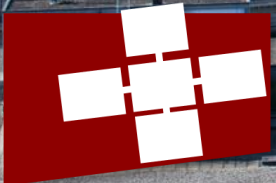


# Active Access: A Mechanism for High-Performance Distributed Data-Centric Computations

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

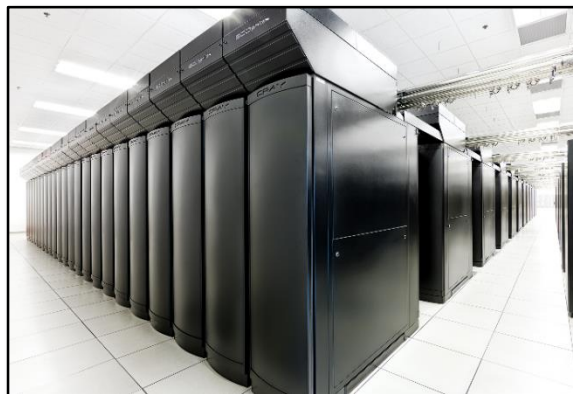
Memory

A

Process q

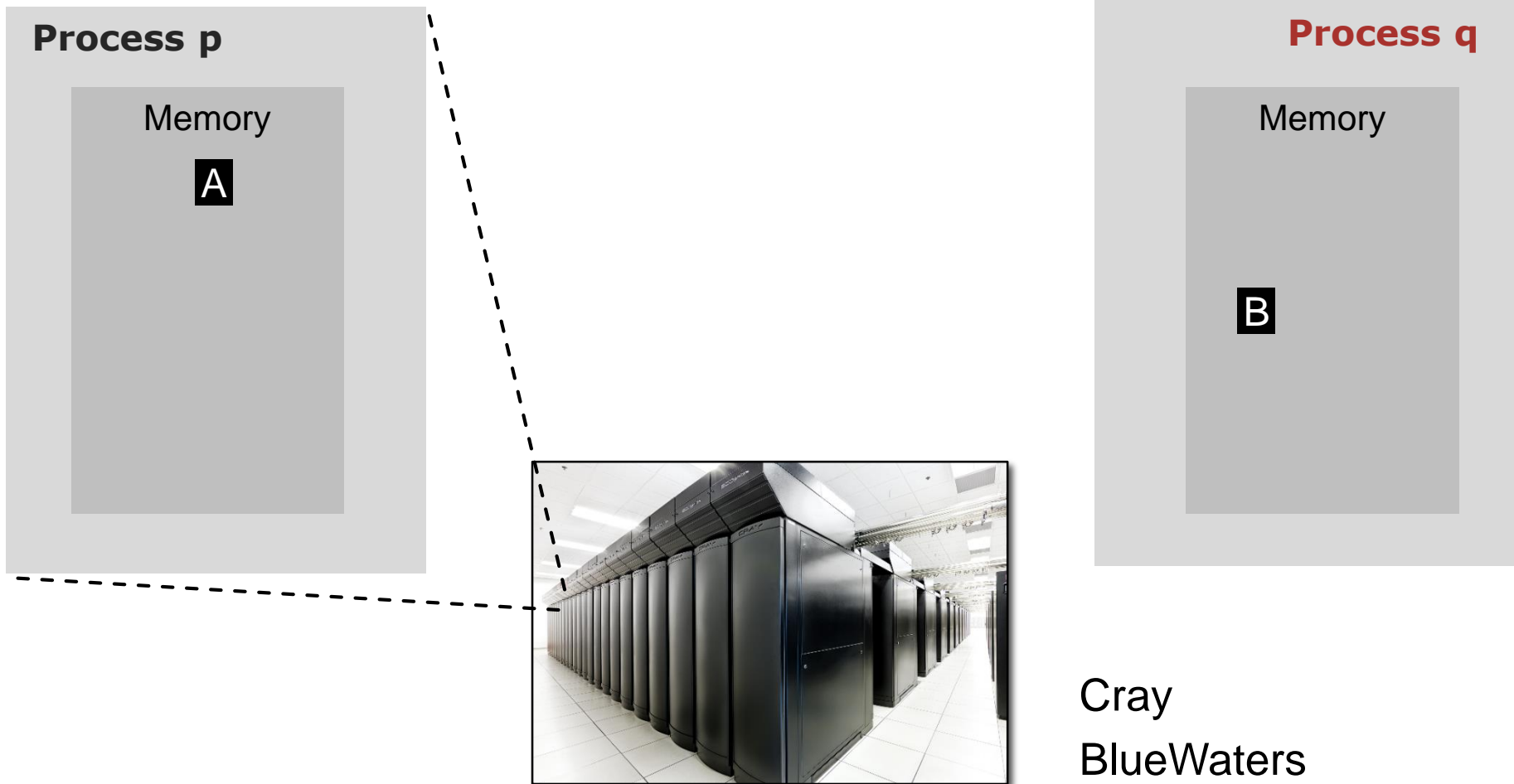
Memory

B

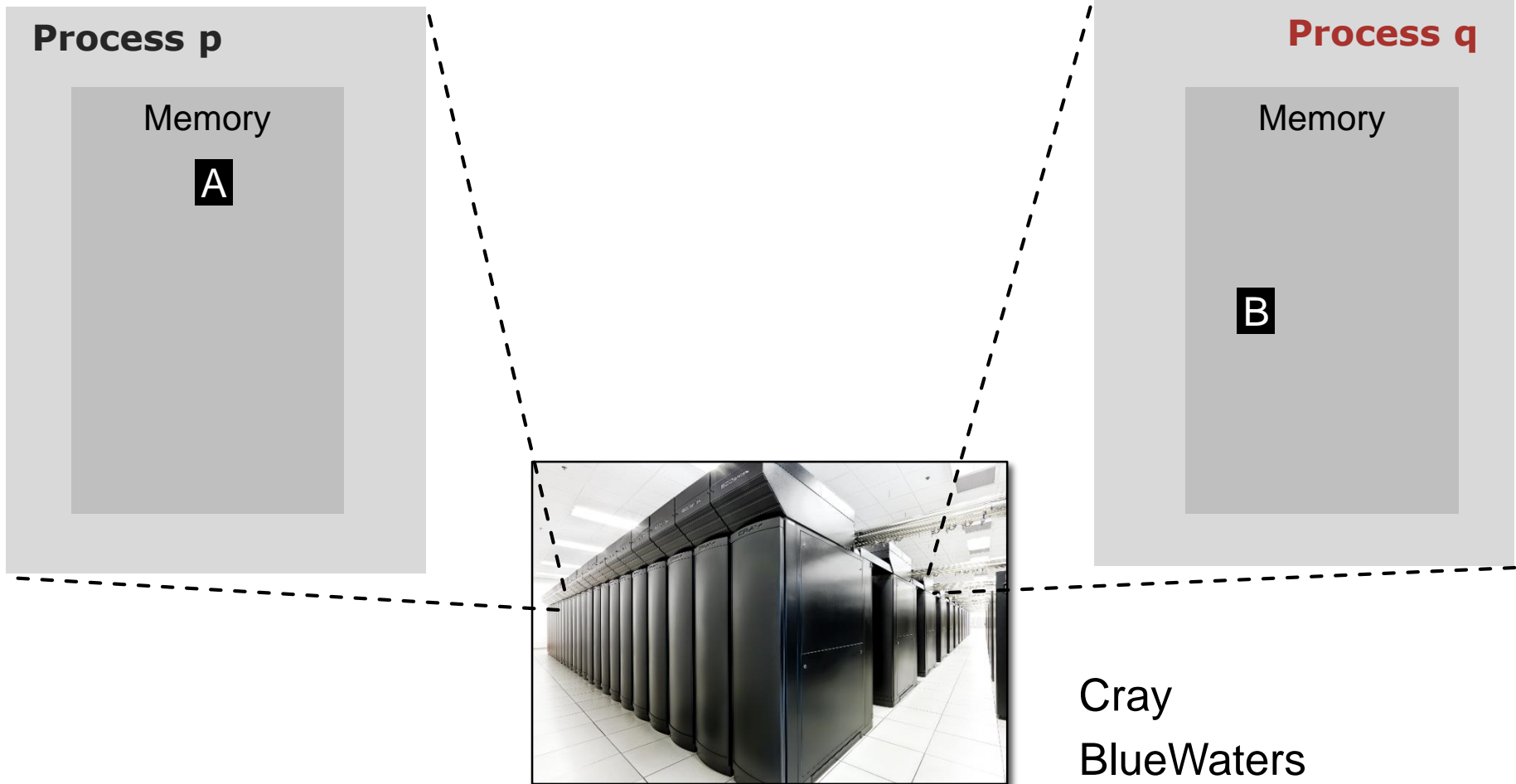


Cray  
BlueWaters

# REMOTE MEMORY ACCESS (RMA) PROGRAMMING



# REMOTE MEMORY ACCESS (RMA) PROGRAMMING



# REMOTE MEMORY ACCESS (RMA) PROGRAMMING

Process p

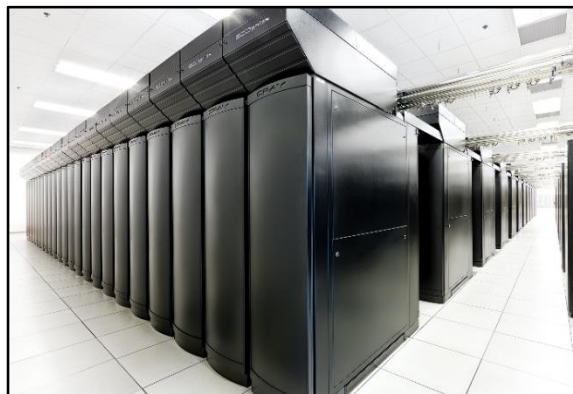
Memory

A

Process q

Memory

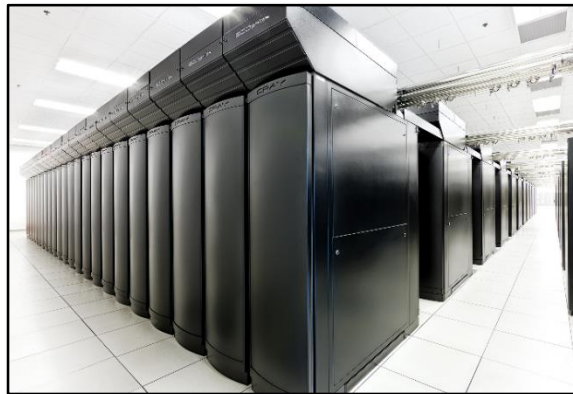
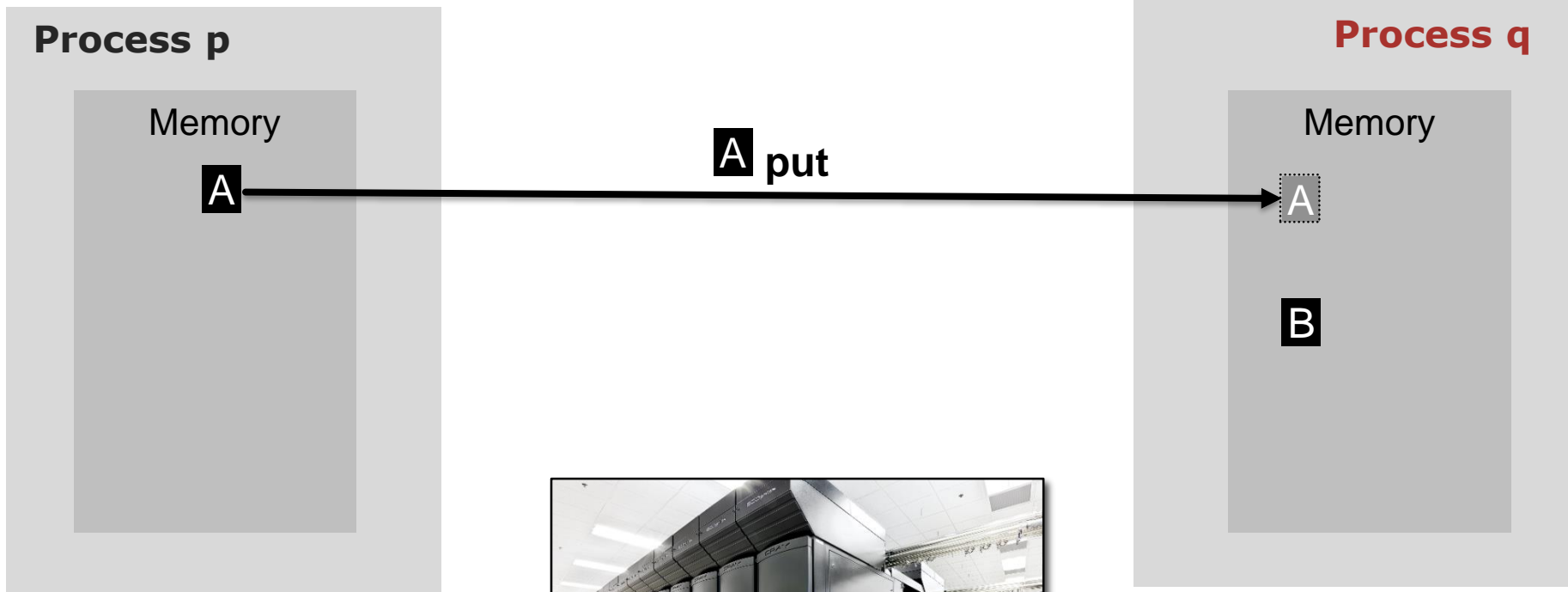
B



Cray  
BlueWaters

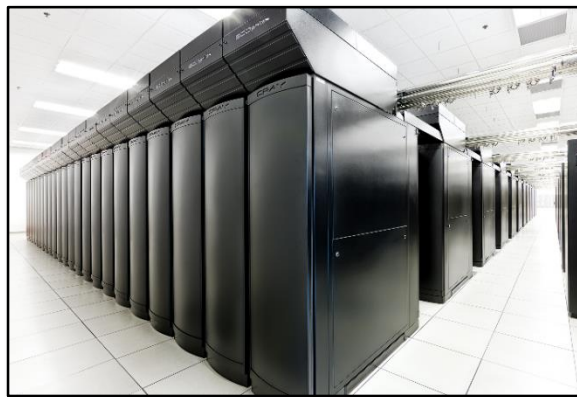
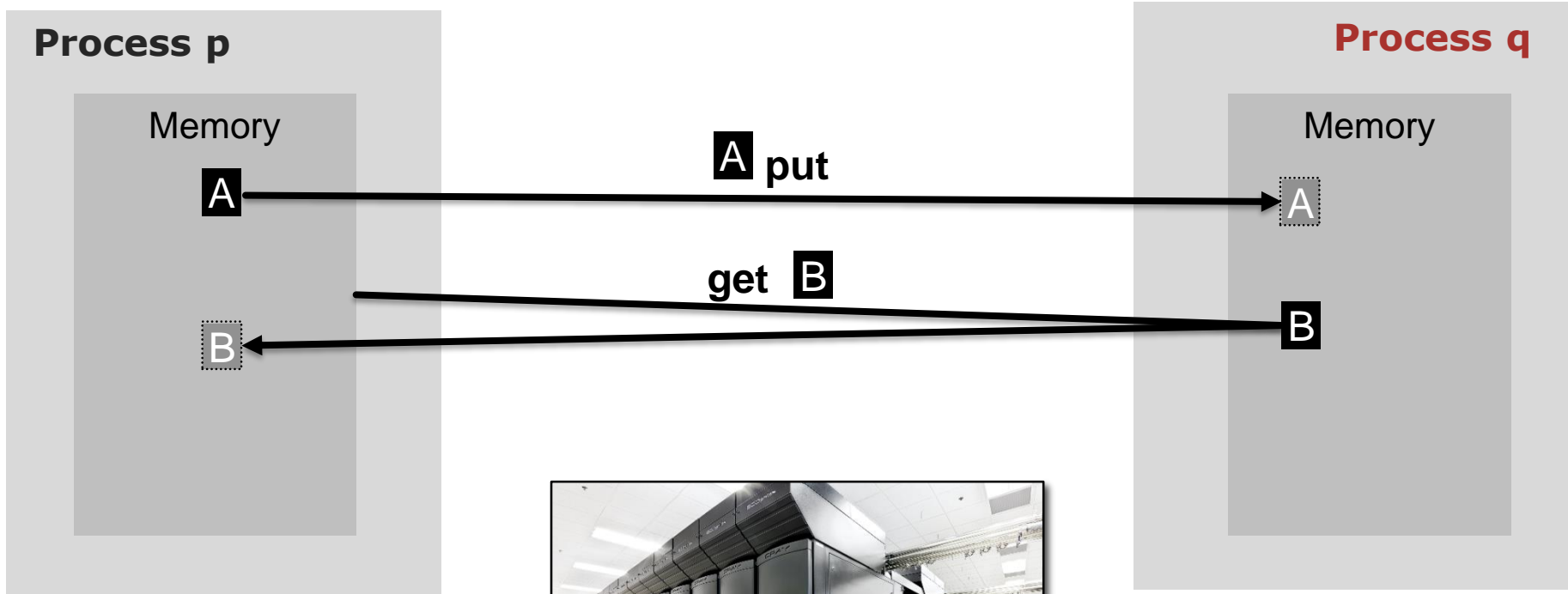


# REMOTE MEMORY ACCESS (RMA) PROGRAMMING



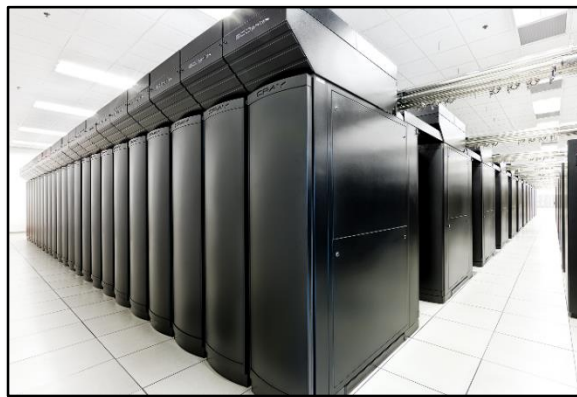
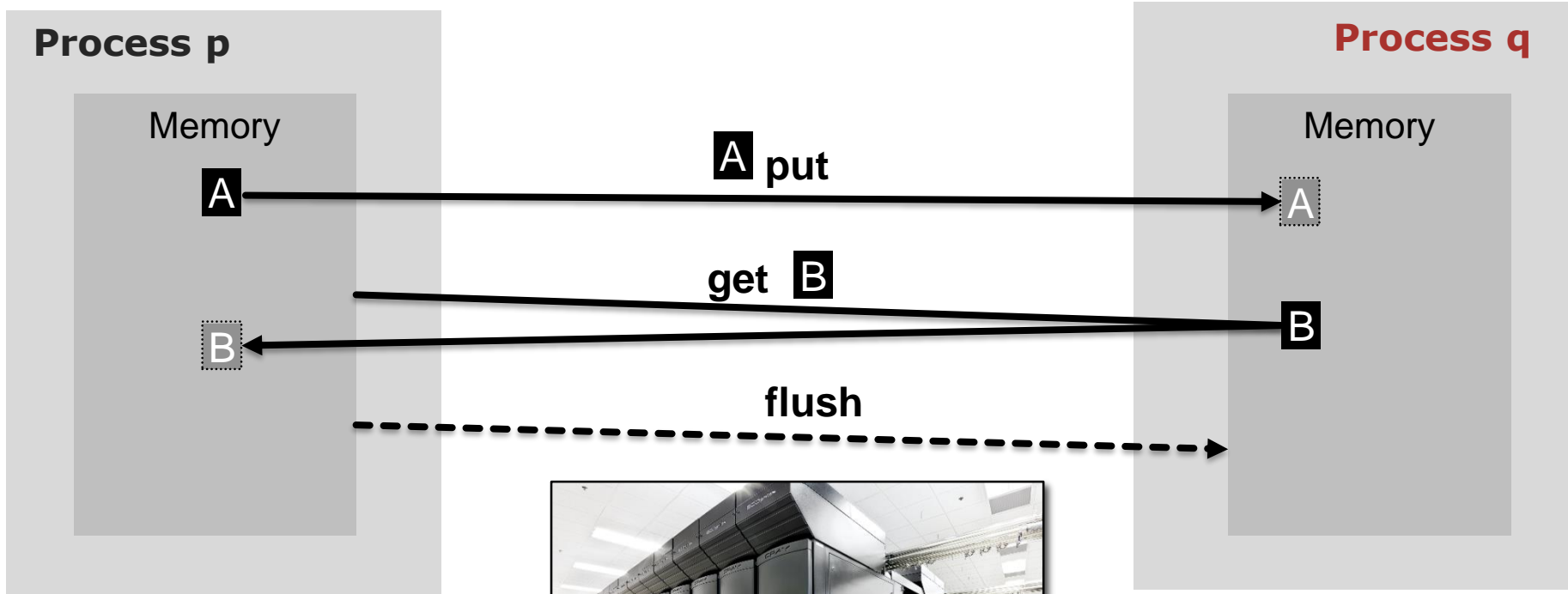
Cray  
 BlueWaters

# REMOTE MEMORY ACCESS (RMA) PROGRAMMING



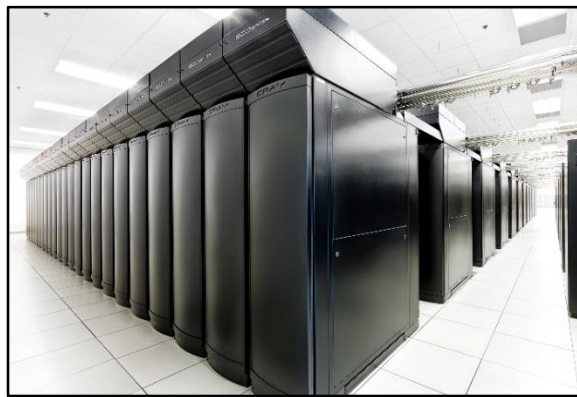
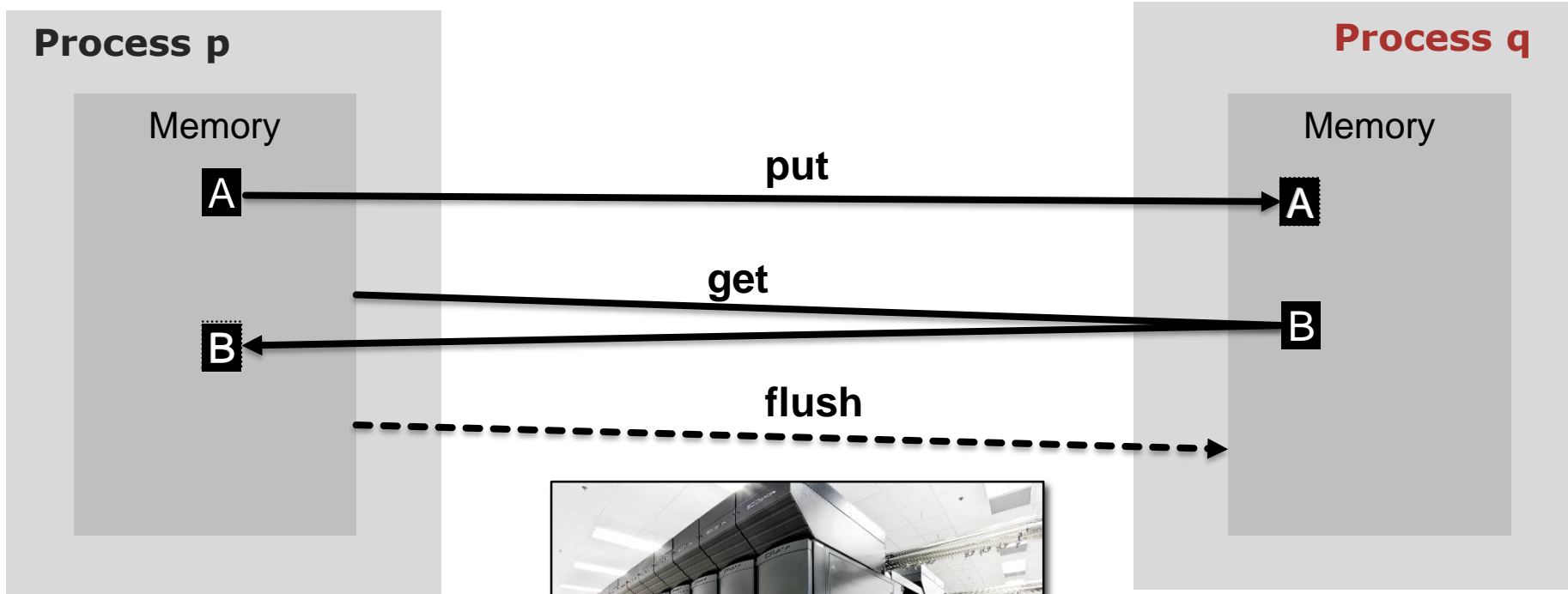
Cray  
 BlueWaters

# REMOTE MEMORY ACCESS (RMA) PROGRAMMING



Cray  
BlueWaters

# REMOTE MEMORY ACCESS (RMA) PROGRAMMING



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)



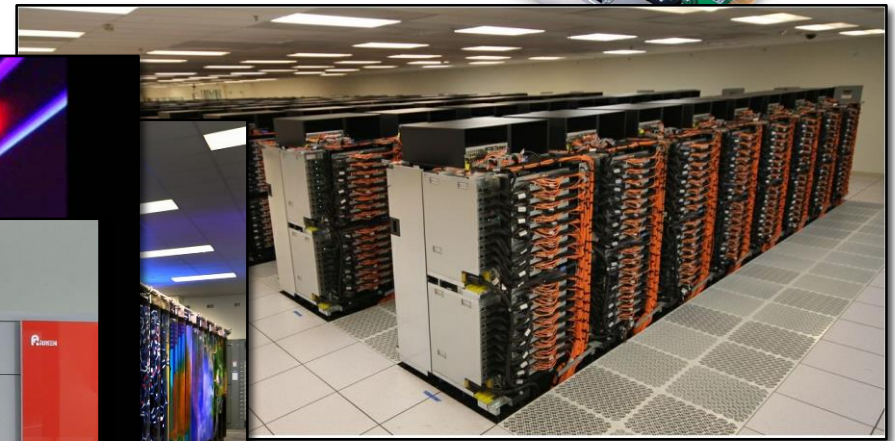
# REMOTE MEMORY ACCESS PROGRAMMING

- 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





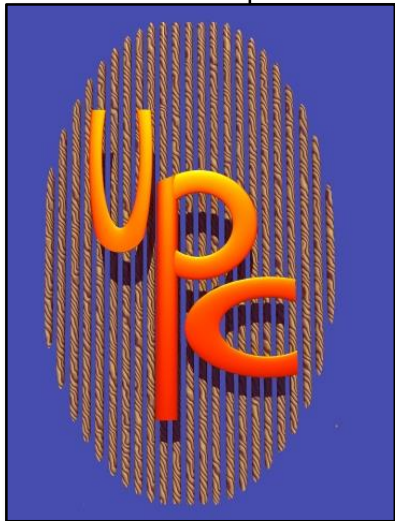
# REMOTE MEMORY ACCESS PROGRAMMING

- Supported by many HPC libraries and languages



# REMOTE MEMORY ACCESS PROGRAMMING

- Supported by many HPC libraries and languages



# REMOTE MEMORY ACCESS PROGRAMMING

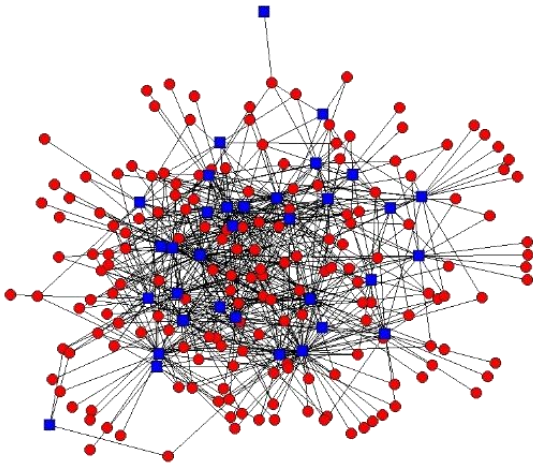
- Enables significant speedups over message passing in many types of applications, e.g.:

[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  $\sim 1.5$  for communication patterns in irregular workloads

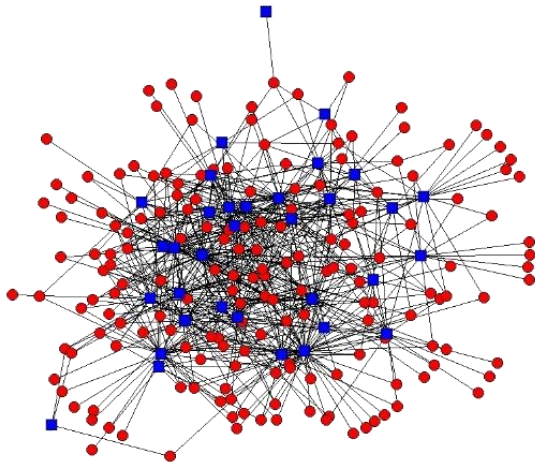
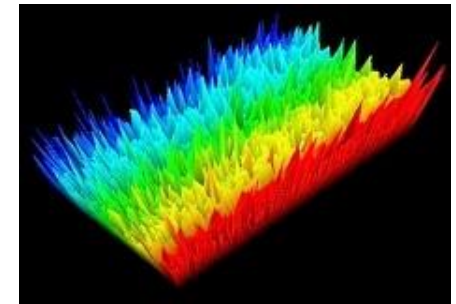


[1] R. Gerstenberger et al. Enabling Highly-Scalable Remote Memory Access Programming with MPI-3 One-Sided. SC13

[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  $\sim 1.5$  for communication patterns in irregular workloads
  - Speedup of  $\sim 1.4-2$  in physics computations



$$\frac{1}{\sqrt{2}} |\text{cat}\rangle + \frac{1}{\sqrt{2}} |\text{dog}\rangle$$

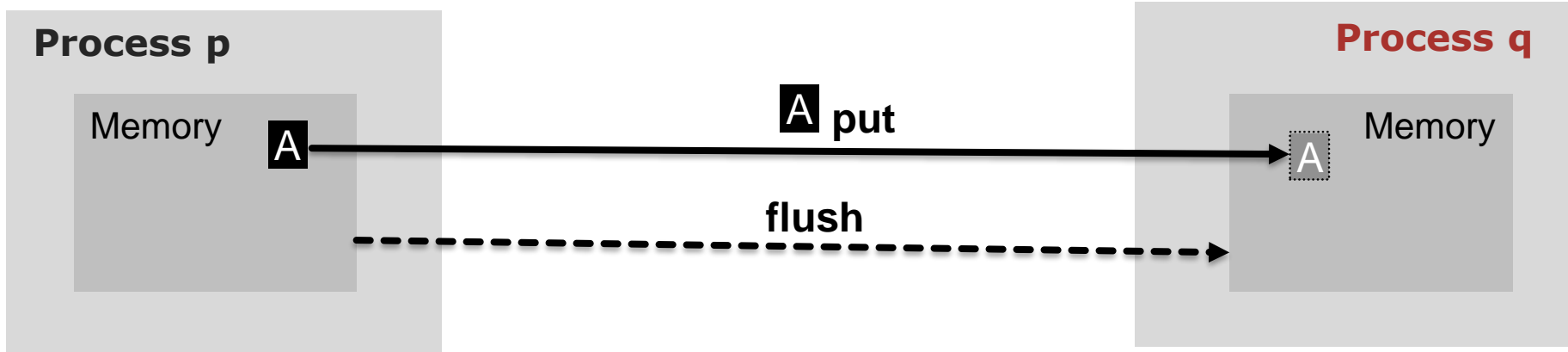
[1] R. Gerstenberger et al. Enabling Highly-Scalable Remote Memory Access Programming with MPI-3 One-Sided. SC13

[2] D. Petrovic et al., High-performance RMA-based broadcast on the Intel SCC. SPAA'12



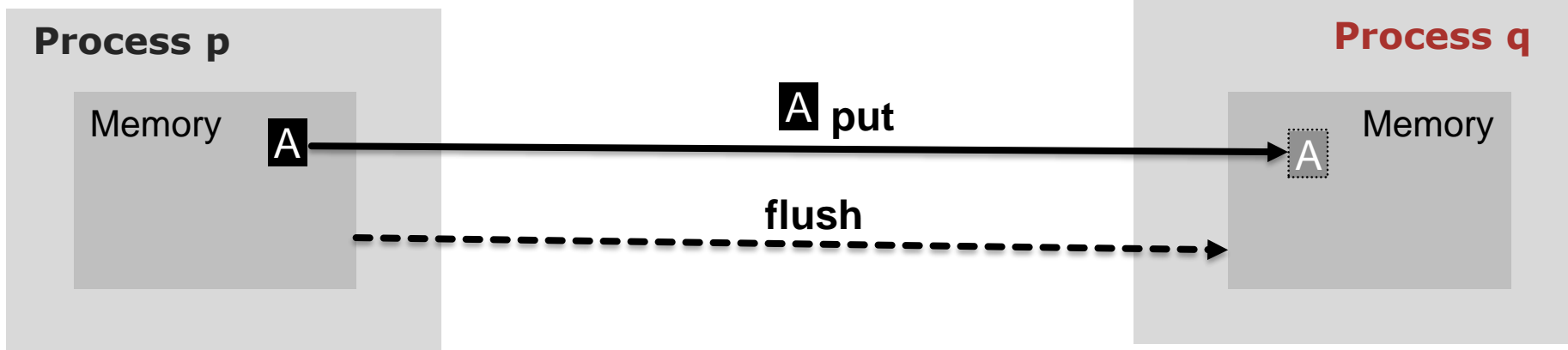
# RMA vs. MESSAGE PASSING

RMA:



# RMA vs. MESSAGE PASSING

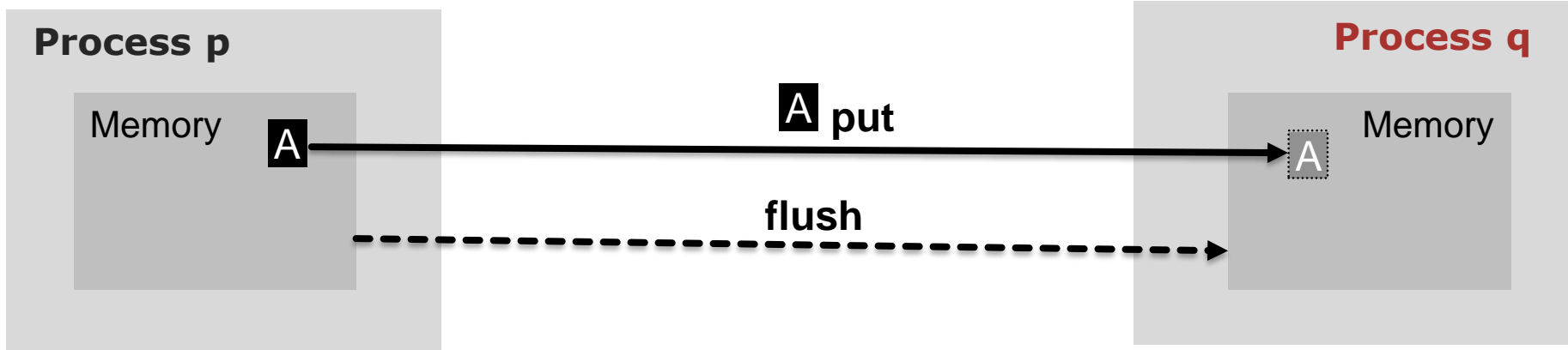
RMA:



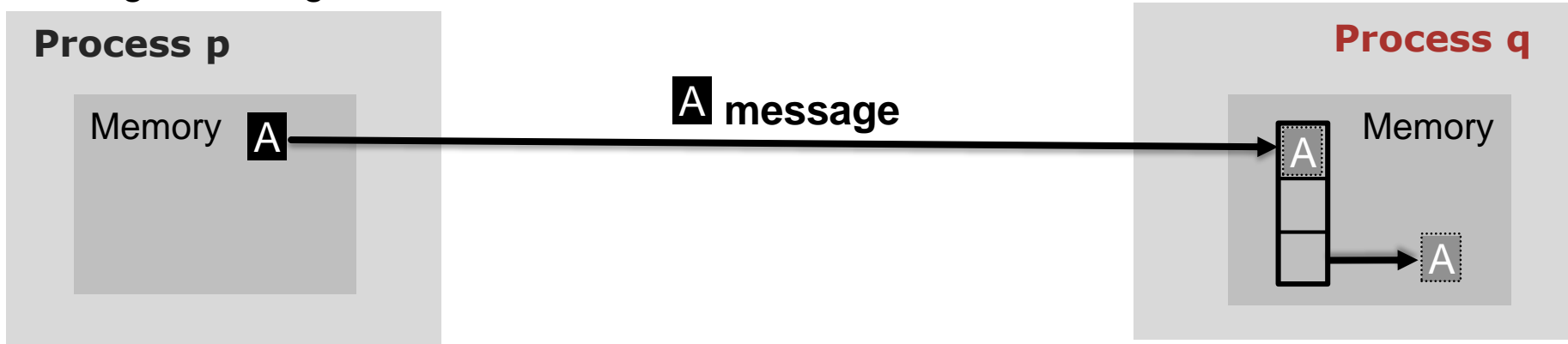
Message Passing:

# RMA vs. MESSAGE PASSING

RMA:



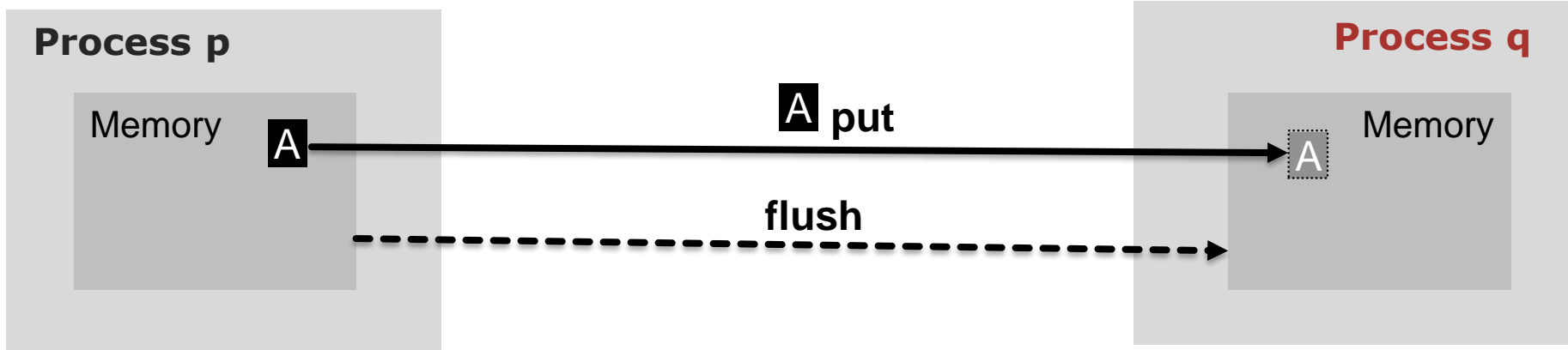
Message Passing:



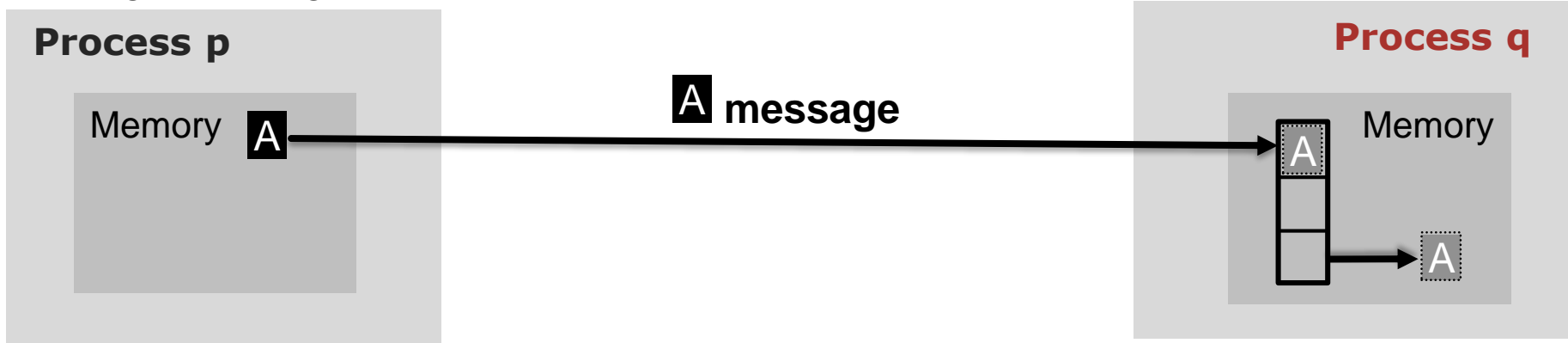
# RMA vs. MESSAGE PASSING

- Communication in RMA is one-sided

RMA:



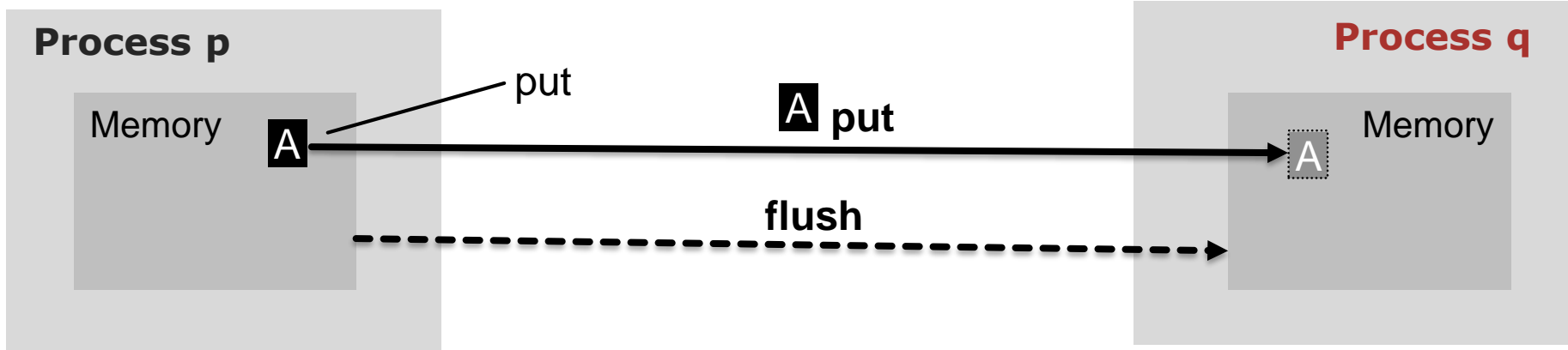
Message Passing:



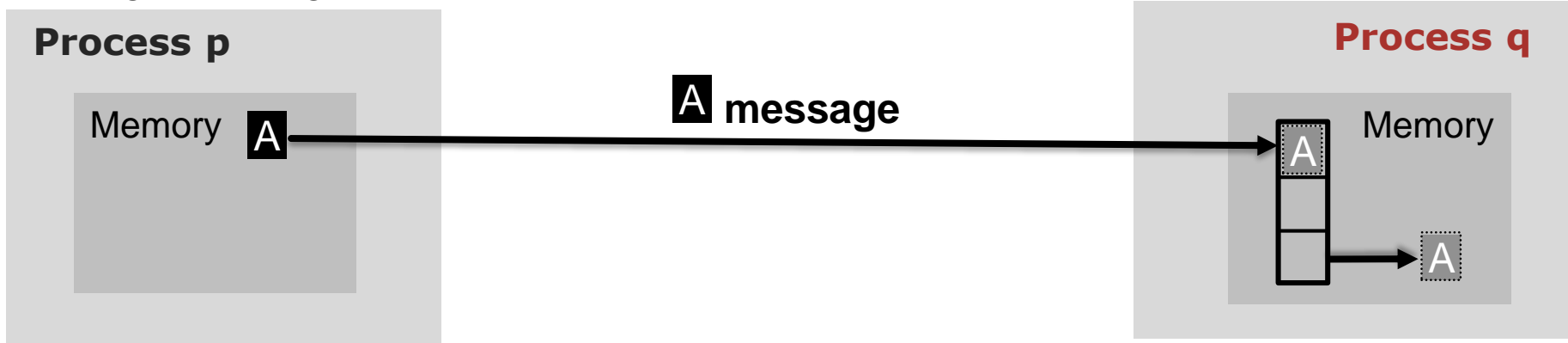
# RMA vs. MESSAGE PASSING

- Communication in RMA is one-sided

RMA:



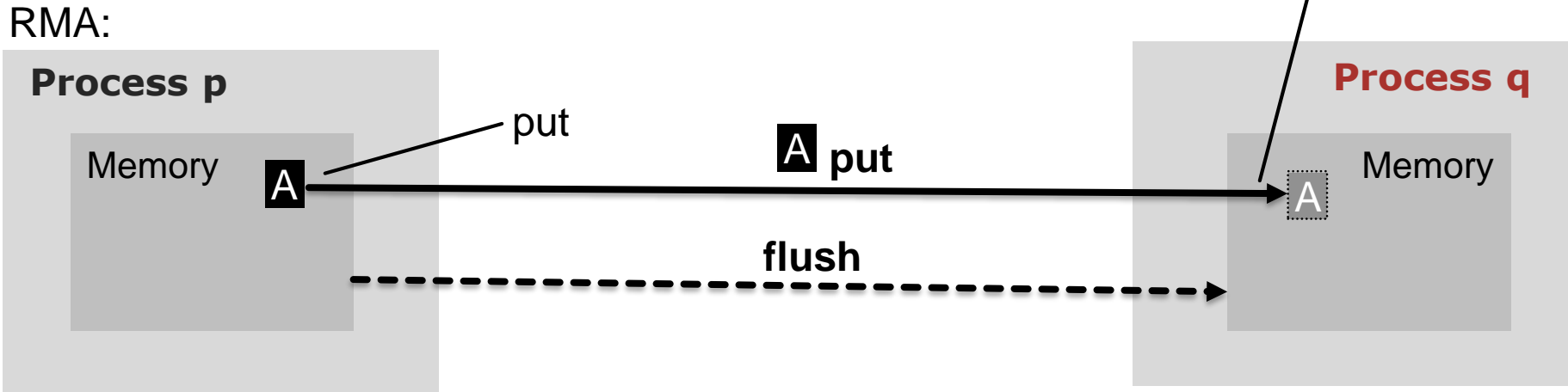
Message Passing:



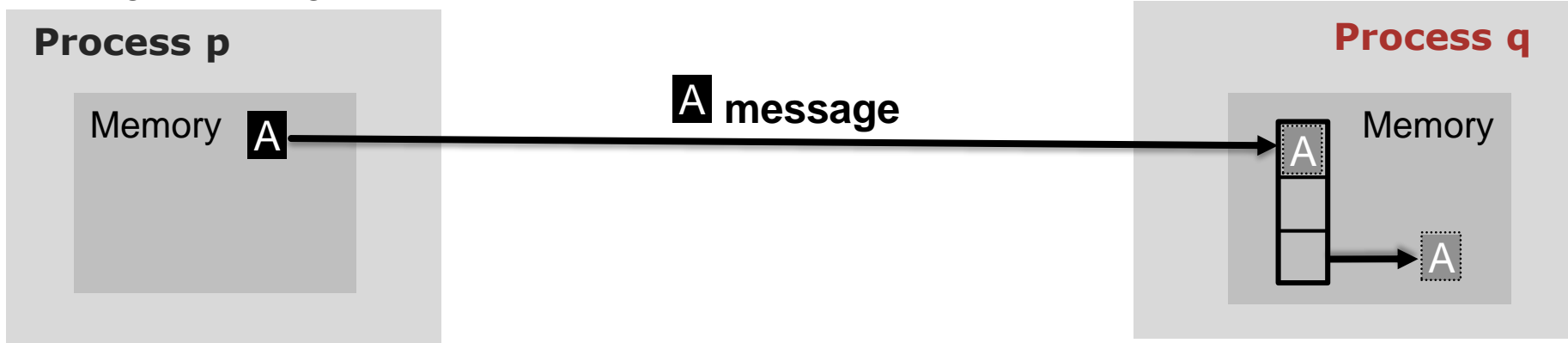


# RMA vs. MESSAGE PASSING

- Communication in RMA is one-sided

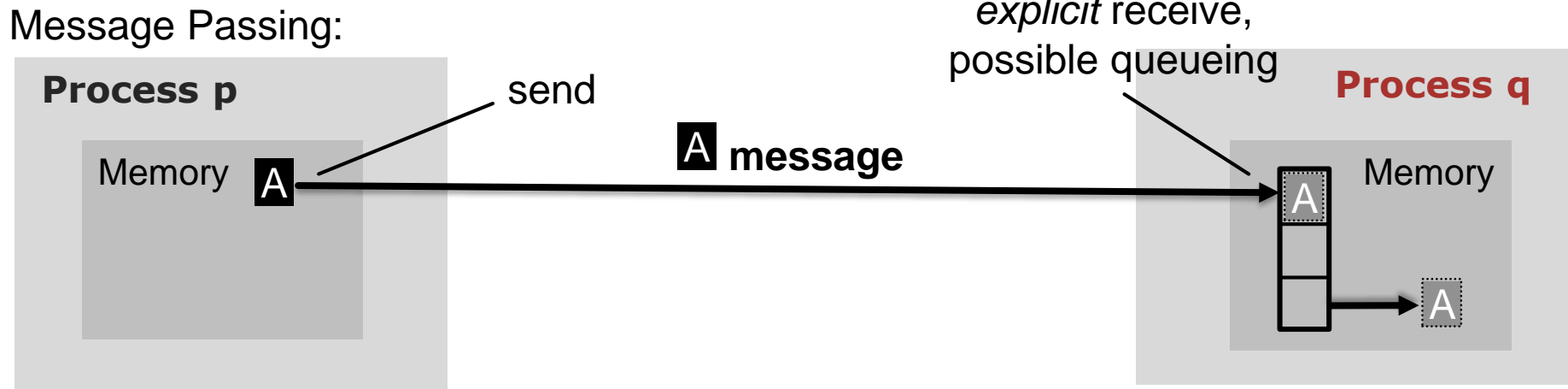
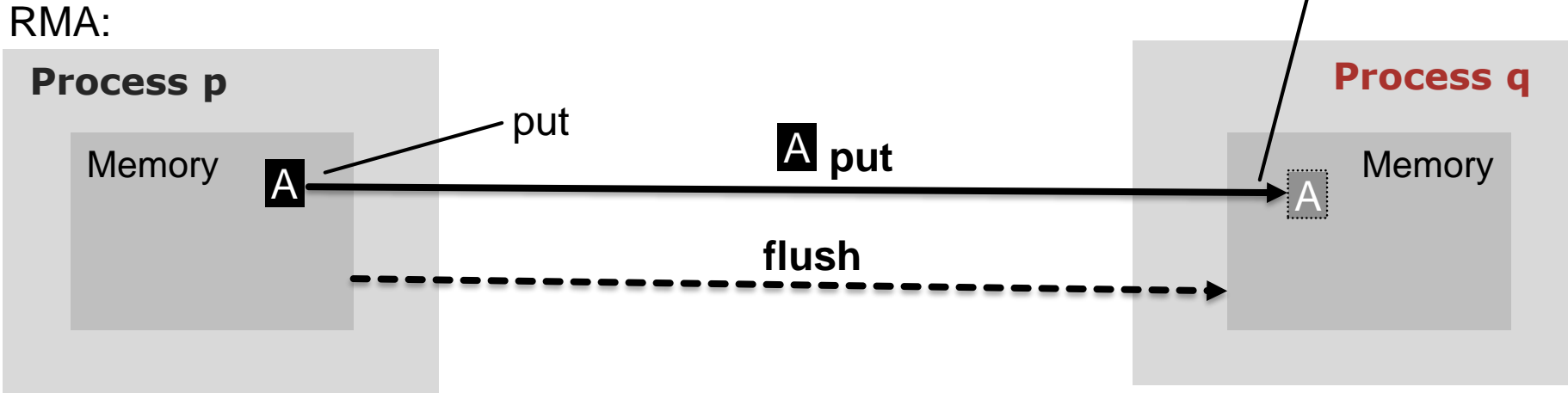


Message Passing:



# RMA vs. MESSAGE PASSING

- Communication in RMA is one-sided



# REMOTE MEMORY ACCESS PROGRAMMING

# REMOTE MEMORY ACCESS PROGRAMMING

- Is it ideal?

# REMOTE MEMORY ACCESS PROGRAMMING

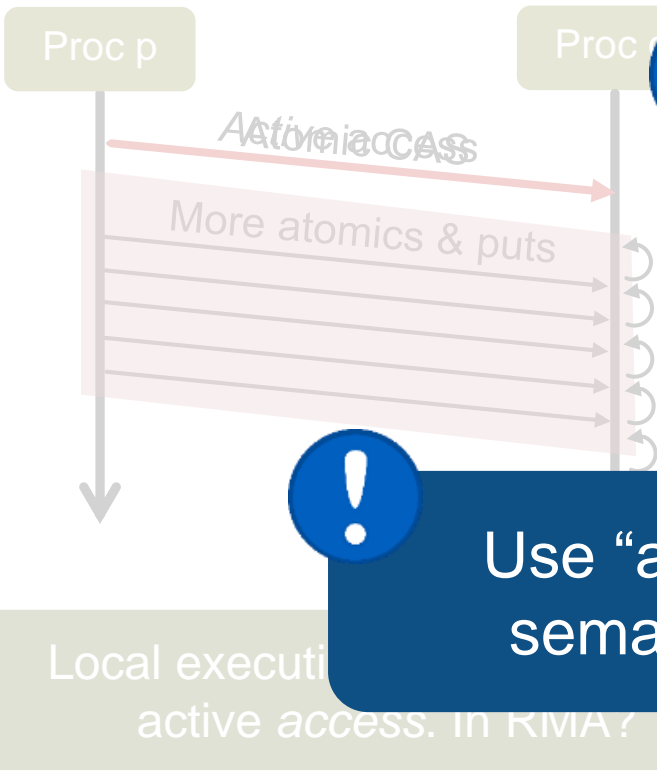
- Is it ideal? 

# REMOTE MEMORY ACCESS PROGRAMMING

- Is it possible to enable it?
- Can we use it for distributed hashtable...

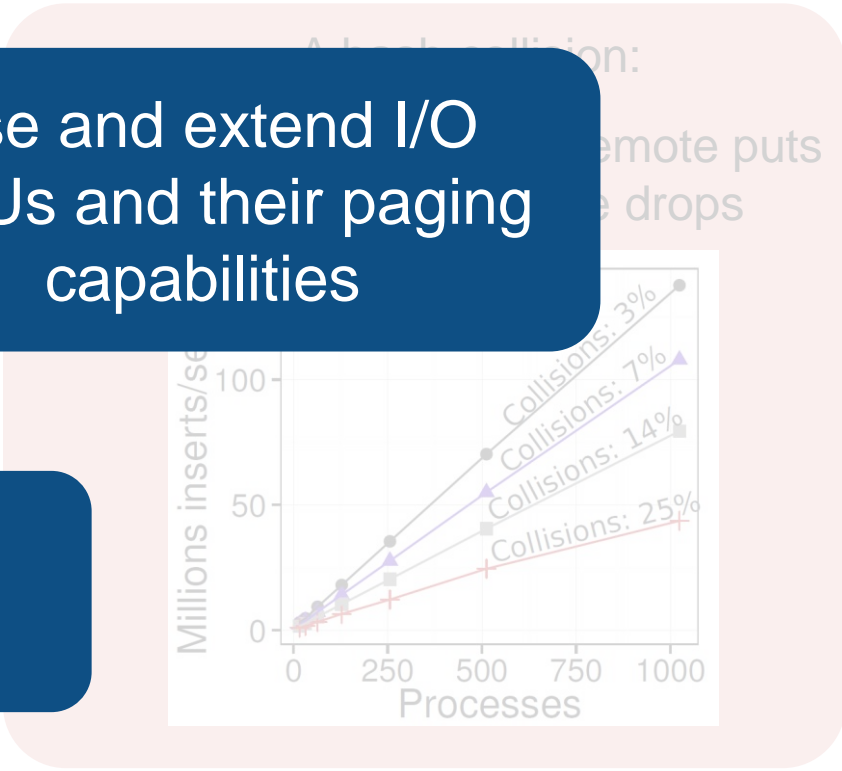
**How to enable it?**

No hash collision:  
 → 1 remote atomic  
 → Up to 5x speedup over MP [1]



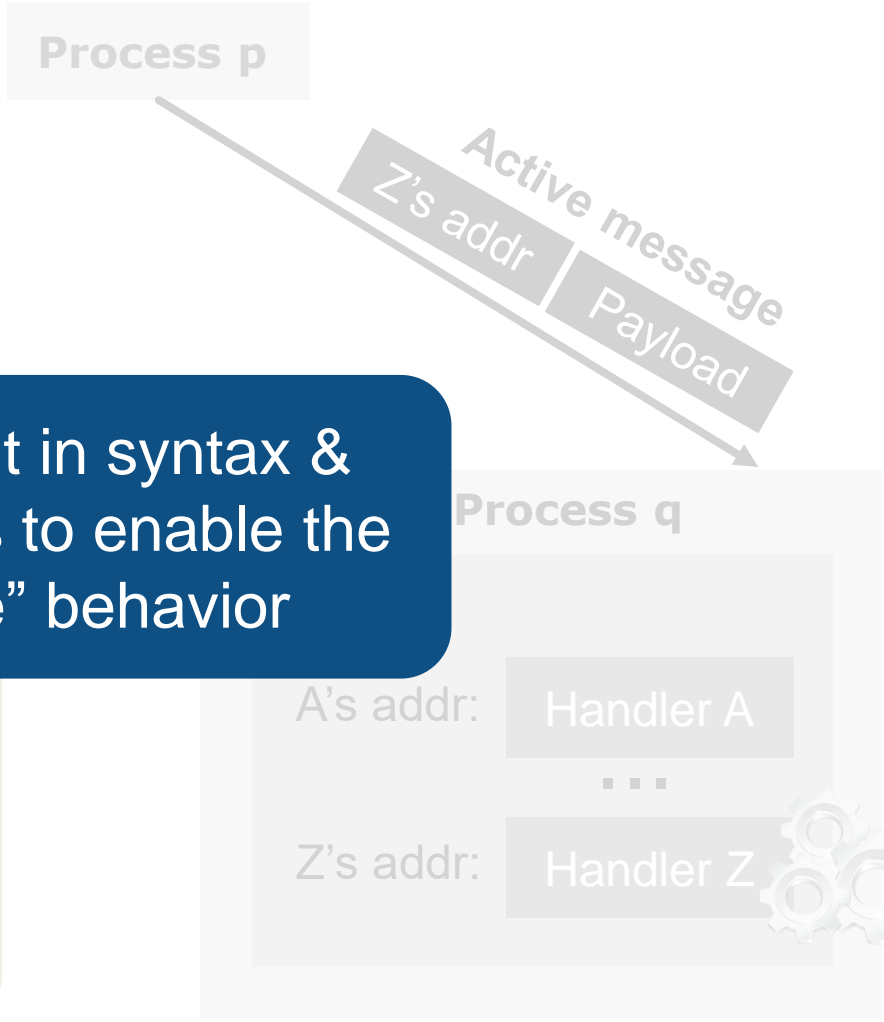
**Use and extend I/O MMUs and their paging capabilities**

**Use "active" semantics**



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

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



**!** We use it in syntax & semantics to enable the “active” behavior

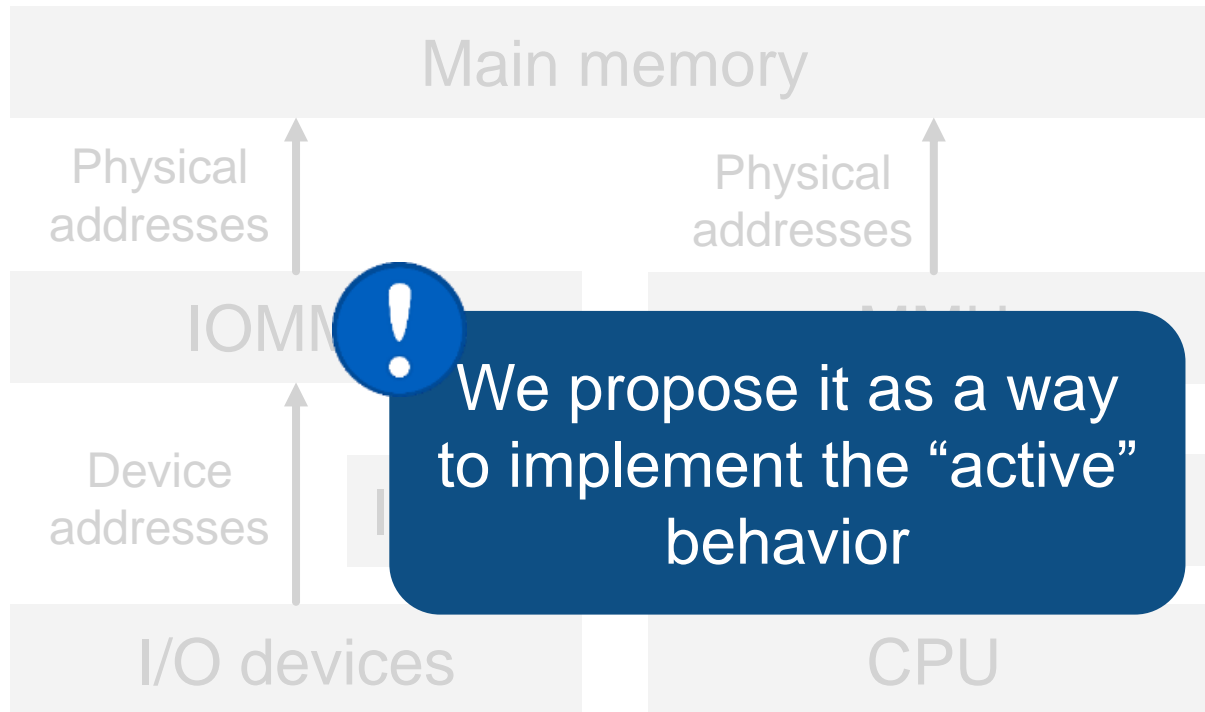
We need *active* puts/gets.

- 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.

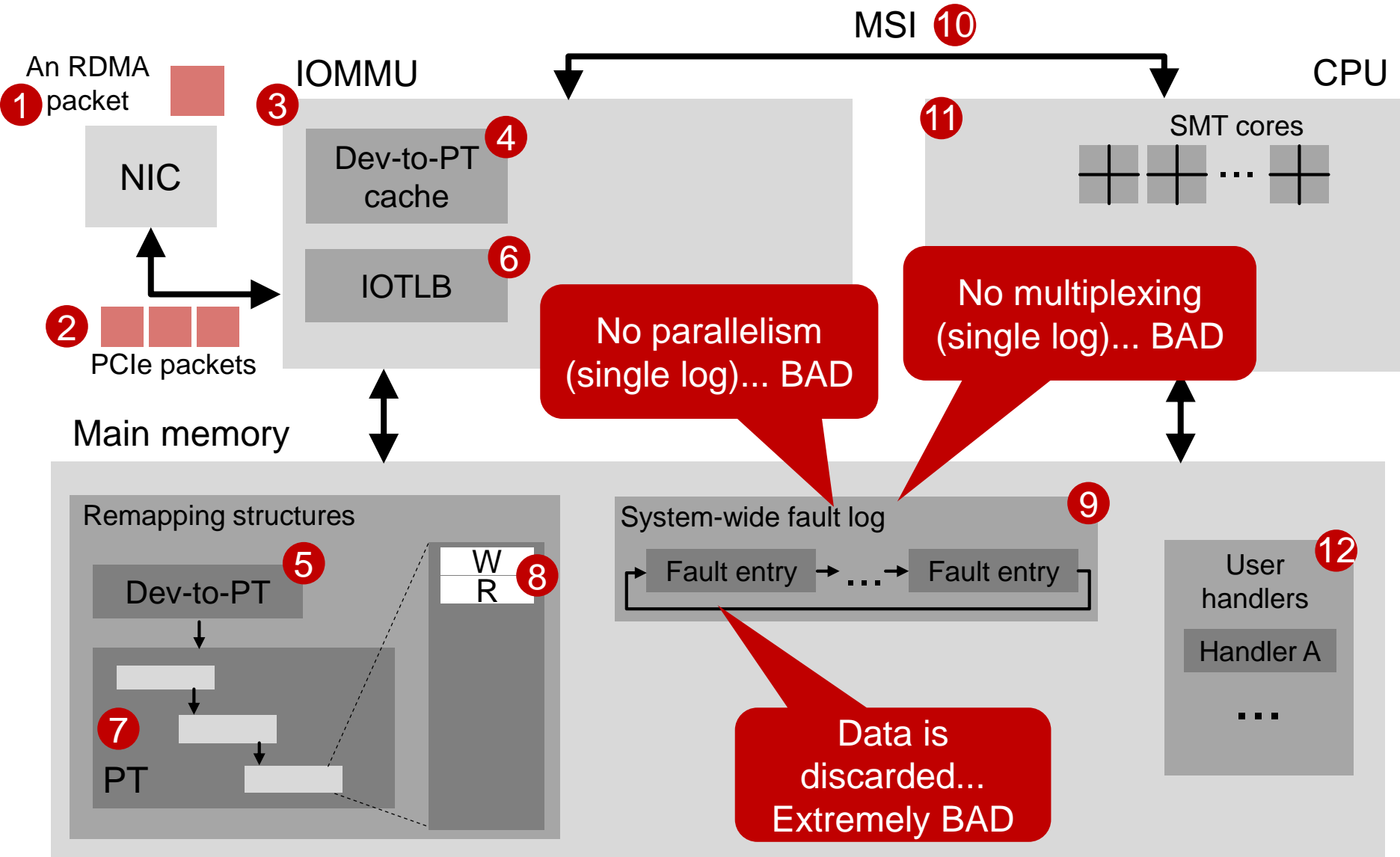


# USE INPUT/OUTPUT MEMORY MANAGEMENT UNITS

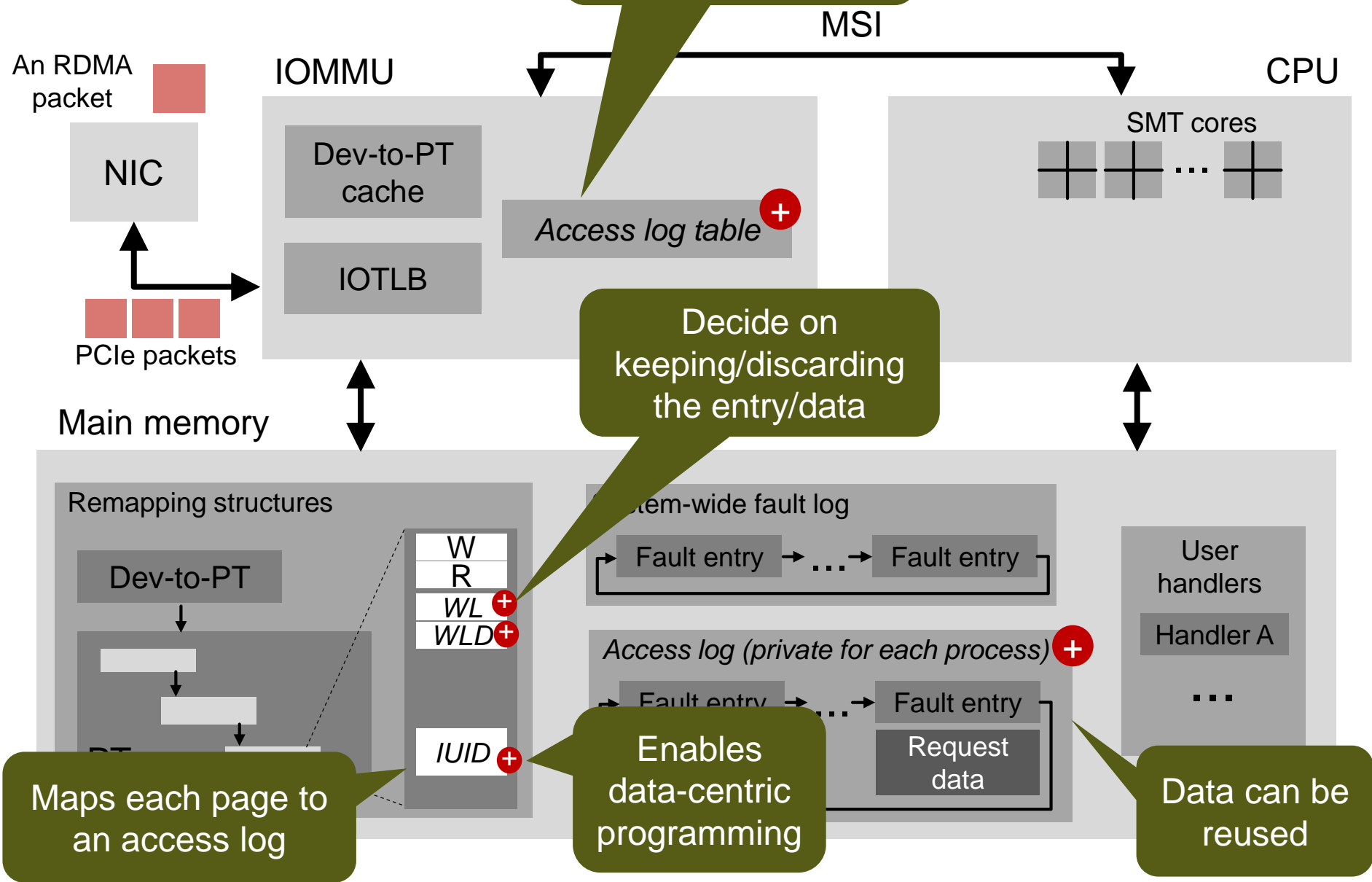


We could use it somehow. But...

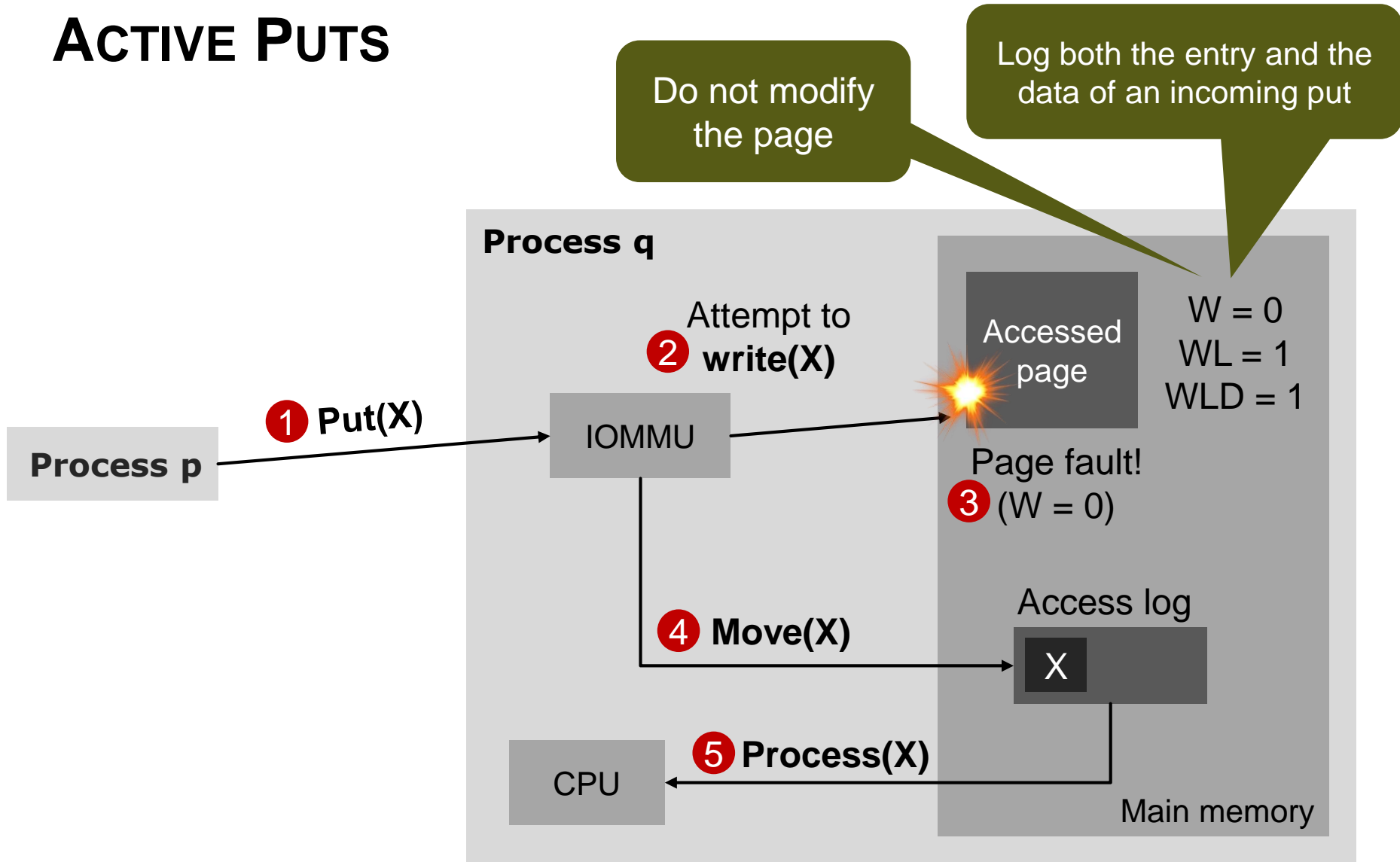
# IOMMUS AND RMA



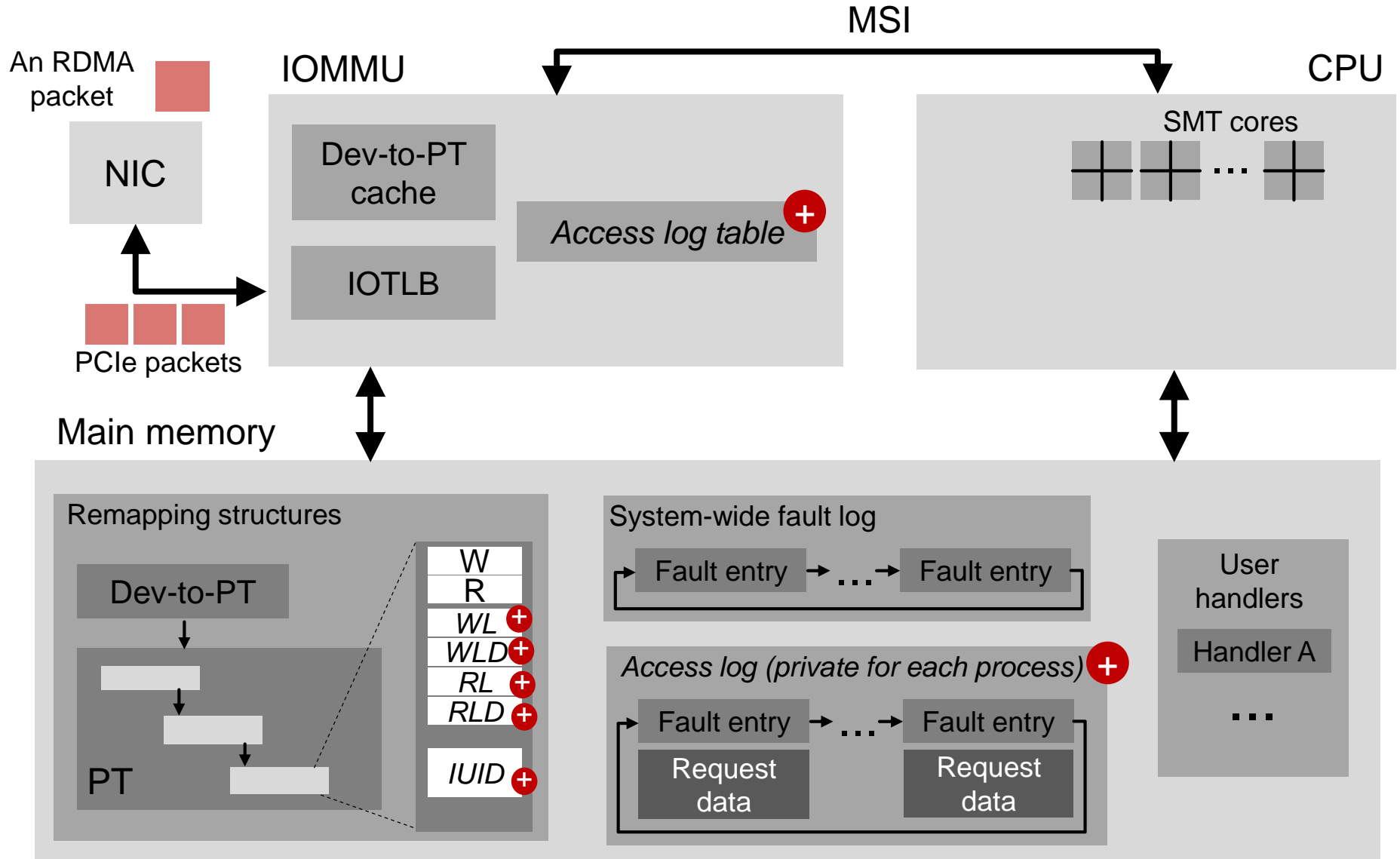
# ACTIVE PUTS



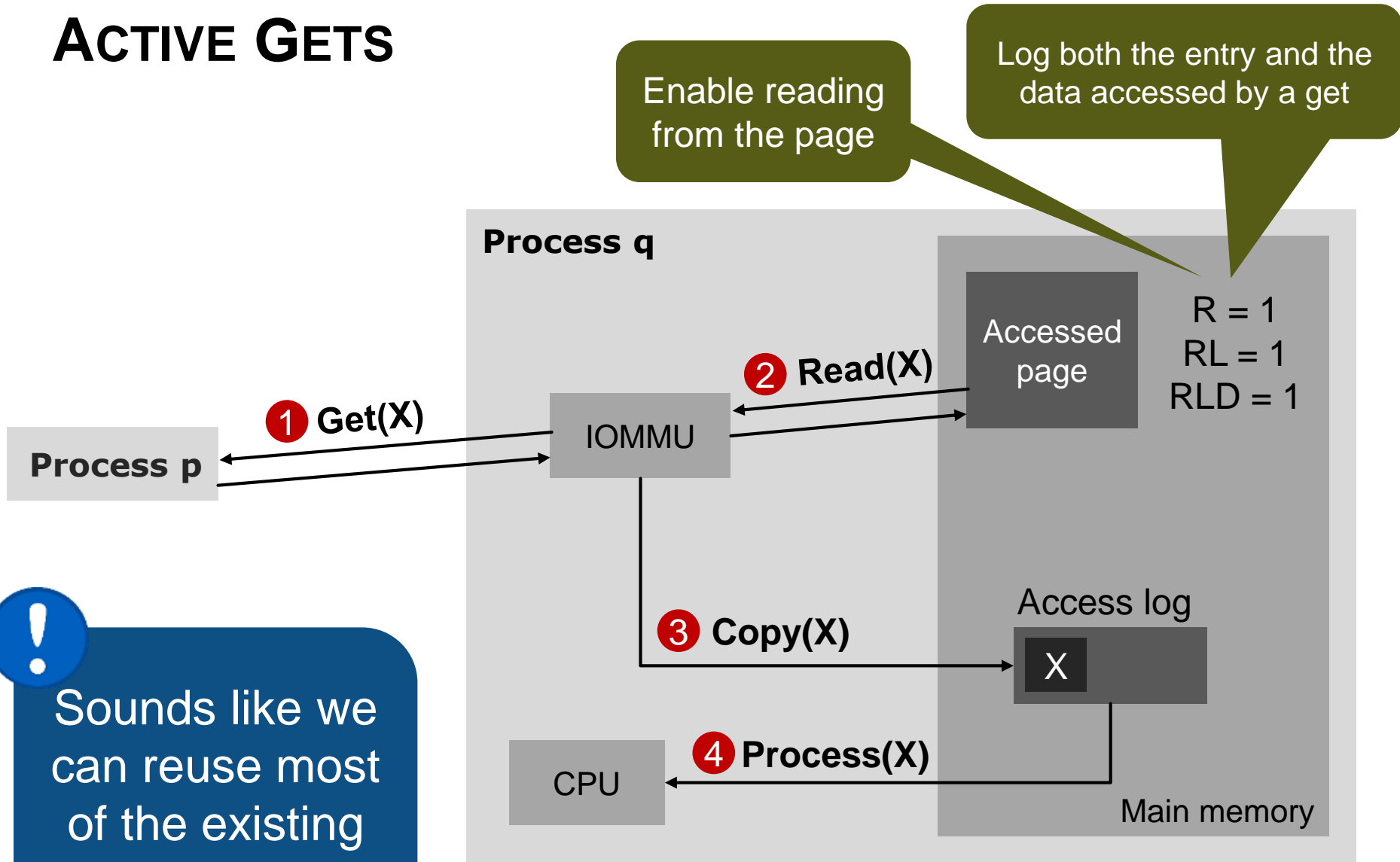
# ACTIVE PUTS



# ACTIVE GETS

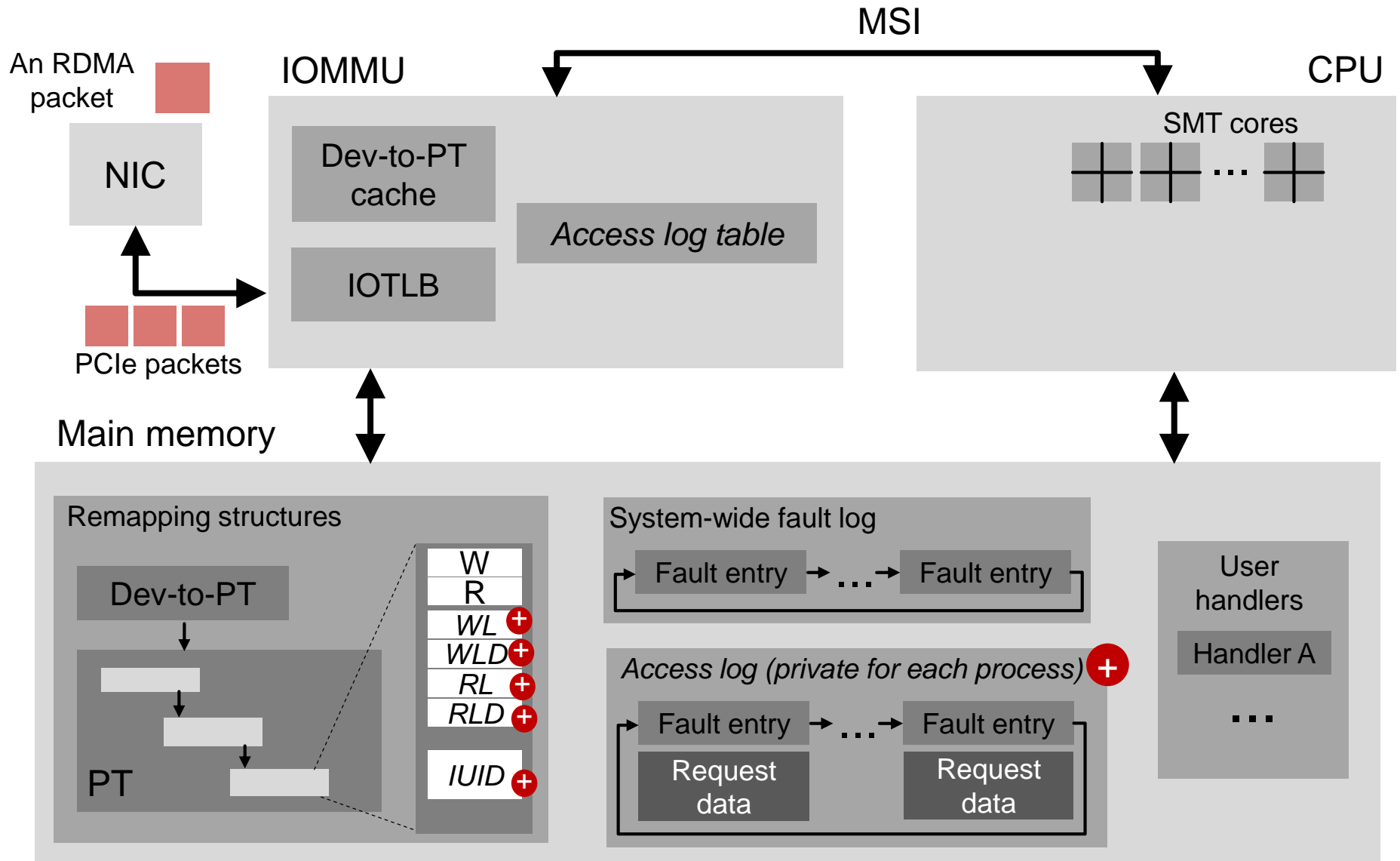


# ACTIVE GETS



! Sounds like we can reuse most of the existing stuff!

# INTERACTIONS WITH THE CPU





# INTERACTIONS WITH THE CPU



Are we done?

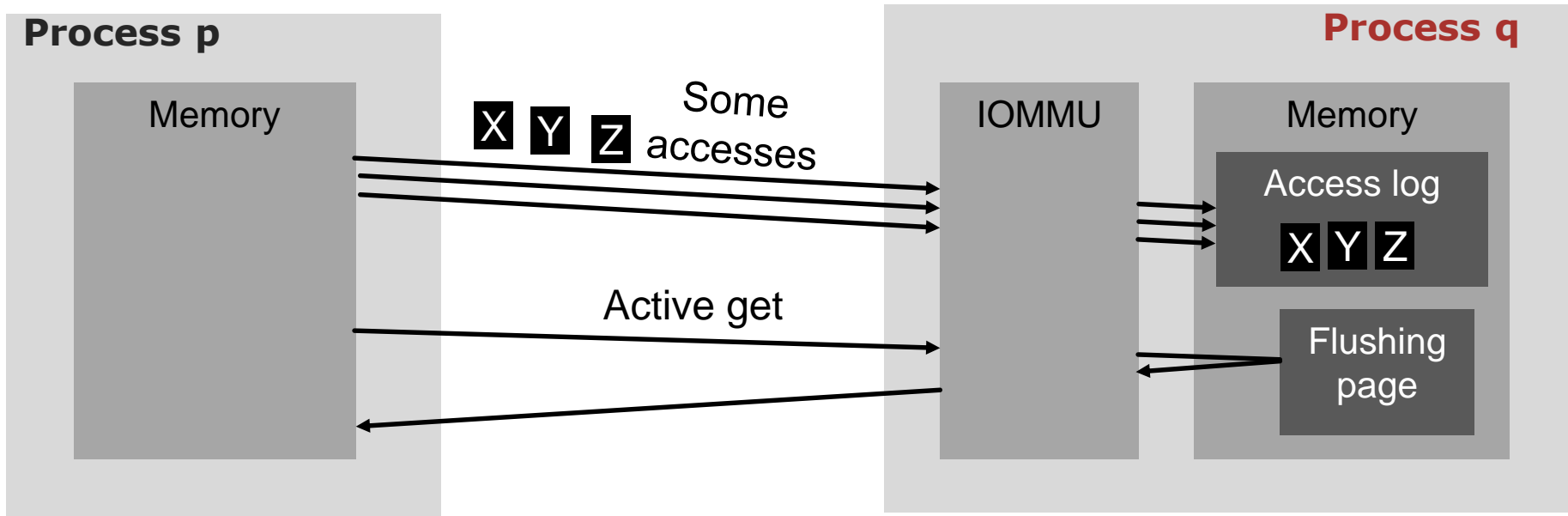
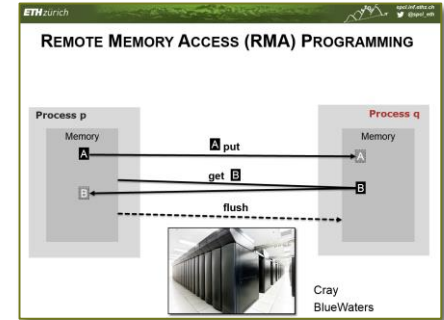
- Interrupts
- Polling
- Direct notifications via software



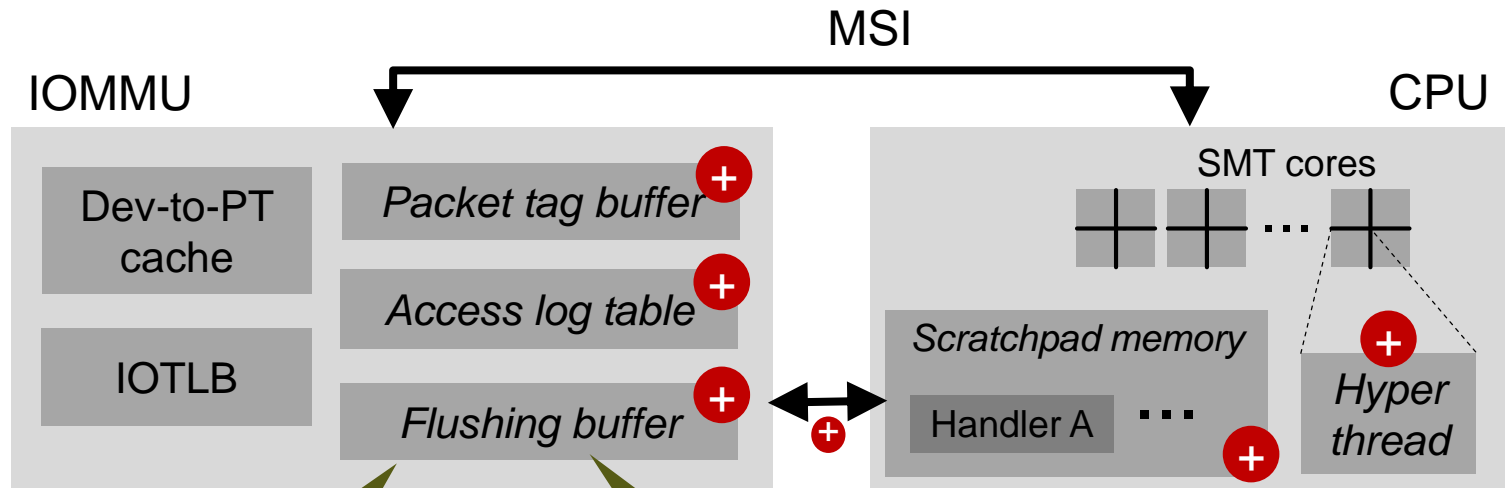
Well...

# 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*



# CONSISTENCY



Contains the addresses of flushing pages

Maps flushing pages to IUIDs and access logs

# Let's summarize...

### 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

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

AM++ [2]  
 GASNet [3]

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. Wilcock et al. AM++: A generalized active message framework. PACT'10.  
 [3] D. Bonachea. GASNet Specification, v1.1. Berkeley Technical Report, 2002.

## Active Messages

### IOMMUs

#### USE INPUT/OUTPUT MEMORY MANAGEMENT UNITS

+

## Active Puts/Gets

### ACTIVE PUTS

An RDMA packet  
 NIC  
 Dev-to-PT caches  
 IOTLB  
 Main memory  
 Remapping structures  
 Dev-to-PT  
 PT  
 RMD  
 CPU  
 SMT cores

### ACTIVE PUTS

Do not modify the page  
 Log both the entry and the data of an incoming put

Process p  
 Put(X)  
 IOMMU  
 Attempt to write(X)  
 Accessed page  
 Page fault (W = 0)  
 Move(X)  
 CPU  
 Process(X)  
 Main memory

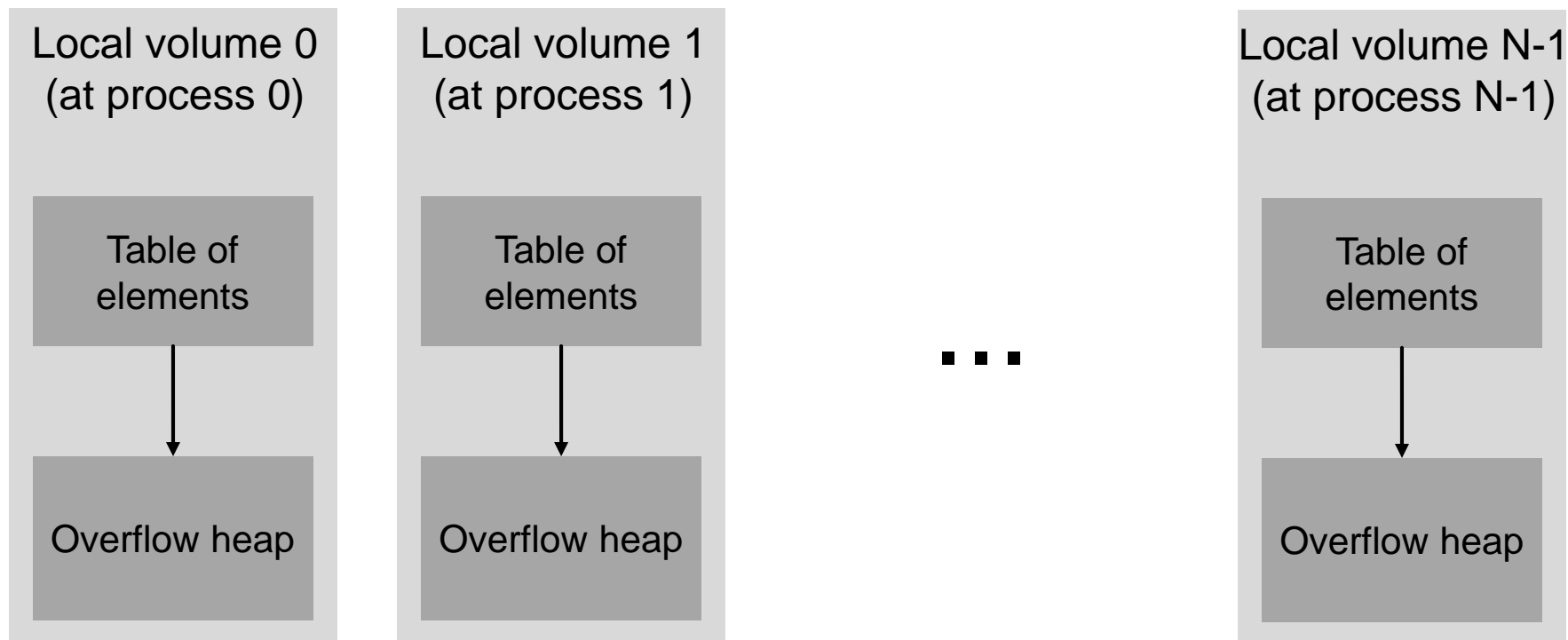
# How can we use it?

# ACTIVE ACCESS USE-CASES

## DISTRIBUTED HASHTABLE

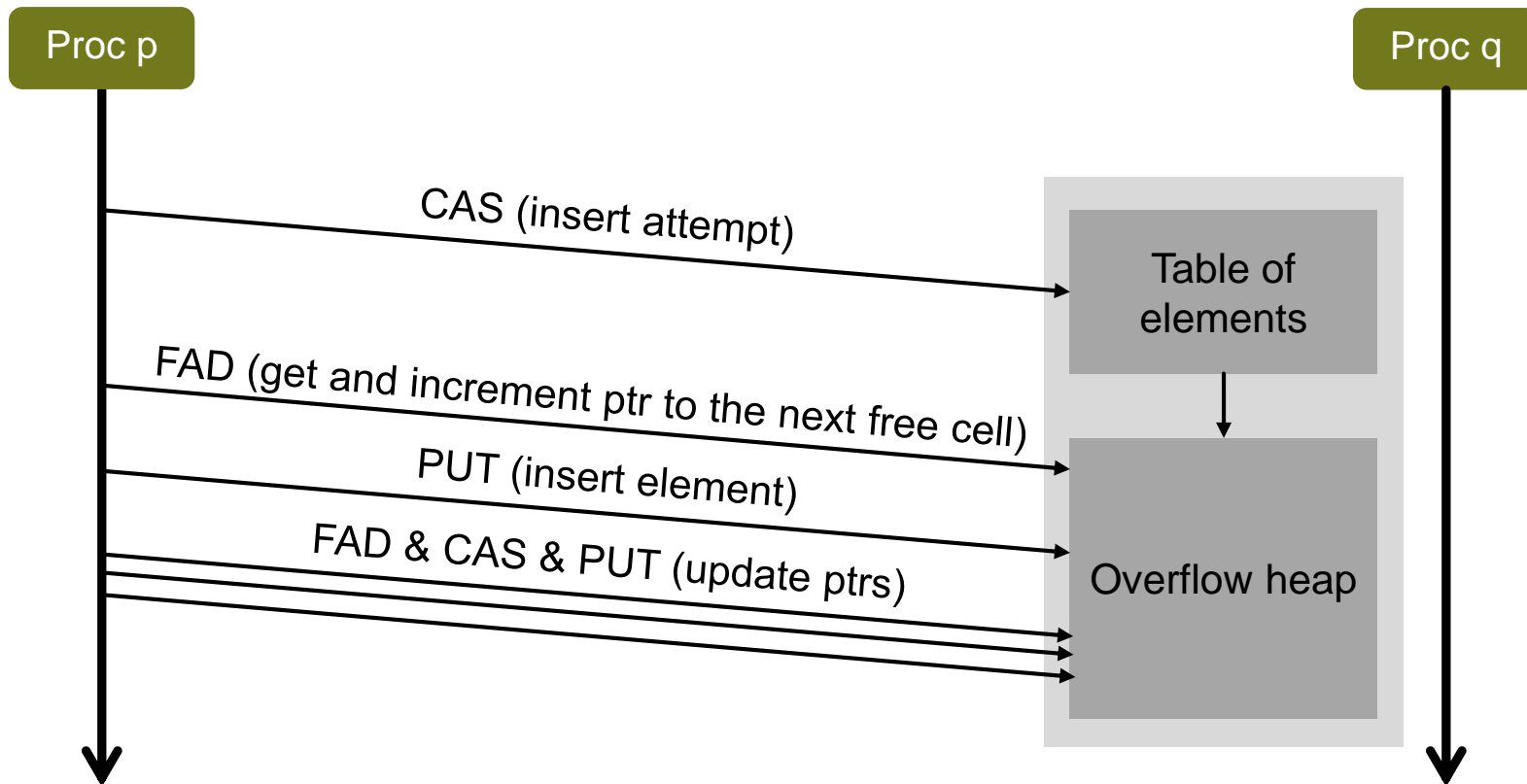


- Used to construct key-value stores (e.g., Memcached [1])



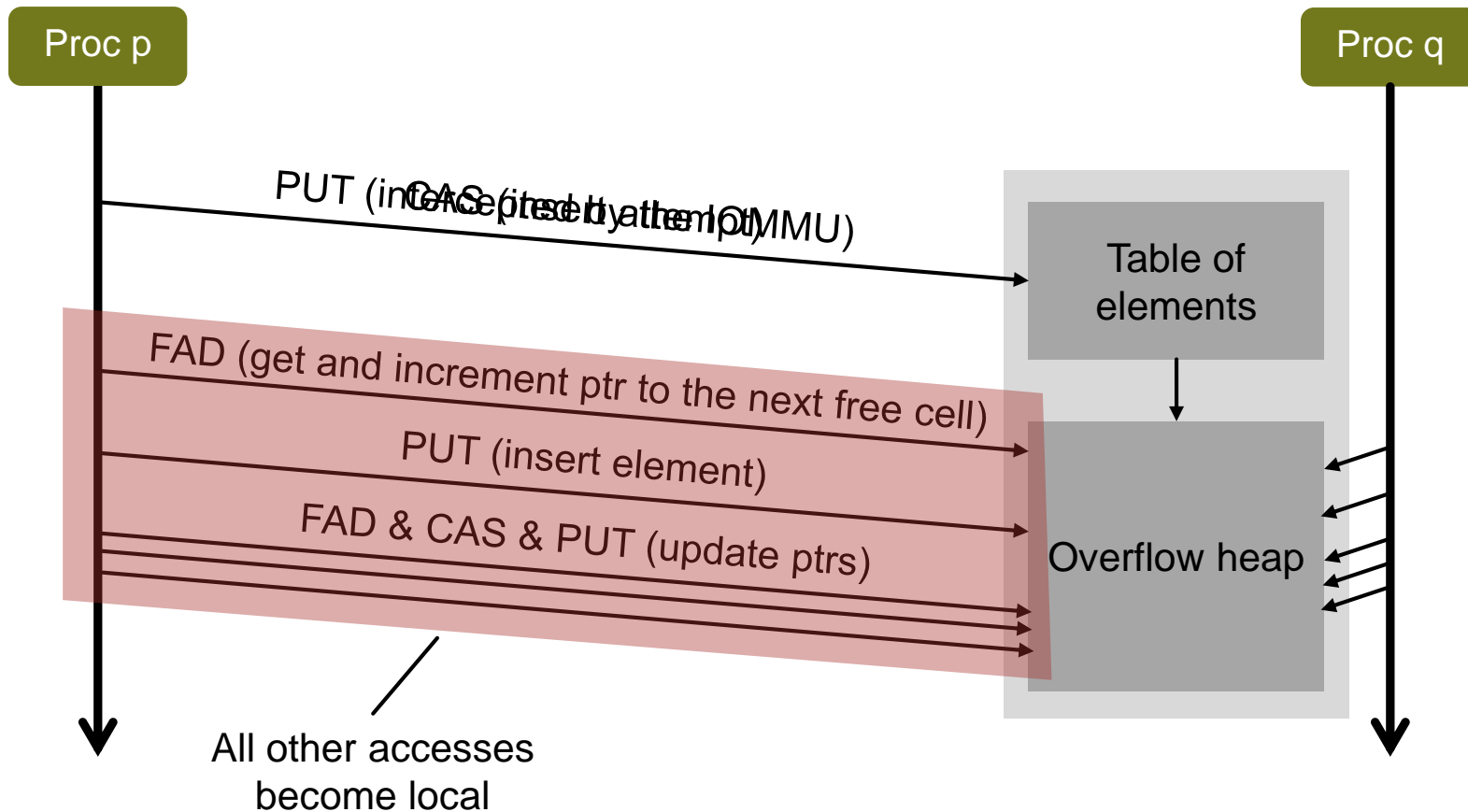
# ACTIVE ACCESS USE-CASES

## DISTRIBUTED HASHTABLE: INSERTS (RMA)



# ACTIVE ACCESS USE-CASES

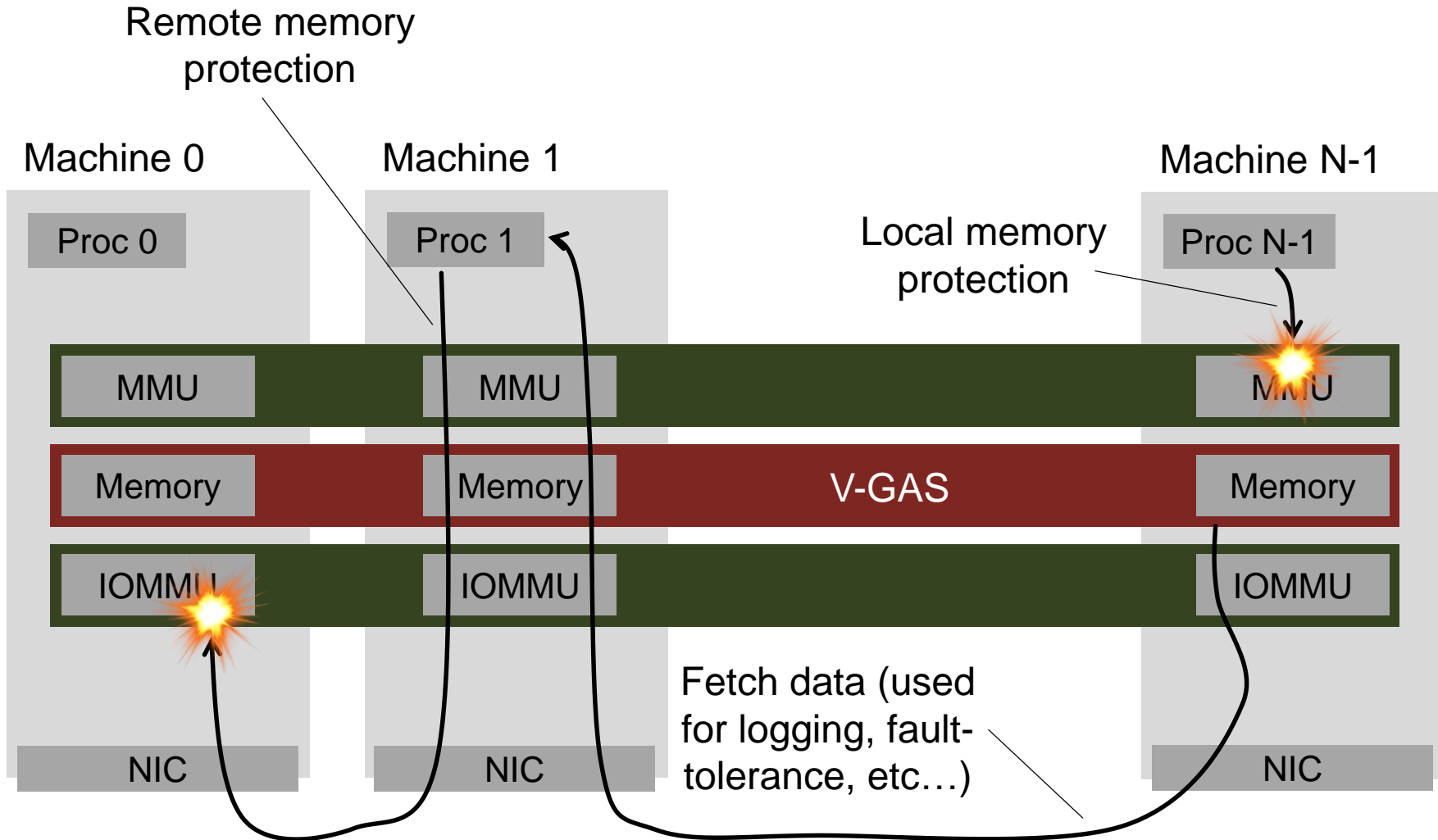
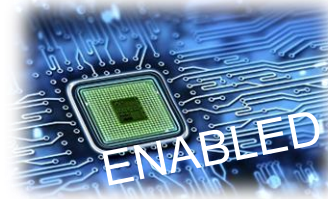
## DISTRIBUTED HASHTABLE: INSERTS (AA)





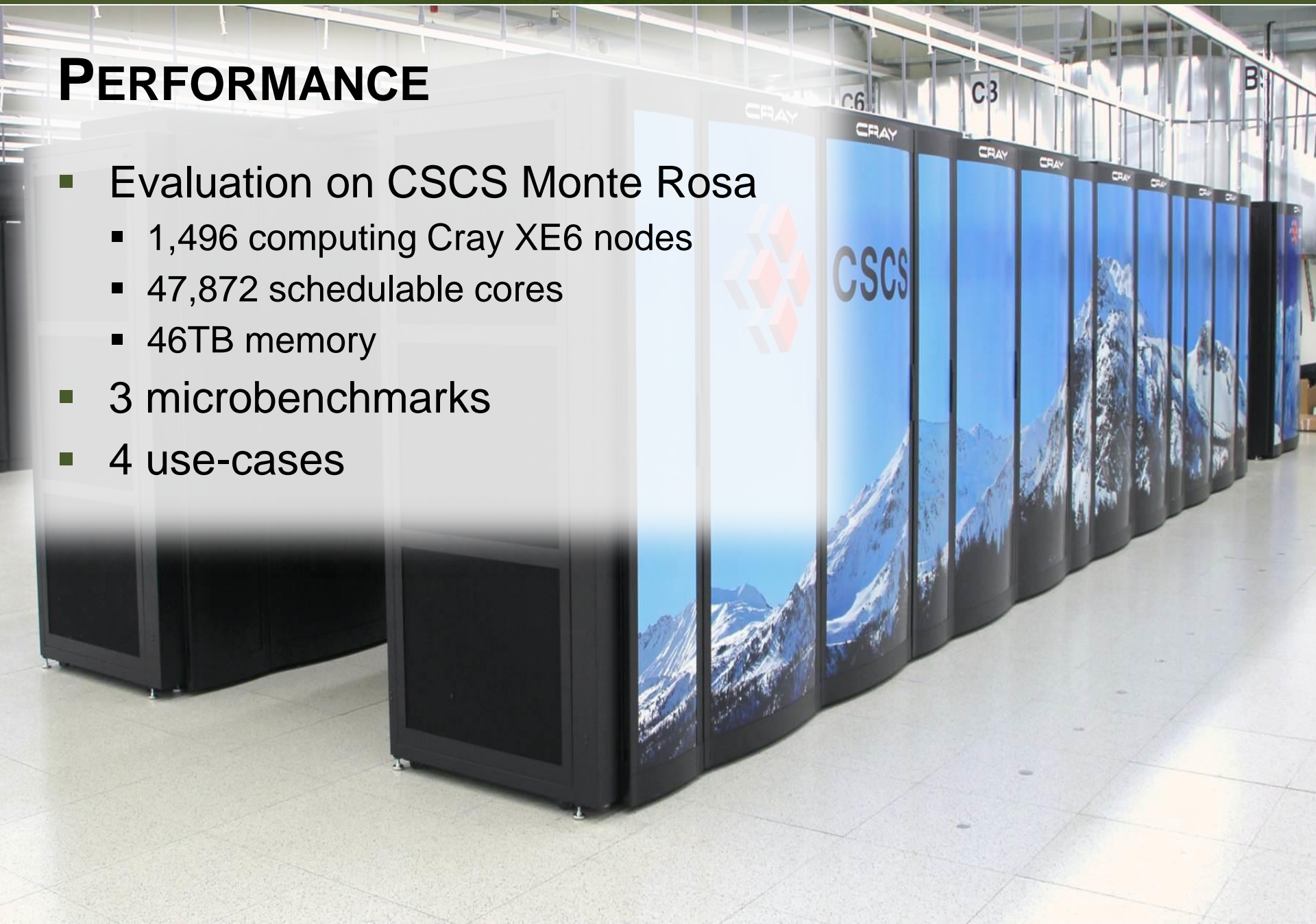
# ACTIVE ACCESS USE-CASES

## VIRTUAL GLOBAL ADDRESS SPACE (V-GAS)



# PERFORMANCE

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



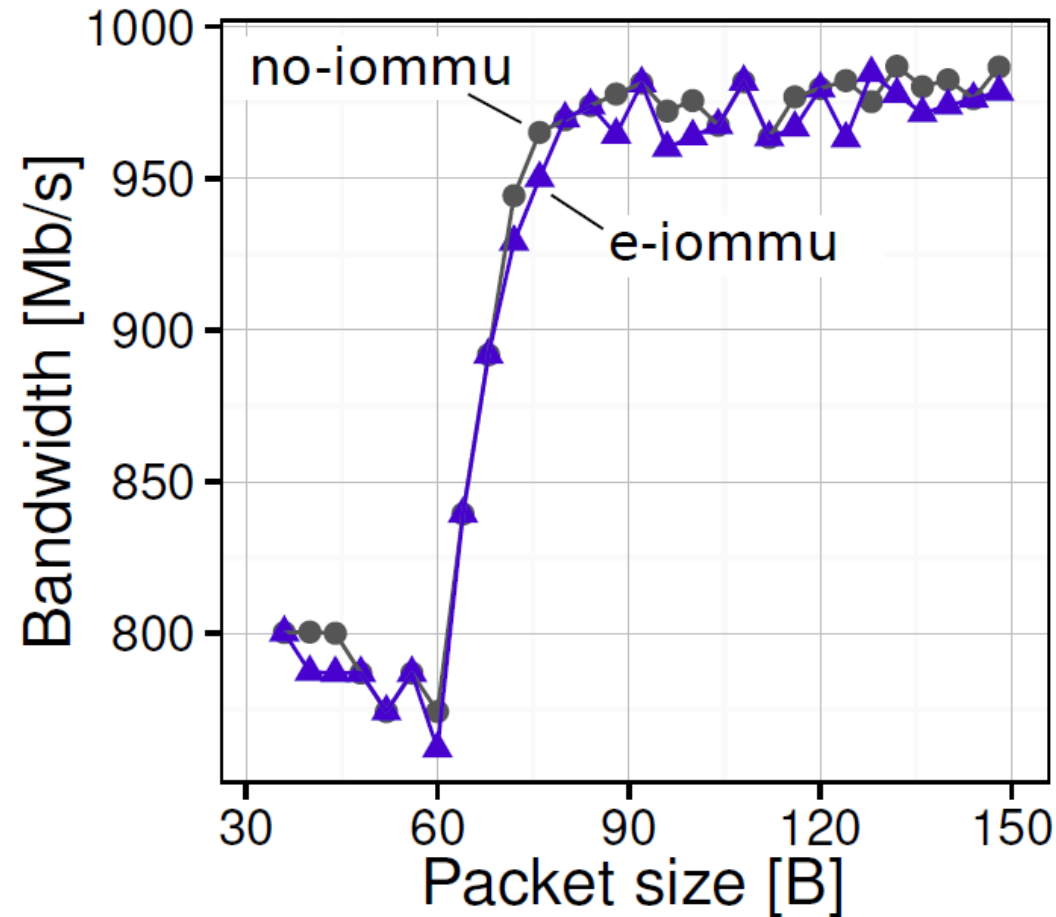
# PERFORMANCE: MICROBENCHMARKS

## RAW DATA TRANSFER

- 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

Active Access

AA-Poll

AA-Int

AA-SP

Active Messages

AM AM-Onload

AM-Exp AM-Ints

RMA



DMAPP



Cell



RoCE



DCMF LAPI  
PAMI

Myricom MX

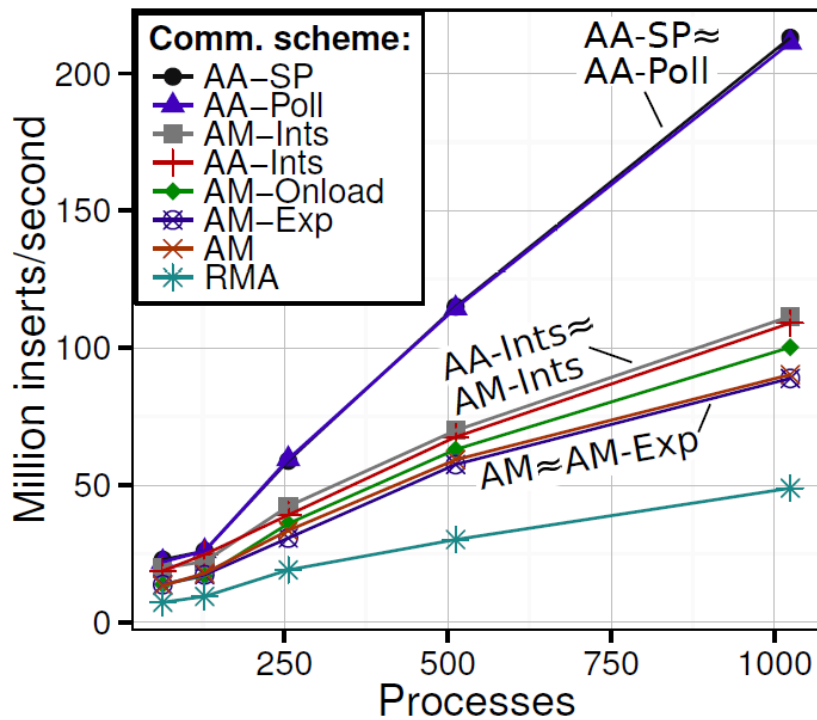
AM+ +GASNet

# PERFORMANCE: LARGE-SCALE CODES

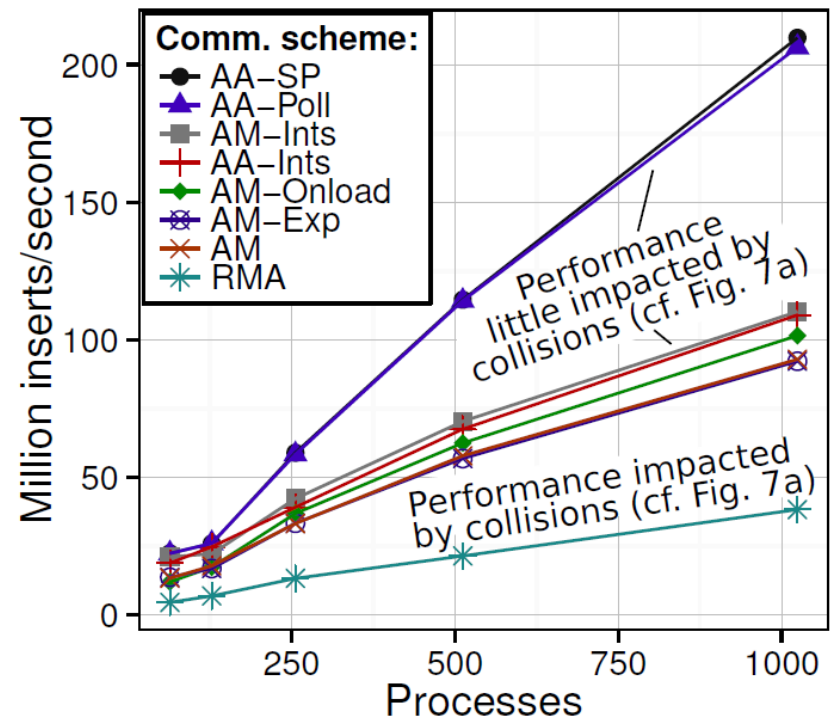
## DISTRIBUTED HASHTABLE



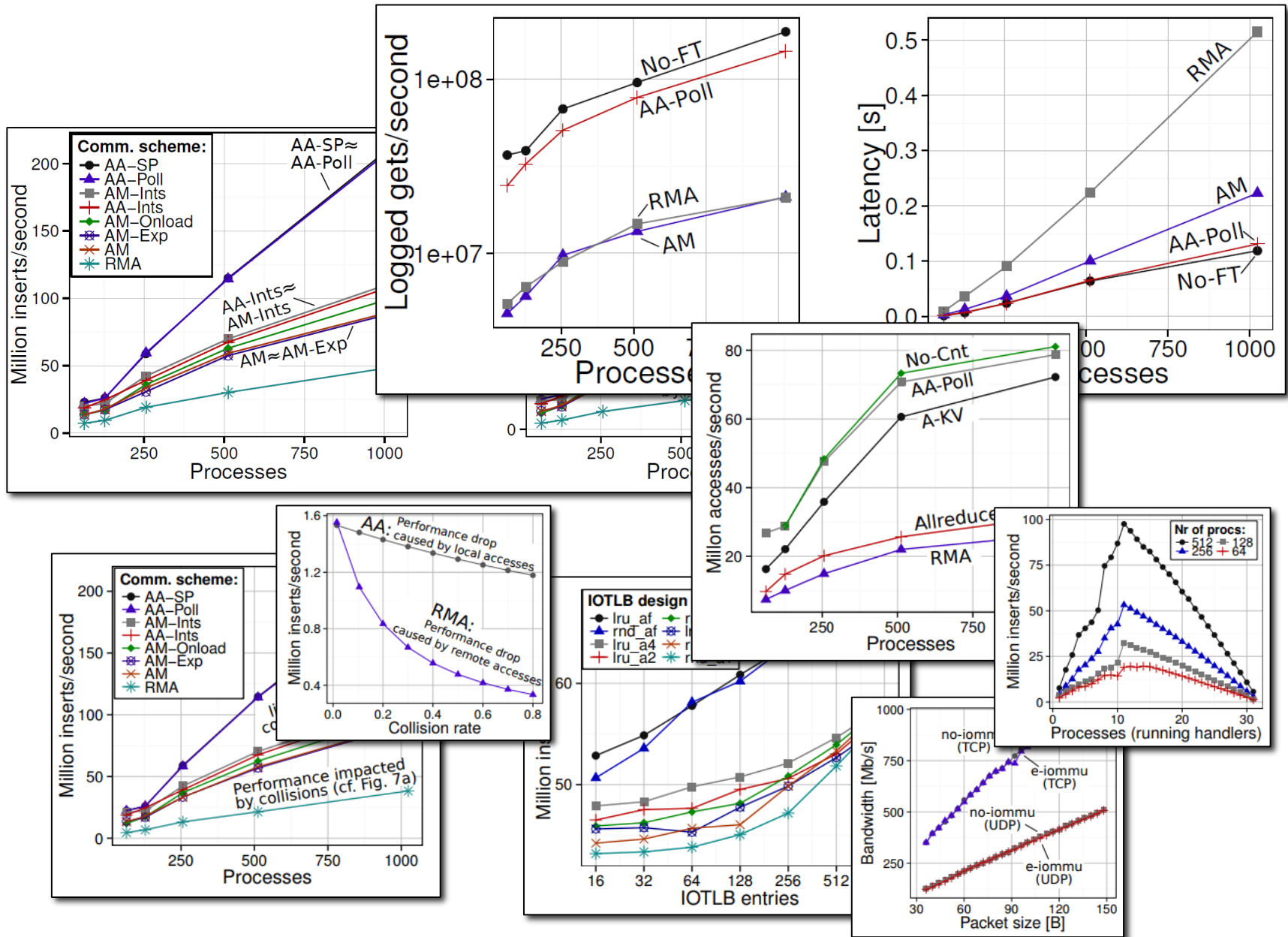
Collisions: 5%



Collisions: 25%



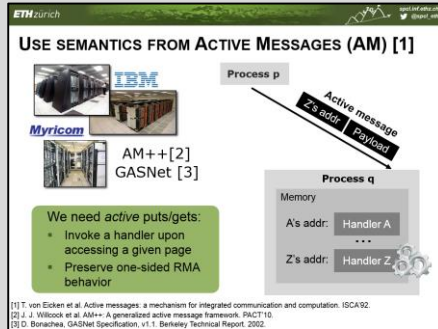






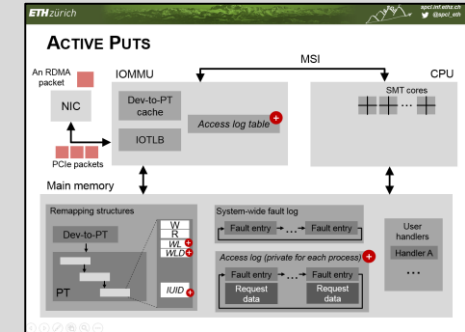
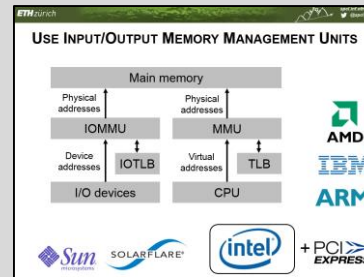
# CONCLUSIONS

## Active Access



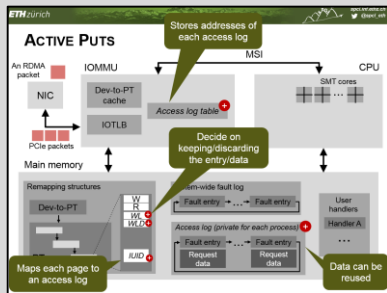
Alleviates RMA's problems with AMs while preserving one-sided semantics

Uses commodity & common IOMMUs

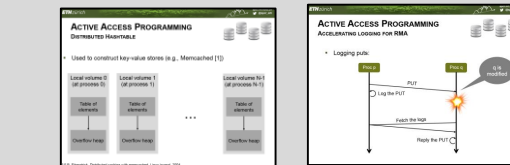


Extends paging capabilities in a distributed environment

## Data-centric programming

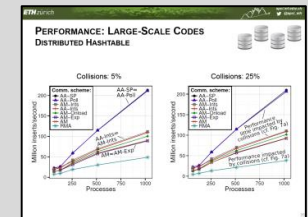
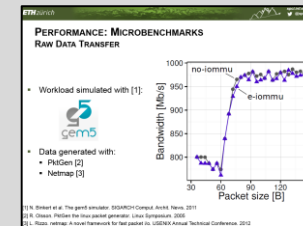


Addresses of pages guide the execution of handlers

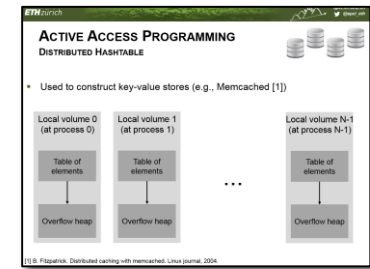
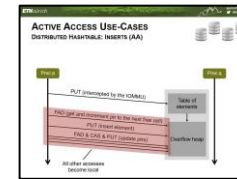
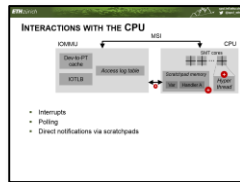
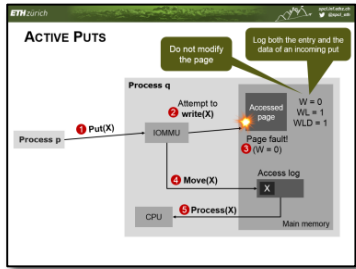
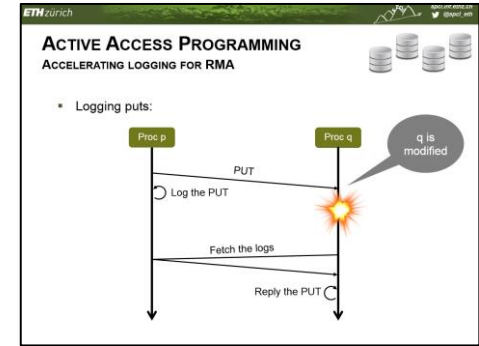
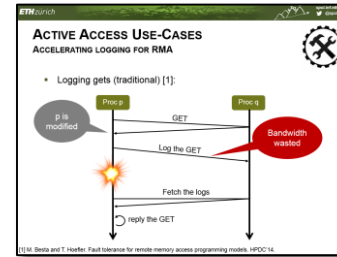
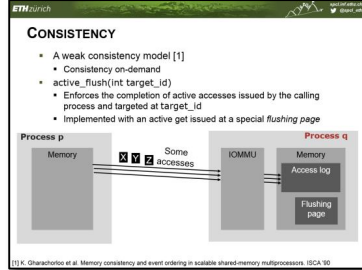
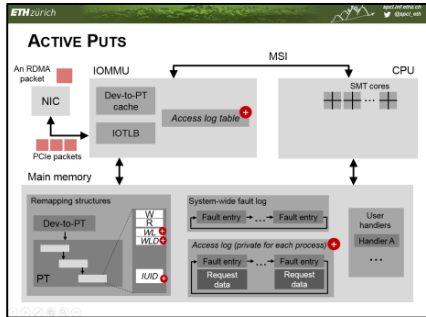


Hashtables, logging schemes, counters, V-GAS, checkpointing...

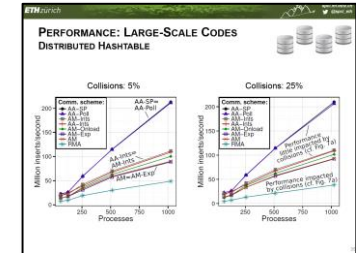
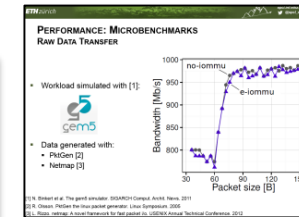
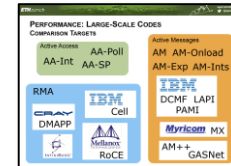
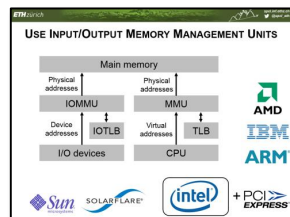
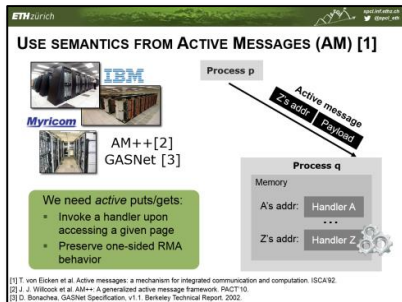
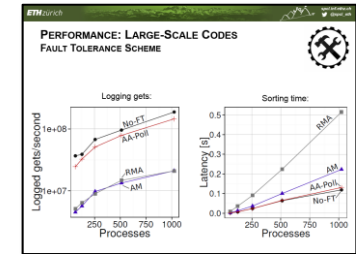
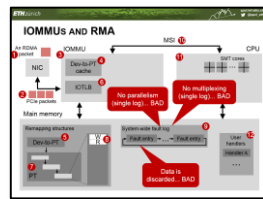
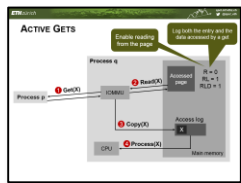
## Performance



Accelerates various distributed codes



# Thank you for your attention



[1] T. van Eicken et al. Active messages: a mechanism for integrated communication and computation. ISCA'02.  
 [2] J. J. Wilcock et al. AM++: A generalized active message framework. PACT'10.  
 [3] D. Bonachea. GASNet Specification, v1.1. Berkeley Technical Report, 2002.

[1] S. Brueck et al. The great server: IBM/STI Cray XE6. March 2007.  
 [2] S. Brueck. Netmap: The network programming framework. Linux Symposium, 2008.  
 [3] S. Brueck. Netmap: A user-friendly framework for packet I/O. Linux Symposium, Technical Session, 2009.

# ACTIVE ACCESS USE-CASES

## ACCELERATING LOGGING FOR RMA



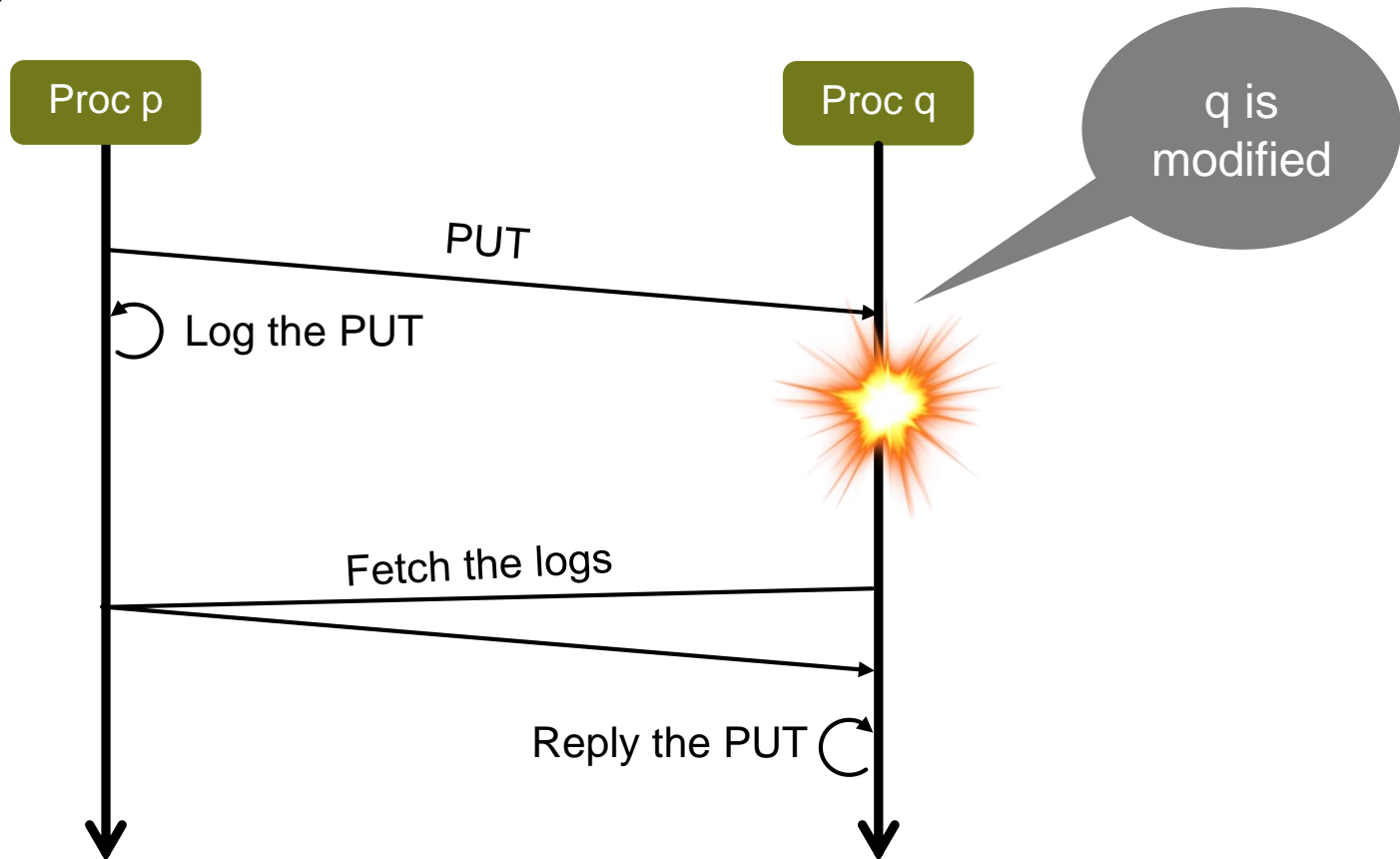
- Logging – a popular mechanism for fault-tolerance.
- Remote communication (puts/gets) is logged.
- Upon a process crash, it is restored and uses the logs to replay its previous actions.
- Logs are stored in volatile memories.

# ACTIVE ACCESS USE-CASES

## ACCELERATING LOGGING FOR RMA



- Logging puts:

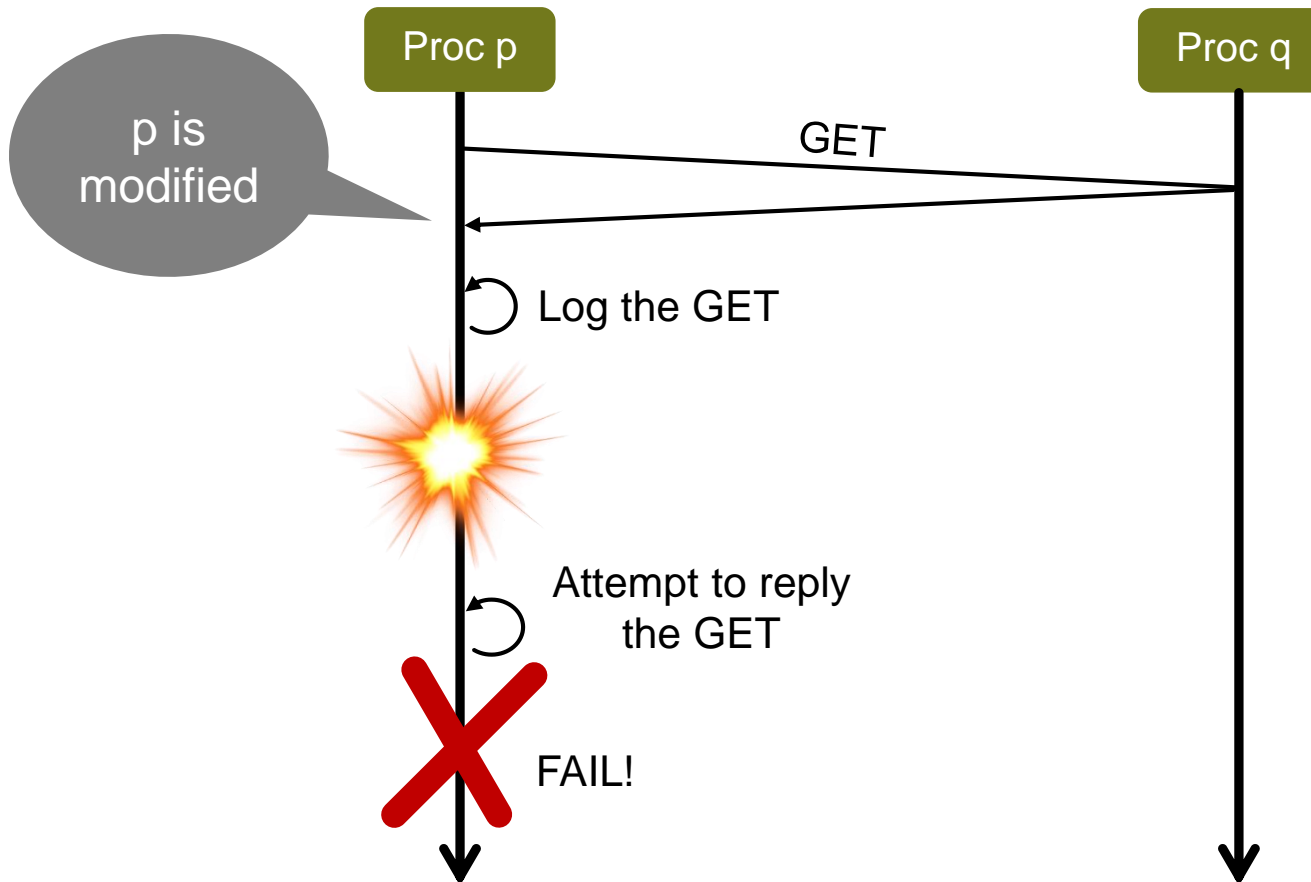


# ACTIVE ACCESS USE-CASES

## ACCELERATING LOGGING FOR RMA



- Logging gets (naive):

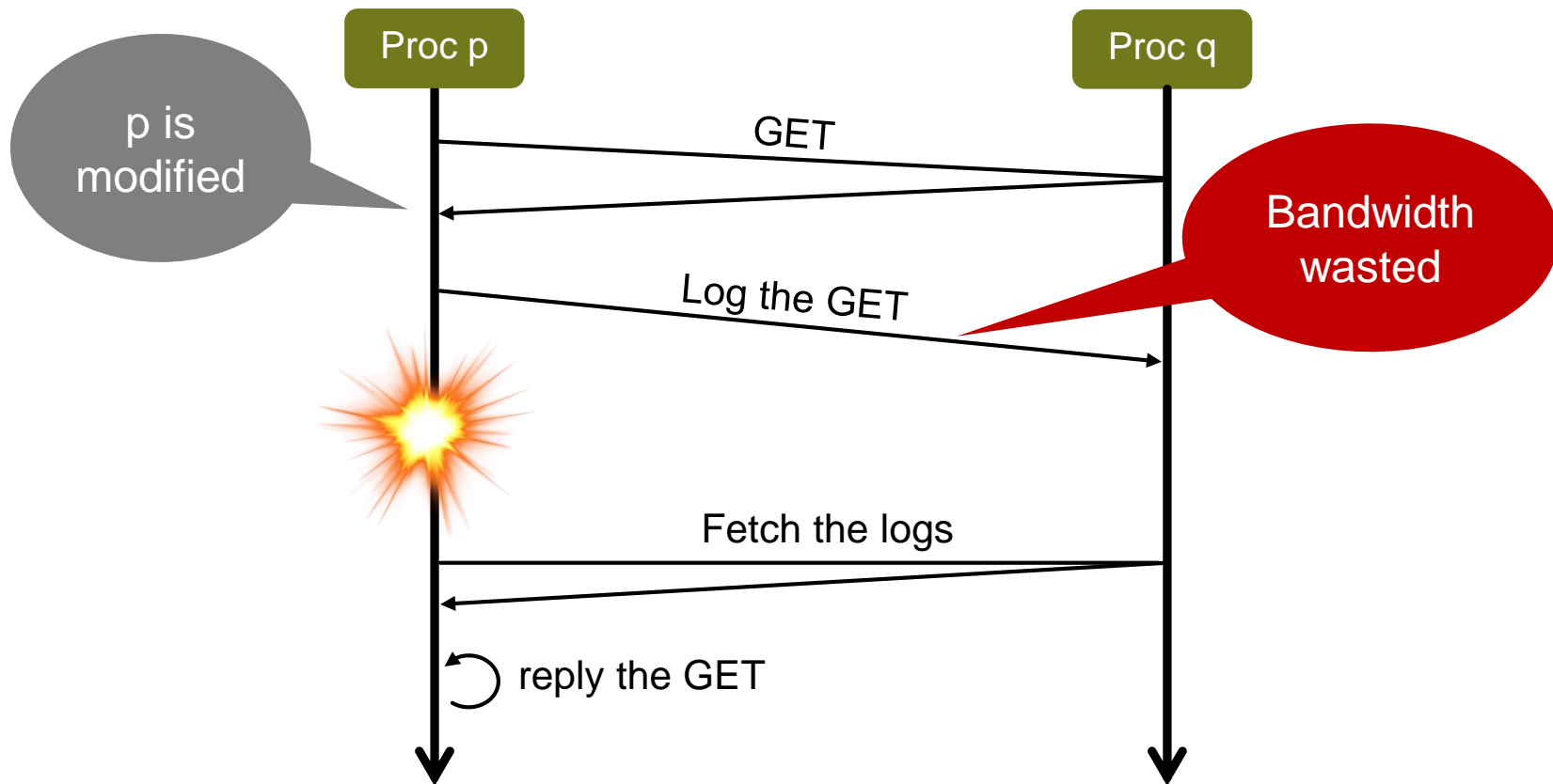


# ACTIVE ACCESS USE-CASES

## ACCELERATING LOGGING FOR RMA



- Logging gets (traditional) [1]:



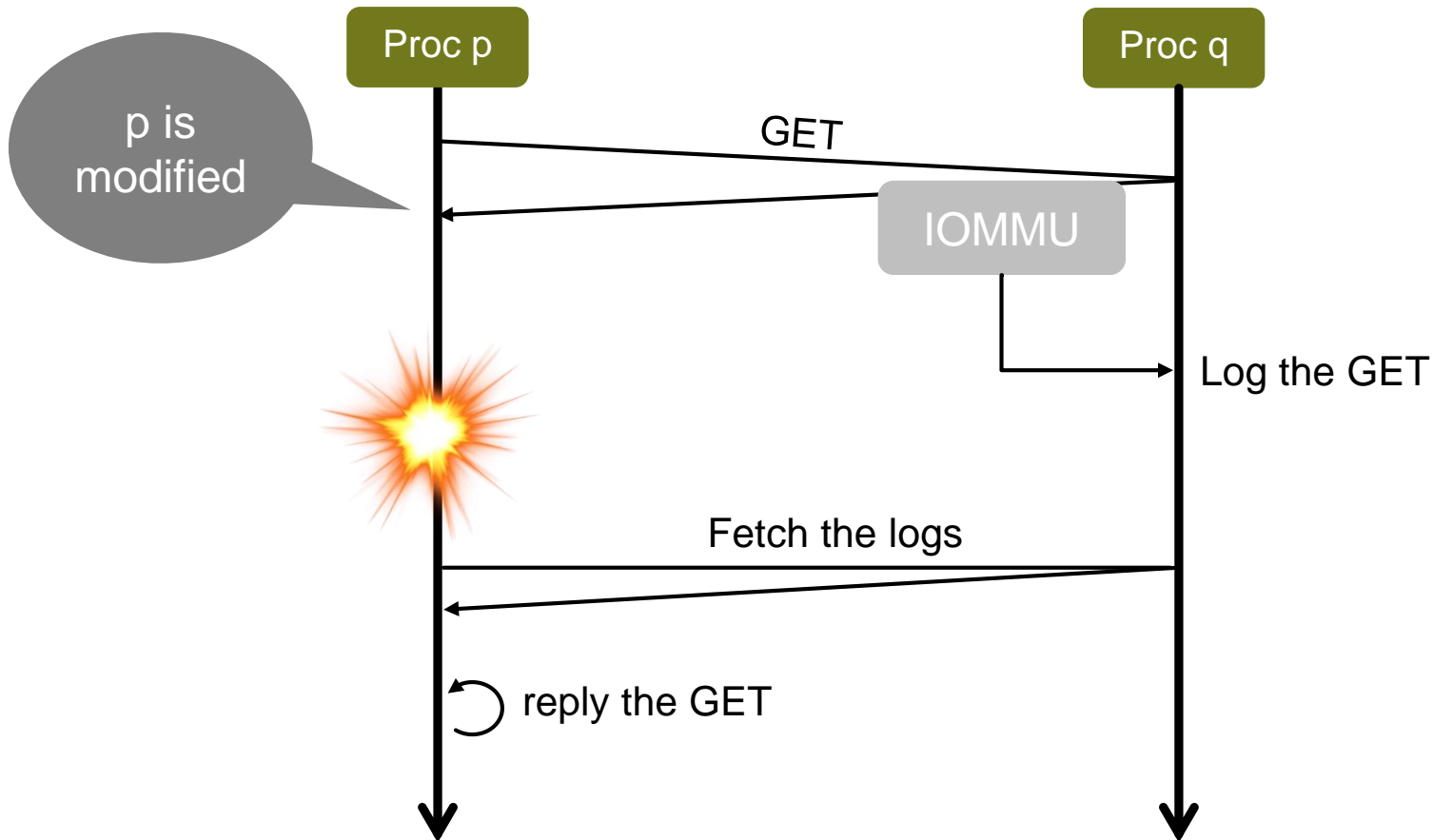
[1] M. Besta and T. Hoefler. Fault tolerance for remote memory access programming models. HPDC'14.

# ACTIVE ACCESS USE-CASES

## ACCELERATING LOGGING FOR RMA



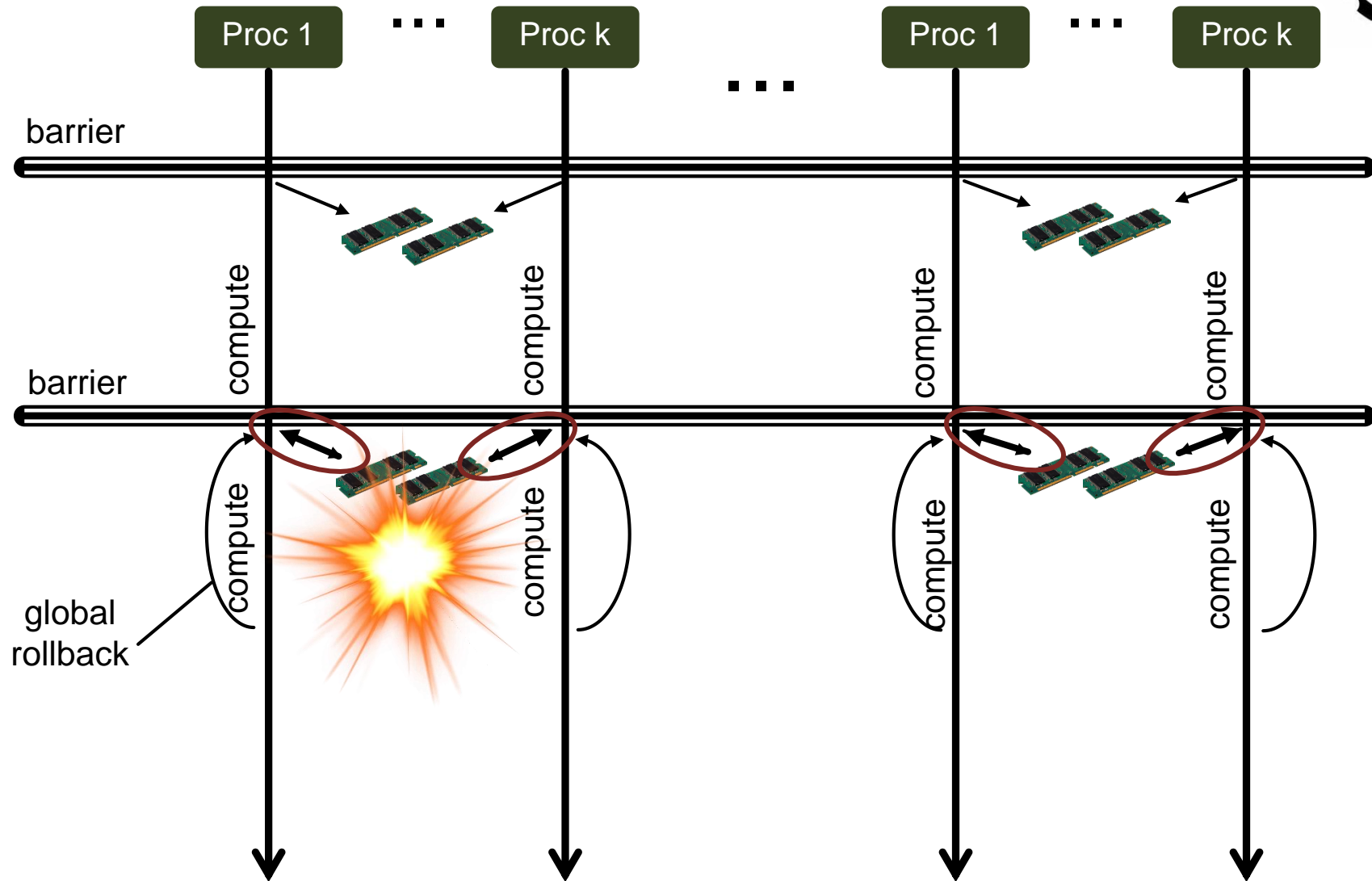
- Logging gets (AA):



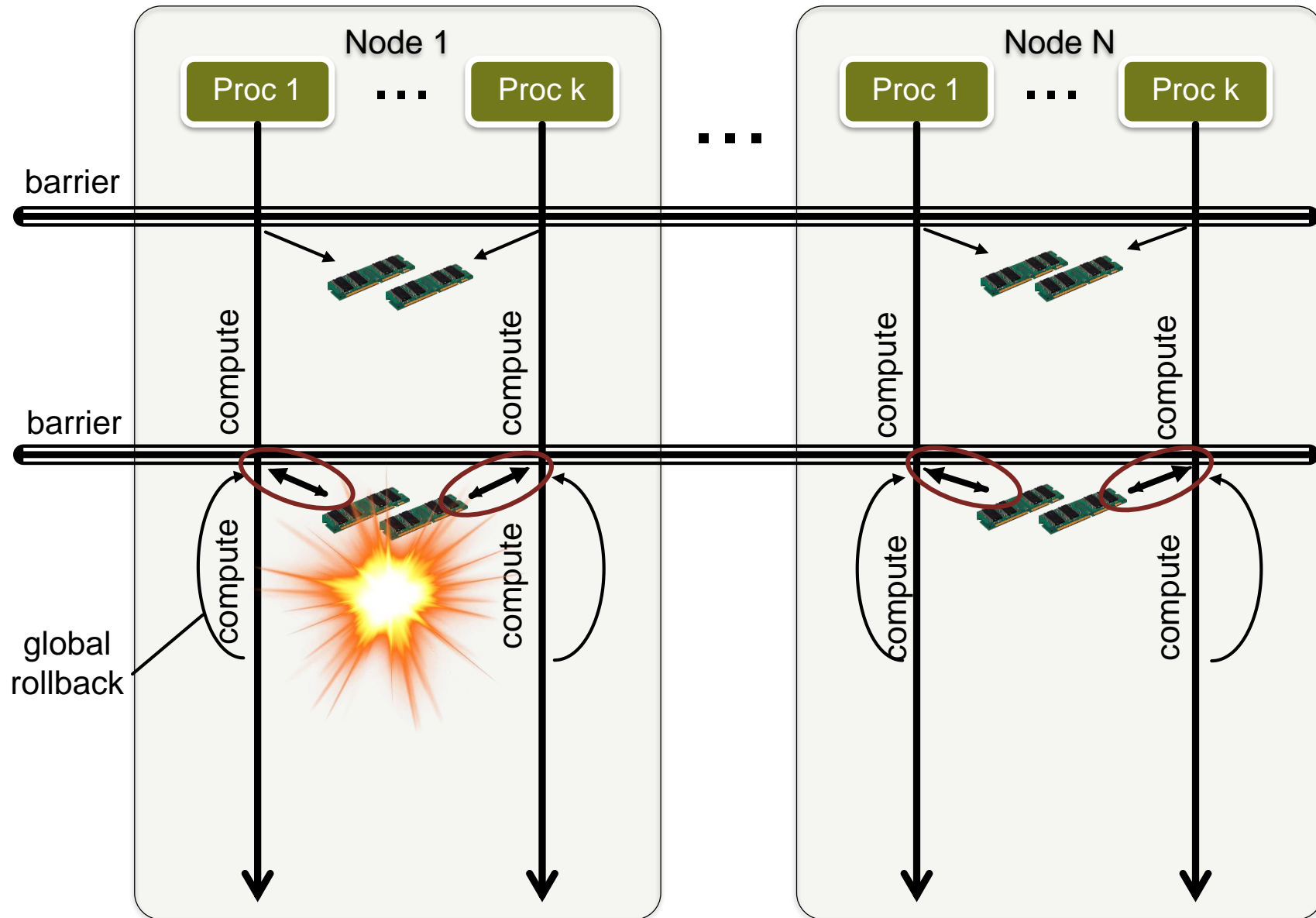


# ACTIVE ACCESS USE-CASES

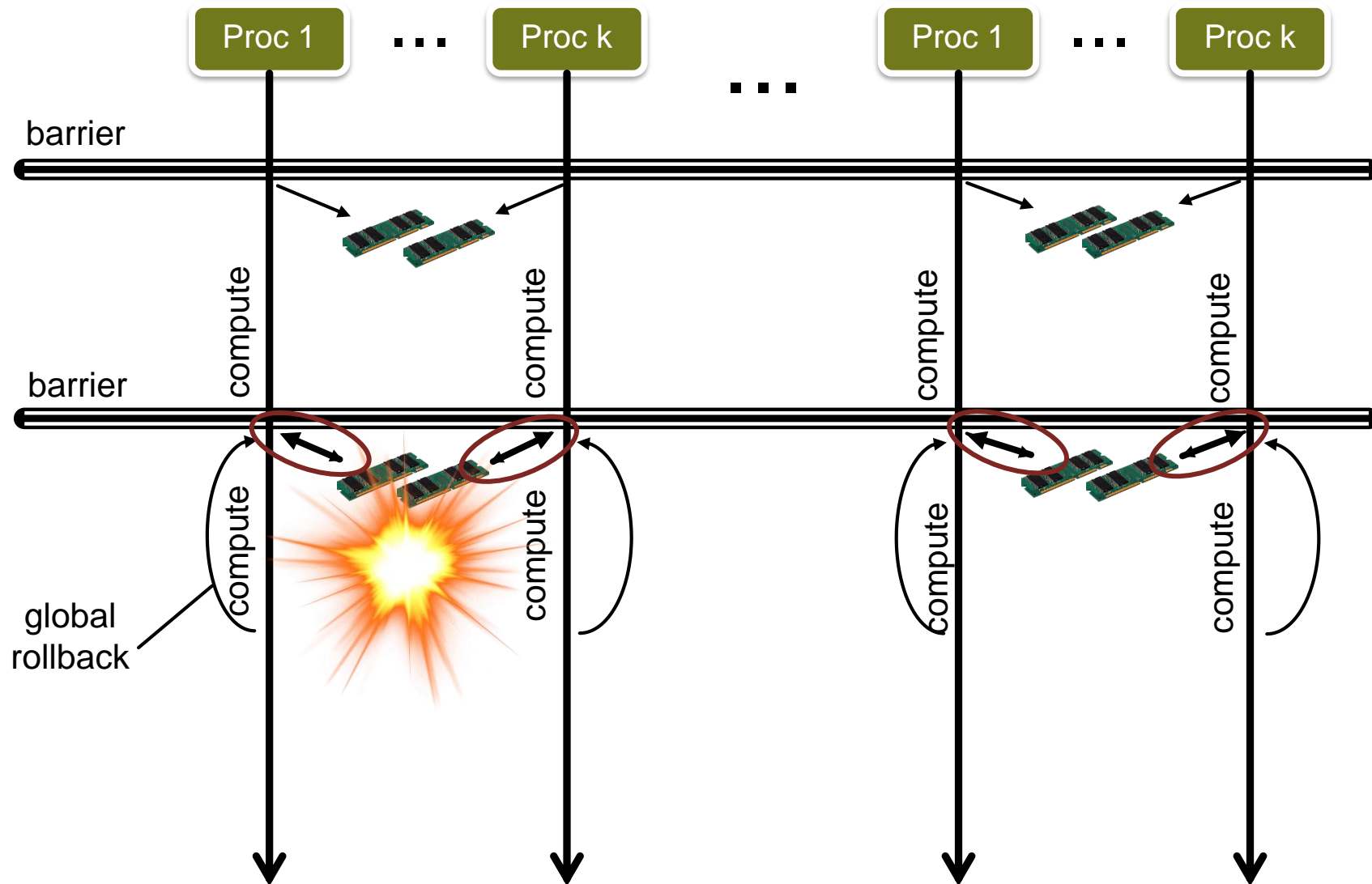
## INCREMENTAL CHECKPOINTING FOR RMA



# COORDINATED CHECKPOINTING (MP)



# COORDINATED CHECKPOINTING (MP)

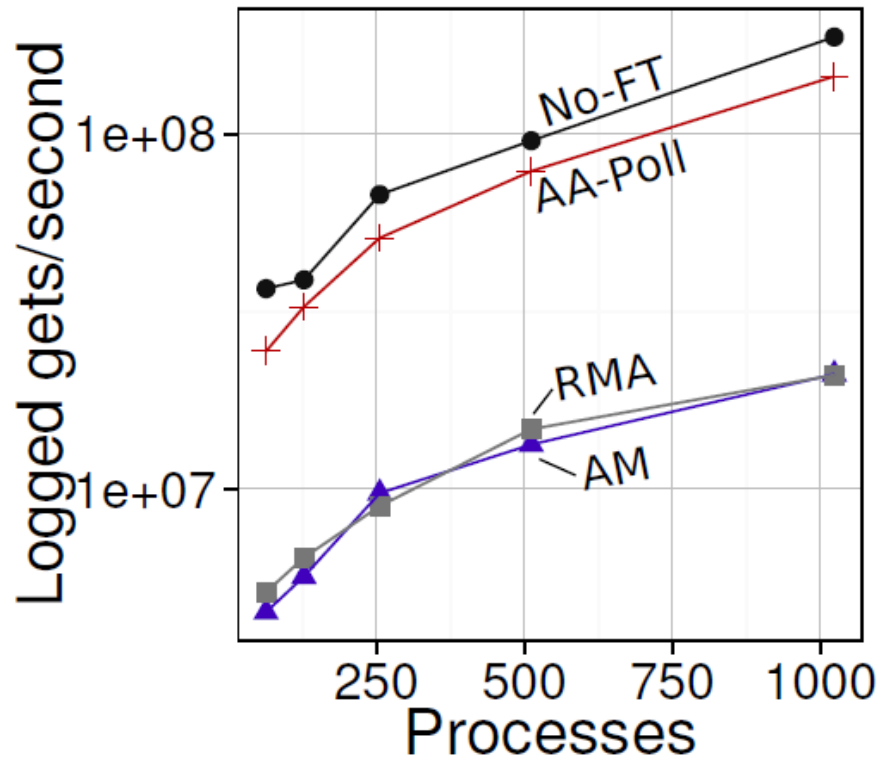


# PERFORMANCE: LARGE-SCALE CODES

## FAULT TOLERANCE SCHEME



Logging gets:



Sorting time:

