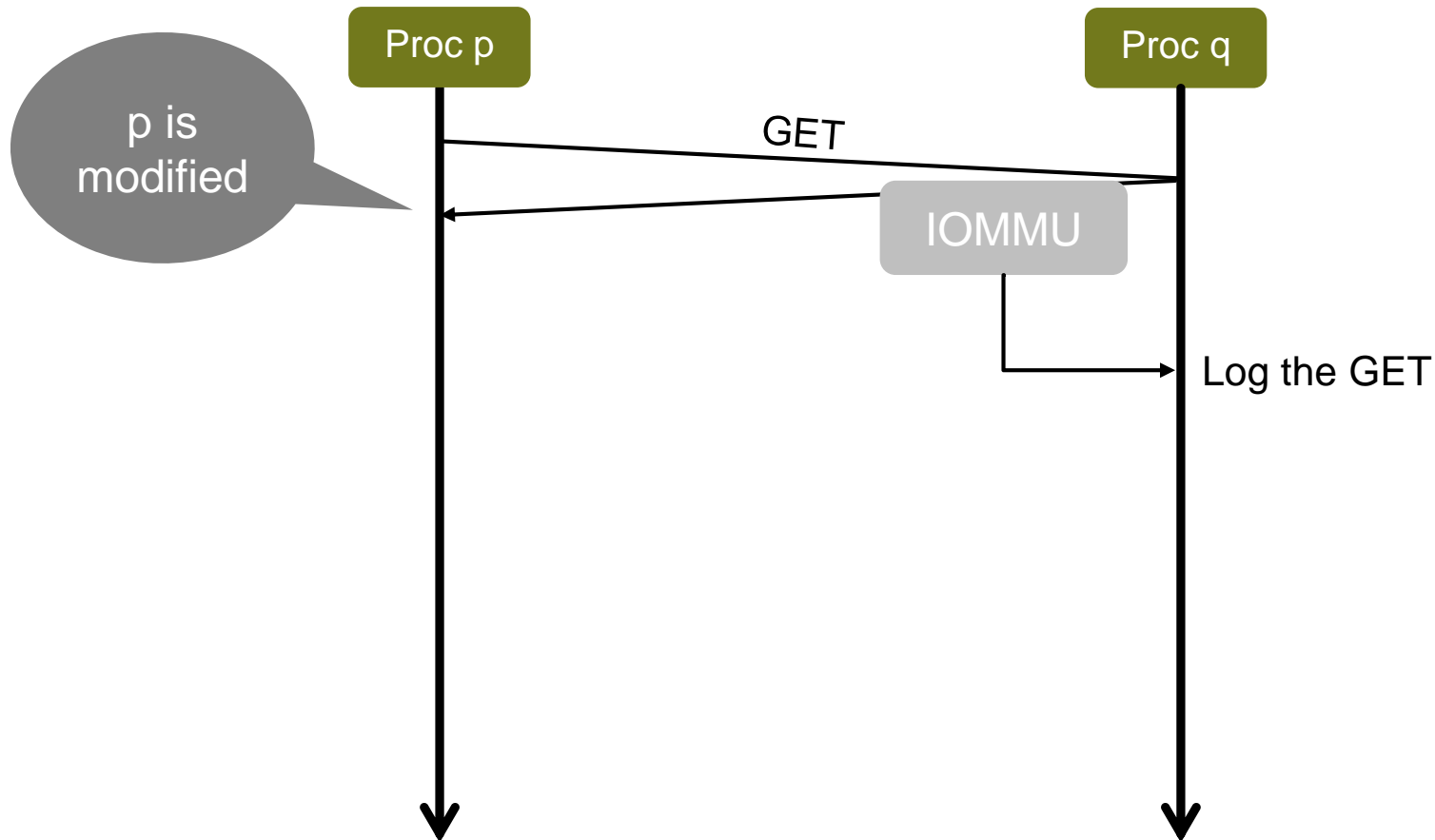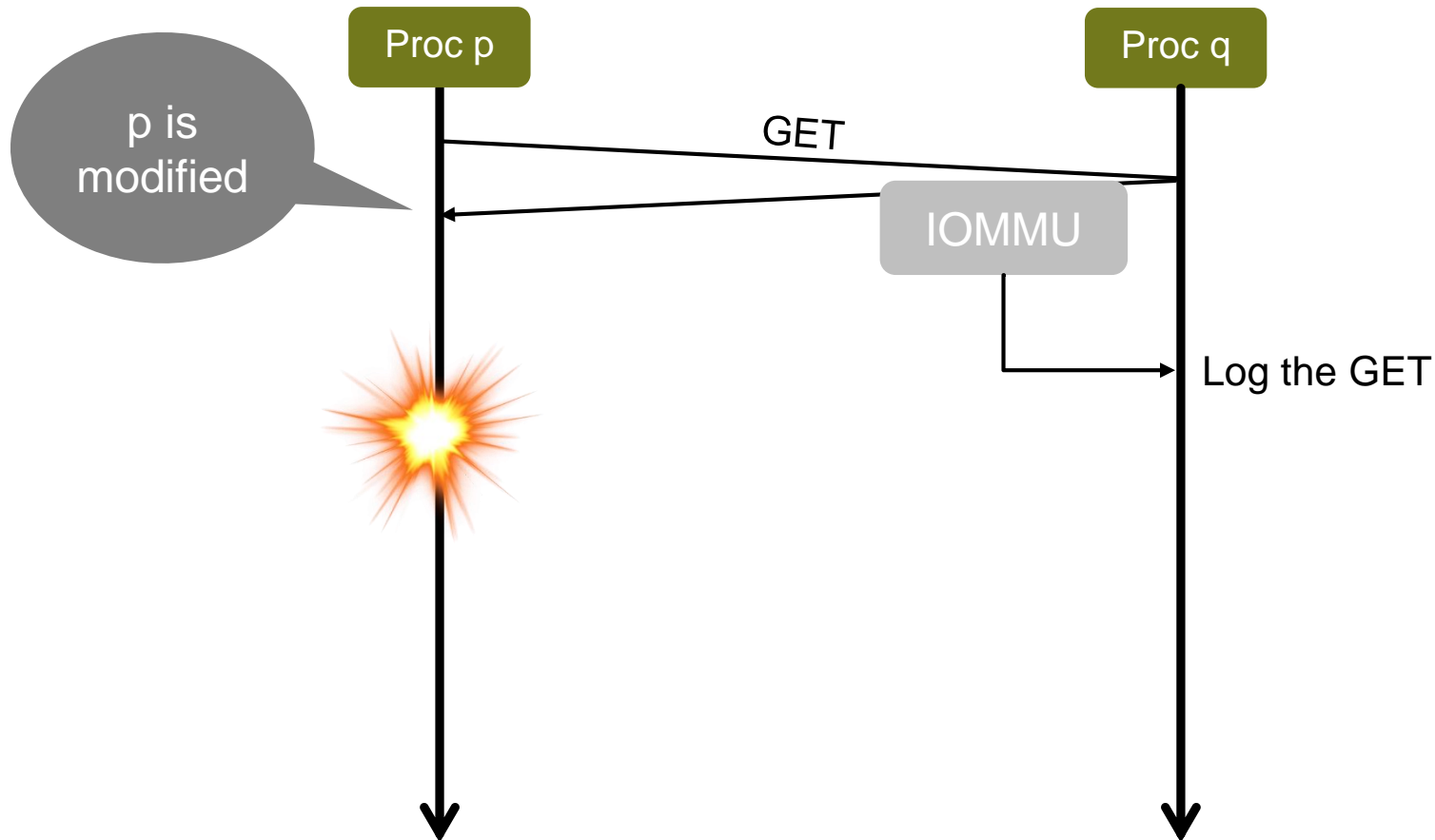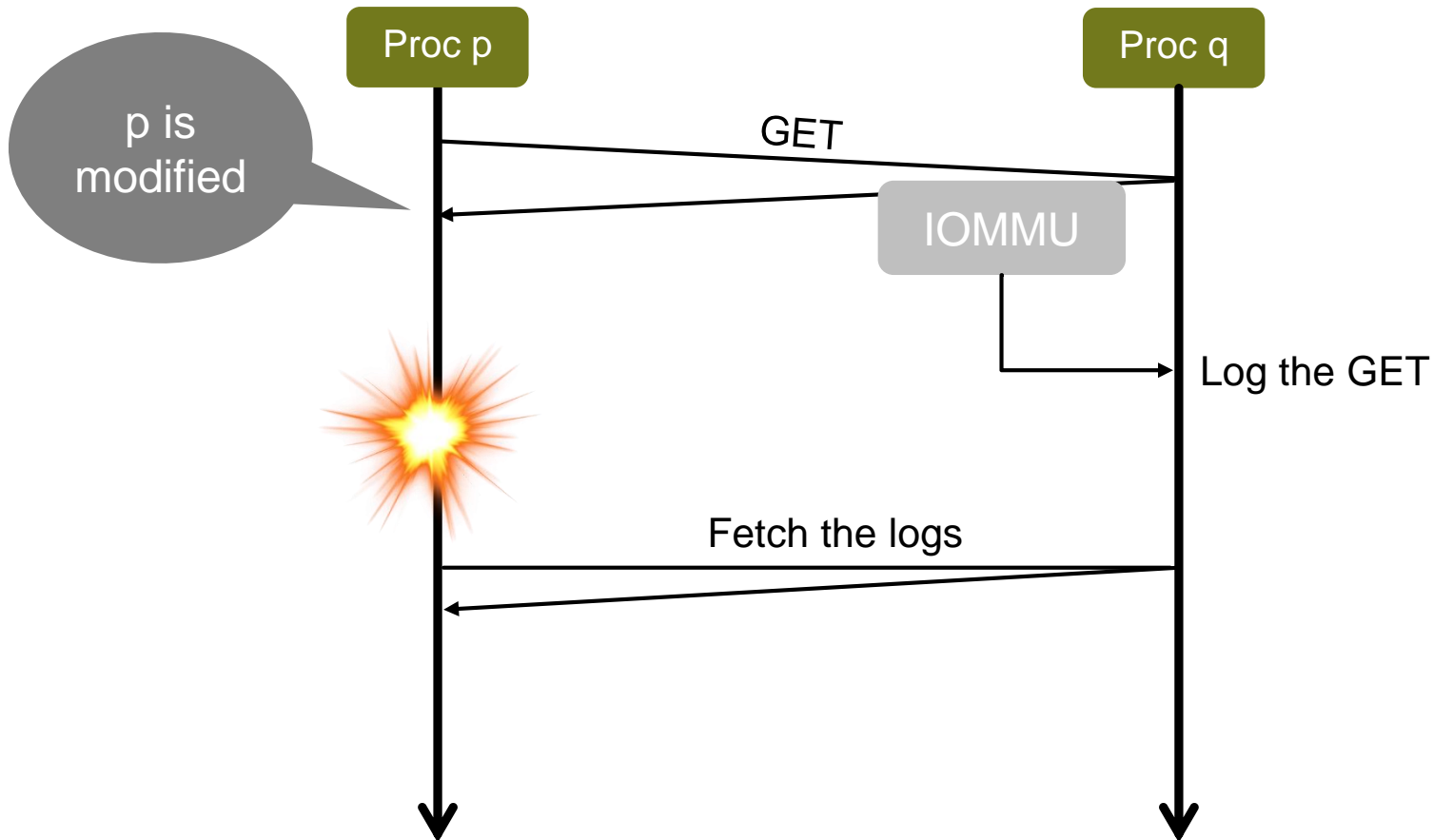## ACCELERATING LOGGING FOR RMA

# ACTIVE ACCESS USE-CASES

## ACCELERATING LOGGING FOR RMA

- Logging puts:

# ACTIVE ACCESS USE-CASES

## ACCELERATING LOGGING FOR RMA

- Logging puts:

# ACTIVE ACCESS USE-CASES

## ACCELERATING LOGGING FOR RMA

- Logging puts:

# ACTIVE ACCESS USE-CASES
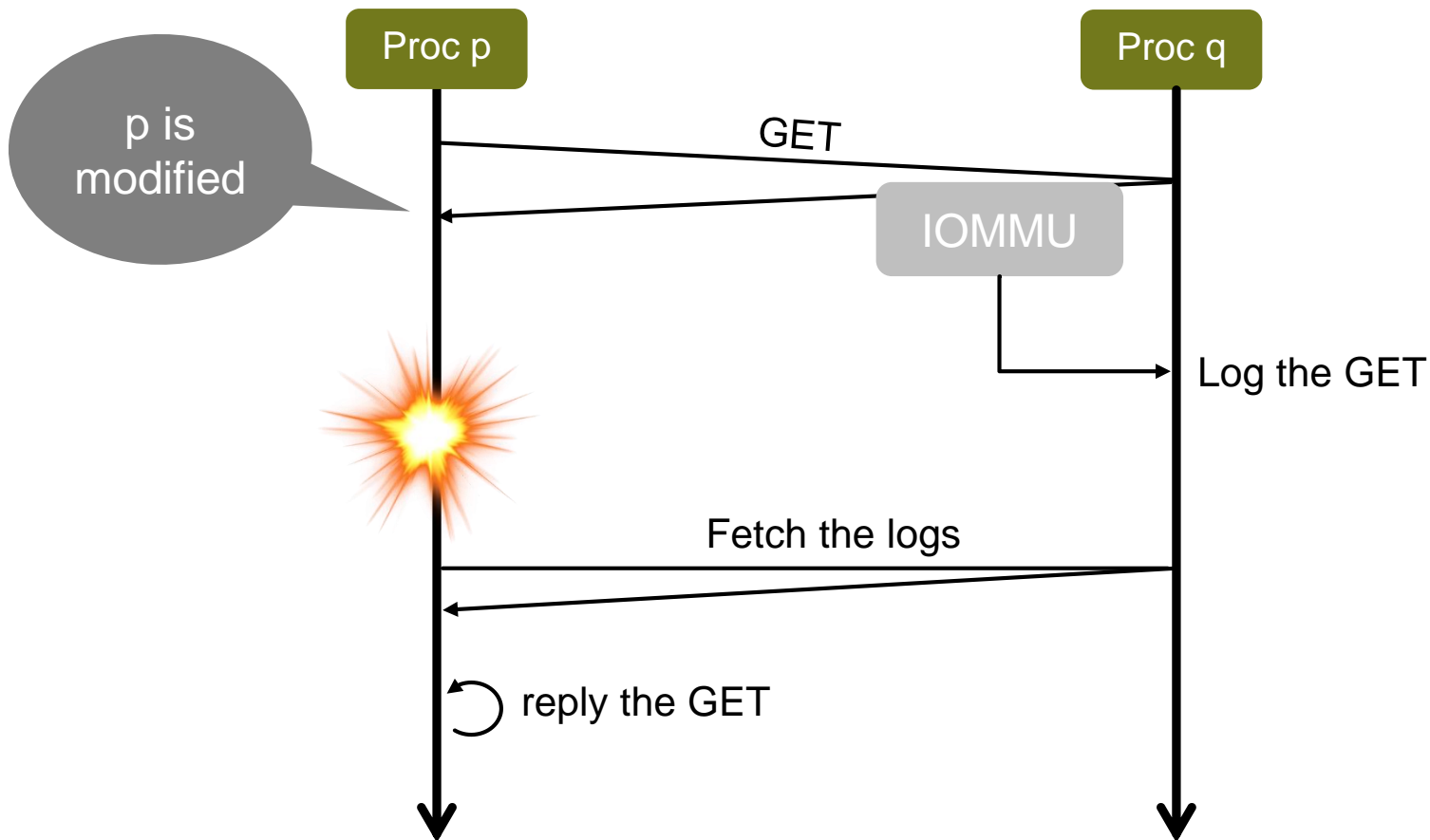## ACCELERATING LOGGING FOR RMA

- Logging puts:

# ACTIVE ACCESS USE-CASES
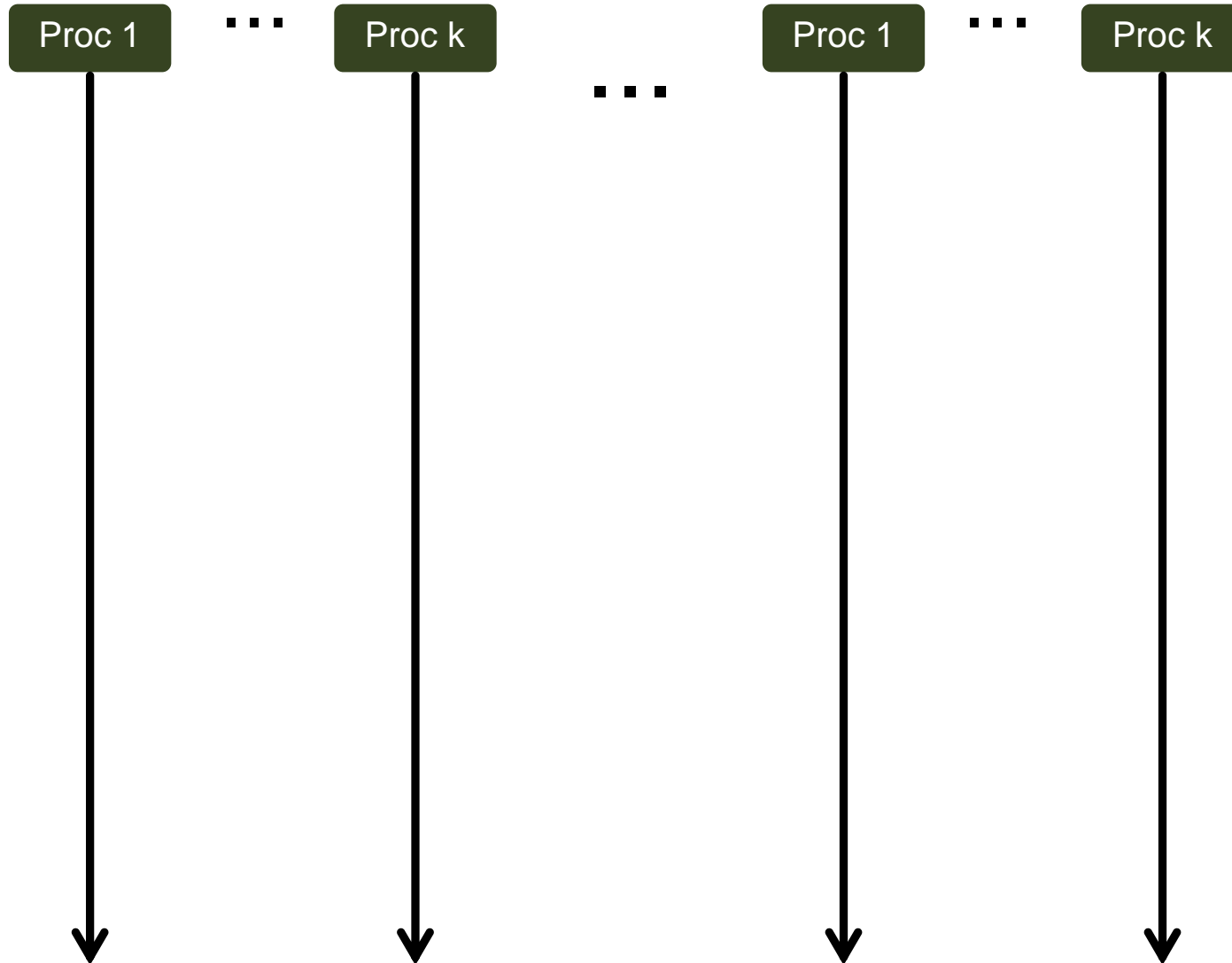## ACCELERATING LOGGING FOR RMA

- Logging puts:

# ACTIVE ACCESS USE-CASES
## ACCELERATING LOGGING FOR RMA

- Logging puts:

# ACTIVE ACCESS USE-CASES
## ACCELERATING LOGGING FOR RMA

- Logging puts:

# ACTIVE ACCESS USE-CASES
## ACCELERATING LOGGING FOR RMA

- Logging puts:

# ACTIVE ACCESS USE-CASES
## ACCELERATING LOGGING FOR RMA

- Logging gets (naive):

# ACTIVE ACCESS USE-CASES
## ACCELERATING LOGGING FOR RMA

- Logging gets (naive):

# ACTIVE ACCESS USE-CASES
## ACCELERATING LOGGING FOR RMA

- Logging gets (naive):

# ACTIVE ACCESS USE-CASES

## ACCELERATING LOGGING FOR RMA

- Logging gets (naive):

# ACTIVE ACCESS USE-CASES
## ACCELERATING LOGGING FOR RMA

- Logging gets (naive):

# ACTIVE ACCESS USE-CASES
## ACCELERATING LOGGING FOR RMA

- Logging gets (naive):

# ACTIVE ACCESS USE-CASES
## ACCELERATING LOGGING FOR RMA

- Logging gets (naive):

# ACTIVE ACCESS USE-CASES
## ACCELERATING LOGGING FOR RMA

- Logging gets (traditional) [1]:



[1] M. Besta and T. Hoefler. Fault tolerance for remote memory access programming models. HPDC'14.

# ACTIVE ACCESS USE-CASES
## ACCELERATING LOGGING FOR RMA

- Logging gets (traditional) [1]:



[1] M. Besta and T. Hoefler. Fault tolerance for remote memory access programming models. HPDC'14.

# ACTIVE ACCESS USE-CASES
## ACCELERATING LOGGING FOR RMA

- Logging gets (traditional) [1]:



[1] M. Besta and T. Hoefler. Fault tolerance for remote memory access programming models. HPDC'14.

# ACTIVE ACCESS USE-CASES
## ACCELERATING LOGGING FOR RMA

- Logging gets (traditional) [1]:



[1] M. Besta and T. Hoefler. Fault tolerance for remote memory access programming models. HPDC'14.

# ACTIVE ACCESS USE-CASES
## ACCELERATING LOGGING FOR RMA

- Logging gets (traditional) [1]:

[1] M. Besta and T. Hoefler. Fault tolerance for remote memory access programming models. HPDC'14.

# ACTIVE ACCESS USE-CASES
## ACCELERATING LOGGING FOR RMA

- Logging gets (traditional) [1]:



[1] M. Besta and T. Hoefler. Fault tolerance for remote memory access programming models. HPDC'14.

# ACTIVE ACCESS USE-CASES
## ACCELERATING LOGGING FOR RMA

- Logging gets (AA):

# ACTIVE ACCESS USE-CASES
## ACCELERATING LOGGING FOR RMA

- Logging gets (AA):

# ACTIVE ACCESS USE-CASES
## ACCELERATING LOGGING FOR RMA

- Logging gets (AA):

# ACTIVE ACCESS USE-CASES
## ACCELERATING LOGGING FOR RMA

- Logging gets (AA):

# ACTIVE ACCESS USE-CASES
## ACCELERATING LOGGING FOR RMA

- Logging gets (AA):

# ACTIVE ACCESS USE-CASES
## INCREMENTAL CHECKPOINTING FOR RMA

# ACTIVE ACCESS USE-CASES

## INCREMENTAL CHECKPOINTING FOR RMA

# ACTIVE ACCESS USE-CASES

## INCREMENTAL CHECKPOINTING FOR RMA

| Proc 1 | · · · | Proc k | · · · | Proc 1 | · · · | Proc k |

compute

compute

compute

compute

compute

compute

compute

compute

# ACTIVE ACCESS USE-CASES
## INCREMENTAL CHECKPOINTING FOR RMA

# ACTIVE ACCESS USE-CASES

## INCREMENTAL CHECKPOINTING FOR RMA

# ACTIVE ACCESS USE-CASES

## INCREMENTAL CHECKPOINTING FOR RMA

# ACTIVE ACCESS USE-CASES

## INCREMENTAL CHECKPOINTING FOR RMA

# ACTIVE ACCESS USE-CASES

## INCREMENTAL CHECKPOINTING FOR RMA

# ACTIVE ACCESS USE-CASES

## INCREMENTAL CHECKPOINTING FOR RMA

# ACTIVE ACCESS USE-CASES
## INCREMENTAL CHECKPOINTING FOR RMA

# ACTIVE ACCESS USE-CASES
## INCREMENTAL CHECKPOINTING FOR RMA

Proc 1 ... Proc k ... Proc 1 ... Proc k

barrier

compute compute compute compute

barrier

compute compute compute compute

# ACTIVE ACCESS USE-CASES
## INCREMENTAL CHECKPOINTING FOR RMA

# ACTIVE ACCESS USE-CASES
## INCREMENTAL CHECKPOINTING FOR RMA

# ACTIVE ACCESS USE-CASES

## INCREMENTAL CHECKPOINTING FOR RMA

# ACTIVE ACCESS USE-CASES
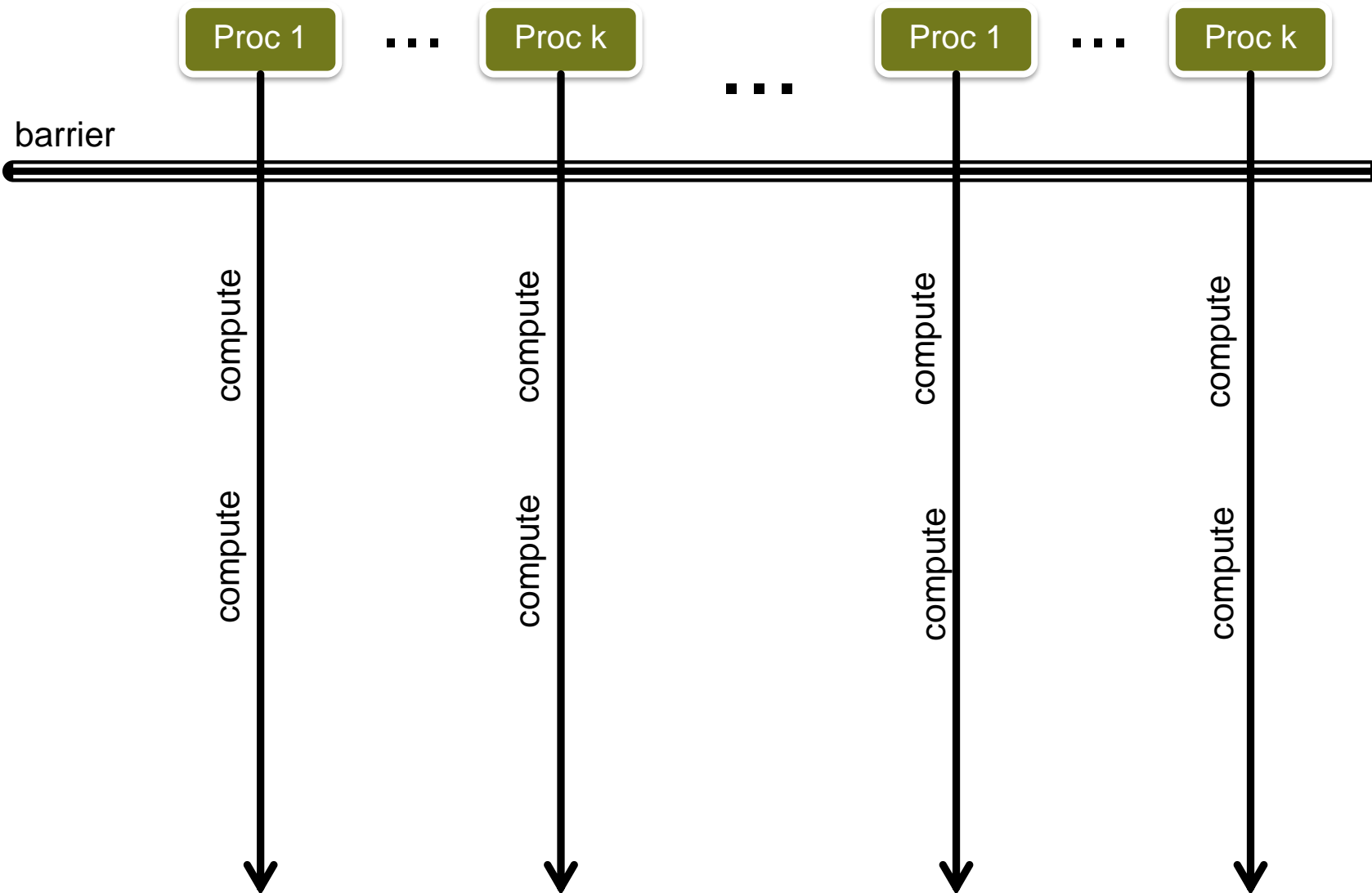## INCREMENTAL CHECKPOINTING FOR RMA

# COORDINATED CHECKPOINTING (MP)

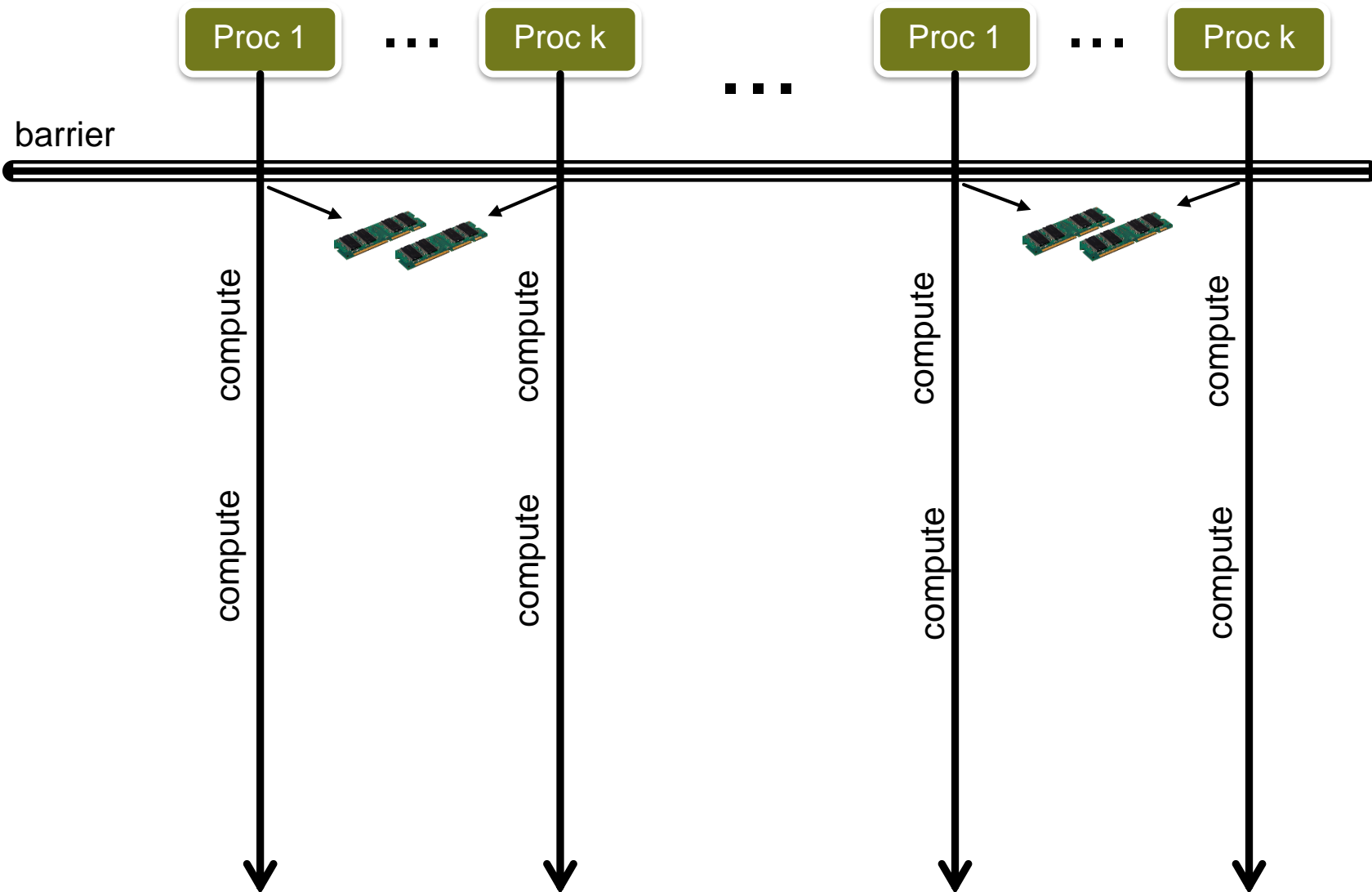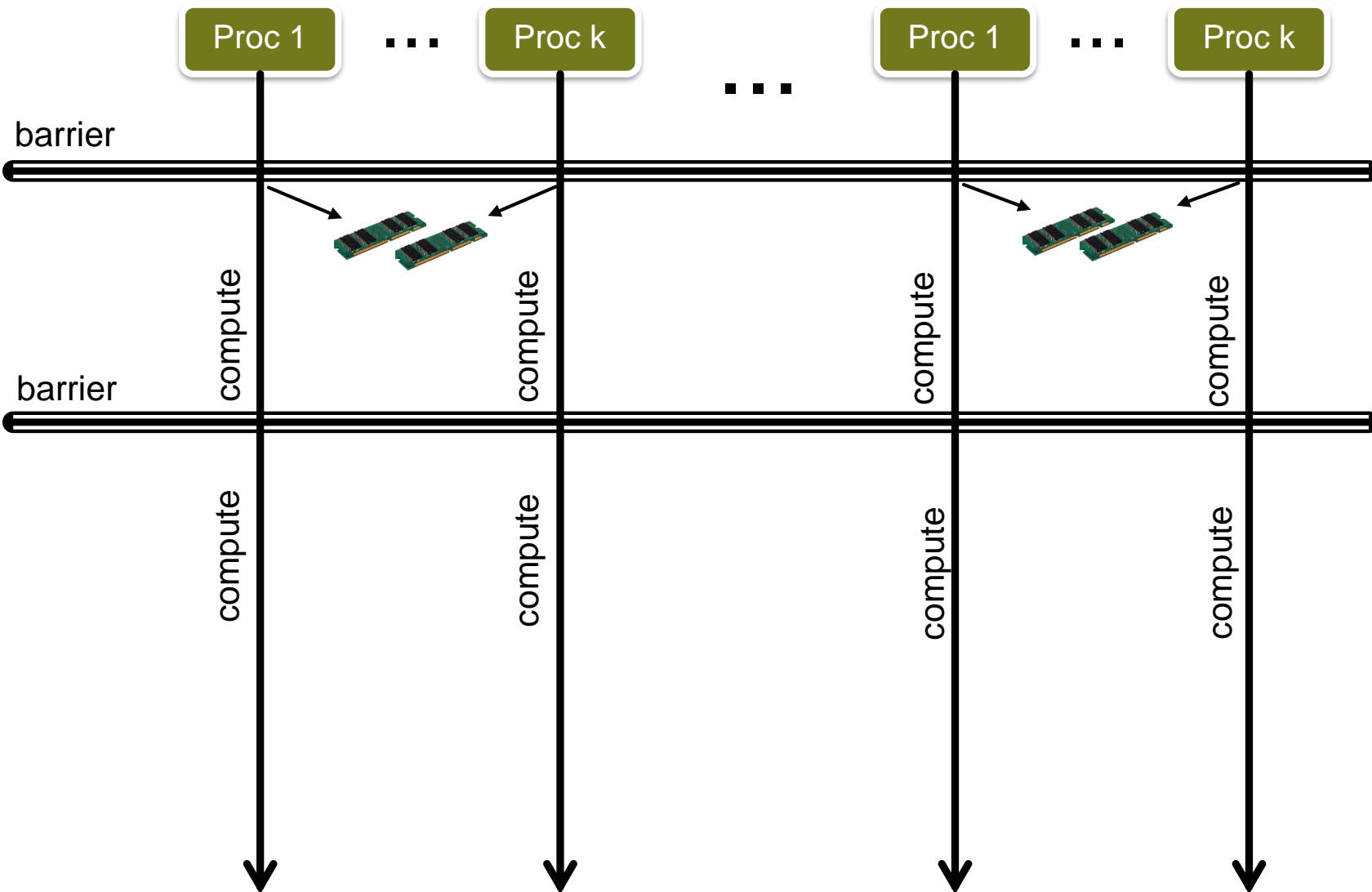# Coordinated Checkpointing (MP)
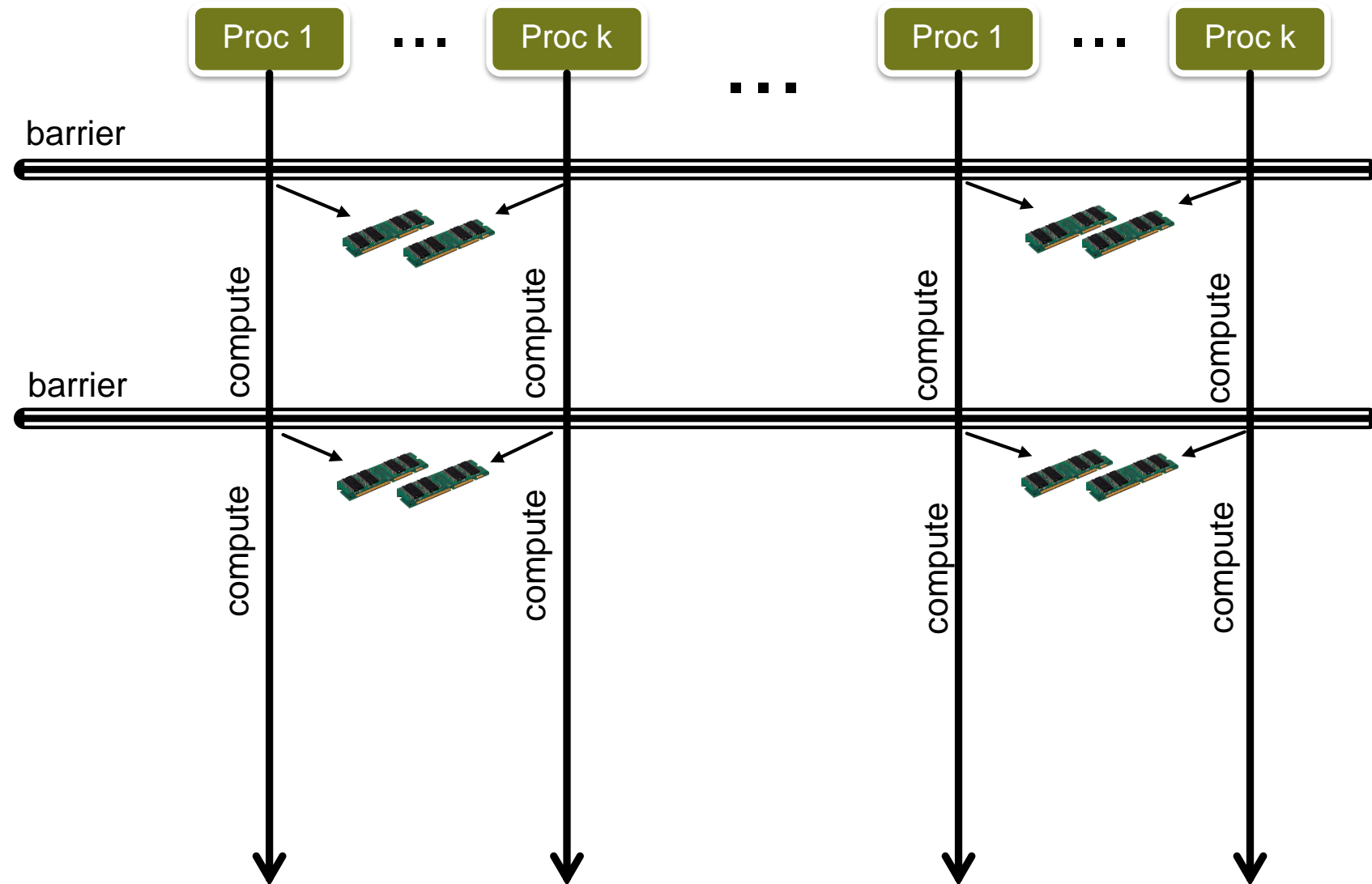
# COORDINATED CHECKPOINTING (MP)

# COORDINATED CHECKPOINTING (MP)

# COORDINATED CHECKPOINTING (MP)

# COORDINATED CHECKPOINTING (MP)

# COORDINATED CHECKPOINTING (MP)

# COORDINATED CHECKPOINTING (MP)

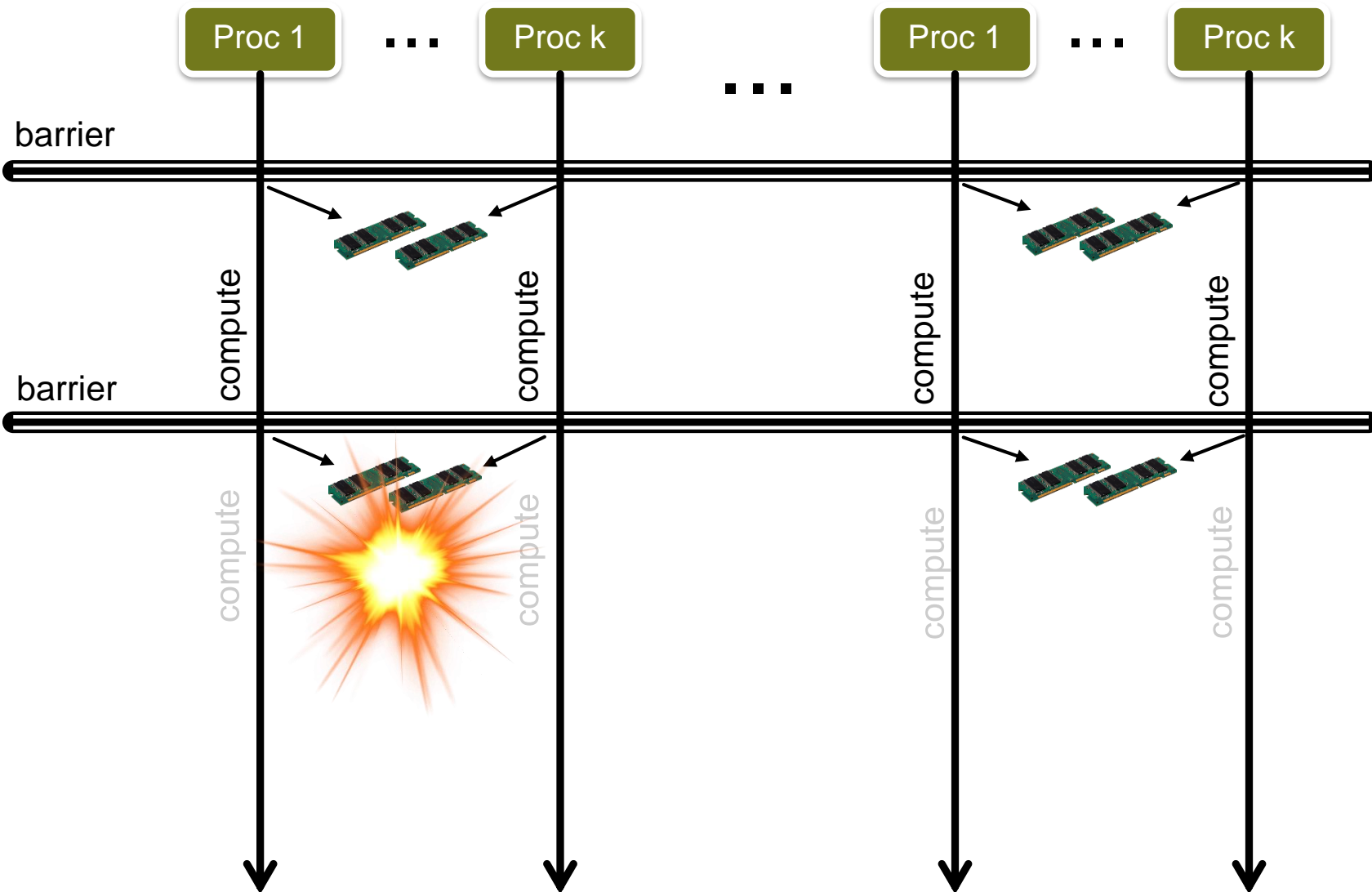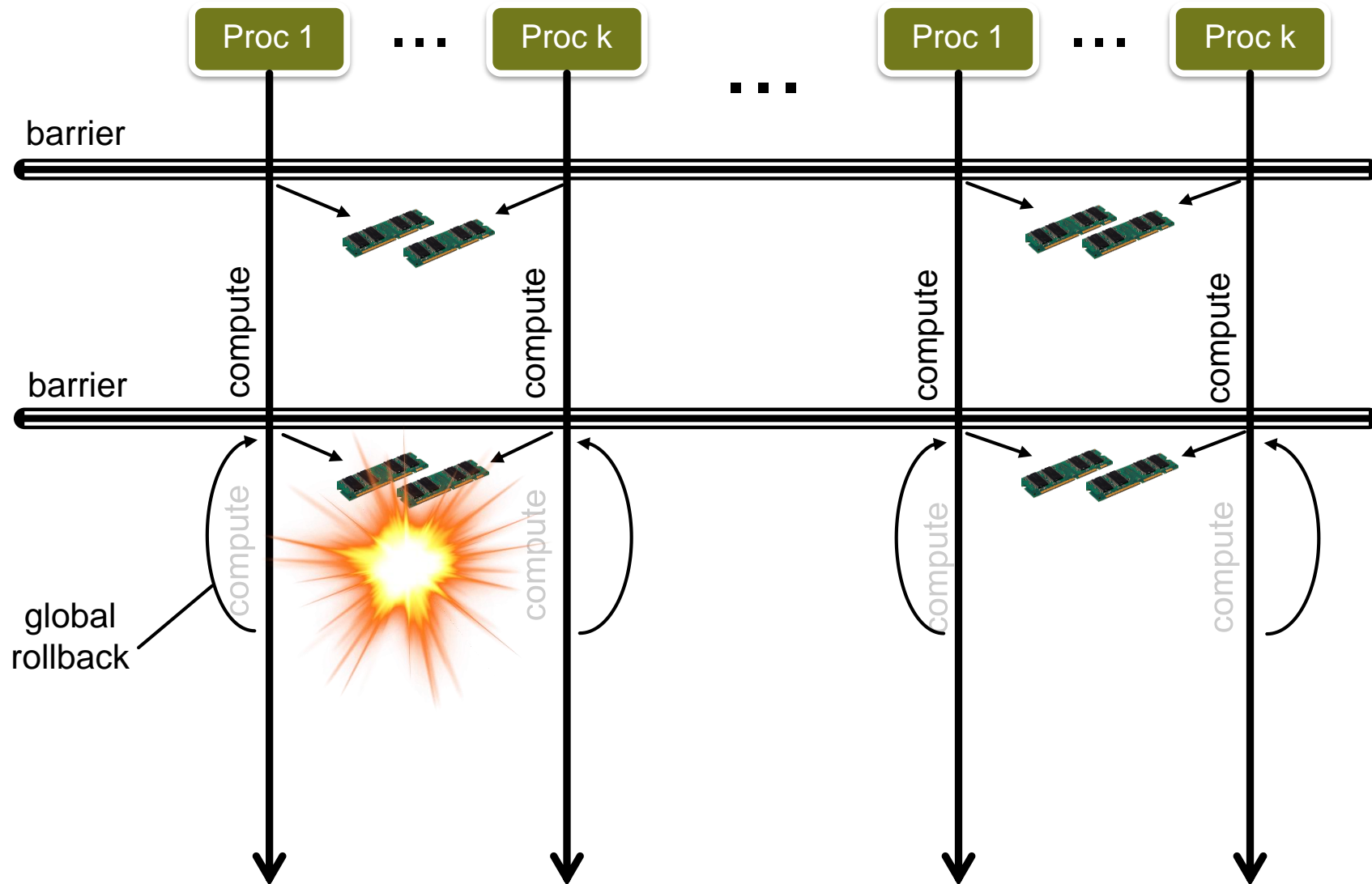# COORDINATED CHECKPOINTING (MP)

# COORDINATED CHECKPOINTING (MP)

# COORDINATED CHECKPOINTING (MP)

# COORDINATED CHECKPOINTING (MP)

# COORDINATED CHECKPOINTING (MP)

# COORDINATED CHECKPOINTING (MP)

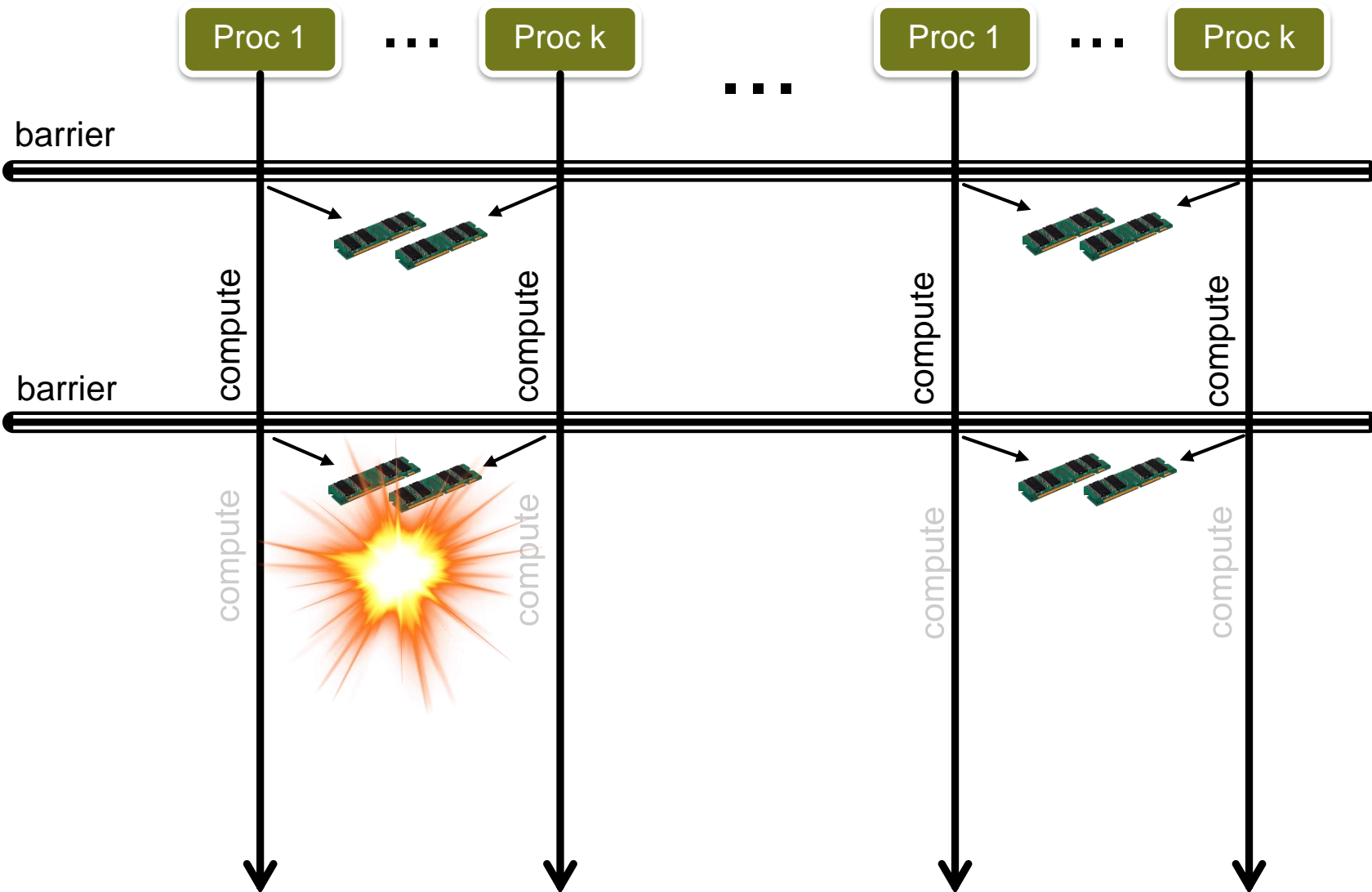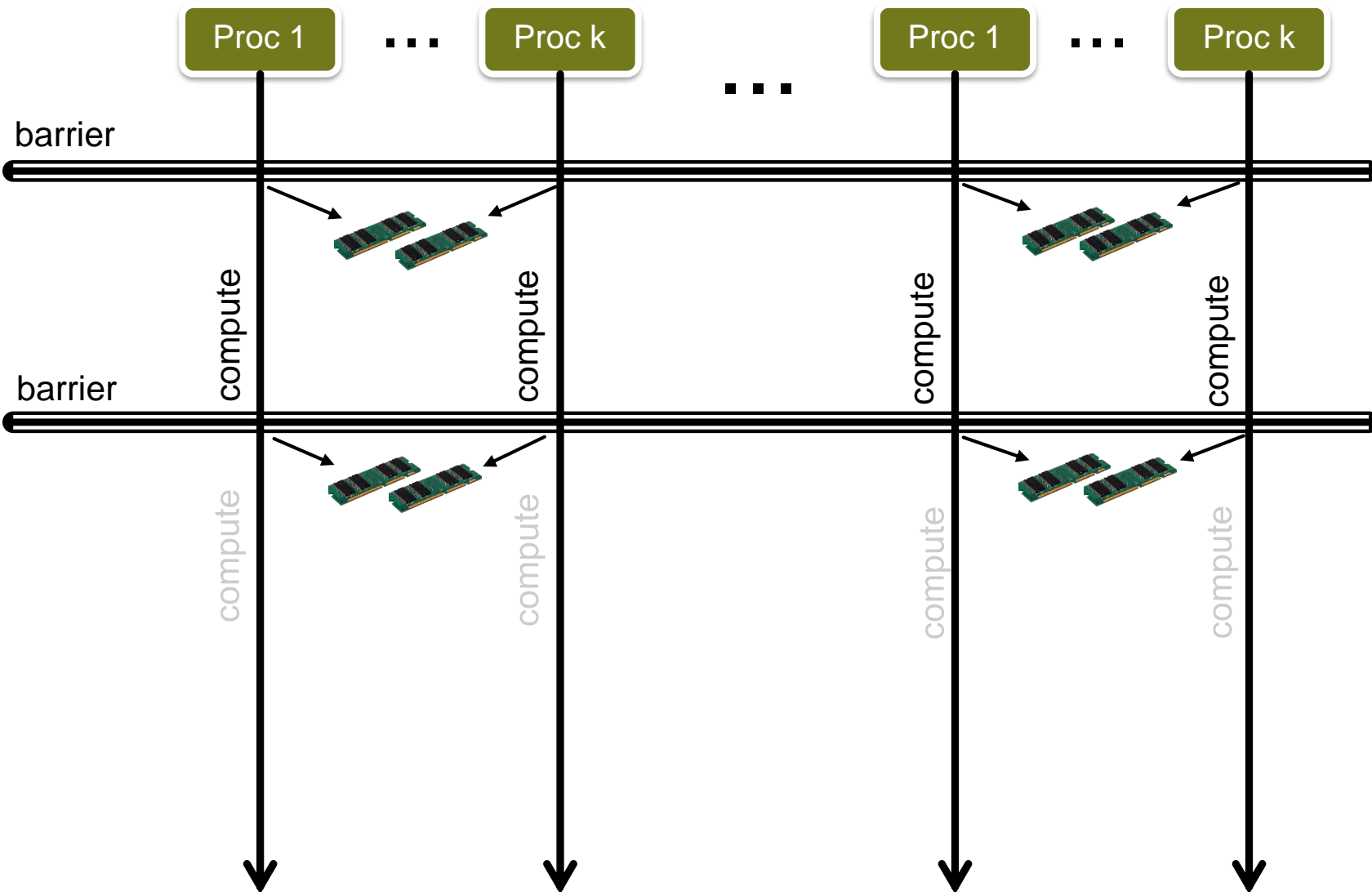# COORDINATED CHECKPOINTING (MP)

# COORDINATED CHECKPOINTING (MP)

# COORDINATED CHECKPOINTING (MP)

# COORDINATED CHECKPOINTING (MP)

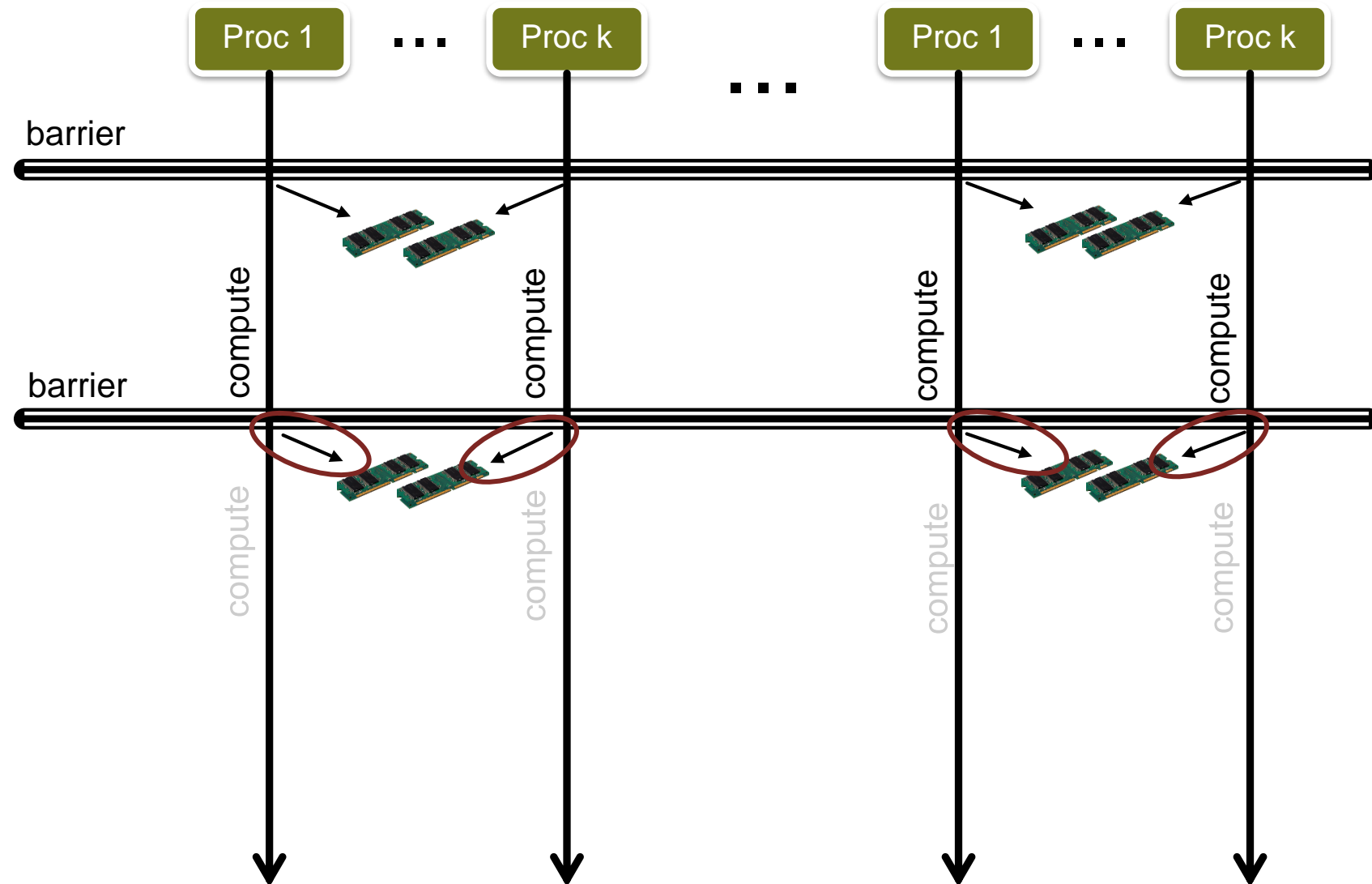# COORDINATED CHECKPOINTING (MP)

# COORDINATED CHECKPOINTING (MP)

# COORDINATED CHECKPOINTING (MP)

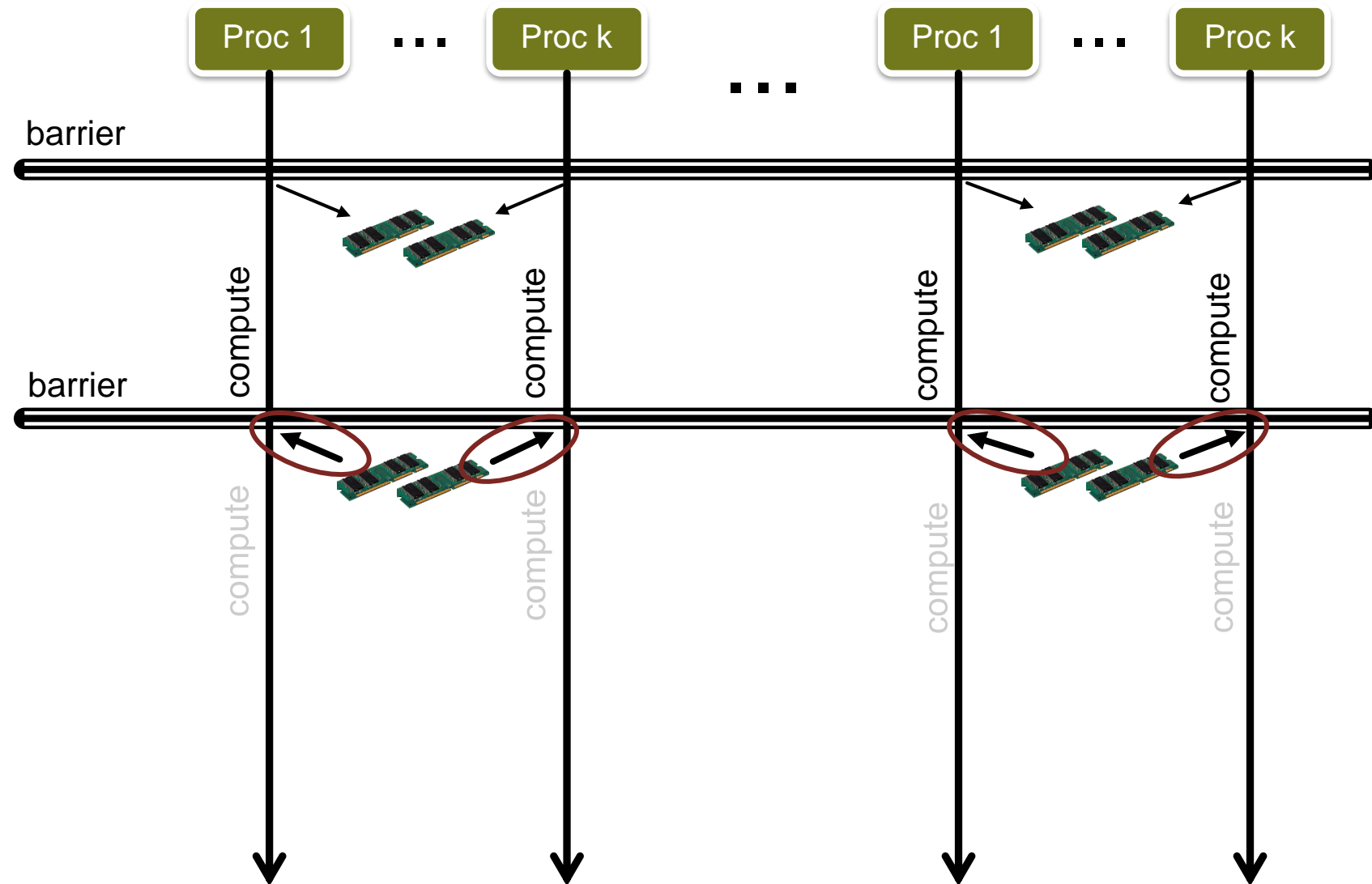# COORDINATED CHECKPOINTING (MP)

# COORDINATED CHECKPOINTING (MP)

# COORDINATED CHECKPOINTING (MP)

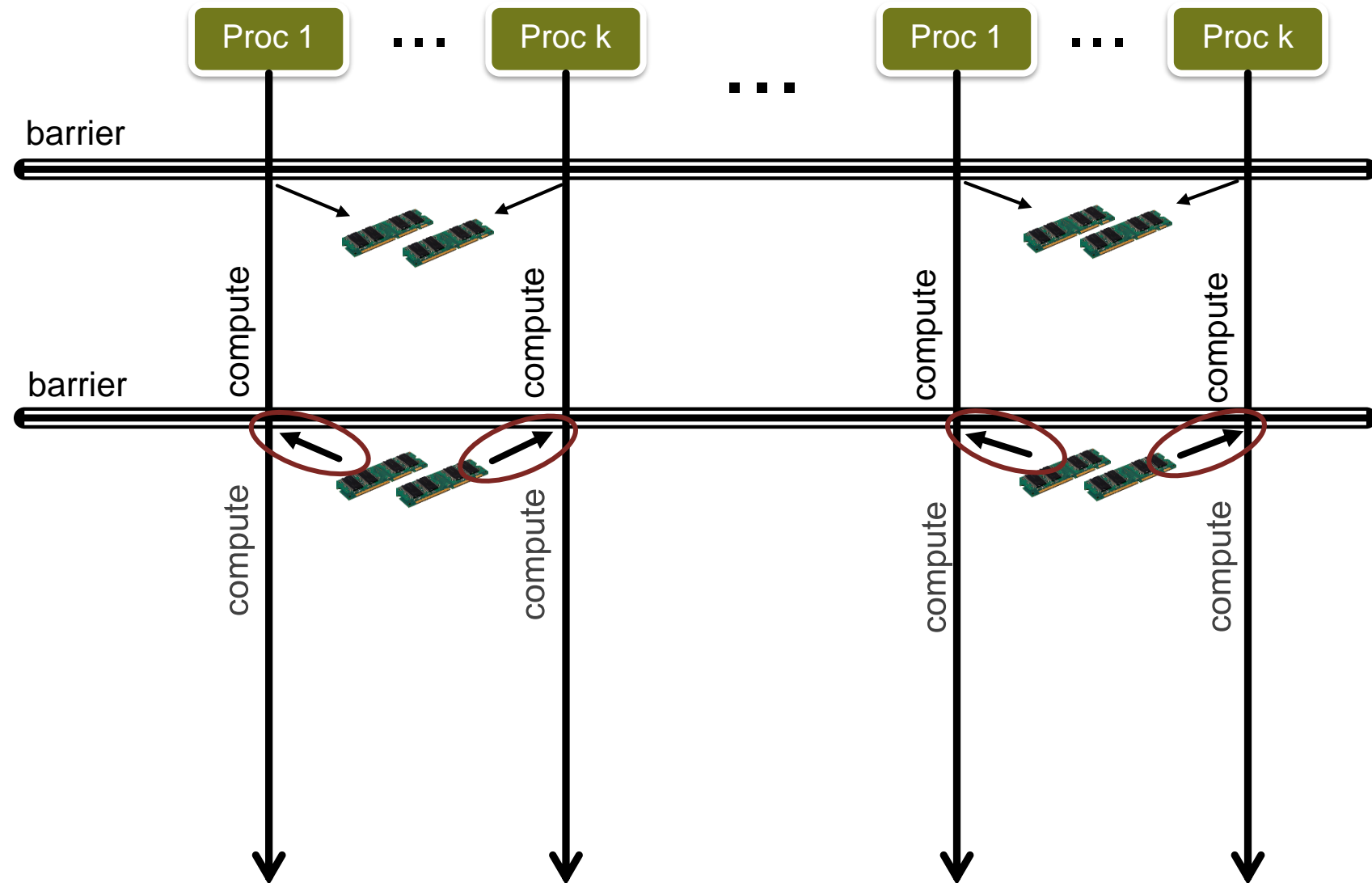# COORDINATED CHECKPOINTING (MP)

# COORDINATED CHECKPOINTING (MP)

# COORDINATED CHECKPOINTING (MP)

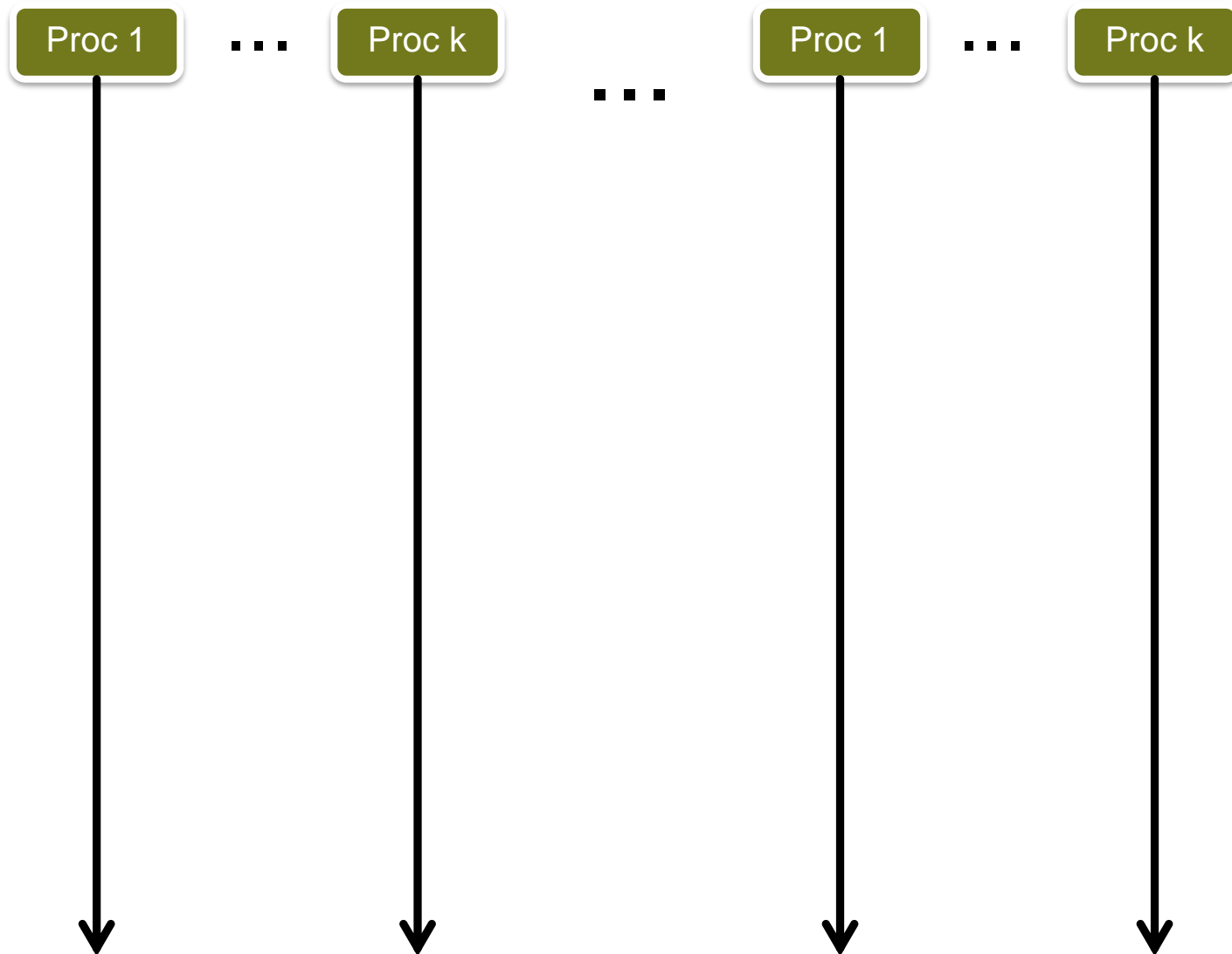# COORDINATED CHECKPOINTING (MP)

# COORDINATED CHECKPOINTING (MP)
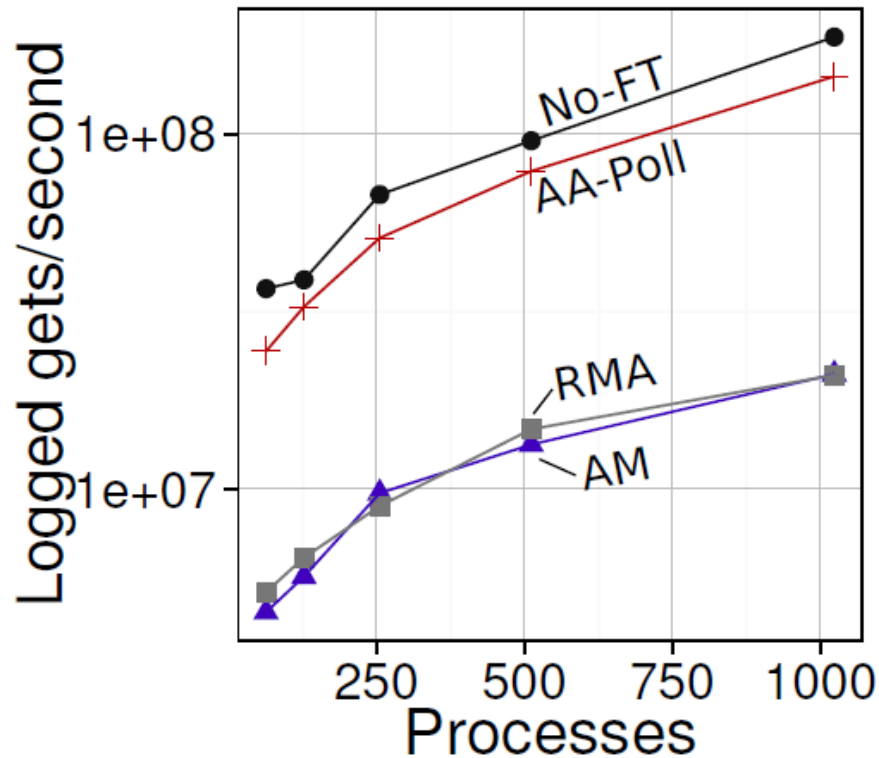
# COORDINATED CHECKPOINTING (MP)

# PERFORMANCE: LARGE-SCALE CODES
## FAULT TOLERANCE SCHEME

Logging gets:

Sorting time: