

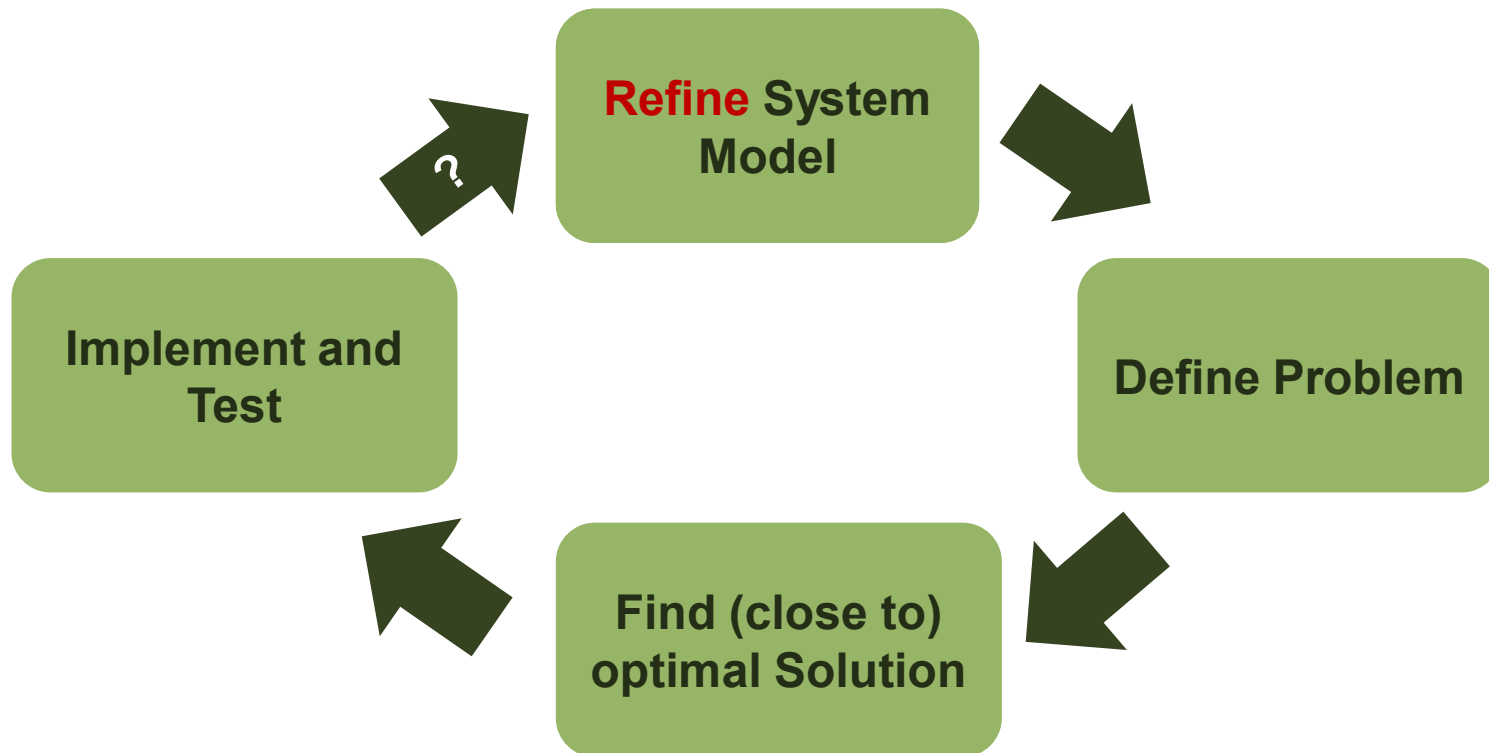
TORSTEN HOEFLER

Remote Memory Access Programming - Tools and Fault Tolerance



Model-based Performance Engineering

- **My dream: provably optimal performance (time and energy)**
 - From problem to machine code

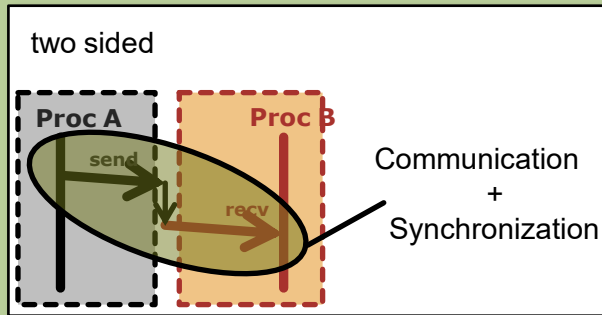


- **A philosophy for system and application design**
 - At different levels of course

State of the Art – Parallel Programming

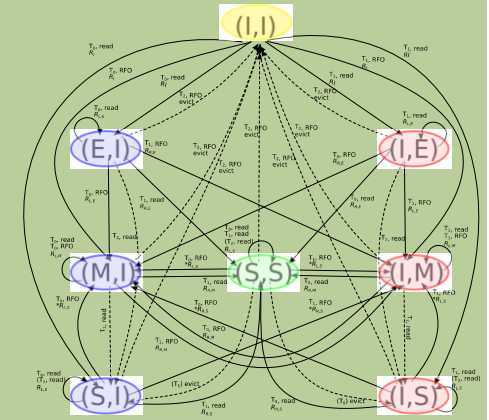
Message Passing (MPI-1)

- De-facto programming model



Coherent Shared Memory

- Hardware support ☺
- Very very hard to optimize [1]

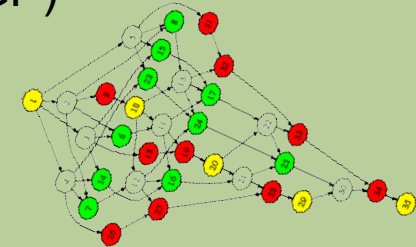


Requirements for a new low-level programming model

1. Messaging needs 100% hardware support (offload) – it's simple after all!
2. Minimal overheads (**tiny**) layer between user and hardware
3. Offer a simple abstract performance model (e.g., LogGP)

DPLASMA

CONSORTIUM FOR SMALL SCALE MODELING
COSMO



MPI-3.0 RMA

- MPI-3.0 supports RMA (“MPI One Sided”)
 - Designed to react to hardware trends
 - Majority of HPC networks support RDMA

**Random datacenter picture
copyrighted by Reuters (yes, they
go after academics with claims for
10 year old images)**

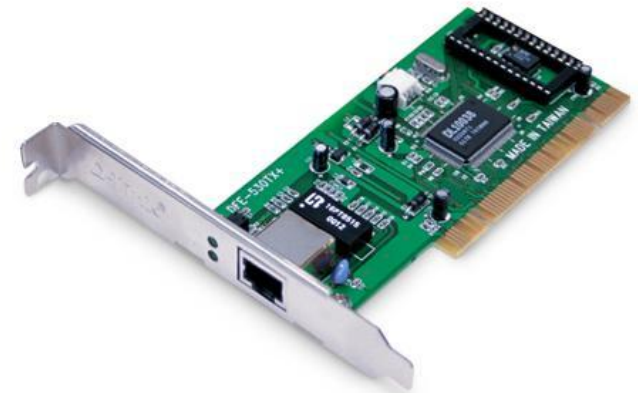


MPI-3.0 RMA

- MPI-3.0 supports RMA (“MPI One Sided”)
 - Designed to react to hardware trends
 - Majority of HPC networks support RDMA

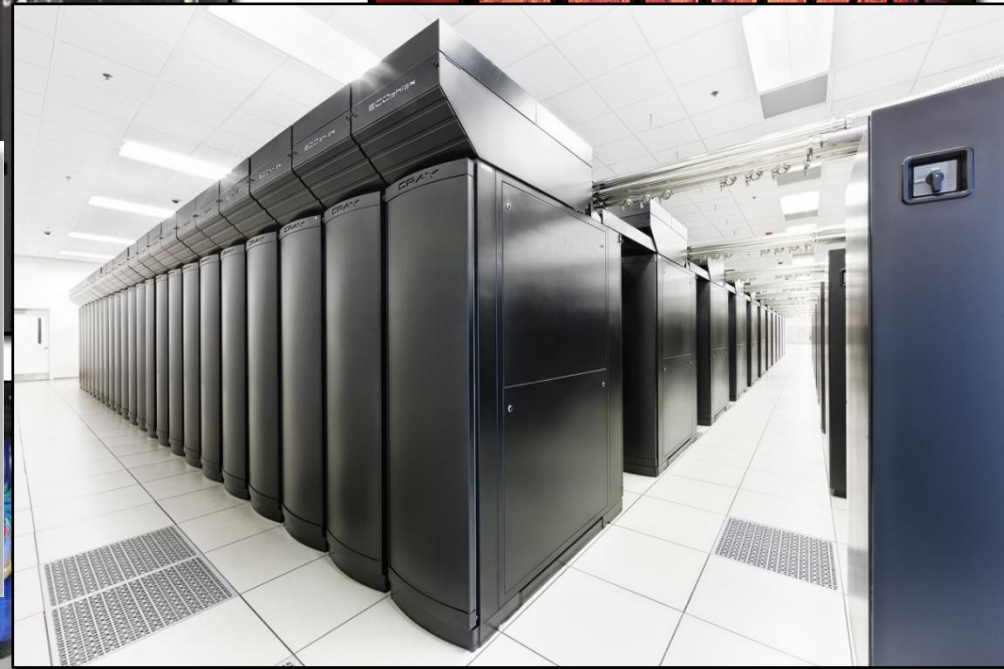


**Random datacenter picture
copyrighted by Reuters (yes, they
go after academics with claims for
10 year old images)**



MPI-3.0 RMA

- MPI-3.0 supports RMA (“MPI One Sided”)
 - Designed to react to hardware trends
 - Majority of HPC networks support RDMA



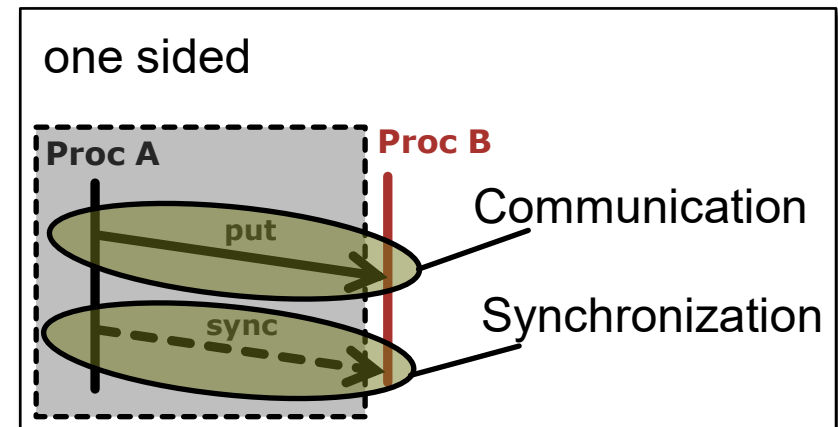
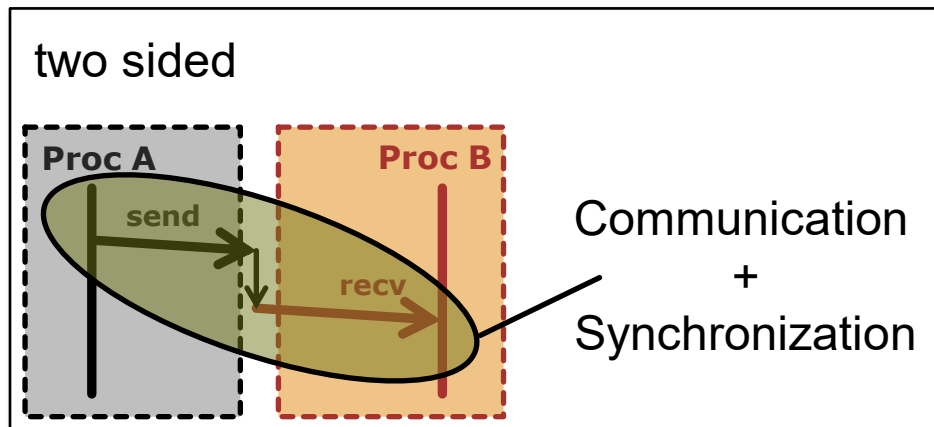
**Random datacenter picture
copyrighted by Reuters (yes, they
go after academics with claims for
10 year old images)**



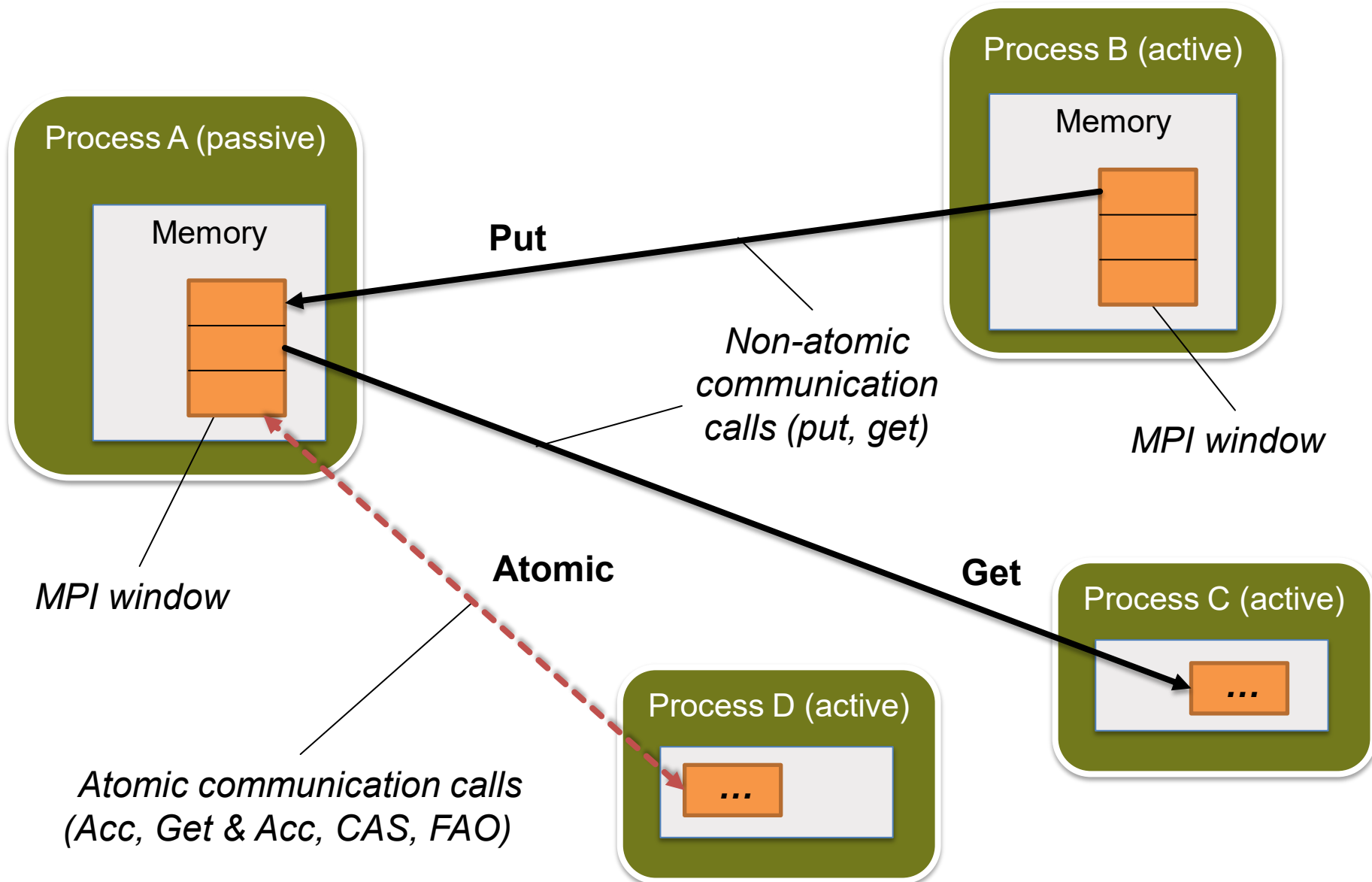
MPI-3.0 RMA

- MPI-3.0 supports RMA (“MPI One Sided”)
 - Designed to react to hardware trends
 - Majority of HPC networks support RDMA
- Communication is „one sided” (no involvement of destination)
- RMA decouples communication & synchronization
 - Different from message passing

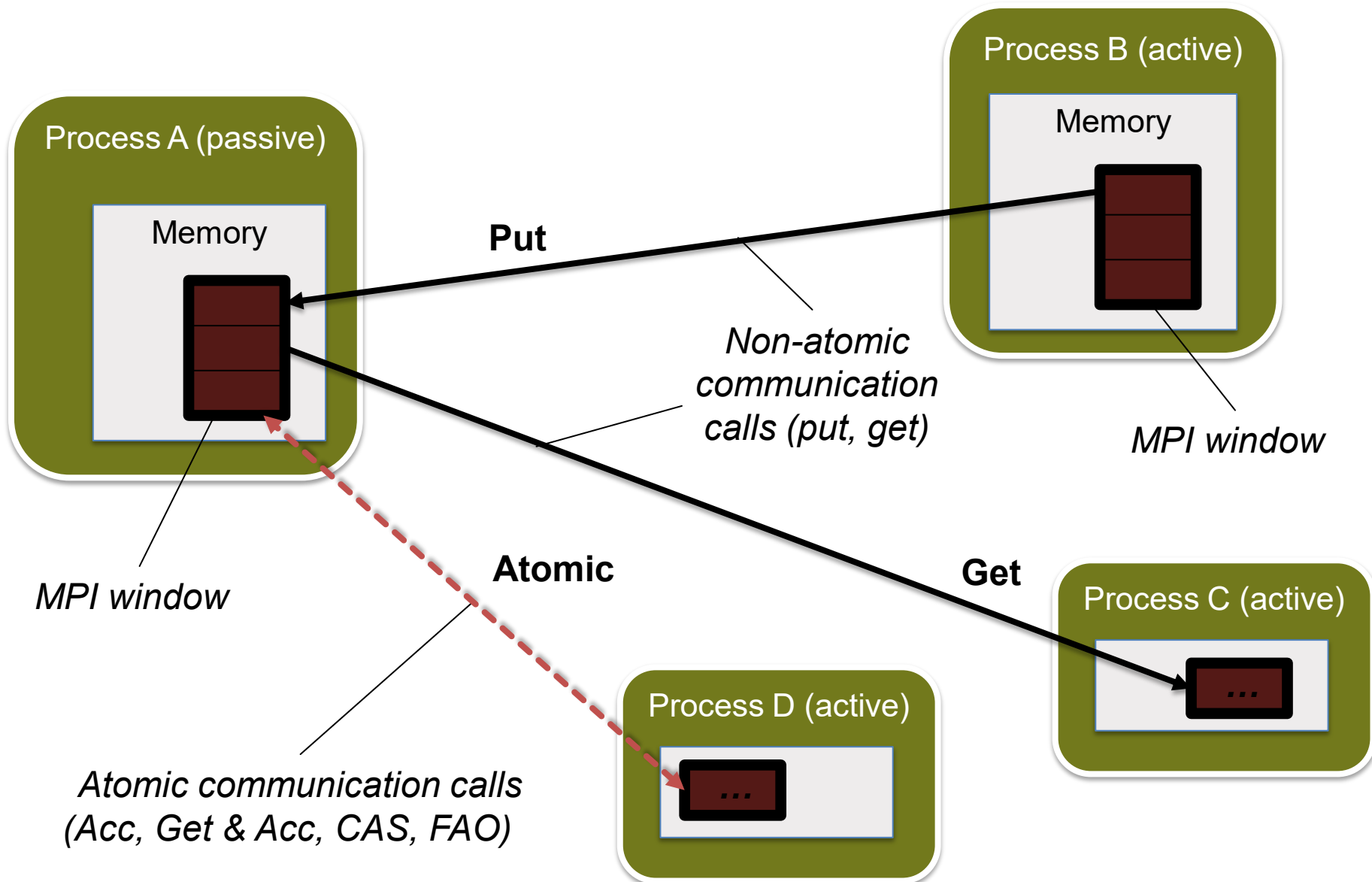
Random datacenter picture
copyrighted by Reuters (yes, they
go after academics with claims for
10 year old images)



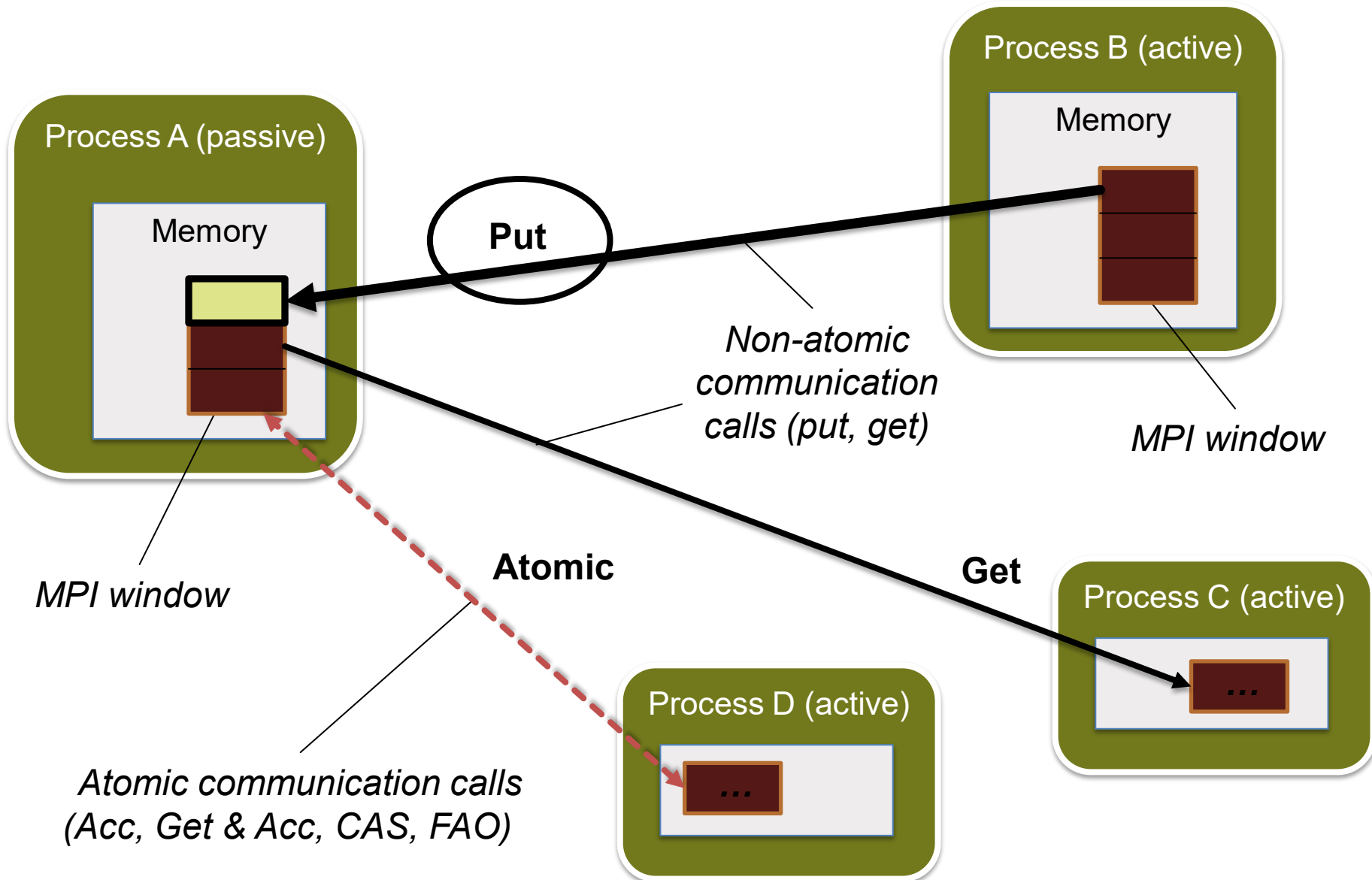
MPI-3 RMA COMMUNICATION OVERVIEW



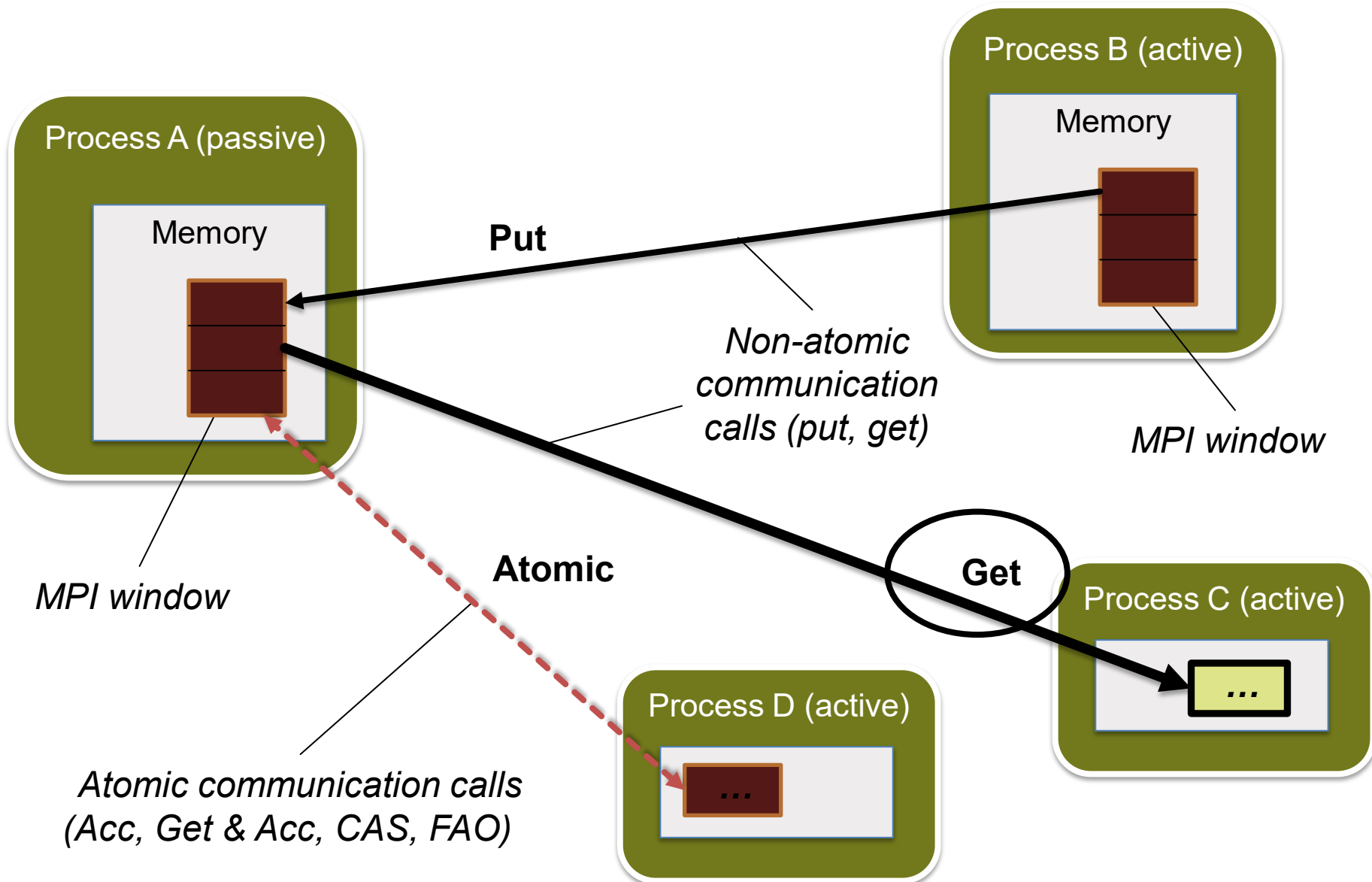
MPI-3 RMA COMMUNICATION OVERVIEW



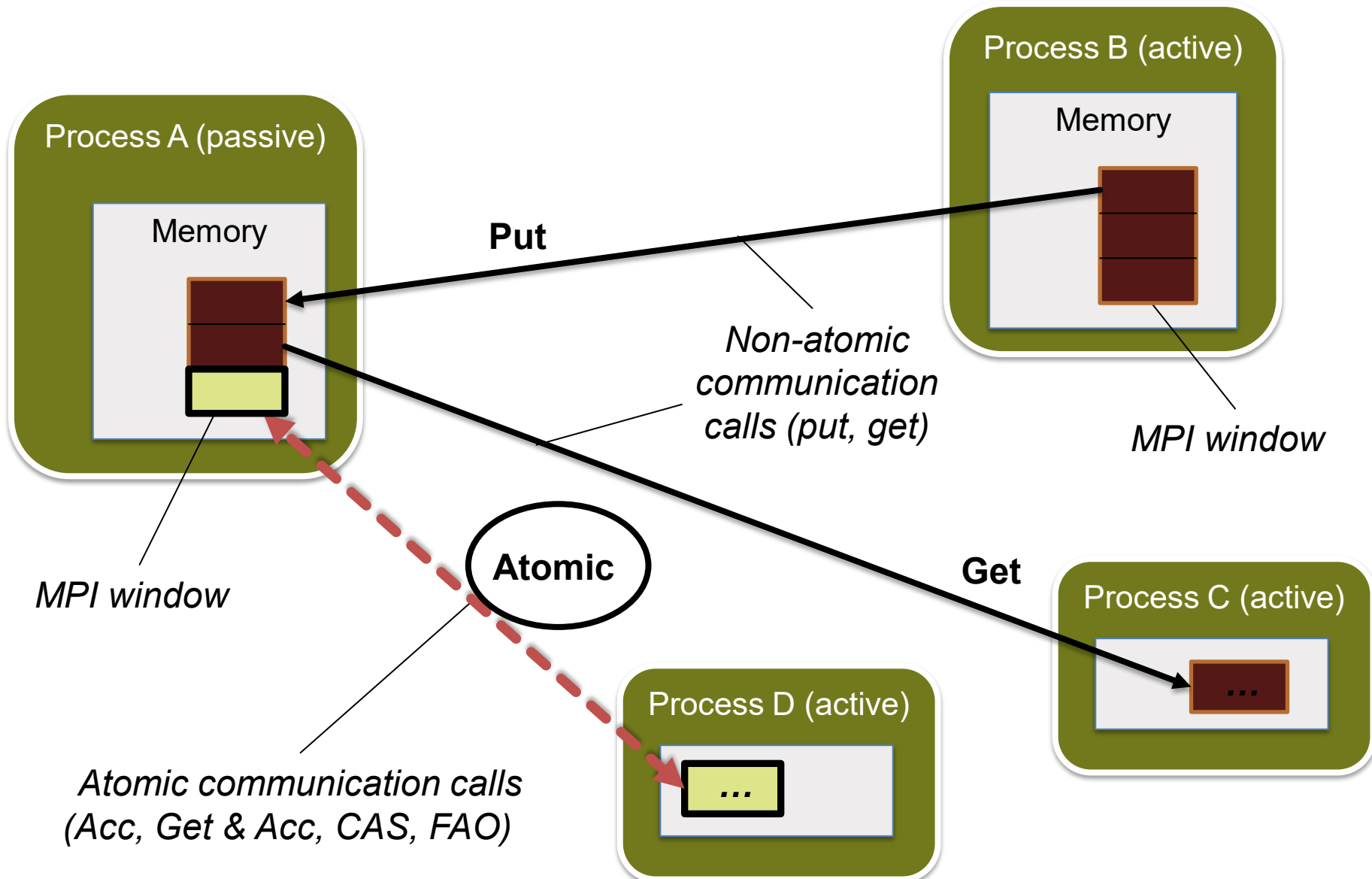
MPI-3 RMA COMMUNICATION OVERVIEW



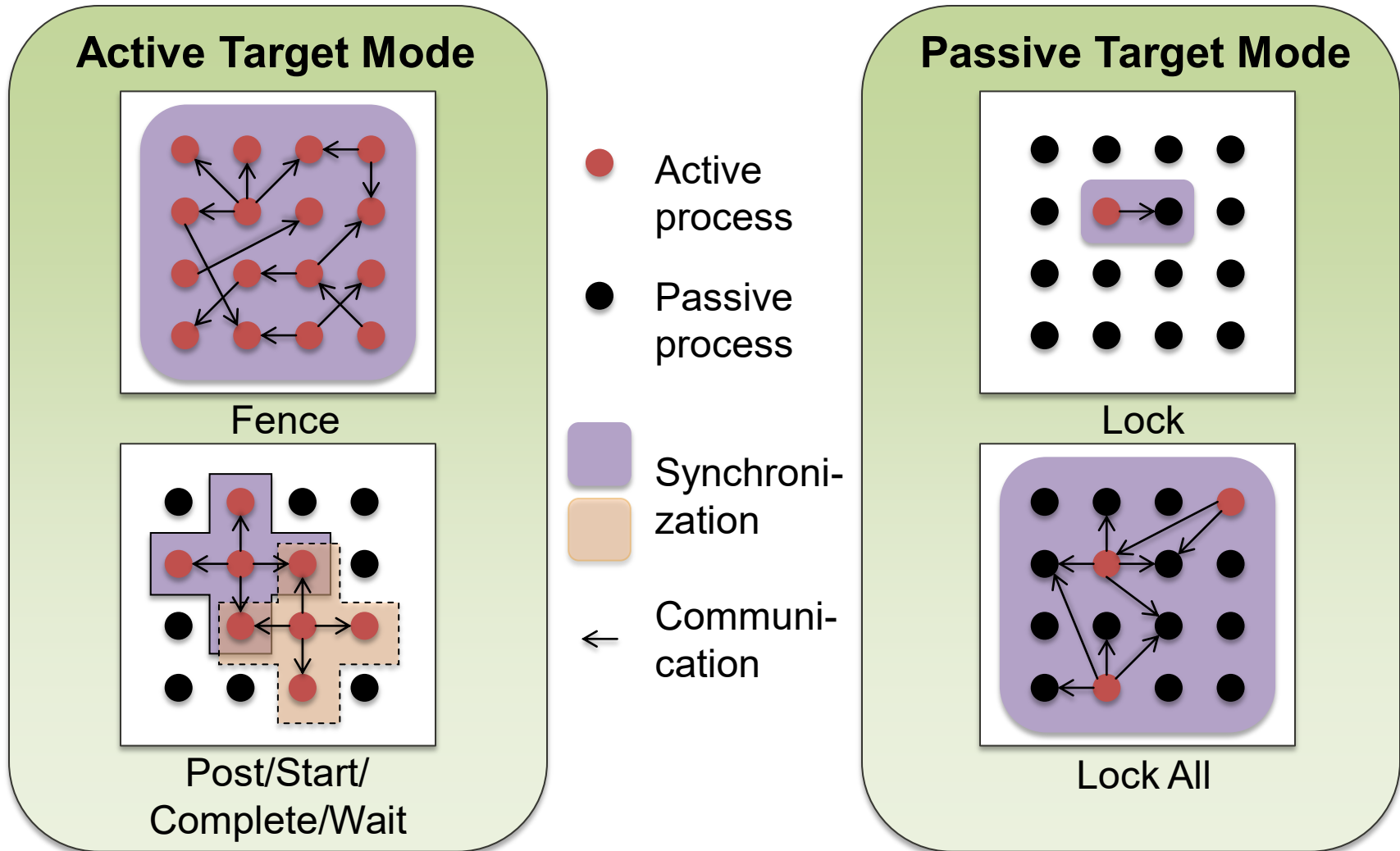
MPI-3 RMA COMMUNICATION OVERVIEW



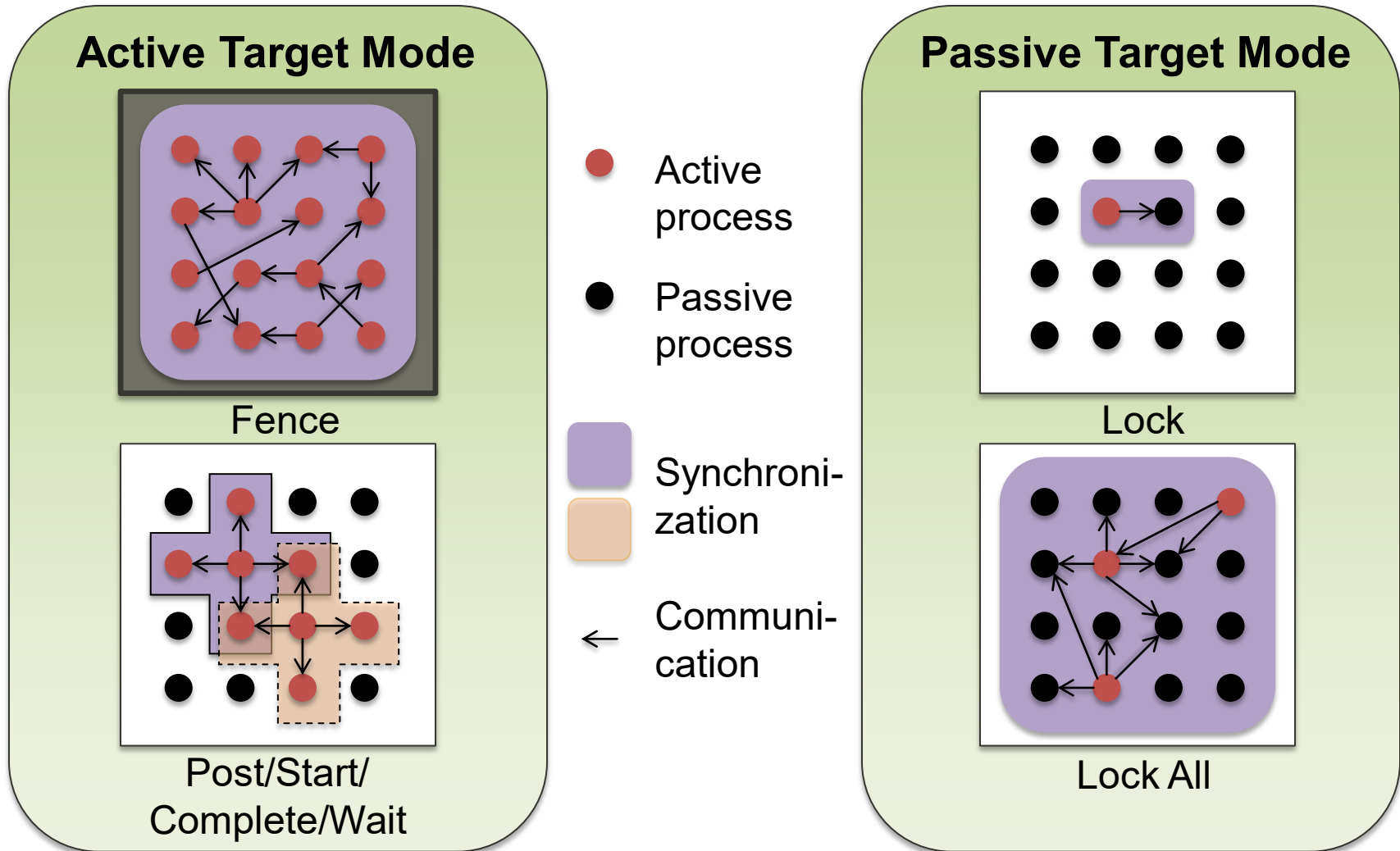
MPI-3 RMA COMMUNICATION OVERVIEW



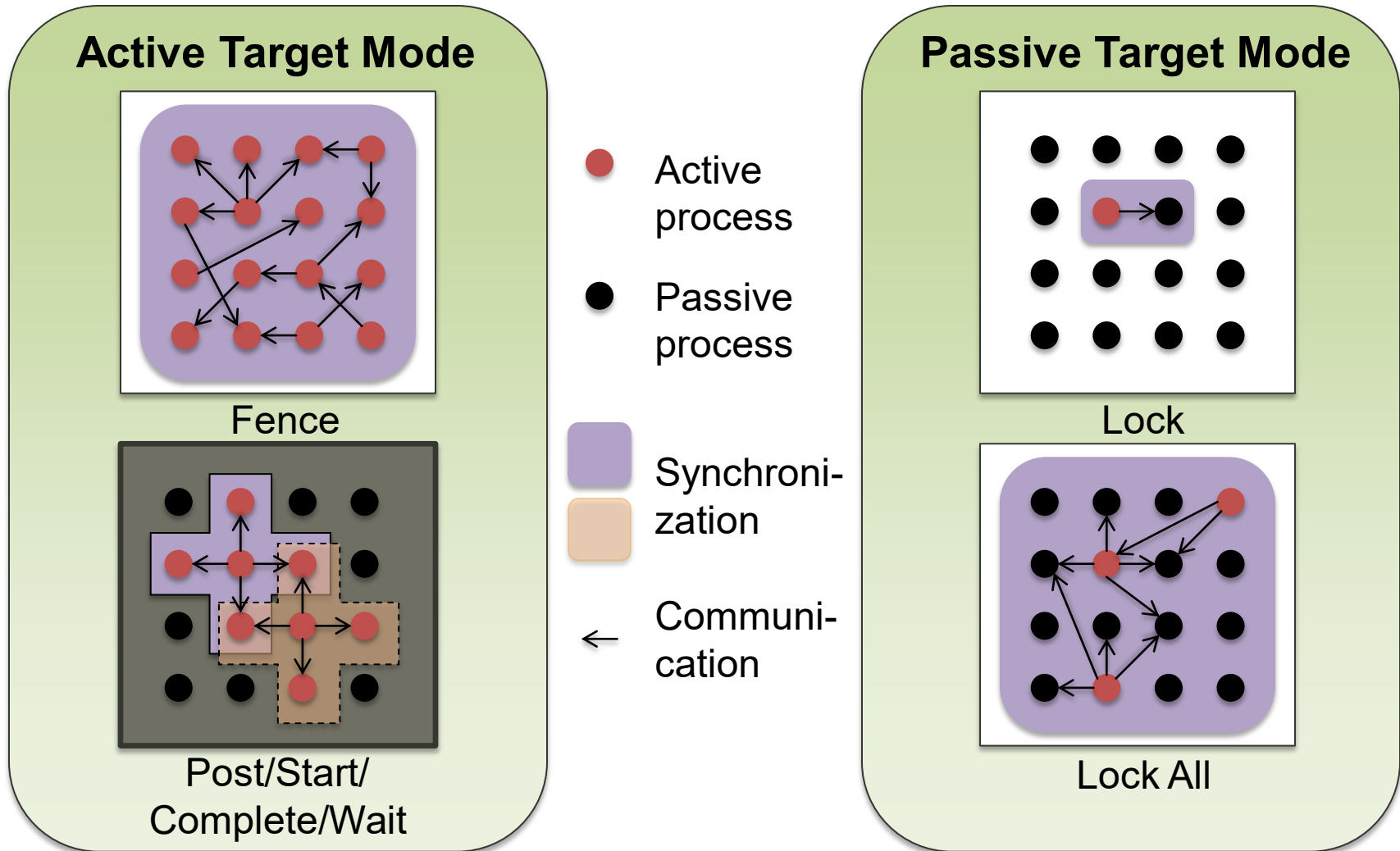
MPI-3.0 RMA SYNCHRONIZATION OVERVIEW



MPI-3.0 RMA SYNCHRONIZATION OVERVIEW

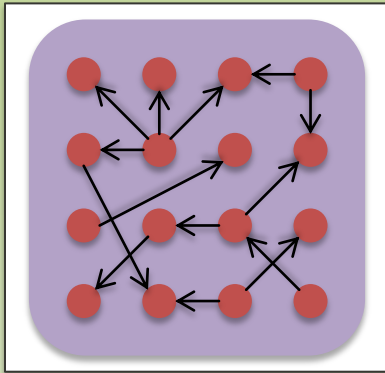


MPI-3.0 RMA SYNCHRONIZATION OVERVIEW

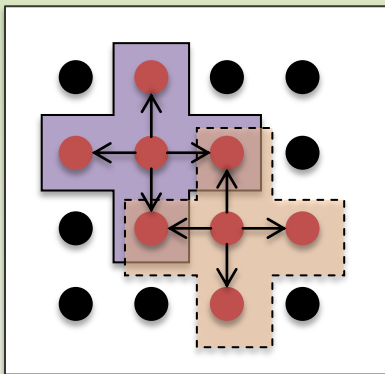


MPI-3.0 RMA SYNCHRONIZATION OVERVIEW

Active Target Mode



Fence



Post/Start/
Complete/Wait

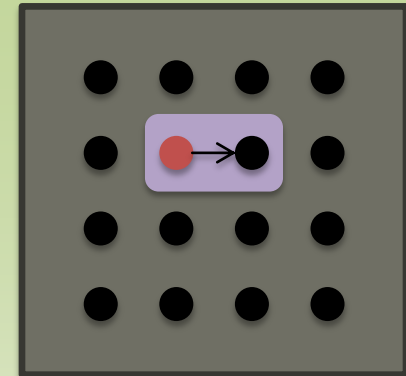
● Active process

● Passive process

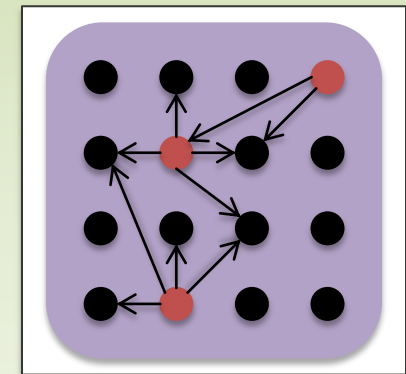
■ Synchroni-
zation

← Communi-
cation

Passive Target Mode



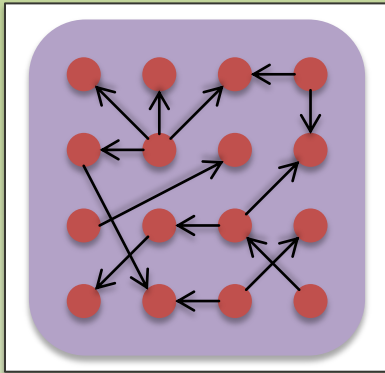
Lock



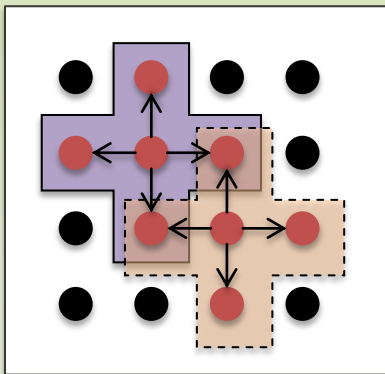
Lock All

MPI-3.0 RMA SYNCHRONIZATION OVERVIEW

Active Target Mode

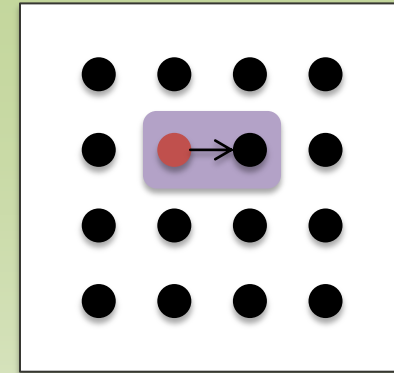


Fence

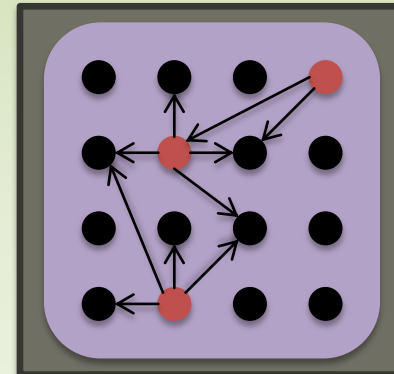


Post/Start/
Complete/Wait

Passive Target Mode



Lock



Lock All

● Active process

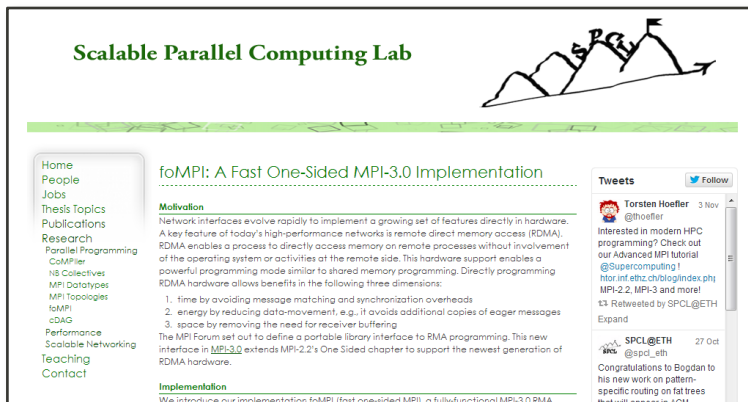
● Passive process

■ Synchroni-
zation

← Communi-
cation

SCALABLE PROTOCOLS & REFERENCE IMPLEMENTATION

- Scalable & generic protocols
 - Can be used on any RDMA network (e.g., OFED/IB)
 - Window creation, communication and synchronization
- foMPI, a fully functional MPI-3 RMA implementation
 - DMAPP: lowest-level networking API for Cray Gemini/Aries systems
 - XPMEM, a portable Linux kernel module



Scalable Parallel Computing Lab

foMPI: A Fast One-Sided MPI-3.0 Implementation

Motivation

Network interfaces evolve rapidly to implement a growing set of features directly in hardware. A key feature of today's high-performance networks is remote direct memory access (RDMA). RDMA enables a process to directly access memory on remote processes without involvement of the operating system or activities at the remote side. This hardware support enables a powerful programming mode similar to shared memory programming. Directly programming RDMA hardware allows benefits in the following three dimensions:

1. time by avoiding message matching and synchronization overheads
2. energy by reducing data-movement, e.g., if avoids additional copies of eager messages
3. space by removing the need for receiver buffering

The MPI Forum set out to define a portable library interface to RMA programming. This new interface in [MPI-3.0](#) extends MPI-2.2's One Sided chapter to support the newest generation of RDMA hardware.

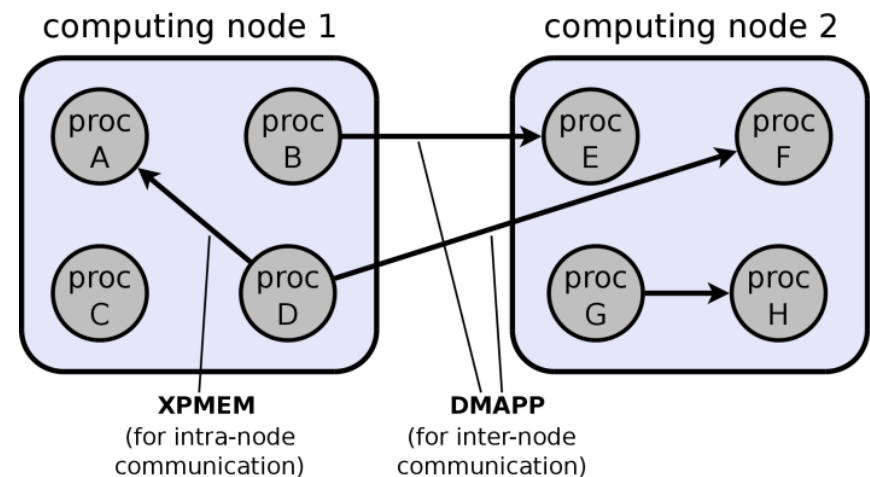
Implementation

We introduce our implementation foMPI (fast one-sided MPI) as fully functional MPI-3.0 RMA

Tweets

Torsten Hoefler @thoefler 3 Nov
Interested in modern HPC programming? Check out our Advanced MPI tutorial @Supercomputing1 http://inf.ethz.ch/blog/index.php/MPI-2.2, MPI-3 and more!
Retweeted by SPCL@ETH Expand

SPCL@ETH @spcl_eth 27 Oct
Congratulations to Bogdan to his new work on pattern-specific routing on fat trees that will speed up MPI



SCALABLE PROTOCOLS & REFERENCE IMPLEMENTATION

- Scalable & generic protocols
 - Can be used on any RDMA network (e.g., OFED/IB)
 - Window creation, communication and synchronization

- foMPI, a fully functional MPI-3 RMA implementation
 - DMAPP: lowest-level networking API for Cray Gemini/Aries systems
 - XPMEM, a portable Linux kernel module

Scalable Parallel Computing Lab

foMPI: A Fast One-Sided MPI-3.0 Implementation

Motivation
 Network interfaces evolve rapidly to implement a growing set of features directly in hardware. A key feature of today's high-performance networks is remote direct memory access (RDMA). RDMA enables a process to directly access memory on remote processes without involvement of the operating system or activities at the remote side. This hardware support enables a powerful programming mode similar to shared memory programming. Directly programming RDMA hardware allows benefits in the following three dimensions:

1. time by avoiding message matching and synchronization overheads
2. energy by reducing data-movement, e.g., if avoids additional copies of eager messages
3. space by removing the need for receiver buffering

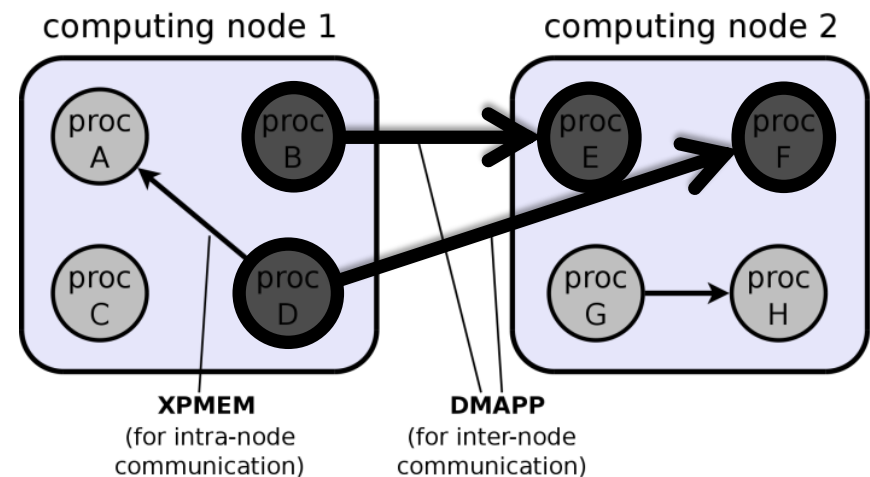
The MPI Forum set out to define a portable library interface to RMA programming. This new interface in [MPI-3.0](#) extends MPI-2.2's One Sided chapter to support the newest generation of RDMA hardware.

Implementation
 We introduce our implementation foMPI (fast one-sided MPI) as fully functional MPI-3.0 RMA

Tweets

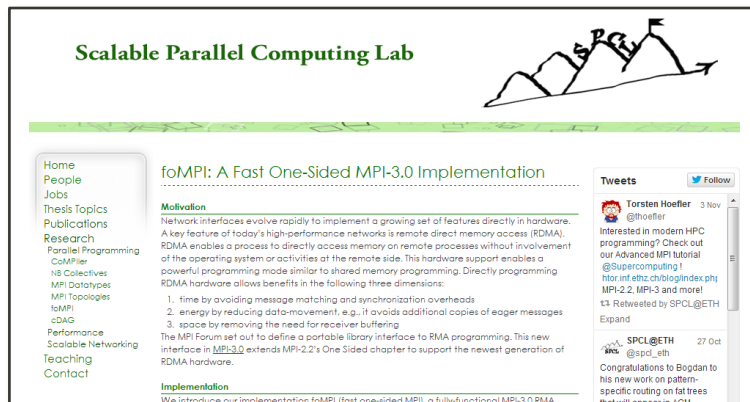
Torsten Hoefler @thoefler 3 Nov
 Interested in modern HPC programming? Check out our Advanced MPI tutorial @Supercomputing1 http://inf.ethz.ch/blog/index.php/MPI-2.2, MPI-3 and more!
 Retweeted by SPCL@ETH Expand

SPCL@ETH @spcl_eth 27 Oct
 Congratulations to Bogdan to his new work on pattern-specific routing on fat trees that will speed up MPI



SCALABLE PROTOCOLS & REFERENCE IMPLEMENTATION

- Scalable & generic protocols
 - Can be used on any RDMA network (e.g., OFED/IB)
 - Window creation, communication and synchronization
- foMPI, a fully functional MPI-3 RMA implementation
 - DMAPP: lowest-level networking API for Cray Gemini/Aries systems
 - XPMEM: a portable Linux kernel module



Scalable Parallel Computing Lab

foMPI: A Fast One-Sided MPI-3.0 Implementation

Motivation

Network interfaces evolve rapidly to implement a growing set of features directly in hardware. A key feature of today's high-performance networks is remote direct memory access (RDMA). RDMA enables a process to directly access memory on remote processes without involvement of the operating system or activities at the remote side. This hardware support enables a powerful programming mode similar to shared memory programming. Directly programming RDMA hardware allows benefits in the following three dimensions:

1. time by avoiding message matching and synchronization overheads
2. energy by reducing data-movement, e.g., if avoids additional copies of eager messages
3. space by removing the need for receiver buffering

The MPI Forum set out to define a portable library interface to RMA programming. This new interface in [MPI-3.0](#) extends MPI-2.2's One Sided chapter to support the newest generation of RDMA hardware.

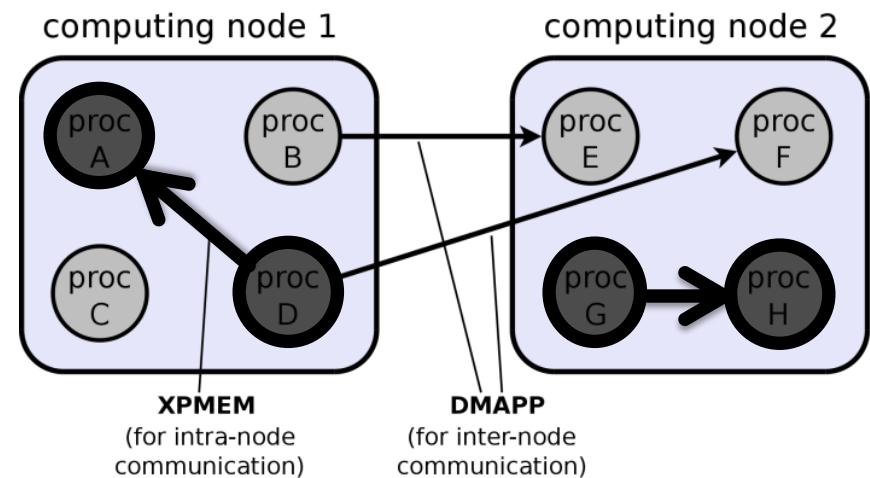
Implementation

We introduce our implementation foMPI (fast one-sided MPI) as fully functional MPI-3.0 RMA

Tweets

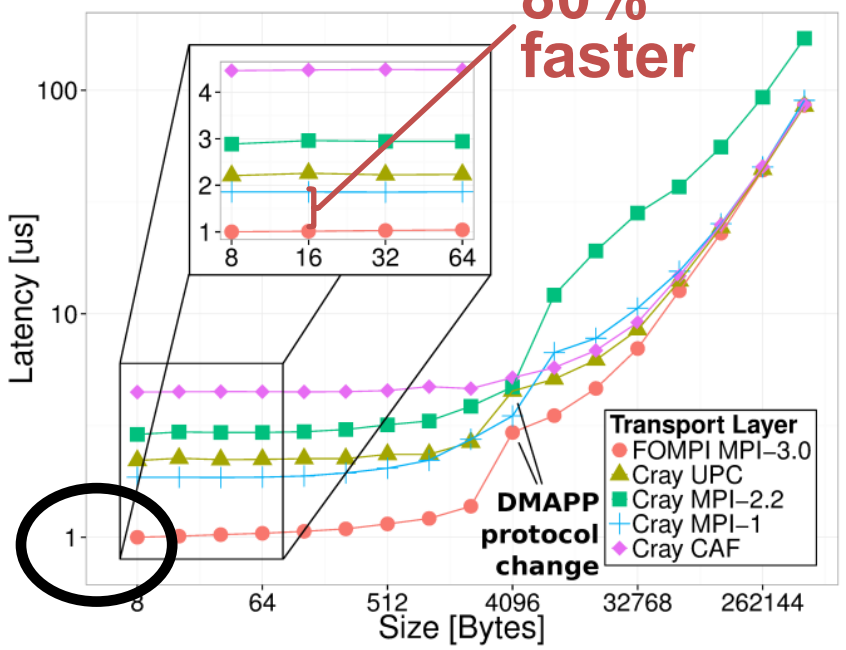
Torsten Hoefler @thoefler 3 Nov
Interested in modern HPC programming? Check out our Advanced MPI tutorial @Supercomputing1 http://inf.ethz.ch/blog/index.php/MPI-2.2, MPI-3 and more!
Retweeted by SPCL@ETH Expand

SPCL@ETH @spcl_eth 27 Oct
Congratulations to Bogdan to his new work on pattern-specific routing on fat trees that will speed up MPI

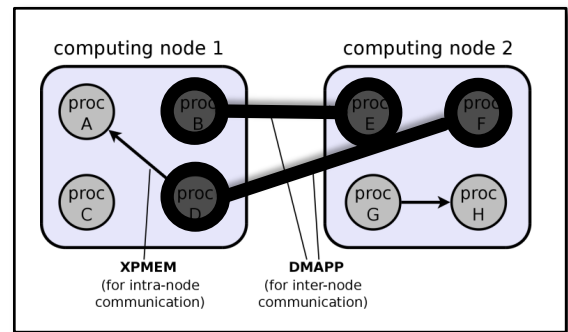
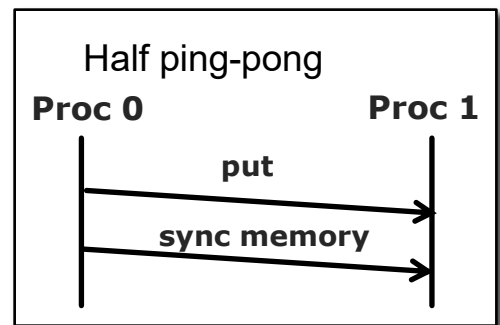
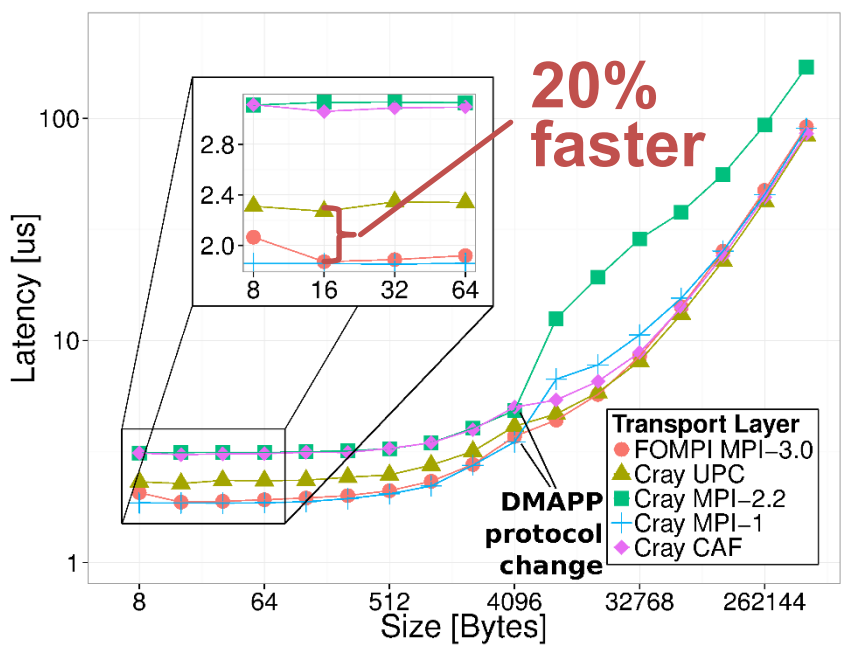


PERFORMANCE INTER-NODE: LATENCY

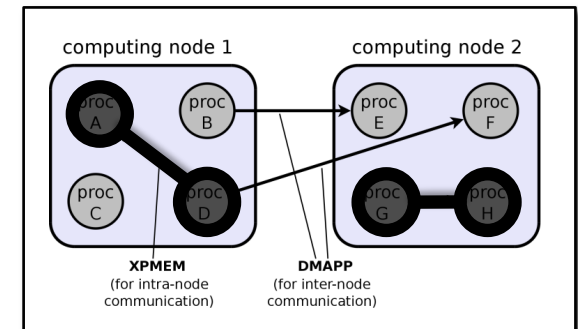
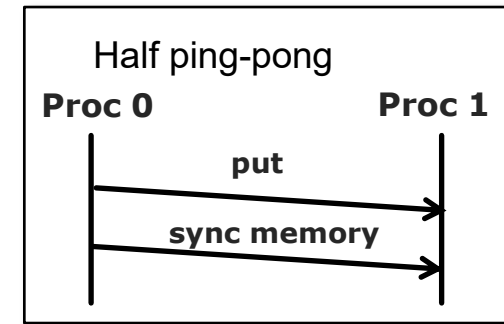
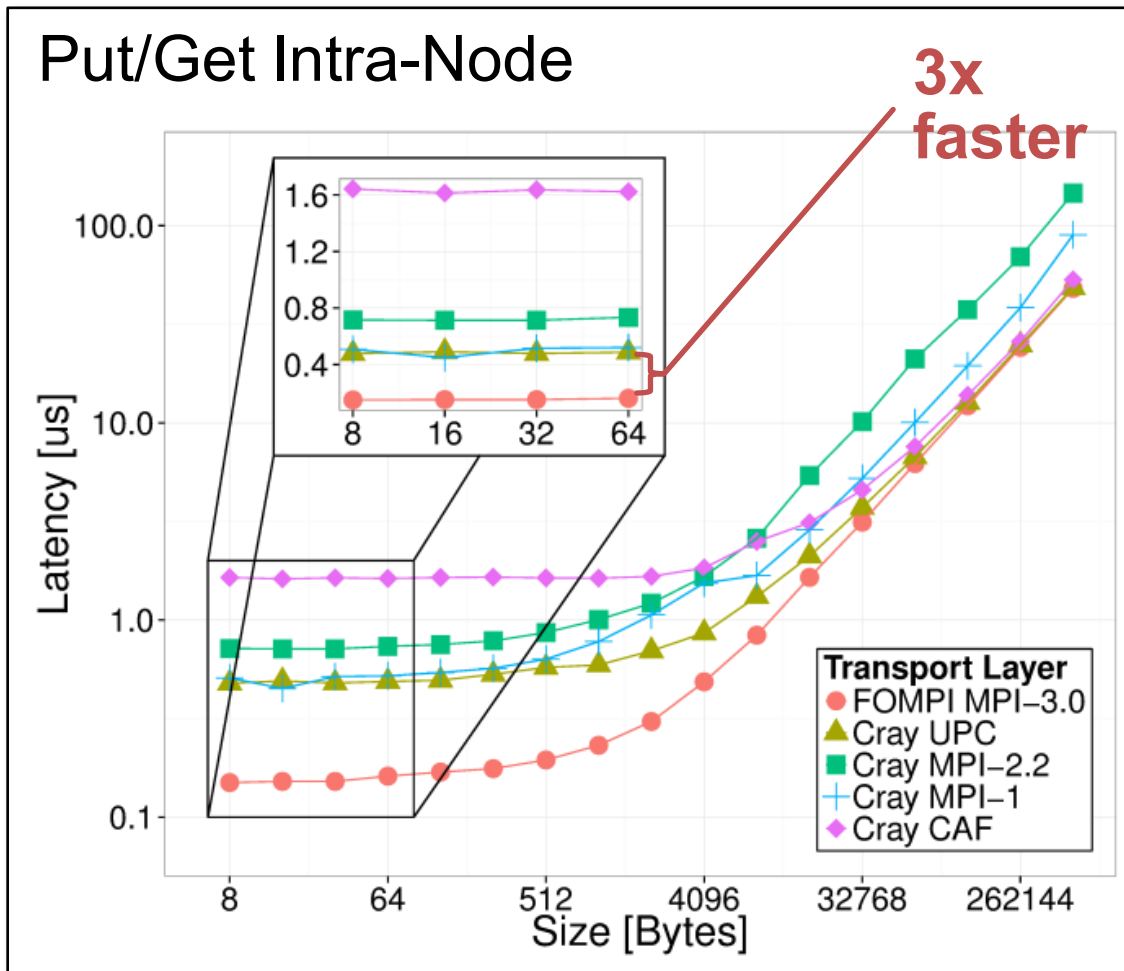
Put Inter-Node



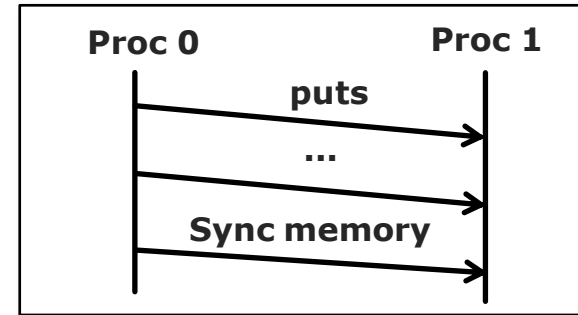
Get Inter-Node



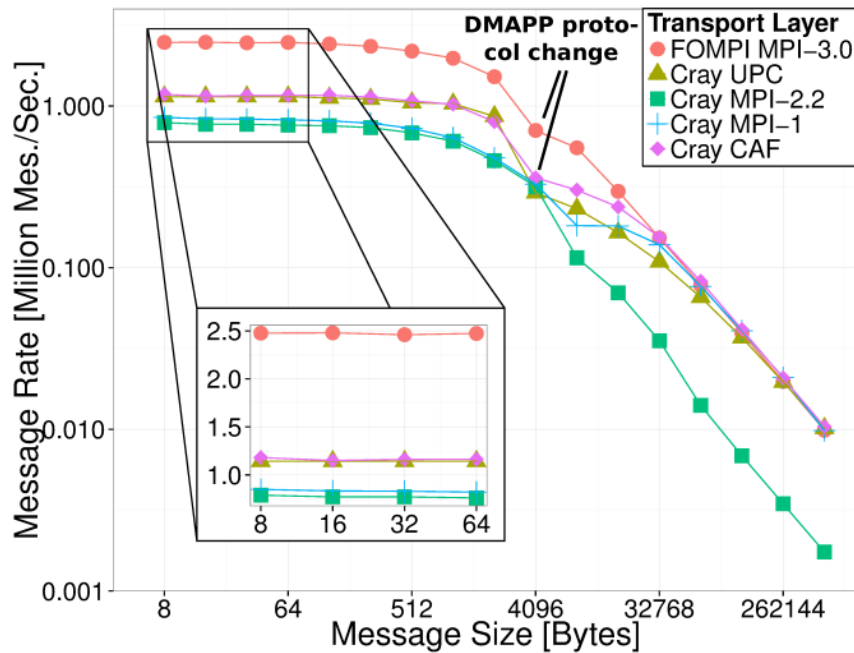
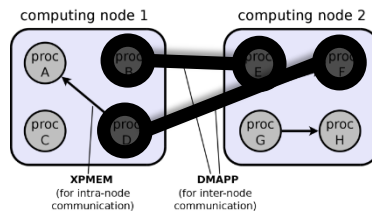
PERFORMANCE INTRA-NODE: LATENCY



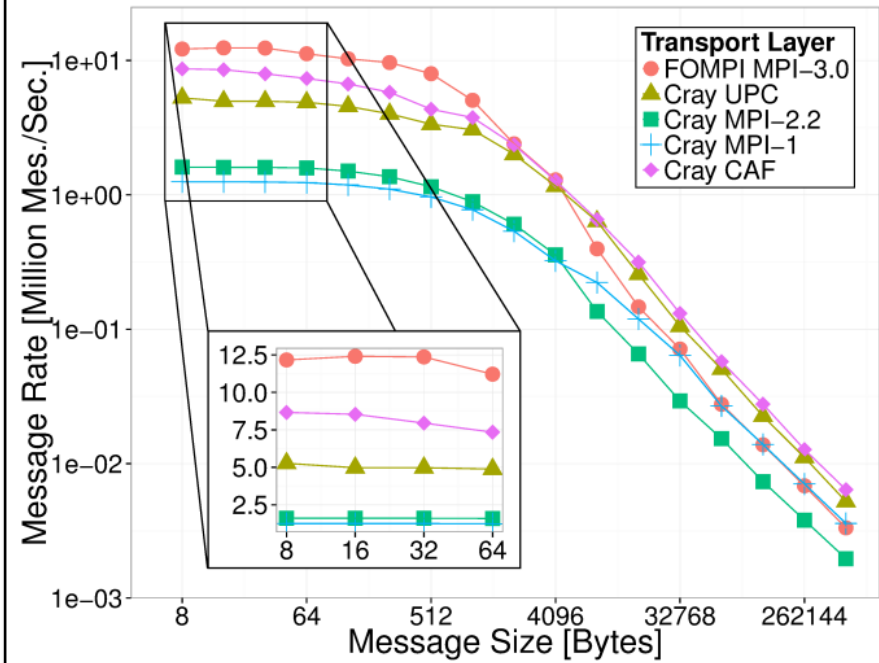
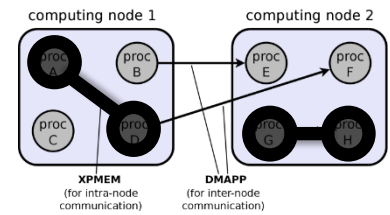
PERFORMANCE: MESSAGE RATE



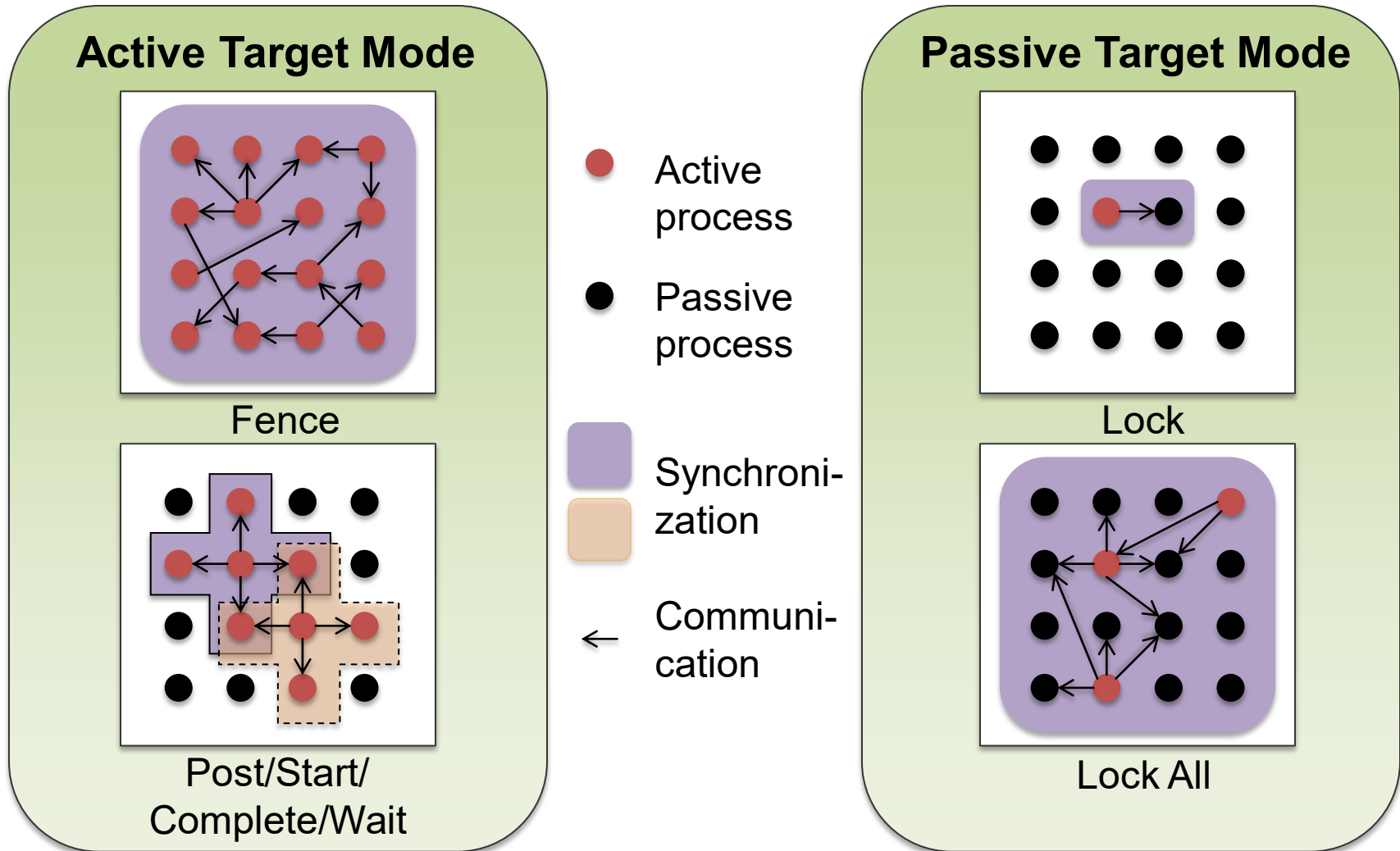
Inter-Node



Intra-Node



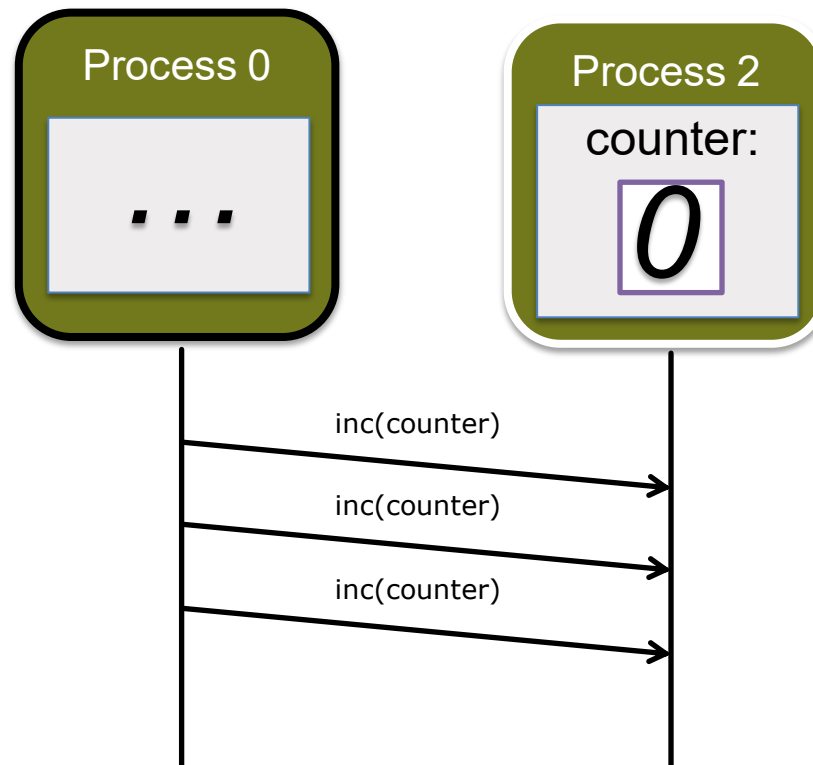
PART 3: SYNCHRONIZATION



FLUSH SYNCHRONIZATION

| | |
|--------------|--------|
| Time bound | $O(1)$ |
| Memory bound | $O(1)$ |

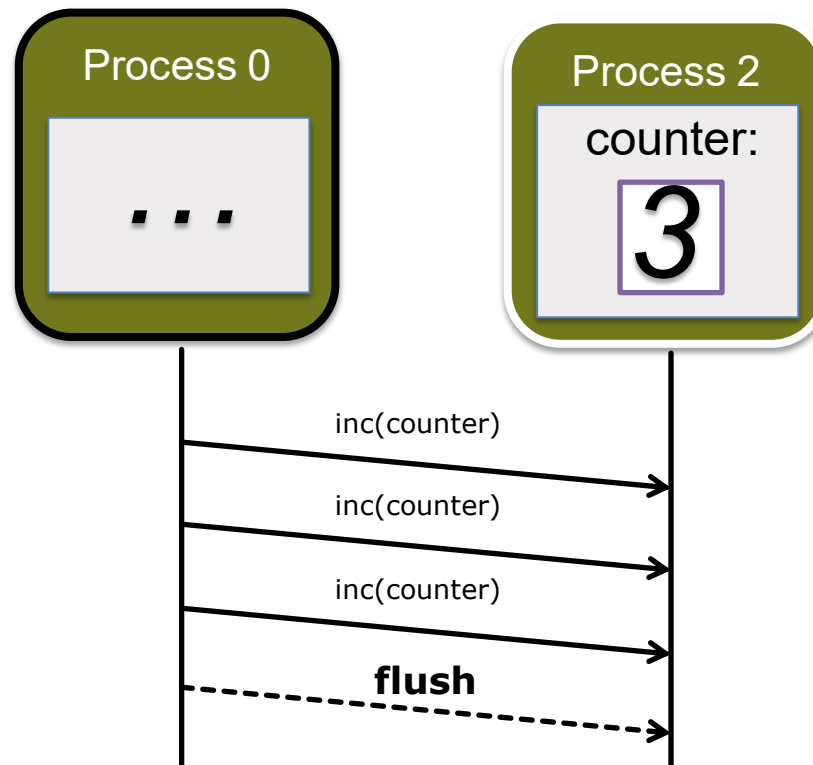
- Guarantees remote completion
- Issues a remote bulk synchronization and an x86 mfence
- One of the most performance critical functions, we add only **78 x86** CPU instructions to the critical path



FLUSH SYNCHRONIZATION

| | |
|--------------|--------|
| Time bound | $O(1)$ |
| Memory bound | $O(1)$ |

- Guarantees remote completion
- Issues a remote bulk synchronization and an x86 mfence
- One of the most performance critical functions, we add only **78 x86** CPU instructions to the critical path



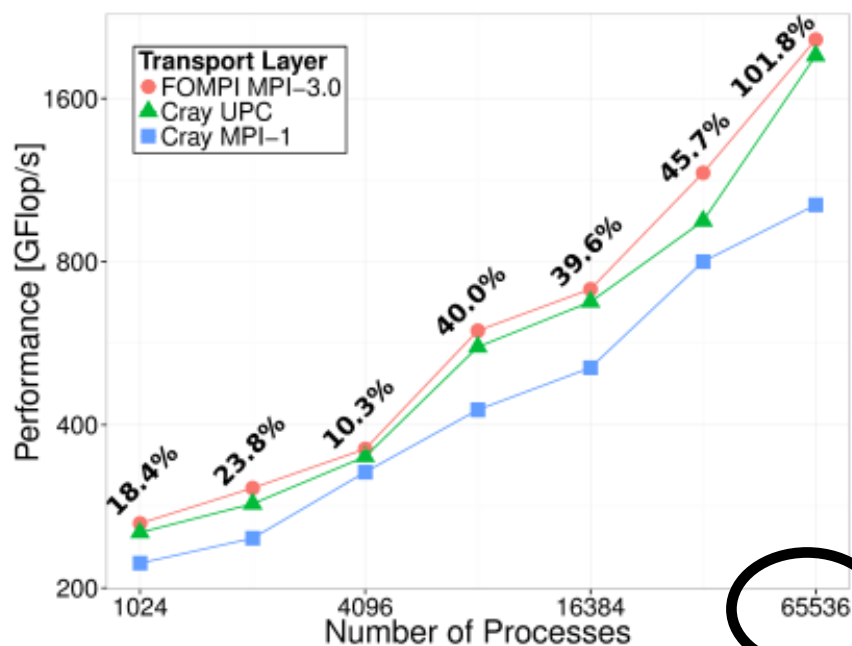
PERFORMANCE

- Evaluation on Blue Waters System
 - 22,640 computing Cray XE6 nodes
 - 724,480 schedulable cores
- All microbenchmarks
- 4 applications
- One nearly full-scale run 😊

PERFORMANCE: APPLICATIONS

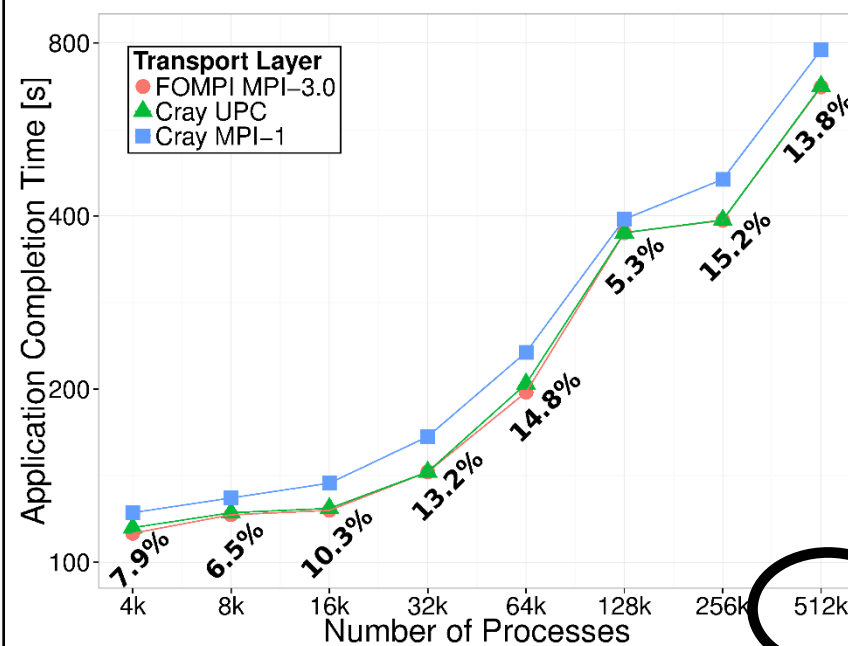
Annotations represent performance gain of foMPI over Cray MPI-1.

NAS 3D FFT [1] Performance



scale
to 65k procs

MILC [2] Application Execution Time



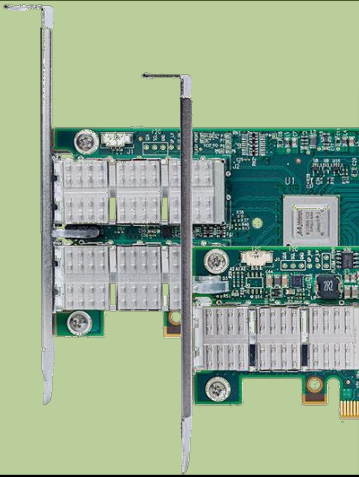
scale
to 512k procs

INTERMEZZO CONCLUSIONS

Requirements for a new low-level programming model

1. Messaging needs 100% hardware support (offload) – it's simple after all!
2. Minimal overheads (**tiny**) layer between user and hardware
3. Offer a simple abstract performance model (e.g., LogGP)

1) Messaging needs hardware support



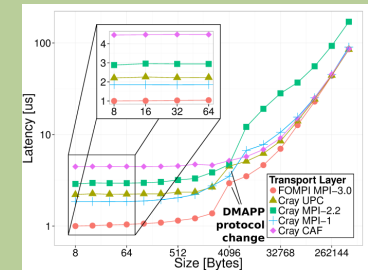
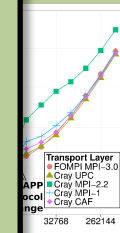
2) Minimal overheads

1. Only 78 x86 CPU instructions

optimal collective communication

3) Simple performance model

| PERFORMANCE MODELLING | |
|---|--|
| Performance functions for synchronization protocols | |
| Fence | $P_{fence} = 2,9\mu s \cdot \log_2(p)$ |
| PSCW | $P_{start} = 0,7\mu s, P_{wait} = 1,8\mu s$ $P_{post} = P_{complete} = 350ns \cdot k$ |
| Locks | $P_{lock,excl} = 5,4\mu s$ $P_{lock,shrd} = P_{lock,all} = 2,7\mu s$ $P_{unlock} = P_{unlock,all} = 0,4\mu s$ $P_{flush} = 76ns$ $P_{sync} = 17ns$ |
| Performance functions for communication protocols | |
| Put/get | $P_{put} = 0,16ns \cdot s + 1\mu s$ $P_{get} = 0,17ns \cdot s + 1,9\mu s$ |
| Atomics | $P_{acc,sum} = 28ns \cdot s + 2,4\mu s$ $P_{acc,min} = 0,8ns \cdot s + 7,3\mu s$ |



WHAT ARE THE CHALLENGES?

- RMA is supported by many HPC libraries and languages



FAULT TOLERANCE + RMA

15.8h of MTBF (for nodes)
for TSUBAME2.0

FAULT TOLERANCE + RMA

- Fault tolerance is well studied for message passing
- Scarce research exists for fault tolerance for RMA
 - Most mechanisms from MP are not applicable
 - Applicable schemes (e.g., simple checkpointing) may be highly inefficient

Message Passing

Coordinated Checkpointing (CC)

Uncoordinated checkpointing and message logging (UC)

RMA

logging memory accesses vs. messages ?

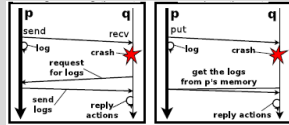
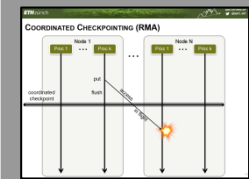
checkpointing in RMA-based applications ?

fault tolerance mechanisms ?

performance ?

OVERVIEW

MP vs. RMA



Uncoord. checkpt. and logging (UC)

Coord. checkpt. (CC)

Generic model

$$D = \langle \mathcal{P}, \mathcal{E}, \mathcal{S}, \dots \rangle, \quad hb \rightarrow \quad co \rightarrow$$

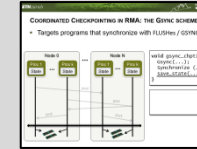
$$a \xrightarrow{cohb} \equiv a \xrightarrow{co} b \wedge a \xrightarrow{hb} b$$

$$po \rightarrow \quad so \rightarrow$$

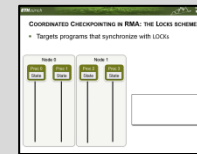
$$PUT(p \rightrightarrows q) \\ GET(p \leftrightharpoons q)$$



CC in RMA

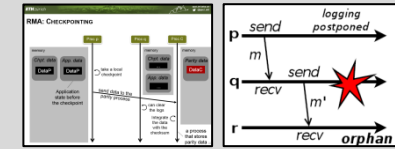


„Gsync” Scheme

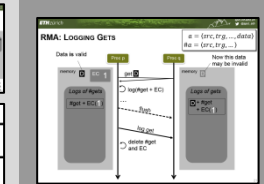
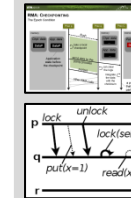


„Locks” Scheme

UC in RMA



Checkpointing



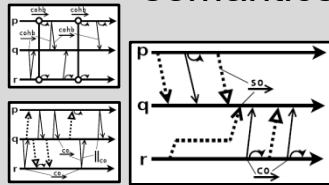
Logging

Recovery in RMA



Basic scheme

Extended RMA semantics



Theory

THEOREM 4.2. *The recovery algorithm 2 preserves the $cohb$ order (referred to as the gsync order).*

Deadlock freedom
Correct recovery

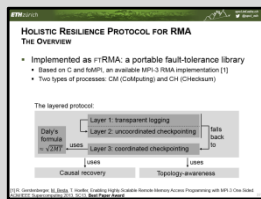
Topology-awareness

$$\langle \mathcal{P}, \mathcal{E}, \mathcal{S}, \mathcal{H}, \mathcal{G}, \xrightarrow{po}, \xrightarrow{so}, \xrightarrow{hb}, \xrightarrow{co}, \mathcal{M} \rangle$$

Model extensions

Decreasing failure prob.

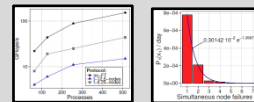
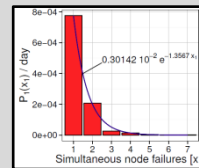
Holistic fault-tolerance library



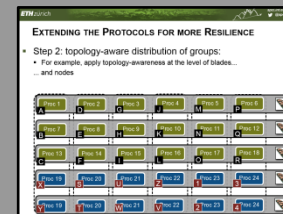
Checkpoints on demand

Optimum CC intervals

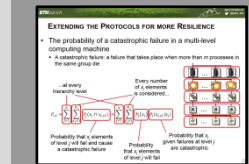
Design and optimizations



Performance



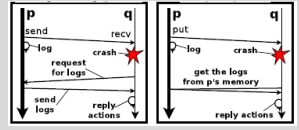
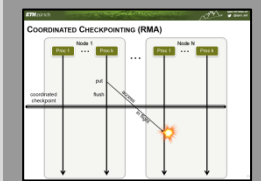
Distribution of processes



$$\frac{D_j \cdot \binom{|G|}{2} \cdot \binom{H_j-2}{x_j-2}}{\binom{H_j}{x_j}}$$

OVERVIEW

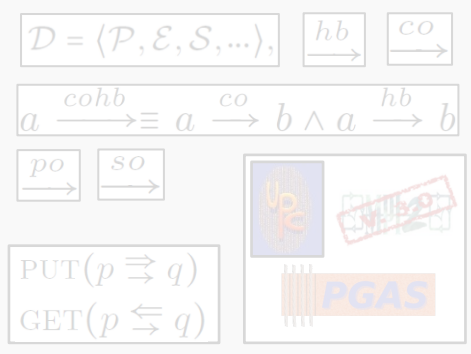
MP vs. RMA



Coord.
checkpt. (CC)

Uncoord.
checkpt. and
logging (UC)

Generic model



CC in RMA

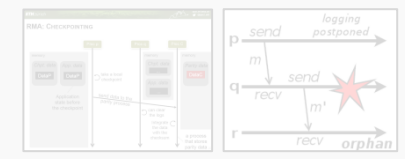


„Gsync”
Scheme



„Locks”
Scheme

UC in RMA



Checkpointing



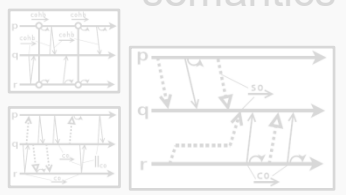
Logging

Recovery in RMA



Basic scheme

Extended RMA
semantics



Theory

THEOREM 4.2. *The recovery rithm 2 preserves the $cohb$ or referred to as the gsync order).*

Deadlock freedom
Correct recovery

Topology-awareness

$$\langle \mathcal{P}, \mathcal{E}, \mathcal{S}, \mathcal{H}, \mathcal{G}, \xrightarrow{po}, \xrightarrow{so}, \xrightarrow{hb}, \xrightarrow{co}, \mathcal{M} \rangle$$

Model extensions

Decreasing
failure prob.

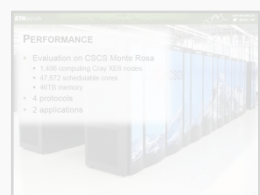
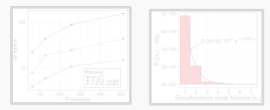
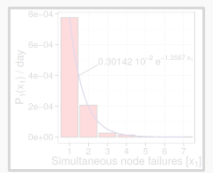
Holistic fault-tolerance library



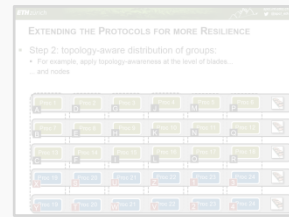
Design and
optimizations

Checkpoints
on demand

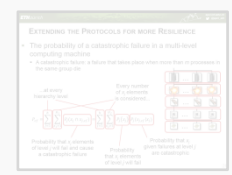
Optimum
CC intervals



Performance



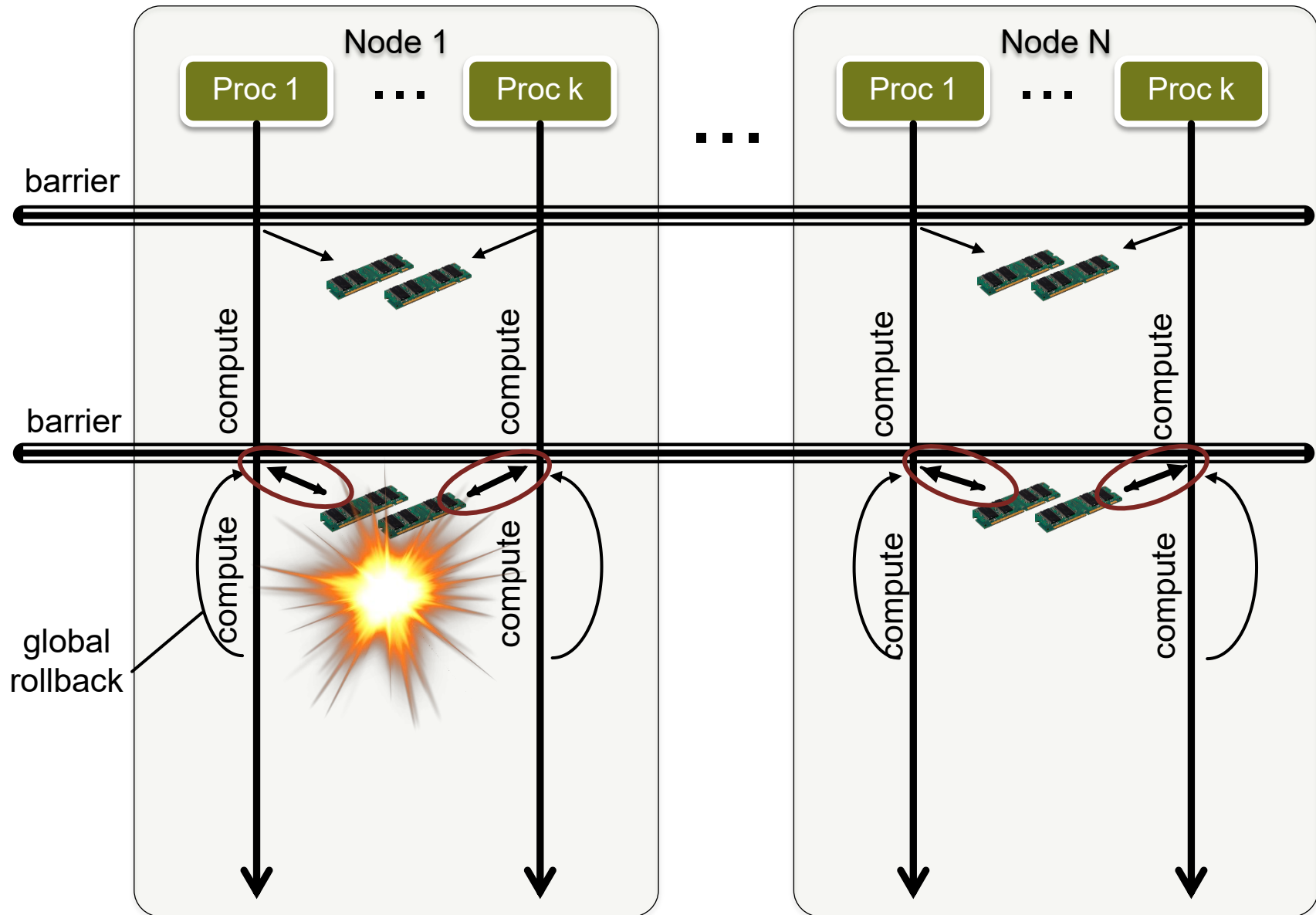
Distribution of
processes



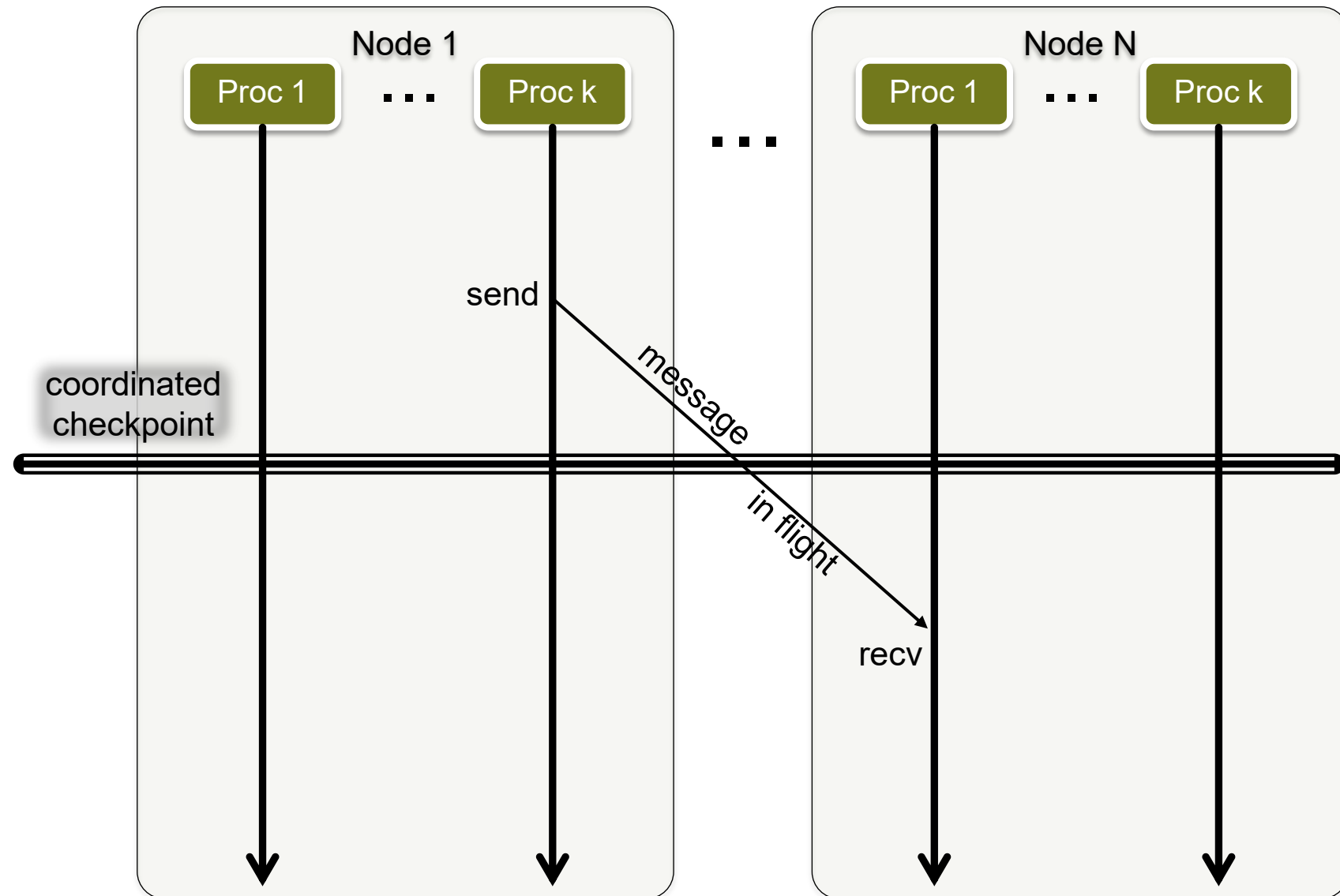
$$\frac{D_j \cdot \binom{G_j}{2} \cdot \binom{H_j-2}{x_j-2}}{\binom{H_j}{x_j}}$$



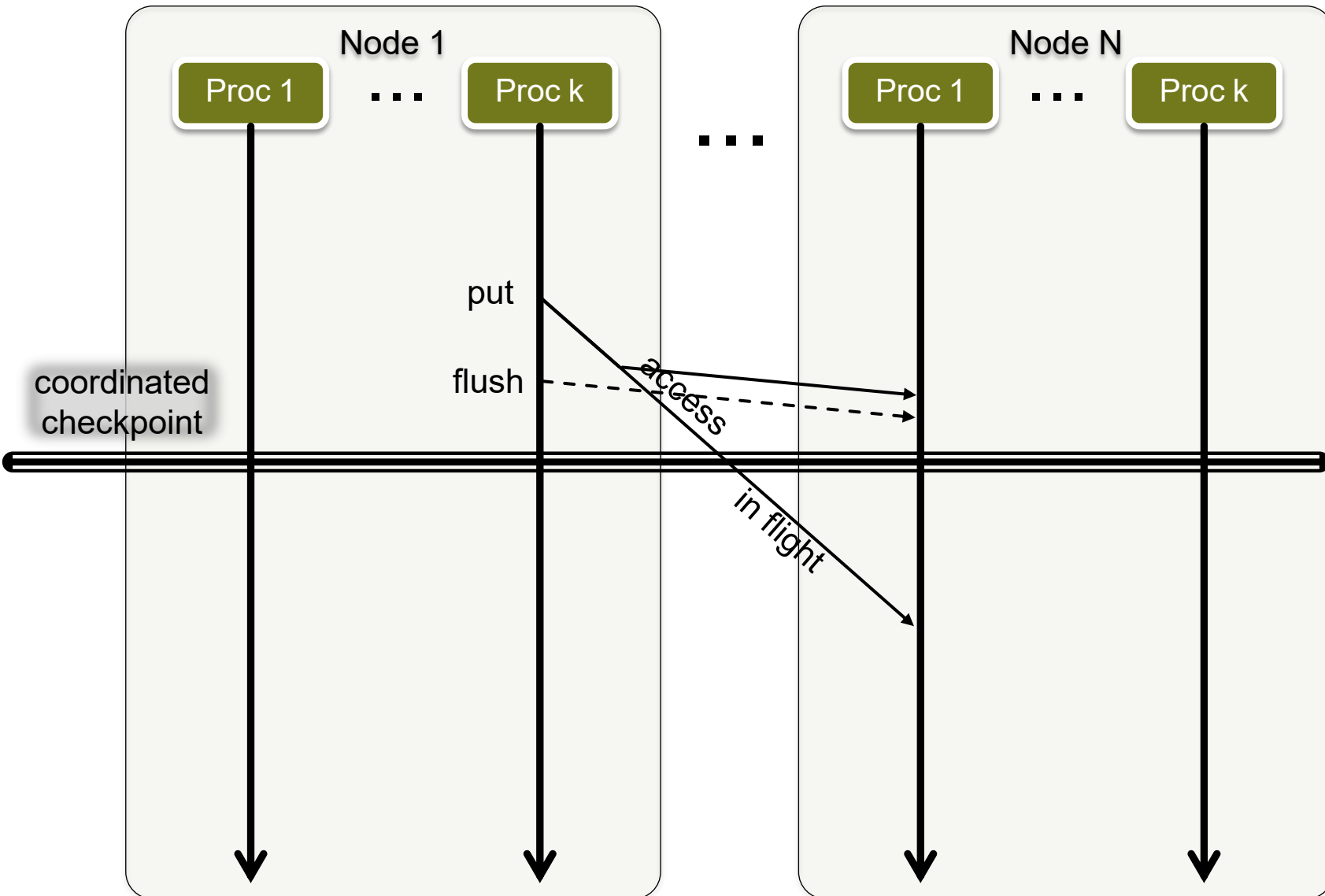
COORDINATED CHECKPOINTING (MP)



COORDINATED CHECKPOINTING (MP)

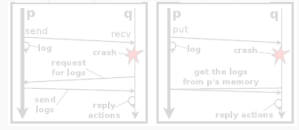
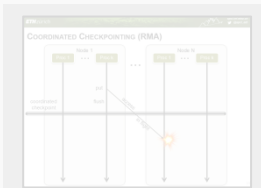


COORDINATED CHECKPOINTING (RMA)



CONTRIBUTIONS

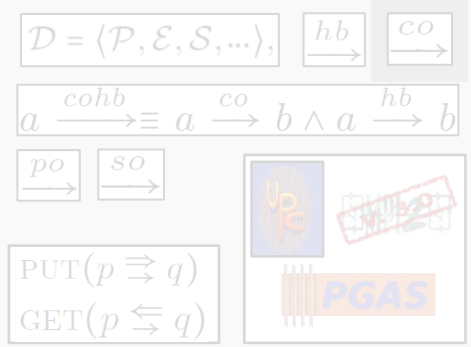
MP vs. RMA



Uncoord.
checkpt. and
logging (UC)

Coord.
checkpt. (CC)

Generic model



CC in RMA

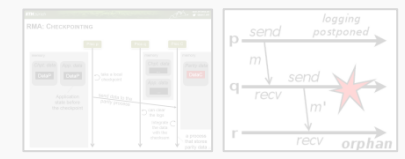


„Gsync”
Scheme

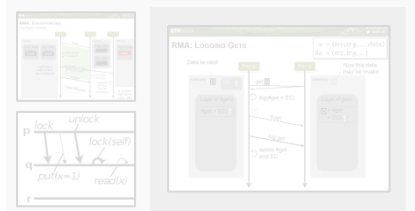


„Locks”
Scheme

UC in RMA

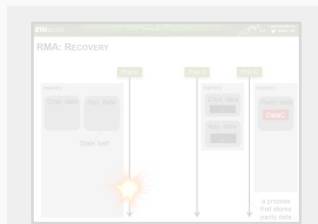


Checkpointing



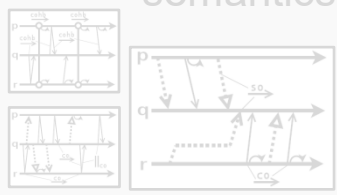
Logging

Recovery in RMA



Basic scheme

Extended RMA
semantics



Theory

THEOREM 4.2. *The recovery rithm 2 preserves the $cohb$ order (referred to as the gsync order).*

Deadlock freedom
Correct recovery

Topology-awareness

$$\langle \mathcal{P}, \mathcal{E}, \mathcal{S}, \mathcal{H}, \mathcal{G}, \xrightarrow{po}, \xrightarrow{so}, \xrightarrow{hb}, \xrightarrow{co}, \mathcal{M} \rangle$$

Model extensions

Decreasing
failure prob.

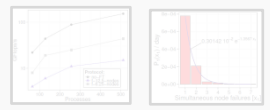
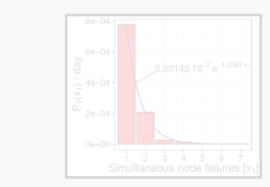
Holistic fault-tolerance library



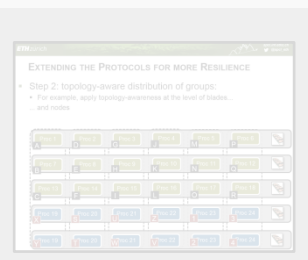
Design and
optimizations

Checkpoints
on demand

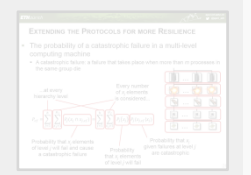
Optimum
CC intervals



Performance

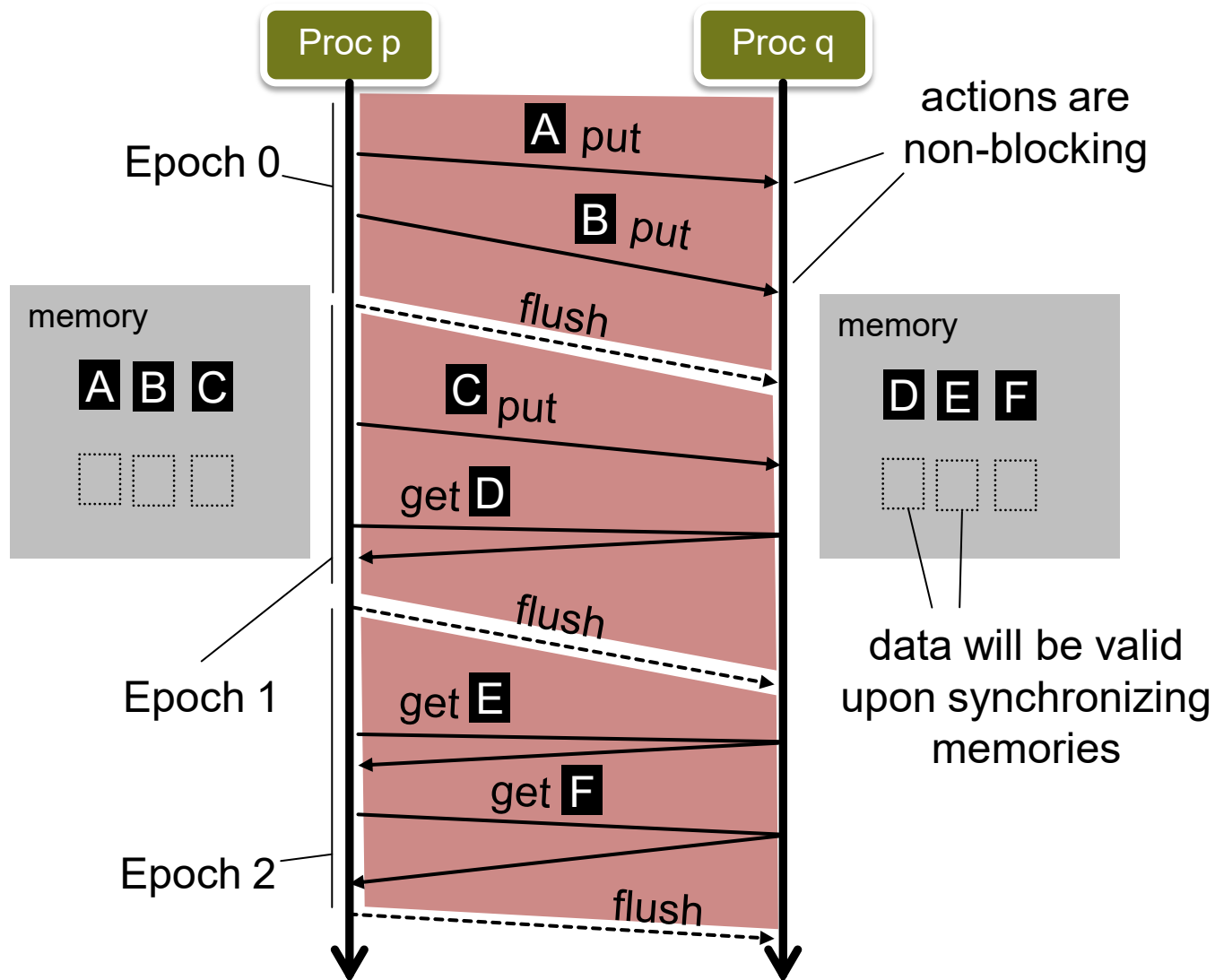


Distribution of
processes

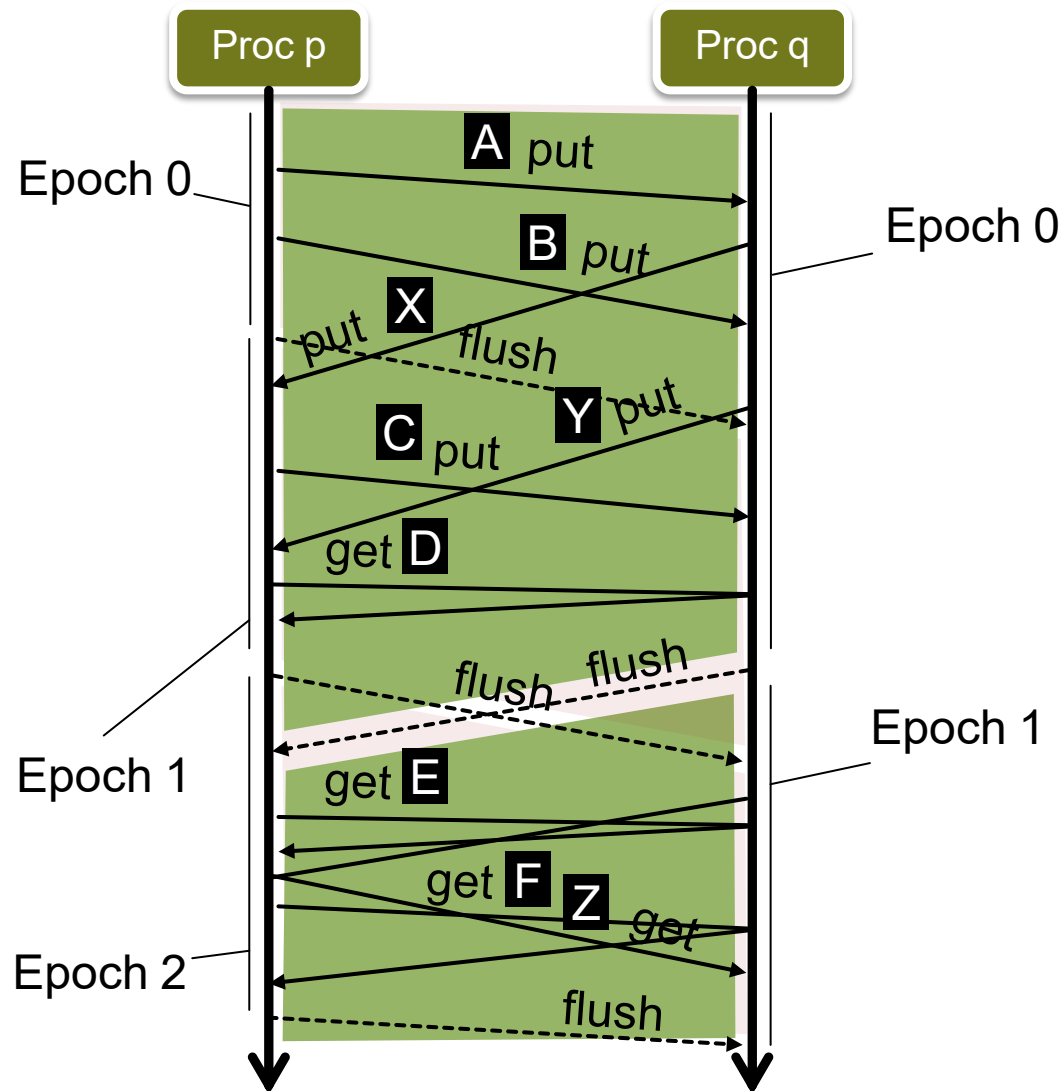


$$\frac{D_j \cdot \binom{|G|}{2} \cdot \binom{H_j - 2}{x_j - 2}}{\binom{H_j}{x_j}}$$

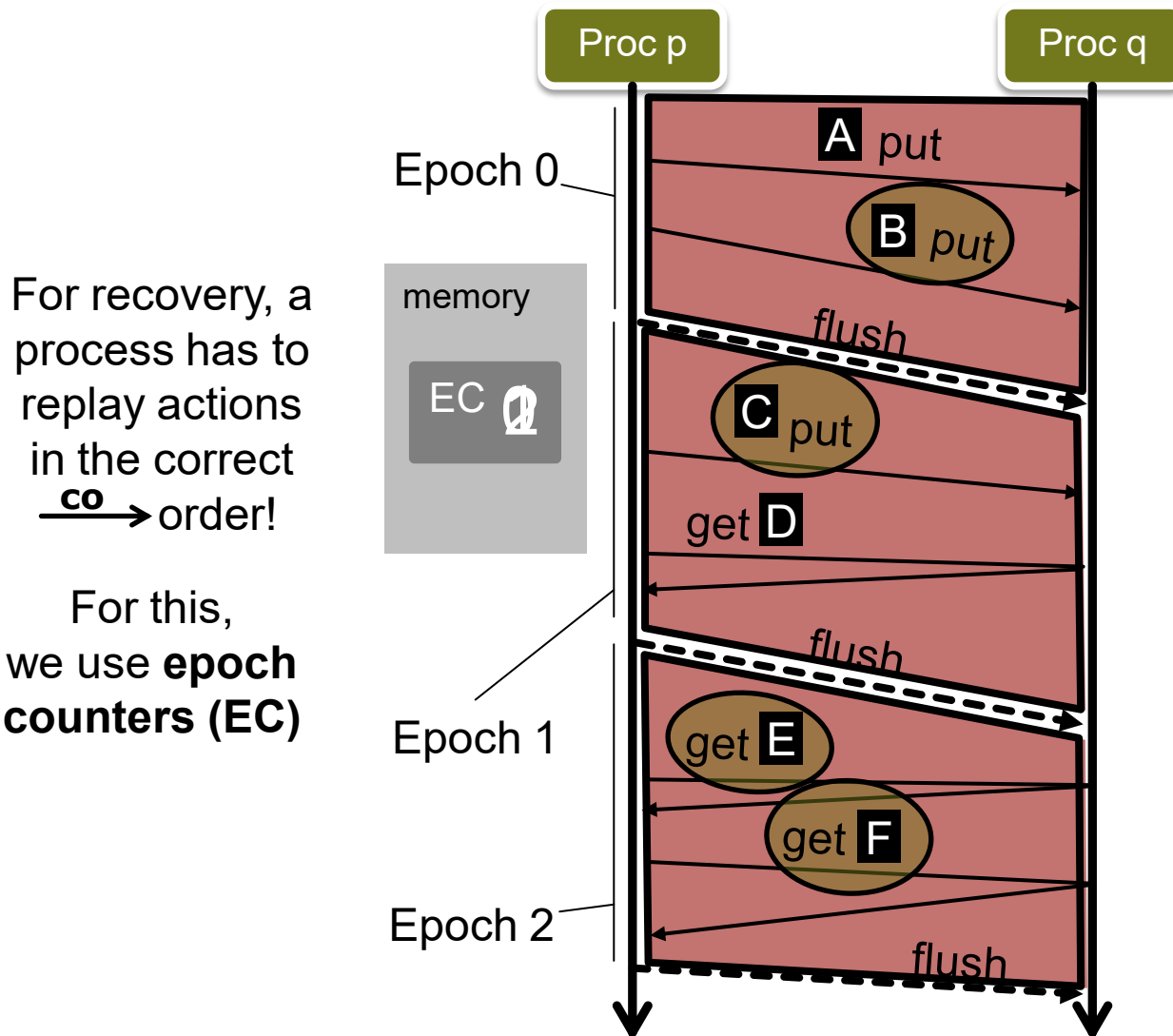
RMA: EPOCHS



RMA: EPOCHS



RMA: THE CONSISTENCY ORDER $\xrightarrow{\text{co}}$

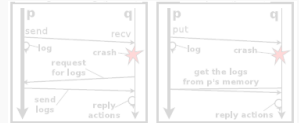
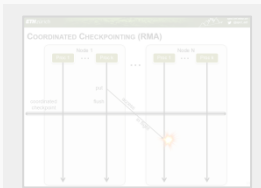


B put $\xrightarrow{\text{co}}$ C put

get E \parallel_{co} get F

CONTRIBUTIONS

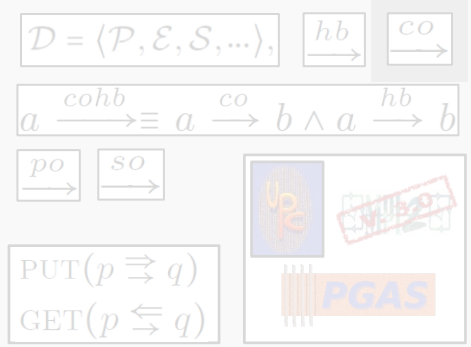
MP vs. RMA



Uncoord.
checkpt. and
logging (UC)

Coord.
checkpt. (CC)

Generic model



CC in RMA

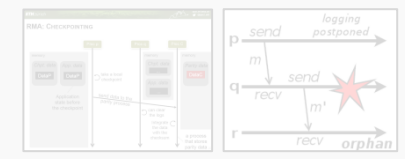


„Gsync”
Scheme



„Locks”
Scheme

UC in RMA

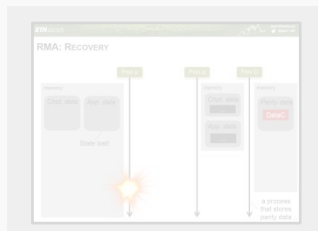


Checkpointing



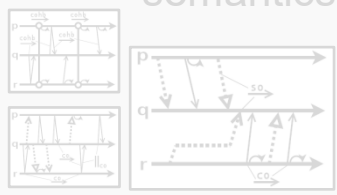
Logging

Recovery in RMA



Basic scheme

Extended RMA
semantics



Theory

THEOREM 4.2. *The recovery rithm 2 preserves the $cohb$ or rferred to as the gsync order).*

Deadlock freedom
Correct recovery

Topology-awareness

$$\langle \mathcal{P}, \mathcal{E}, \mathcal{S}, \mathcal{H}, \mathcal{G}, \xrightarrow{po}, \xrightarrow{so}, \xrightarrow{hb}, \xrightarrow{co}, \mathcal{M} \rangle$$

Model extensions

Decreasing
failure prob.

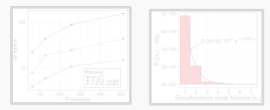
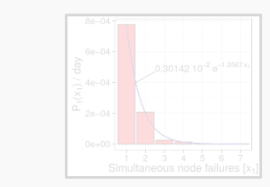
Holistic fault-tolerance library



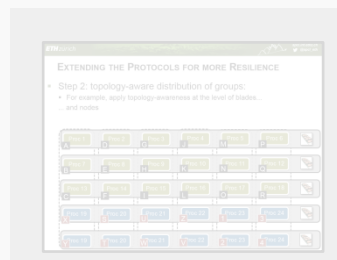
Design and
optimizations

Checkpoints
on demand

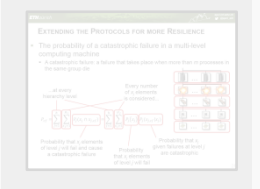
Optimum
CC intervals



Performance

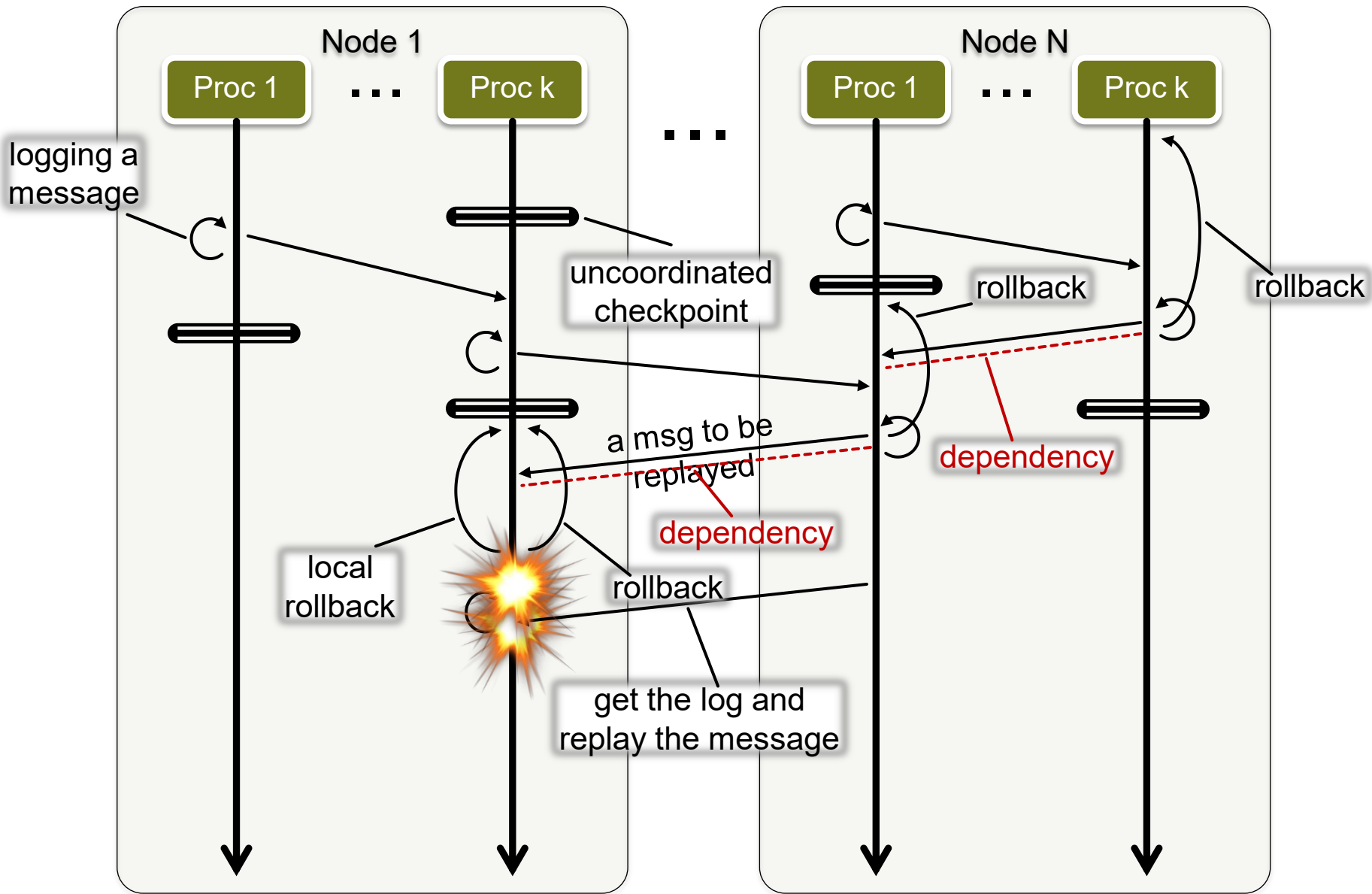


Distribution of
processes



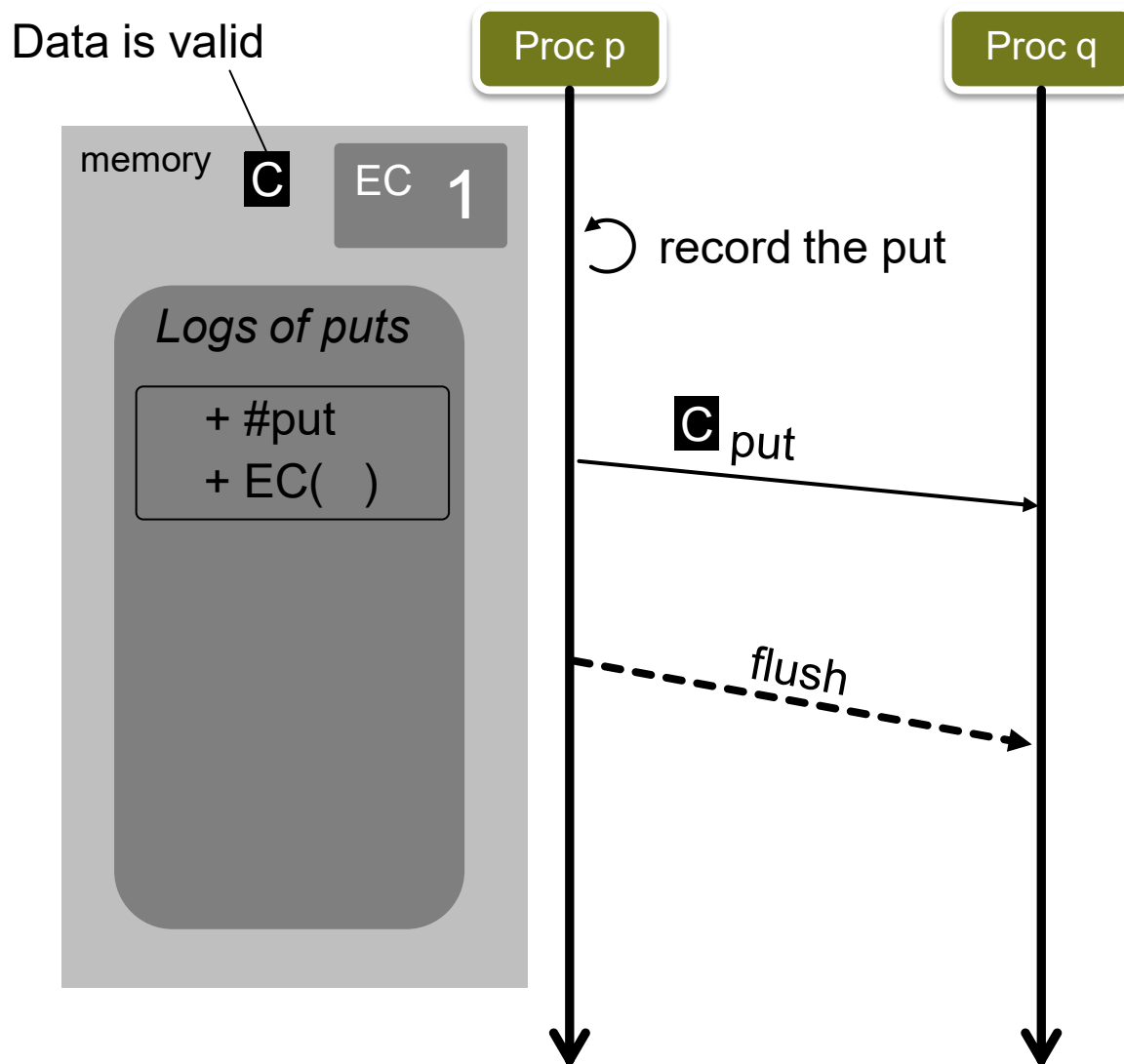
$$\frac{D_j \cdot \binom{|G|}{2} \cdot \binom{H_j - 2}{x_j - 2}}{\binom{H_j}{x_j}}$$

UNCOORDINATED CHECKPOINTING (MP)



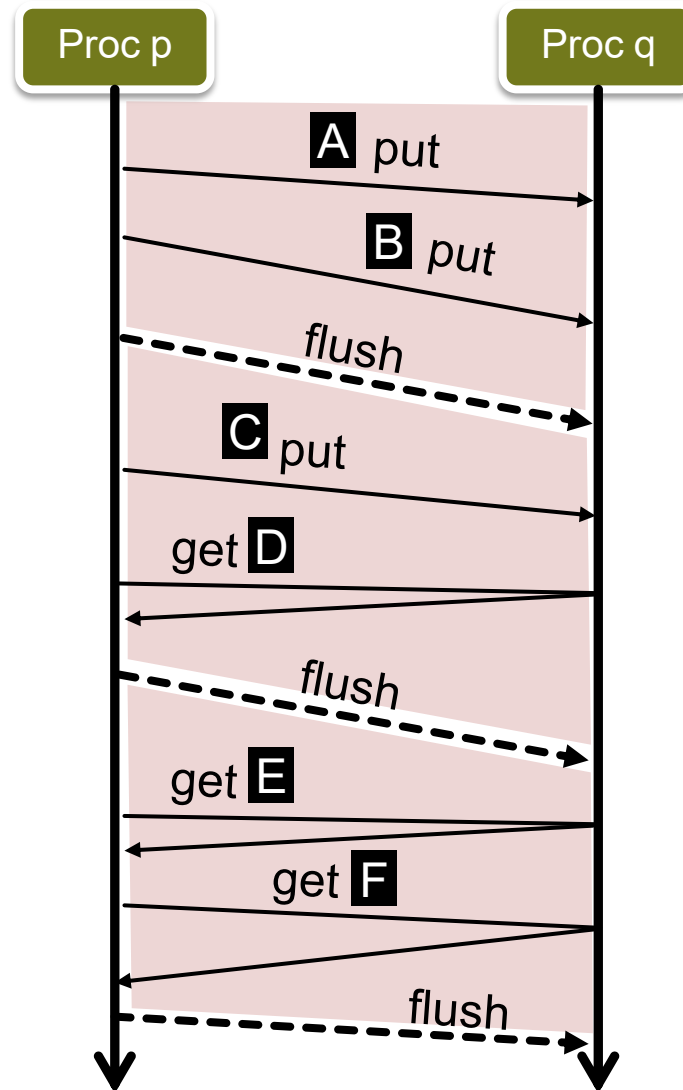
RMA: LOGGING PUTS

$$a = \langle src, trg, \dots, data \rangle$$

$$\#a = \langle src, trg, \dots \rangle$$


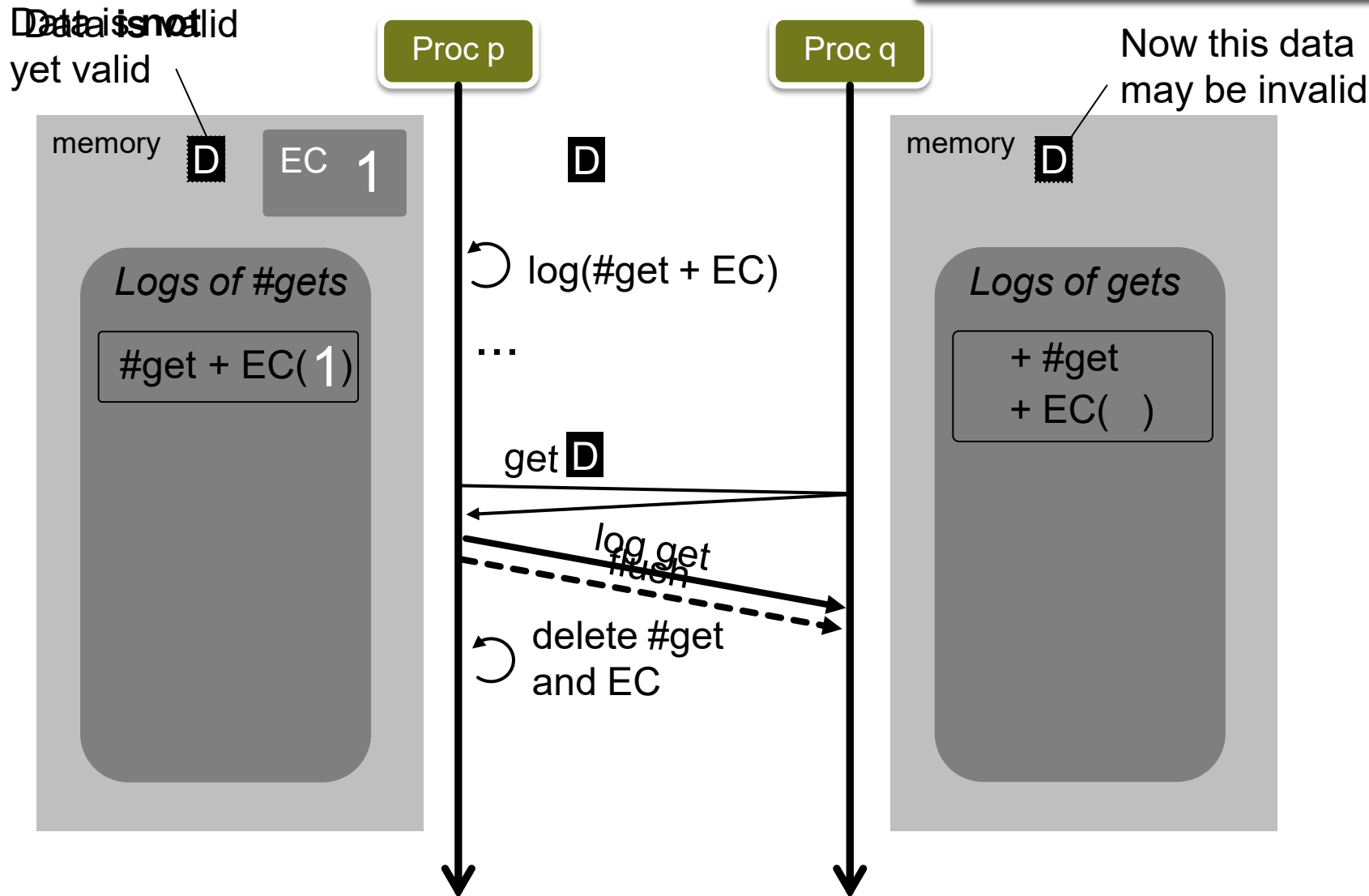
RMA: LOGGING GETS

$$a = \langle src, trg, \dots, data \rangle$$

$$\#a = \langle src, trg, \dots \rangle$$


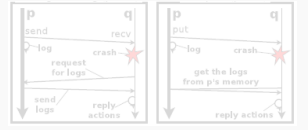
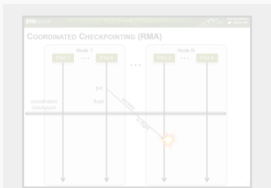
RMA: LOGGING GETS

$a = \langle src, trg, \dots, data \rangle$
 $\#a = \langle src, trg, \dots \rangle$



CONTRIBUTIONS

MP vs. RMA



Coord. checkpt. (CC)

Uncoord. checkpt. and logging (UC)

Generic model

$$\mathcal{D} = \langle \mathcal{P}, \mathcal{E}, \mathcal{S}, \dots \rangle, \quad hb \rightarrow \quad co \rightarrow$$

$$a \xrightarrow{cohb} \equiv a \xrightarrow{co} b \wedge a \xrightarrow{hb} b$$

$$po \rightarrow \quad so \rightarrow$$

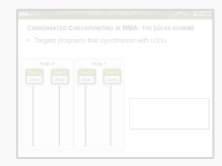
PUT($p \Rightarrow q$)
GET($p \Leftarrow q$)



CC in RMA

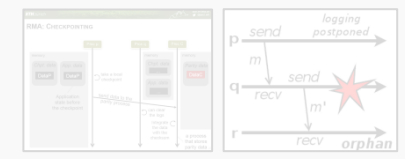


„Gsync” Scheme

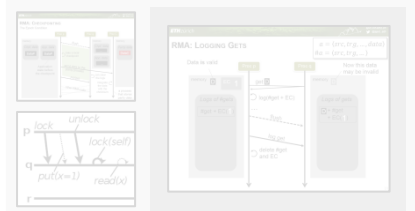


„Locks” Scheme

UC in RMA

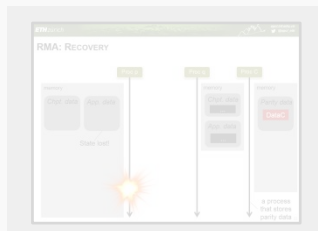


Checkpointing



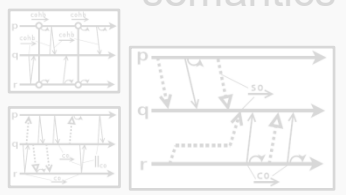
Logging

Recovery in RMA



Basic scheme

Extended RMA semantics



Theory

THEOREM 4.2. The recovery rithm 2 preserves the $cohb \rightarrow$ or ferred to as the gsync order).

Deadlock freedom
Correct recovery

Topology-awareness

$$\langle \mathcal{P}, \mathcal{E}, \mathcal{S}, \mathcal{H}, \mathcal{G}, \xrightarrow{po}, \xrightarrow{so}, \xrightarrow{hb}, \xrightarrow{co}, \mathcal{M} \rangle$$

Model extensions

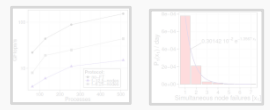
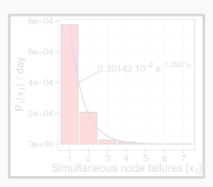
Holistic fault-tolerance library



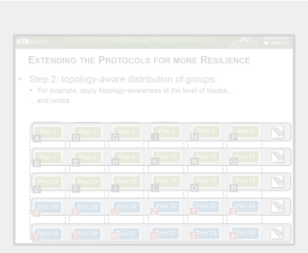
Checkpoints on demand

Optimum CC intervals

Design and optimizations

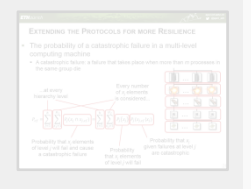


Performance



Distribution of processes

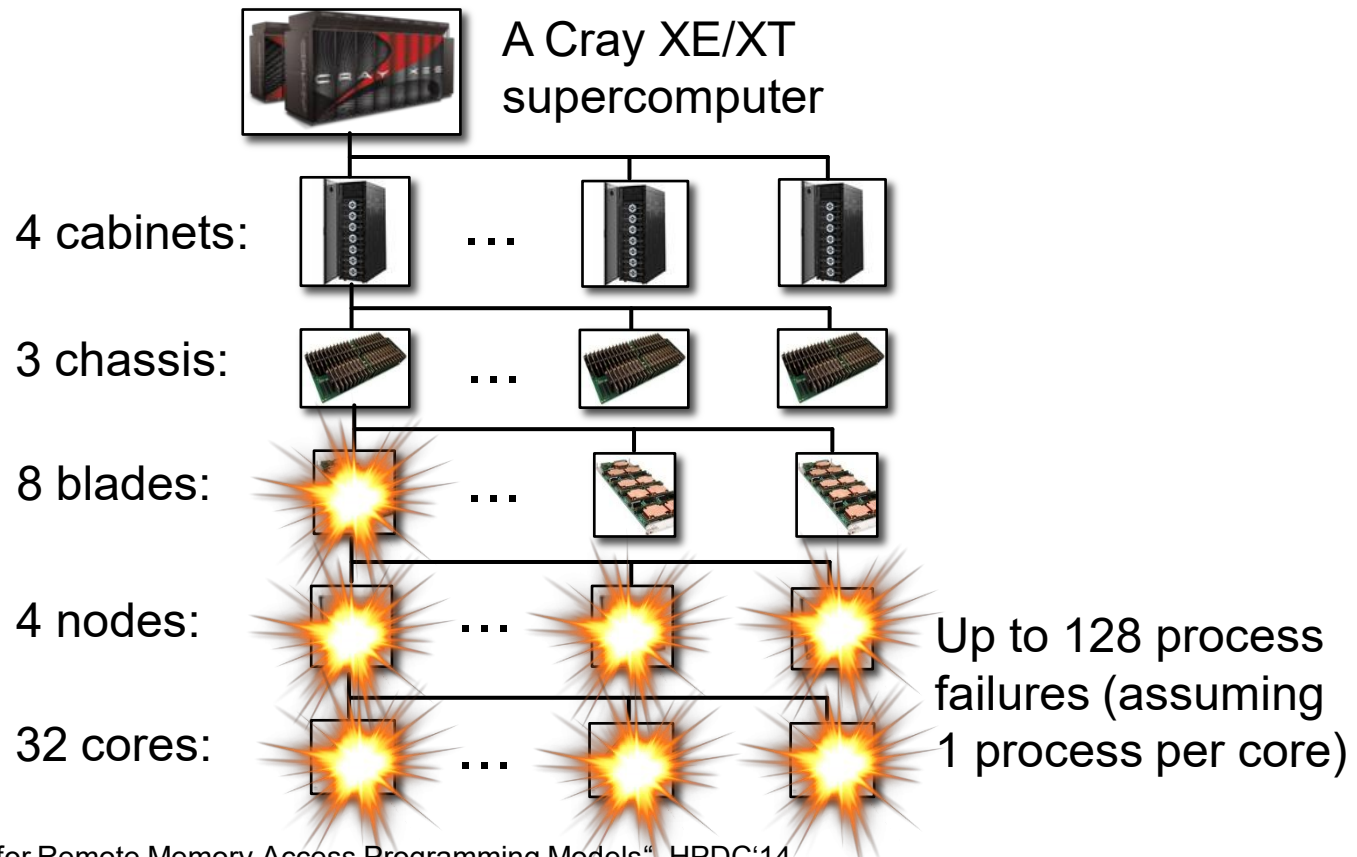
Decreasing failure prob.



$$\frac{D_j \cdot \binom{|G|}{2} \cdot \binom{H_j-2}{x_j-2}}{\binom{H_j}{x_j}}$$

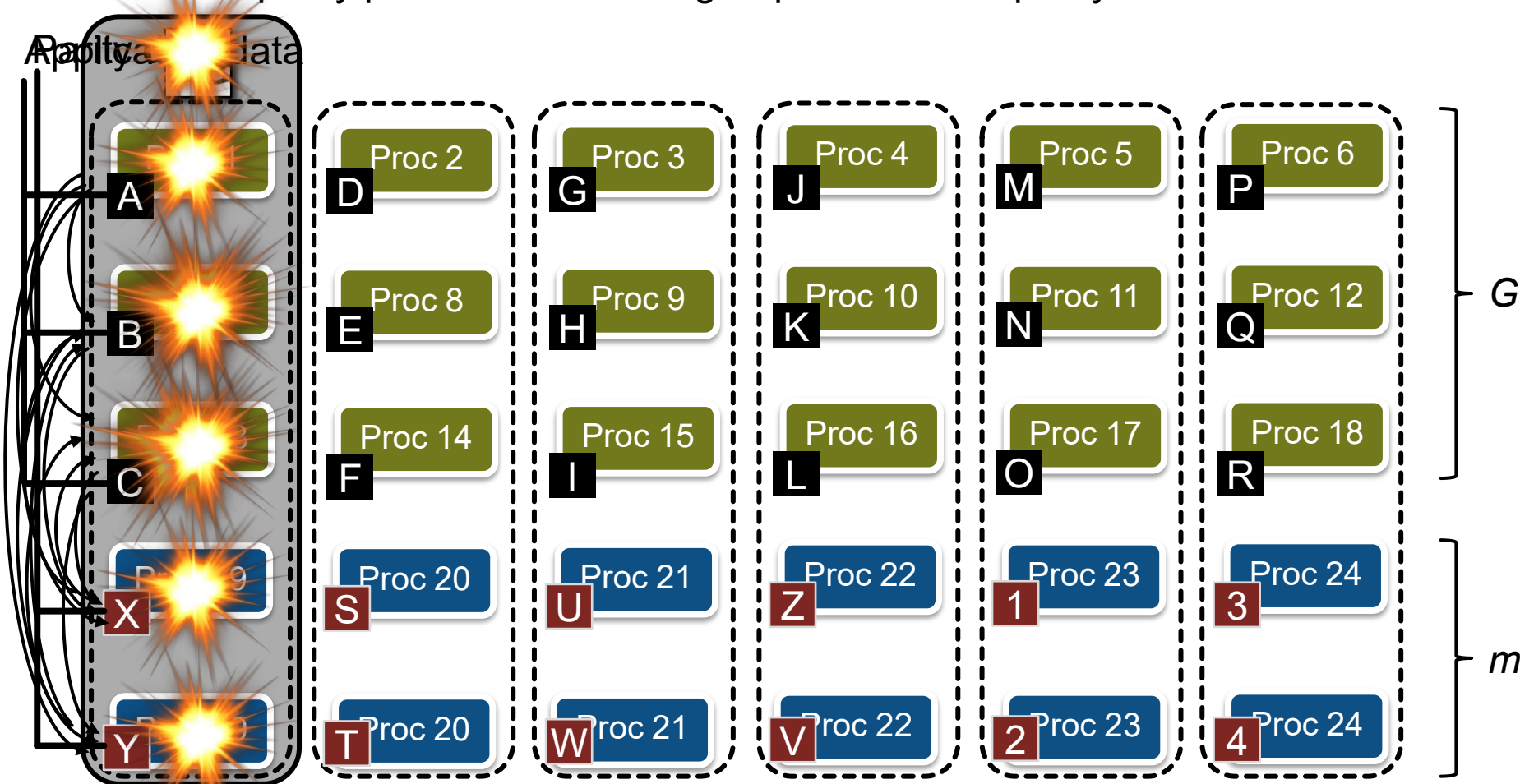
TOPOLOGY-AWARENESS FOR INCREASING RESILIENCE

- Today's supercomputers have a hierarchical layout
- A single hardware crash may kill multiple processes...
- ...protocols will not work and the whole computation is lost ☹️



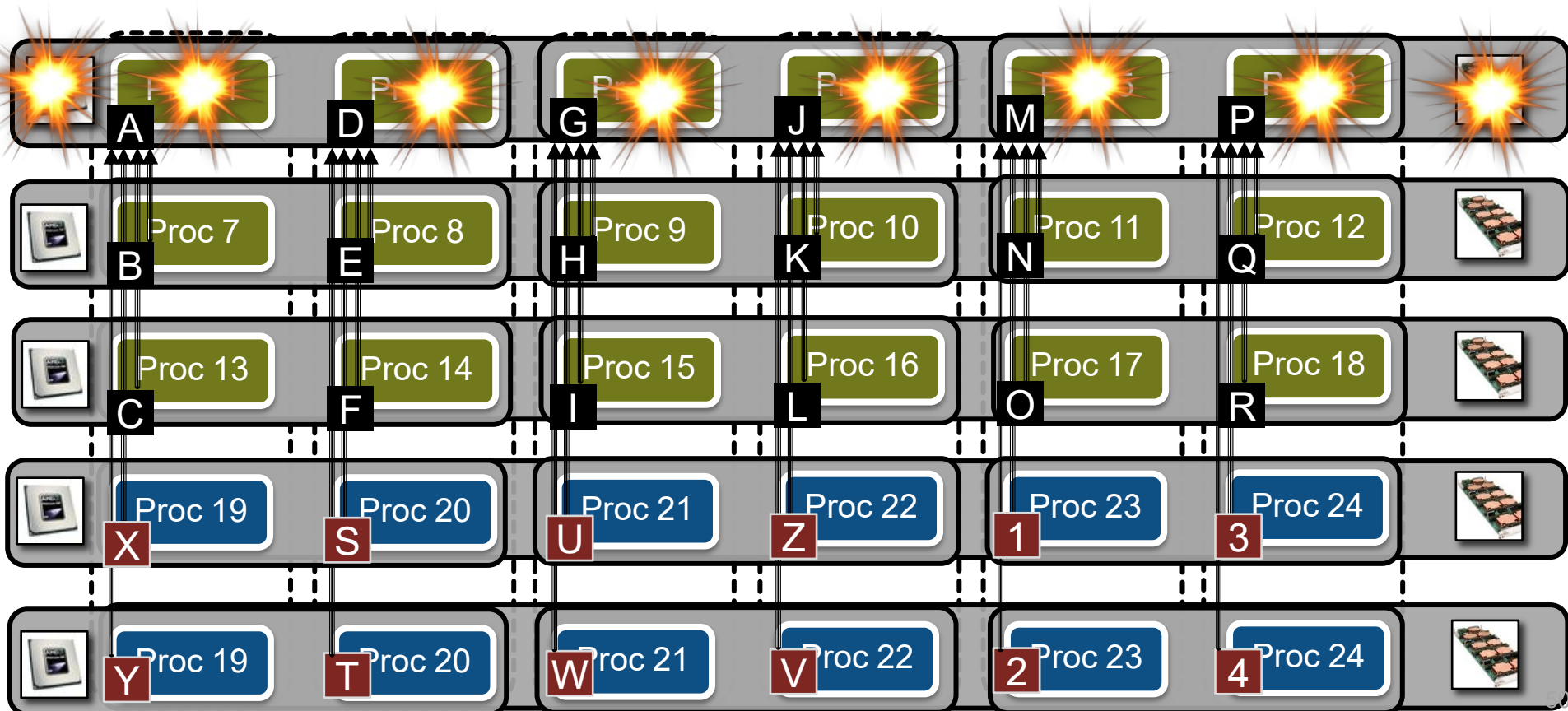
EXTENDING THE PROTOCOLS FOR MORE RESILIENCE

- Step 1: groups of processes:
 - Divide processes into groups of size G each
 - Add m parity processes to each group to store the parity data



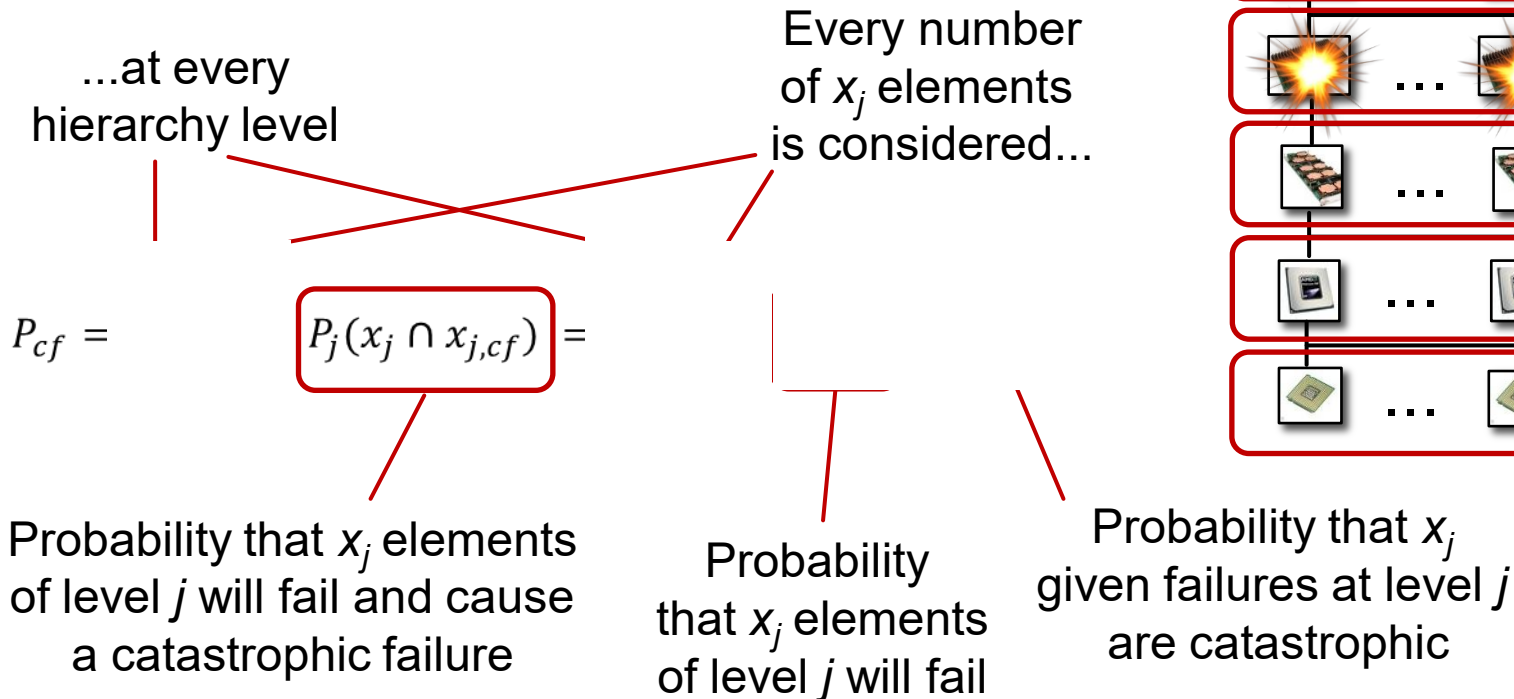
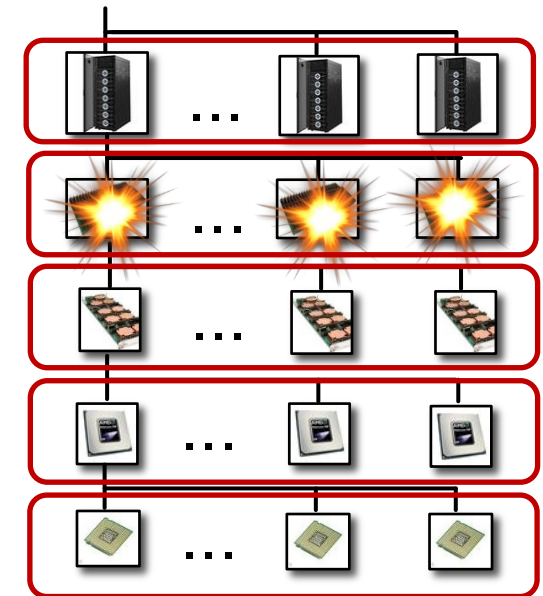
TOPOLOGY-AWARENESS FOR INCREASING RESILIENCE

- Step 2: topology-aware distribution of groups:
 - For example, apply topology-awareness at the level of blades...
... and nodes



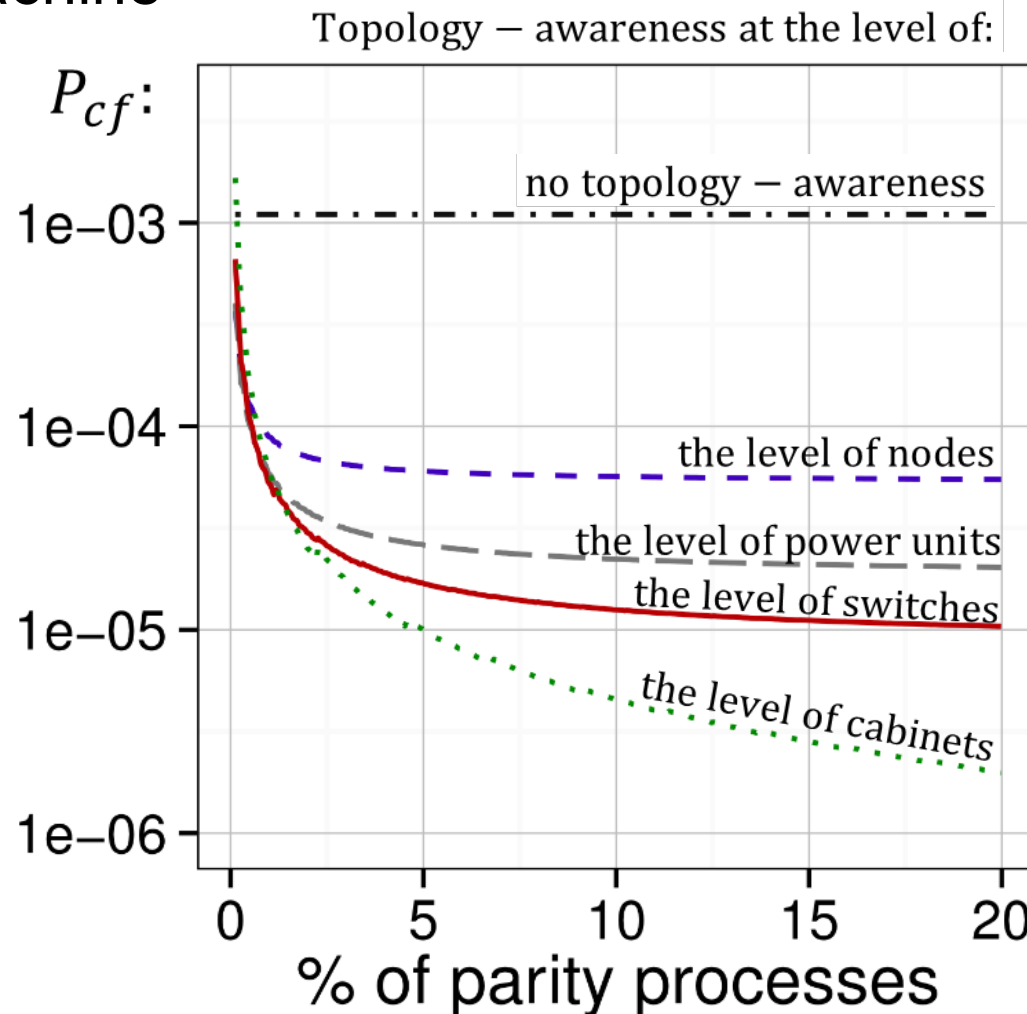
TOPOLOGY-AWARENESS FOR INCREASING RESILIENCE

- The probability of a catastrophic failure in a multi-level computing machine (generalizing [1])
 - A catastrophic failure: a failure that takes place if $> m$ processes in the same group die



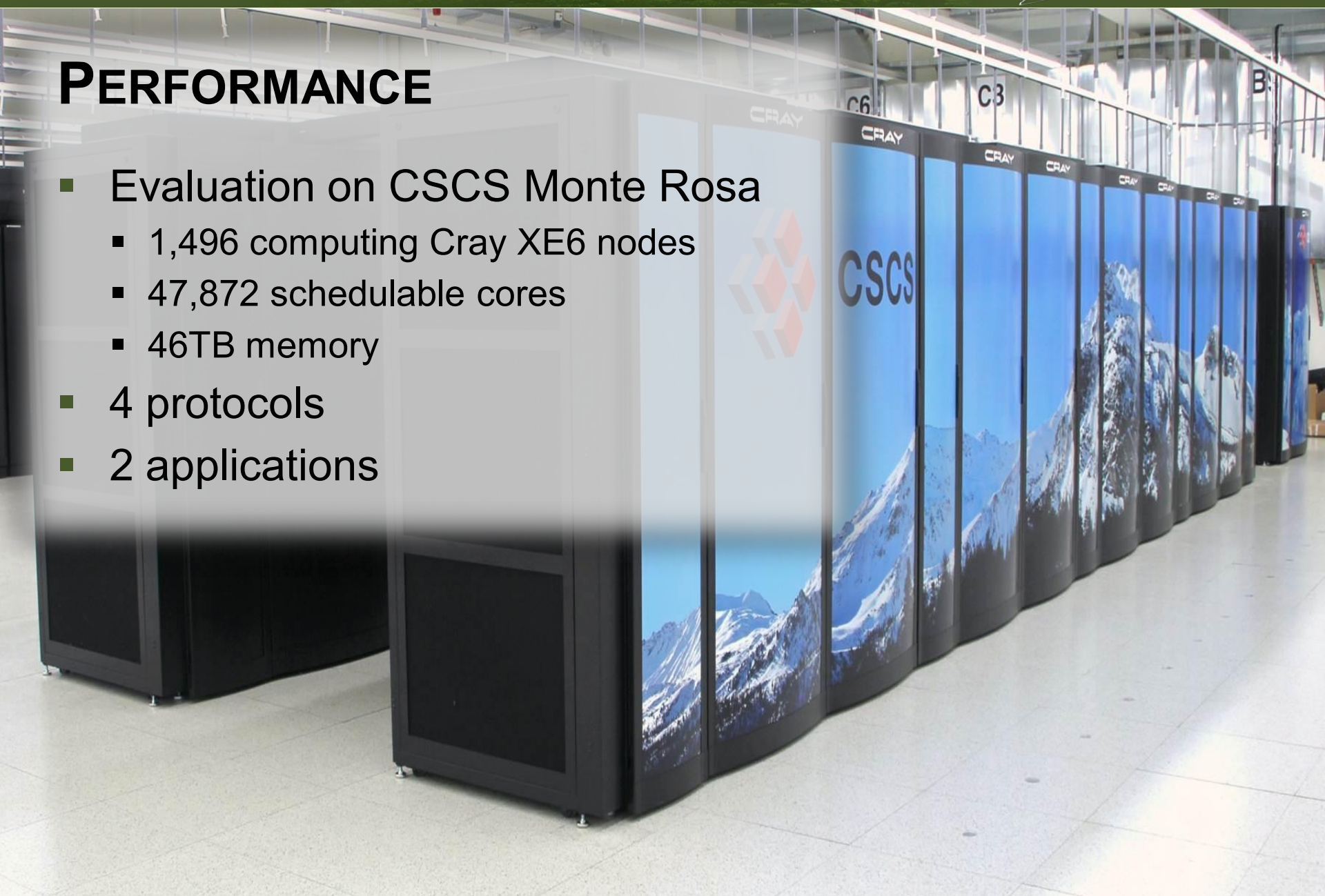
TOPOLOGY-AWARENESS FOR INCREASING RESILIENCE

- The probability of a catastrophic failure in a multi-level computing machine

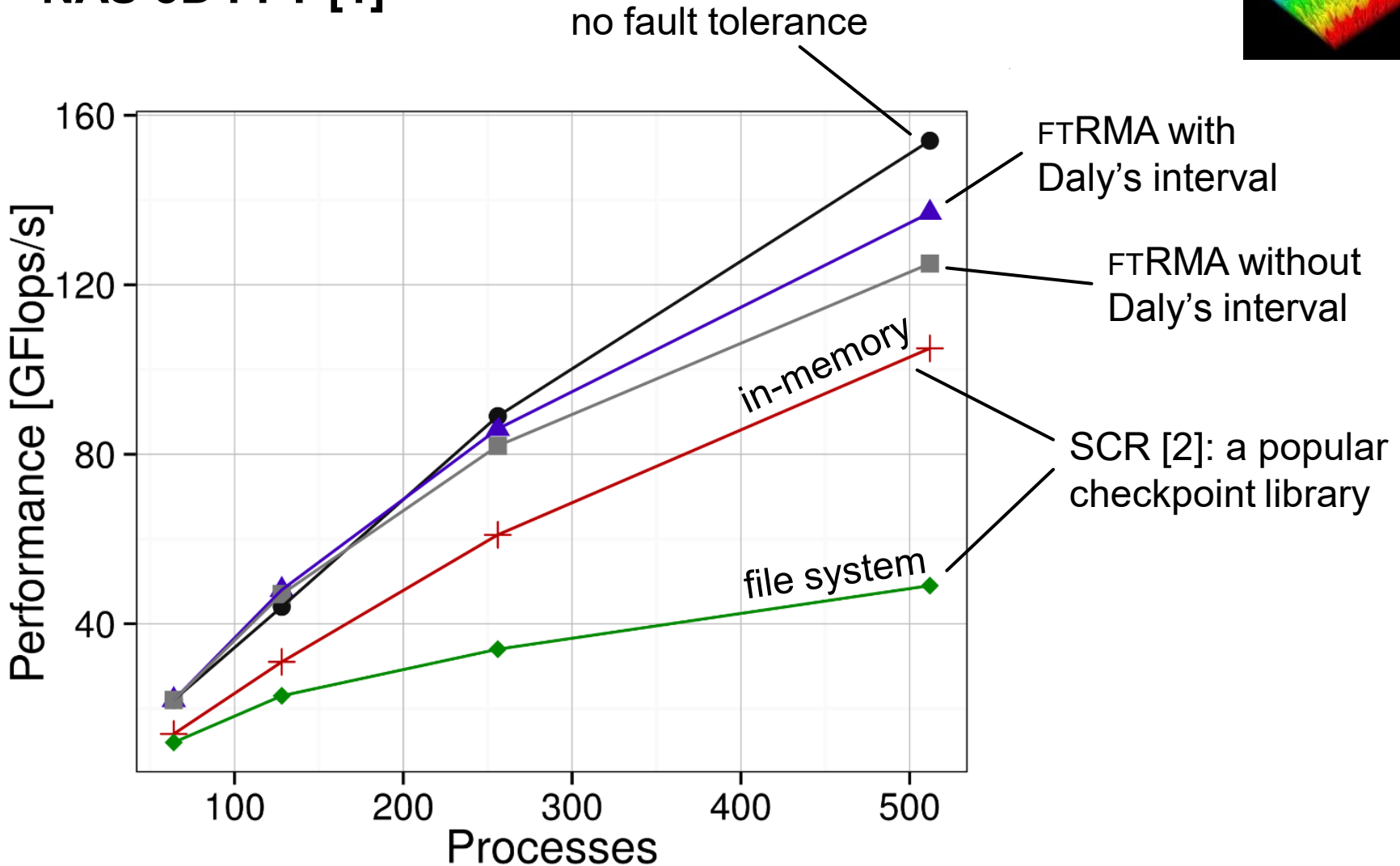
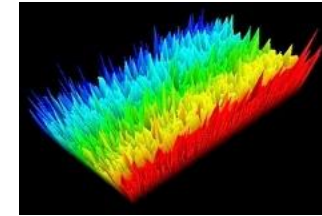


PERFORMANCE

- Evaluation on CSCS Monte Rosa
 - 1,496 computing Cray XE6 nodes
 - 47,872 schedulable cores
 - 46TB memory
- 4 protocols
- 2 applications



PERFORMANCE: COORDINATED CHECKPOINTING NAS 3D FFT [1]

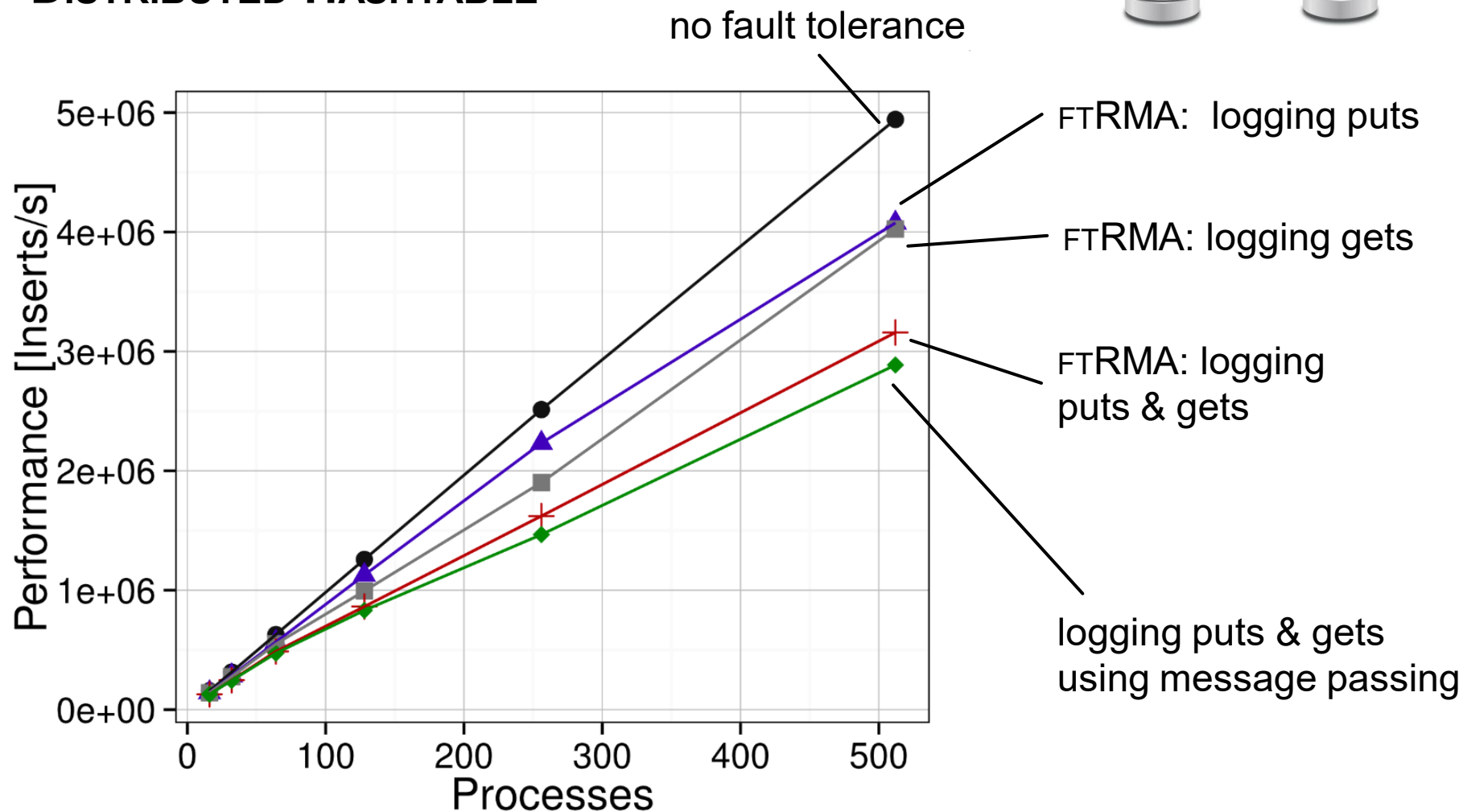


[1] Nishtala et al. Scaling communication-intensive applications on BlueGene/P using one-sided communication and overlap. IPDPS'09

[2] Moody et al. Design, Modeling, and Evaluation of a Scalable Multi-level Checkpointing System. SC10

PERFORMANCE: LOGGING OF PUTS & GETS

DISTRIBUTED HASHTABLE



An insert: 1 get and 1 put are logged
 A collision: 4 gets and 6 puts are logged

Other Recent Results & Open Challenges

■ Networking (topology, routing)

- Besta, TH: “Slim Fly: A Cost Effective Low-Diameter Network Topology”, SC14, best student paper finalist
 - Domke, TH, Matsuoka: “Fail-in-Place Network Design: Interaction between Topology, Routing Algorithm and Failures”, SC14
 - Prisacari, TH, et al.: “Efficient Task Placement and Routing in Dragonfly Networks”, HPDC’14
- ... *optimal routing? Topology mapping? On-chip topologies? ...*

■ Resilience

- Ferreira, et al., TH: “Understanding the Effects of Communication and Coordination on Checkpointing at Scale”, SC14
- ... *optimal checkpointing? Lowest storage vs. runtime overhead? ...*

■ Performance modeling for complex applications

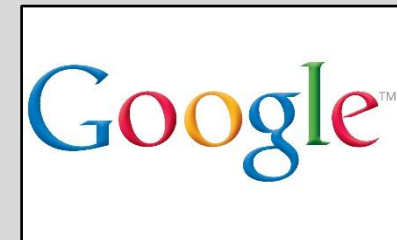
- Bhattacharyya TH: “PEMOGEN: Automatic Adaptive Performance Modeling during Program Runtime”, PACT’14
 - TH, Kwasniewski: “Automatic Complexity Analysis of Explicitly Parallel Programs”, SPAA’14
- ... *designing optimal implementations? Automated modeling for co-design? ...*

ACKNOWLEDGMENTS



Paul Hargrove (and the whole UPC team), the MPI Forum
RMA WG, many others, ...

... and



Thanks for your attention!

