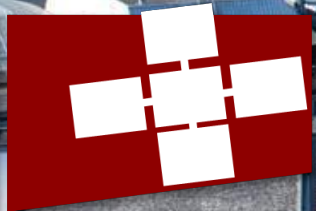


T. HOEFLER

RDMA, Scalable MPI-3 RMA, and Next-Generation Post-RDMA Interconnects

HPC Advisory Council Swiss Workshop, Lugano, Switzerland

WITH HELP OF ROBERT GERSTENBERGER, MACIEJ BESTA, S. DI GIROLAMO, K. TARANOV, R. E. GRANT, R. BRIGHTWELL AND ALL OF SPCL



EuroMPI'19

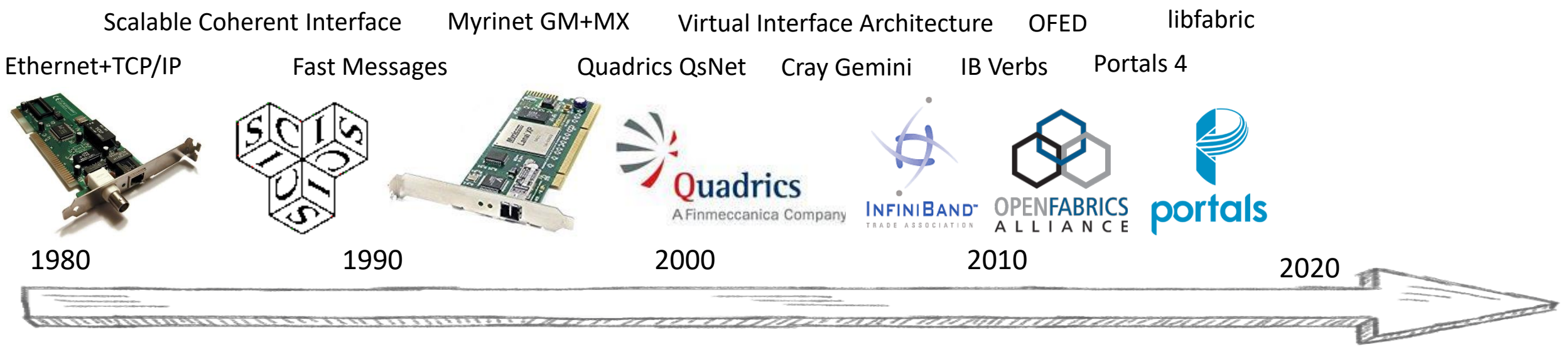
September 11-13 2019

Zurich, Switzerland

<https://eurompi19.inf.ethz.ch>

Submit papers by April 15th!

The Development of High-Performance Networking Interfaces



sockets (active) message based protocol offload remote direct memory access (RDMA)

coherent memory access OS bypass zero copy triggered operations

InfiniBand Trade Association Launches the RoCE Initiative to Advance RDMA over Converged Ethernet Solutions

RoCE delivers significant performance and efficiency gains to cloud, storage, virtualization and hyper-converged infrastructures
businessinsider.com

Microsoft to Drive RDMA Into Datacenters and Clouds
November 18, 2013 by Timothy Prickett Morgan

RDMA over Ethernet - the Rocky road to convergence
17 November 2015 | By Brandon Hoff

TOP 500[®]
 SUPERCOMPUTER SITES June 2017

95 / top-100 systems use RDMA
 >285 / top-500 systems use RDMA

RDMA as MPI-3.0 REMOTE MEMORY ACCESS TRANSPORT

- MPI-3.0 supports RMA (“MPI-3 One Sided”)
 - Designed to react to hardware
 - Majority of HPC networks support



research highlights

Communications of the ACM, October 2018, Vol. 61 No. 10, Pages 106-113

DOI:10.1145/3264413

Enabling Highly Scalable Remote Memory Access Programming with MPI-3 One Sided

By Robert Gerstenberger,* Maciej Besta, and Torsten Hoefler

Abstract

Modern high-performance networks offer remote direct memory access (RDMA) that exposes a process' virtual address space to other processes in the network. The *Message Passing Interface* (MPI) specification has recently been extended with a programming interface called MPI-3 *Remote Memory Access* (MPI-3 RMA) for efficiently exploiting state-of-the-art RDMA features. MPI-3 RMA enables a powerful programming model that alleviates many message passing downsides. In this work, we design and develop bufferless protocols that demonstrate how to implement this interface and support scaling to millions of cores with negligible memory consumption while providing highest performance and minimal overheads. To arm programmers, we provide a spectrum of performance models for RMA functions that enable rigorous mathematical analysis of application performance and facilitate the development of codes that solve given tasks within specified time and energy budgets. We validate the usability of our library and models with several application studies with up to half a million processes. In a wider sense, our work illustrates how to use RMA principles to accelerate computation- and data-intensive codes.

Supercomputers consist of massively parallel nodes, each supporting up to hundreds of hardware threads in a single shared-memory domain. Up to tens of thousands of such nodes can be connected with a high-performance network, providing large-scale distributed-memory parallelism. For example, the Blue Waters machine has >700,000 cores and a peak computational bandwidth of >13 petaflops.

Programming such large distributed computers is far from trivial: an ideal programming model should tame the complexity of the underlying hardware and offer an easy abstraction for the programmer to facilitate the development of high-performance codes. Yet, it should also be able to effectively utilize the available massive parallelism and various heterogeneous processing units to ensure highest scalability and speedups. Moreover, there has been a growing need for the support for *performance modeling*: a rigorous mathematical analysis of application performance. Such formal reasoning facilitates developing codes that solve given tasks within the assumed time and energy budget.

The *Message Passing Interface* (MPI)¹¹ is the *de facto* standard API used to develop applications for distributed-memory supercomputers. MPI specifies message passing as well as remote memory access semantics and offers a rich set of features that facilitate developing highly scalable and portable codes; message passing has been the prevalent model so far. MPI's message passing specification does not prescribe specific ways how to exchange messages and thus enables flex-

1. INTRODUCTION

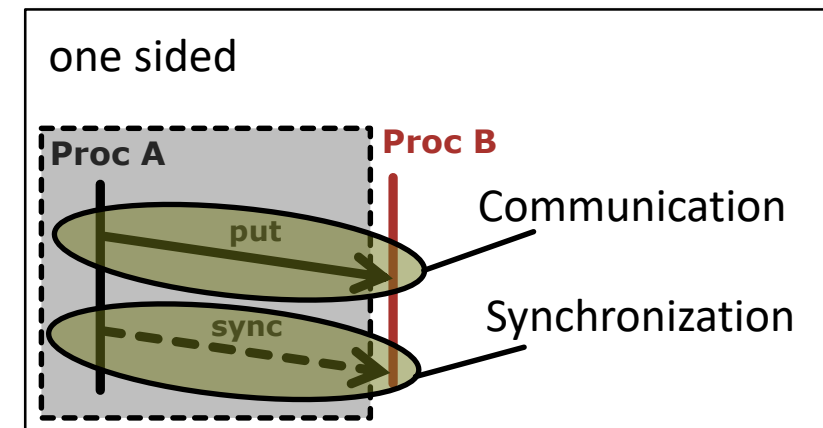
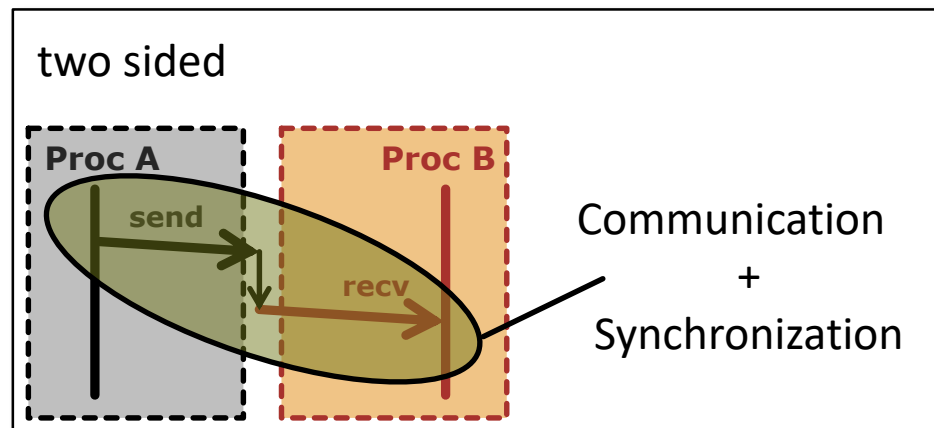
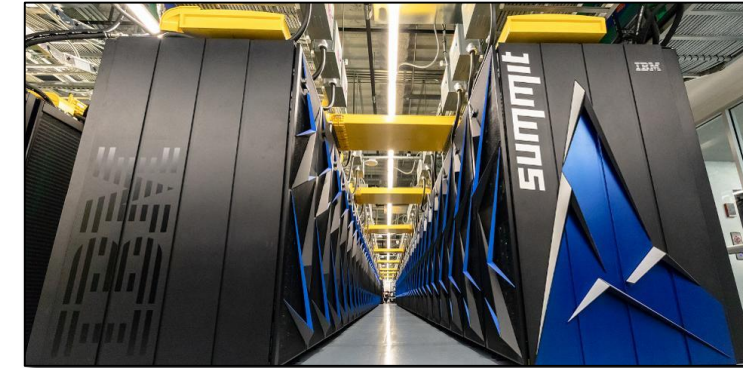
Supercomputers have driven the progress of various soci-

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

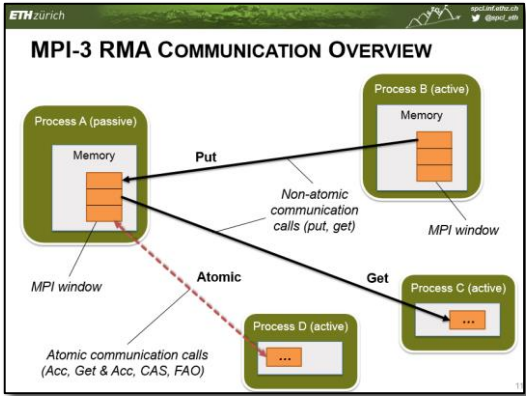
[1] <http://www.mpi-forum.org/docs/mpi-3.0/mpi30-report.pdf>

MPI-3.0 REMOTE MEMORY ACCESS

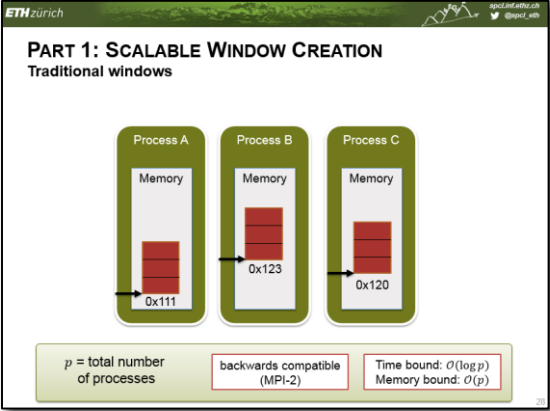
- MPI-3.0 supports RMA (“MPI One Sided”)
 - Designed to react to hardware trends
 - Majority of HPC networks support RDMA
- Communication is „one sided” (no involvement of destination)
- RMA decouples communication & synchronization
 - Different from message passing



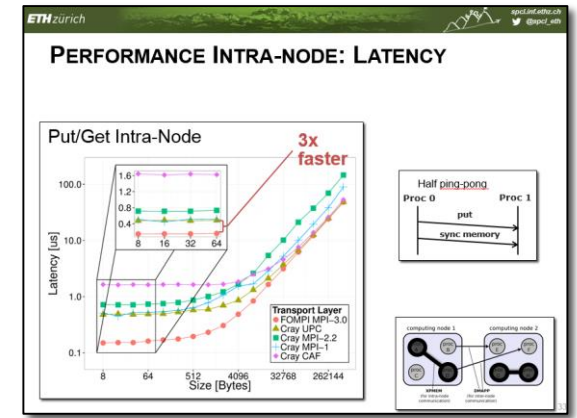
PRESENTATION OVERVIEW



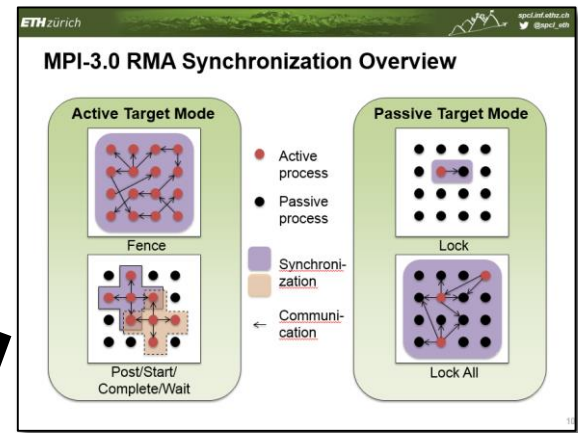
1. Overview of three MPI-3 RMA concepts



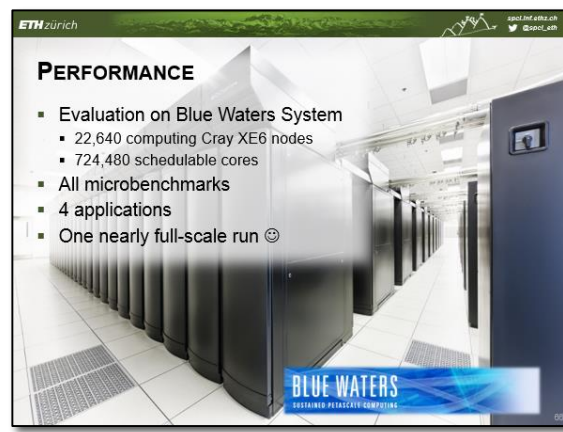
2. MPI window creation



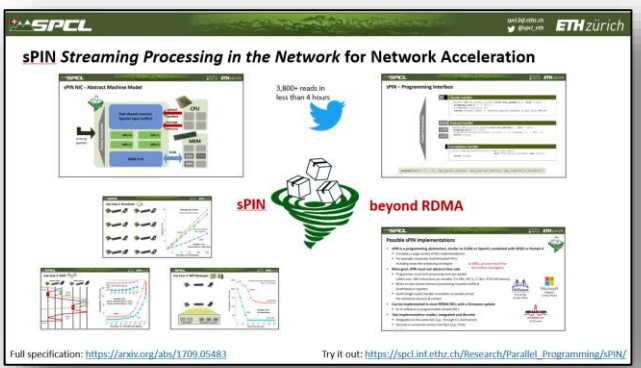
3. Communication



4. Synchronization



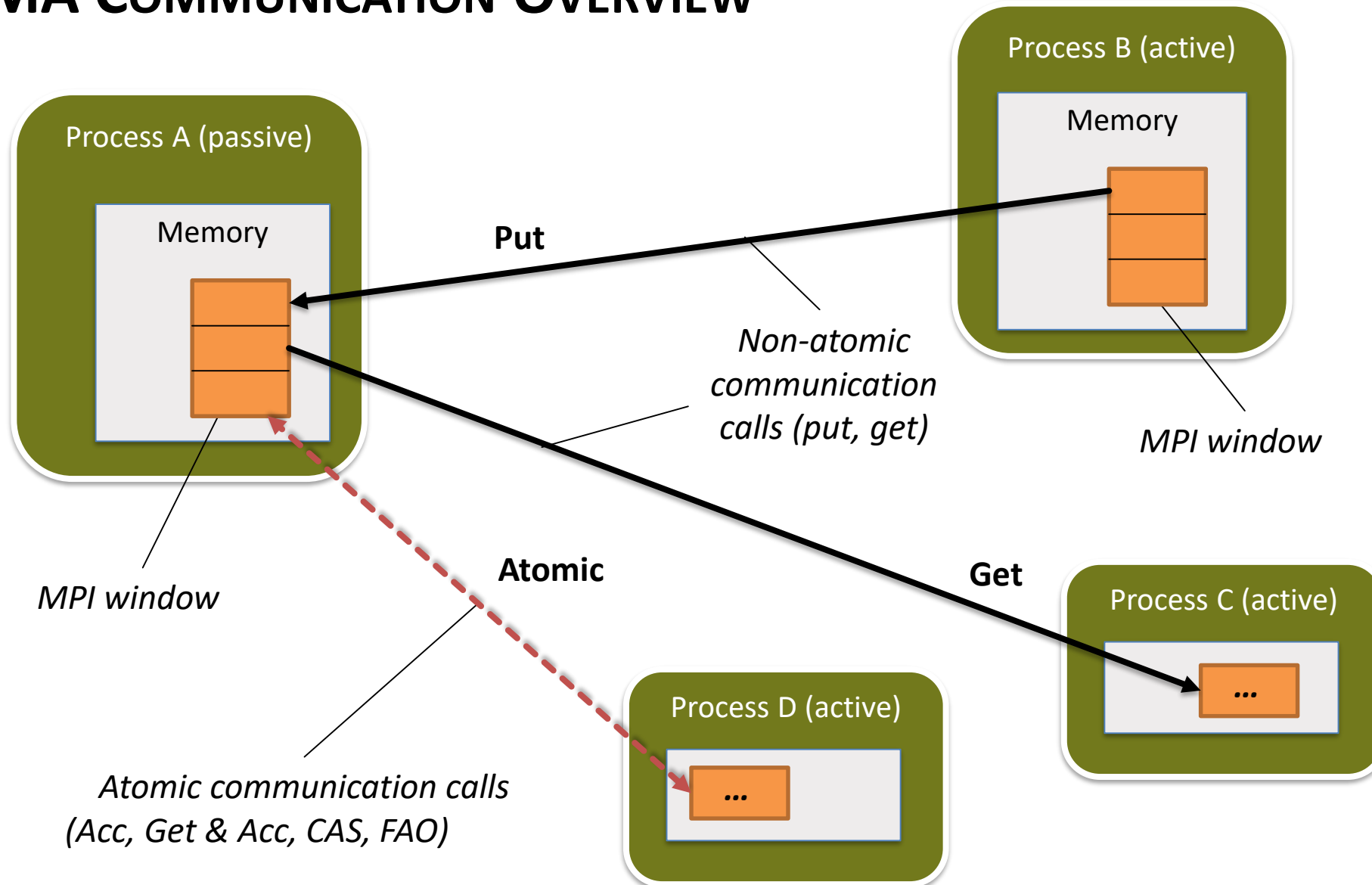
5. Application evaluation



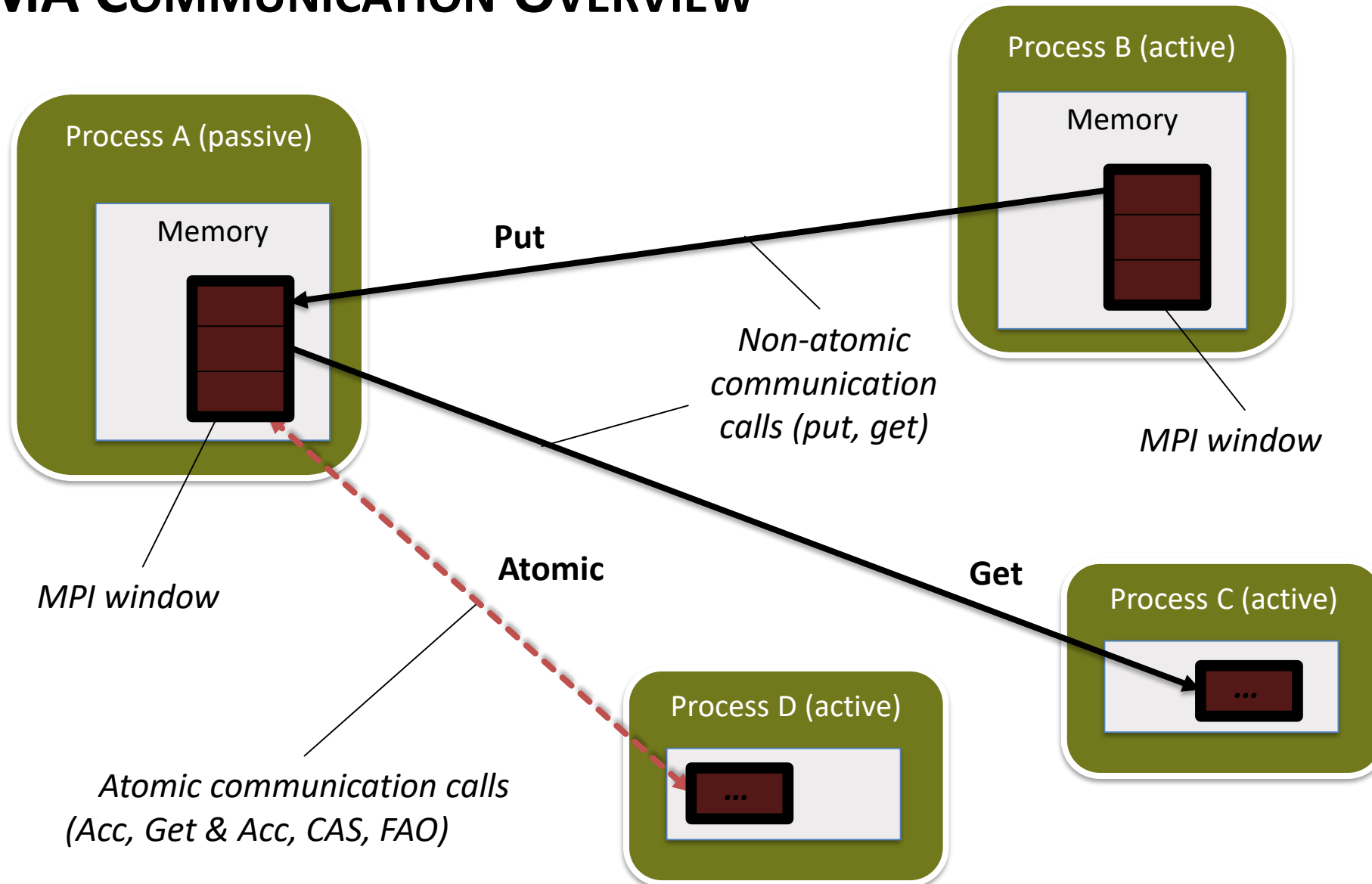
6. Post-RDMA networking

Outlook!

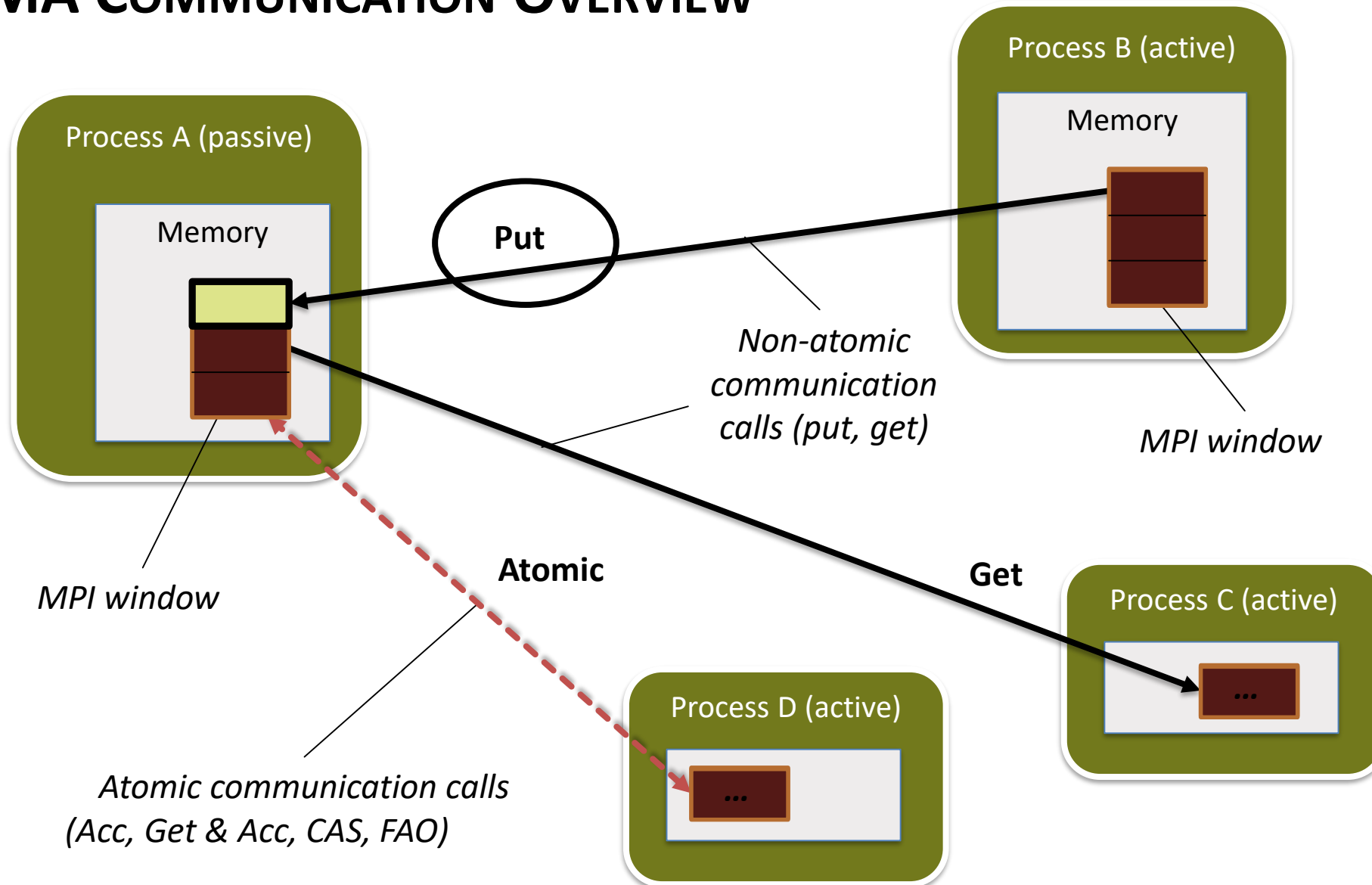
MPI-3 RMA COMMUNICATION OVERVIEW



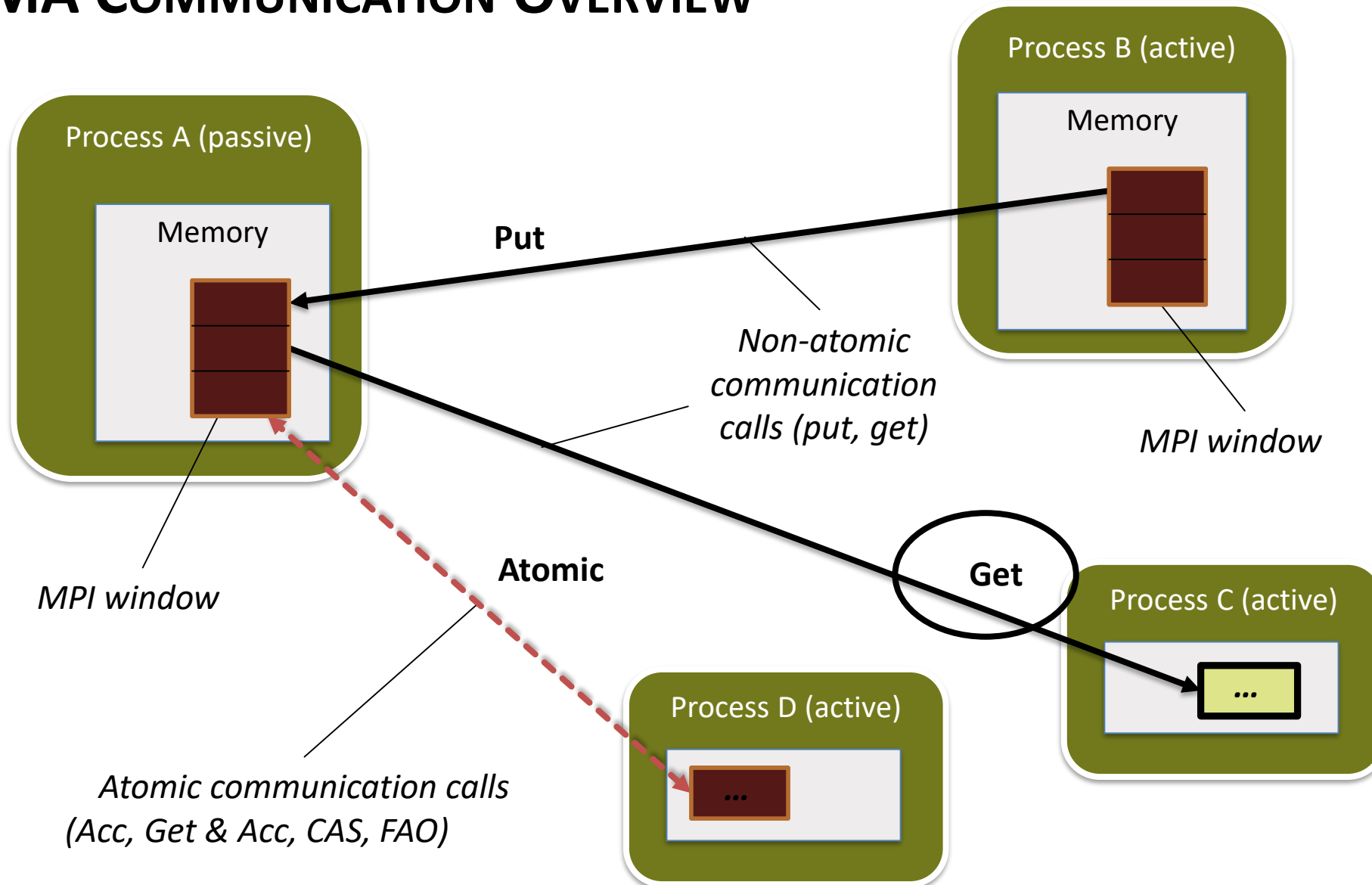
MPI-3 RMA COMMUNICATION OVERVIEW



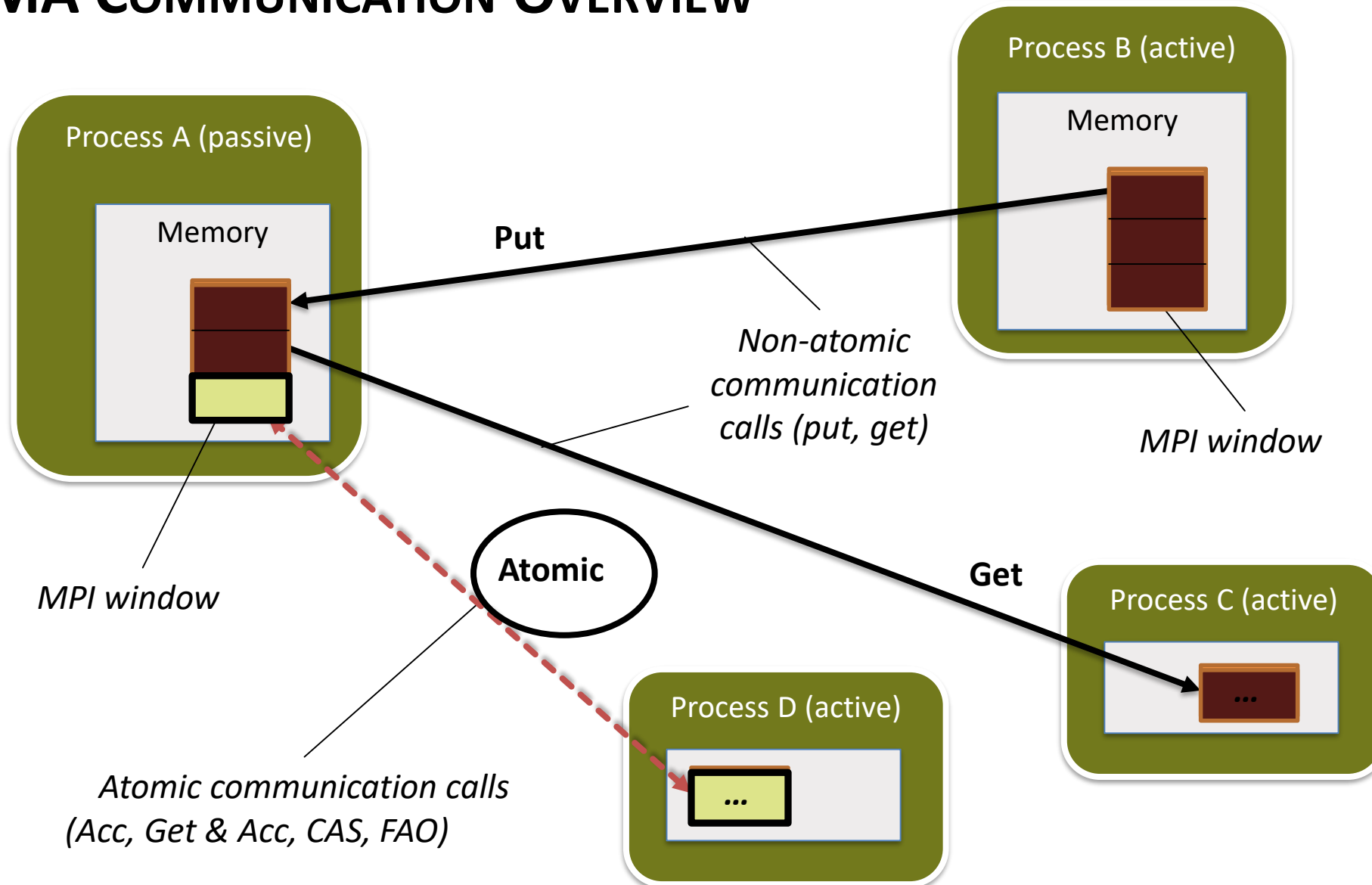
MPI-3 RMA COMMUNICATION OVERVIEW



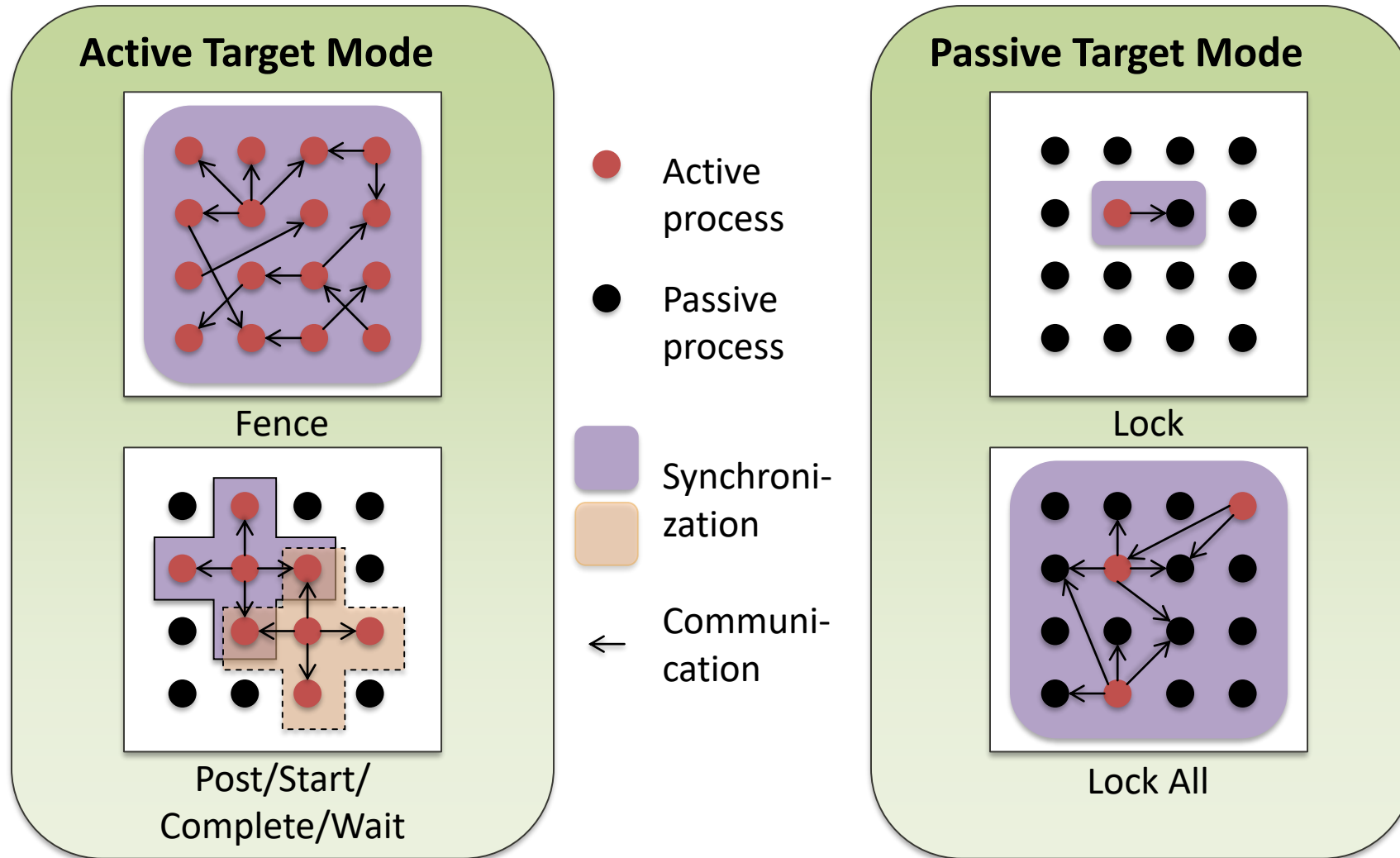
MPI-3 RMA COMMUNICATION OVERVIEW



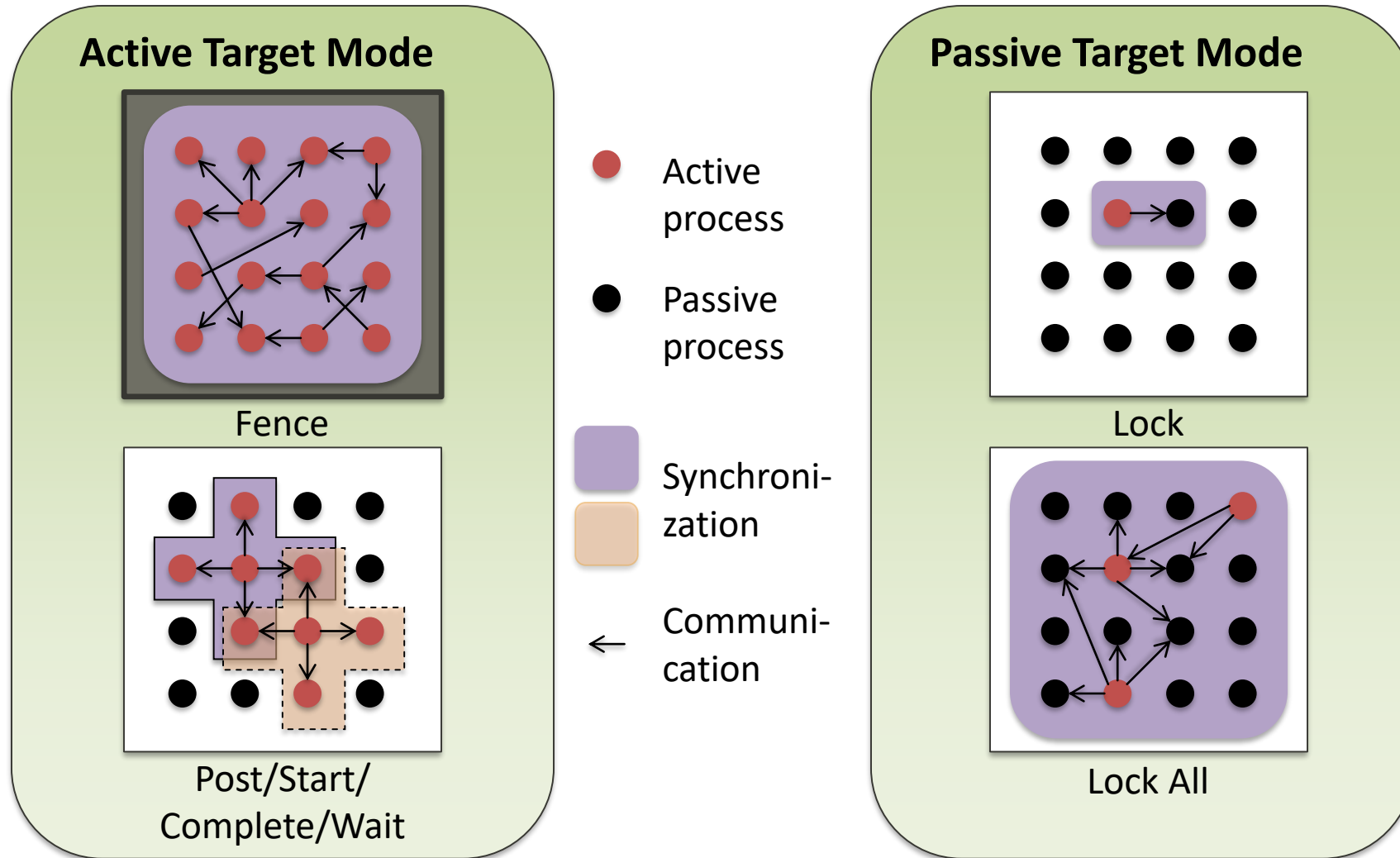
MPI-3 RMA COMMUNICATION OVERVIEW



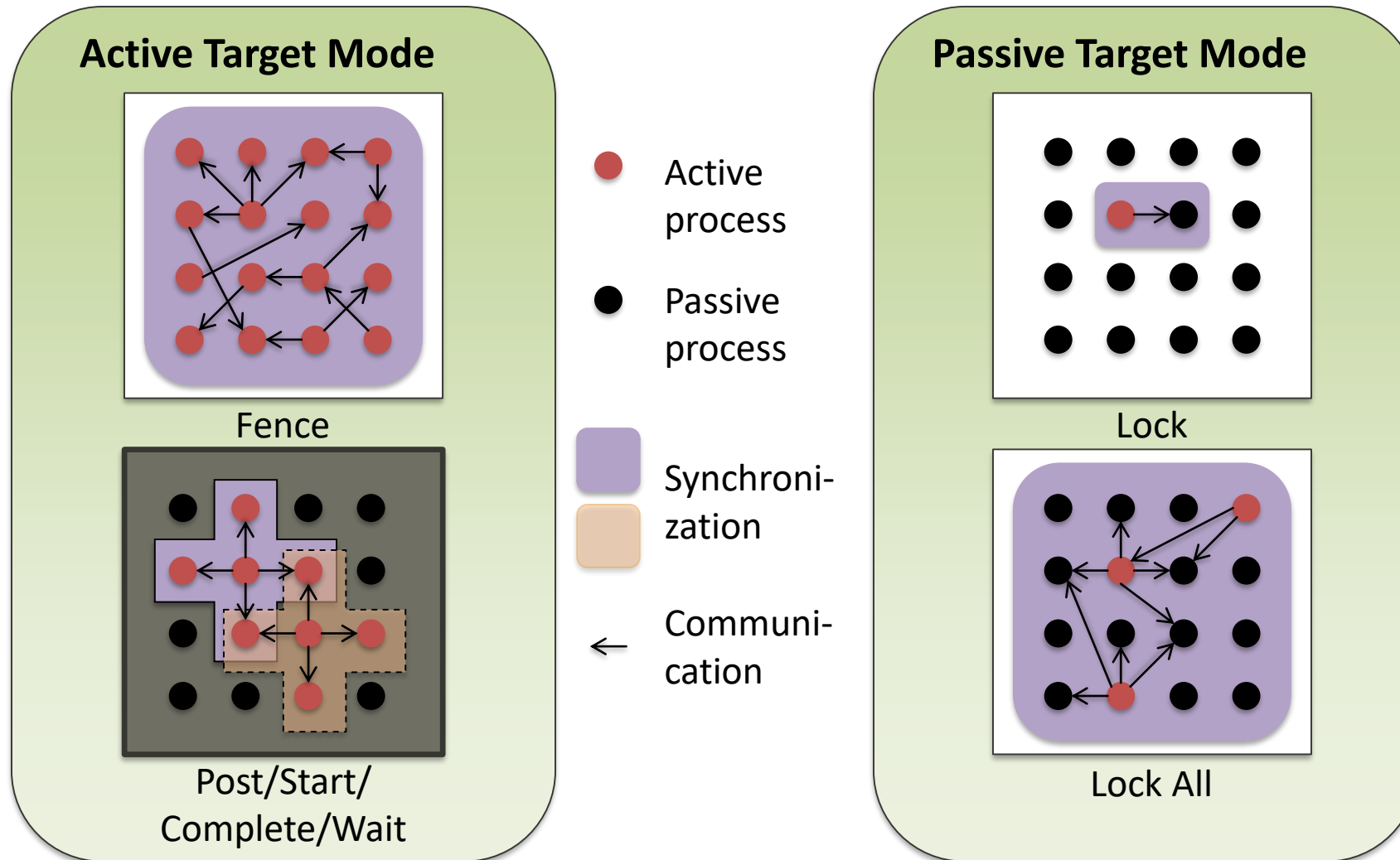
MPI-3.0 RMA Synchronization Overview



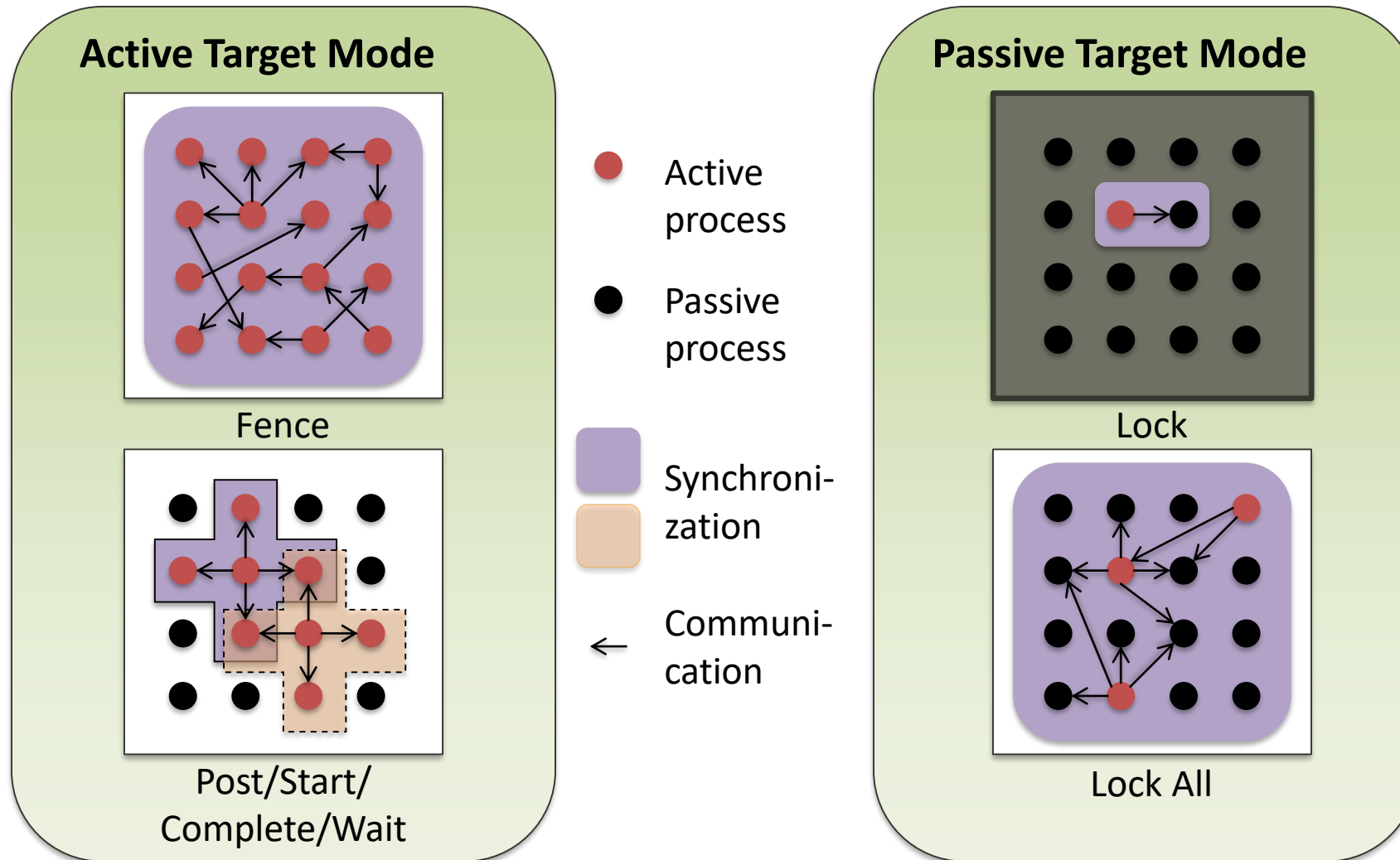
MPI-3.0 RMA Synchronization Overview



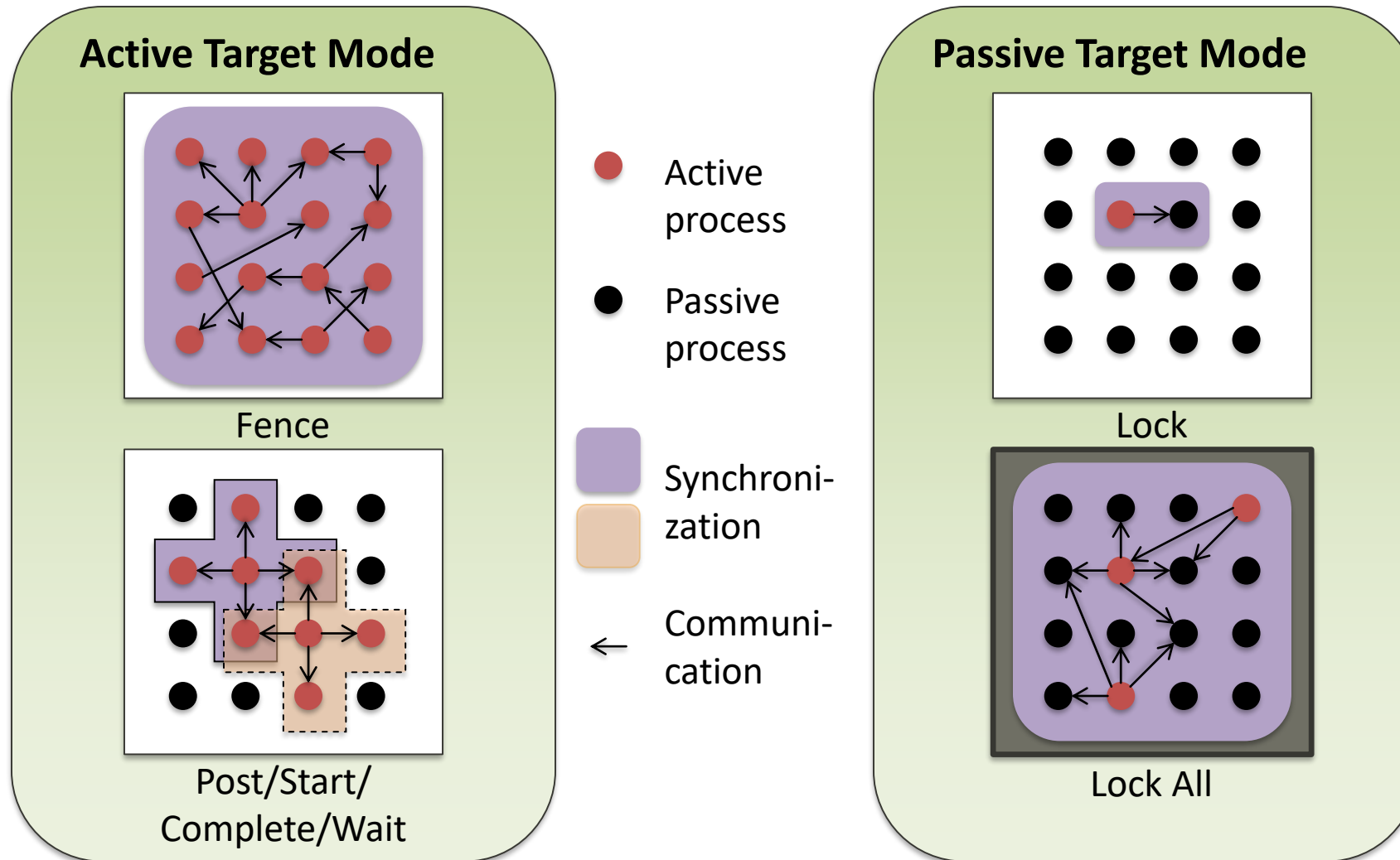
MPI-3.0 RMA Synchronization Overview



MPI-3.0 RMA Synchronization Overview



MPI-3.0 RMA Synchronization Overview



SCALABLE PROTOCOLS & REFERENCE IMPLEMENTATION

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

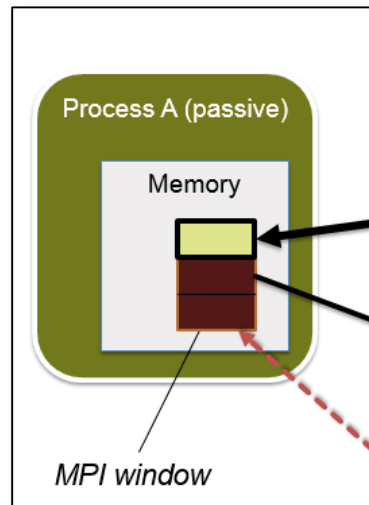
SCALABLE PROTOCOLS & REFERENCE IMPLEMENTATION

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

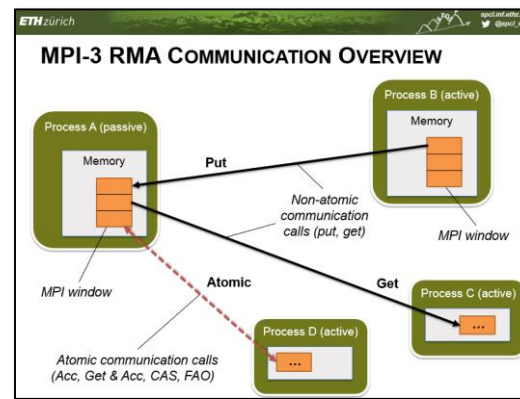


SCALABLE PROTOCOLS & REFERENCE IMPLEMENTATION

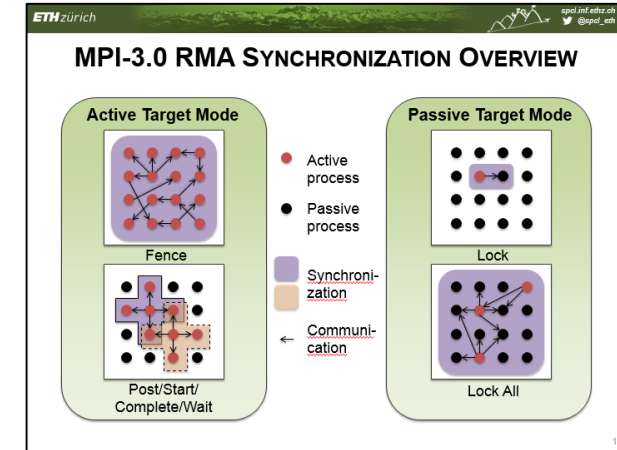
- Scalable & generic protocols
 - Can be used on any RDMA network (e.g., OFED/IB)
 - Window creation, communication and synchronization



Window creation



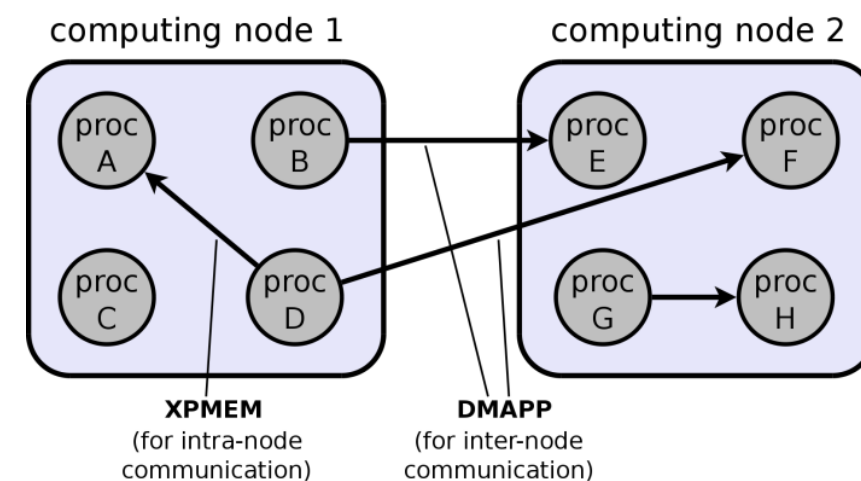
Communication



Synchronization

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:

- time by avoiding message matching and synchronization overheads
- energy by reducing data-movement, e.g., it avoids additional copies of eager messages
- 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), a fully functional MPI-3.0 RMA...

Tweets

Torsten Hoeller 3 Nov
 @thoeller
 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

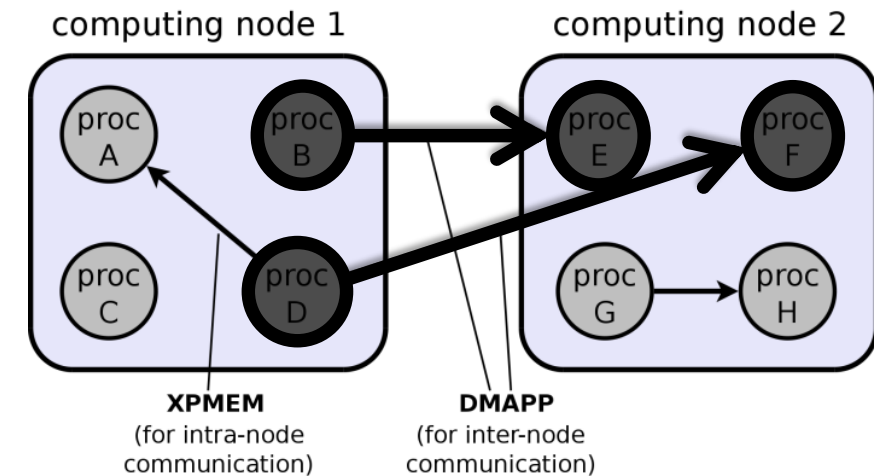
SPCL@ETH 27 Oct
 @spcl_eth
 Congratulations to Bogdan to his new work on pattern-specific routing on fat trees that will increase rdt!

http://spcl.inf.ethz.ch/Research/Parallel_Programming/foMPI

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

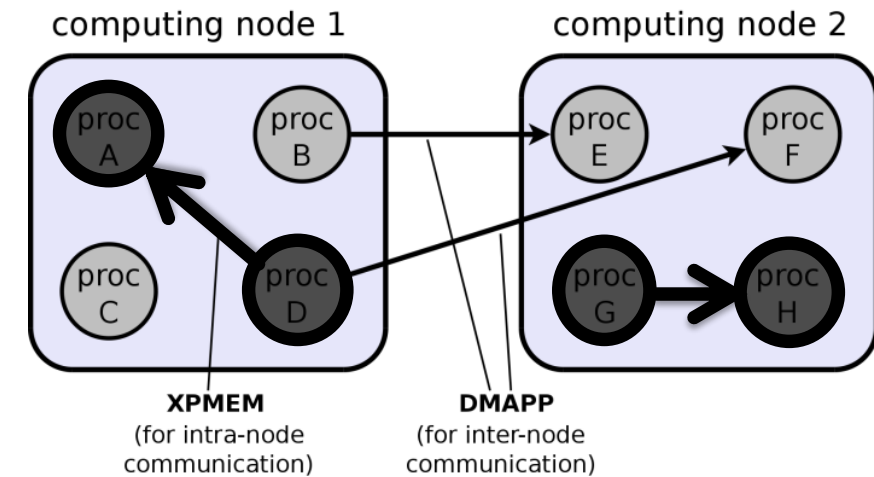


The screenshot shows the website for the Scalable Parallel Computing Lab. The main heading is "foMPI: A Fast One-Sided MPI-3.0 Implementation". The page includes a navigation menu on the left, a main content area with a "Motivation" section, and a "Tweets" section on the right. The "Motivation" section discusses the benefits of RDMA hardware and the development of foMPI. The "Tweets" section shows a tweet from SPCL@ETH dated 27 Oct, congratulating Bogdan on his work on pattern-specific routing on fat trees.

http://spcl.inf.ethz.ch/Research/Parallel_Programming/foMPI

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:

- time by avoiding message matching and synchronization overheads
- energy by reducing data-movement, e.g., it avoids additional copies of eager messages
- 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), a fully-functional MPI-3.0 RMA.

Tweets

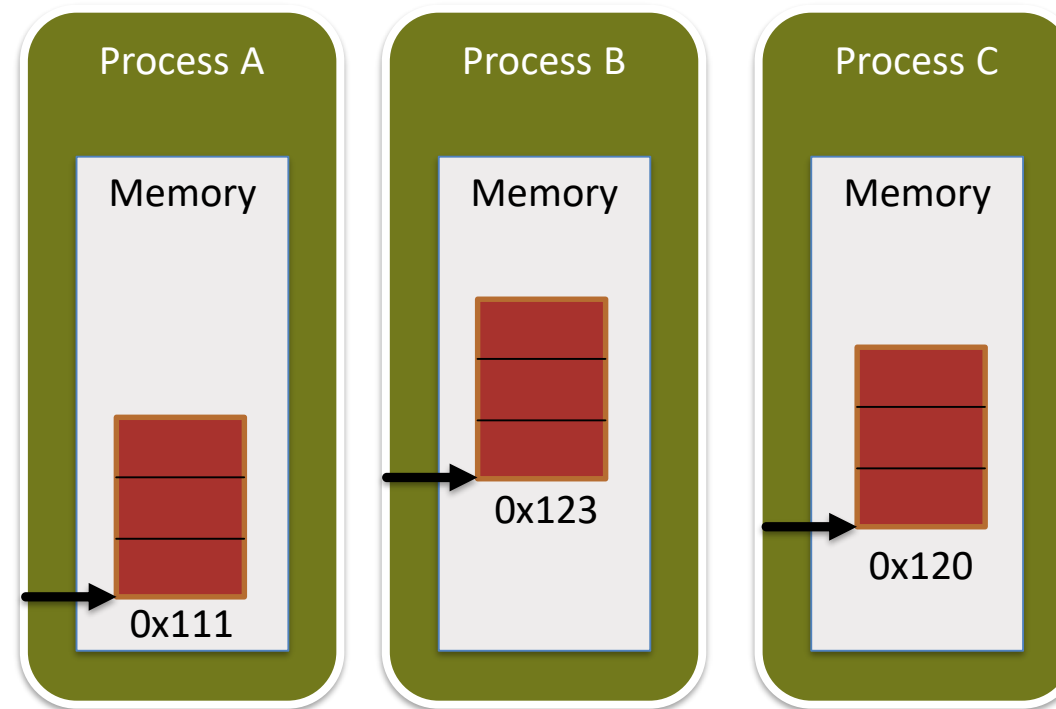
Torsten Hoeller 3 Nov
@thoeller
Interested in modern HPC programming? Check out our Advanced MPI tutorial @Supercomputing1
tor.inf.ethz.ch/blog/index.php/MPI-2.2-MPI-3-and-more/
Retweeted by SPCL@ETH

SPCL@ETH 27 Oct
@spcl_eth
Congratulations to Bogdan to his new work on pattern-specific routing on fat trees that will increase rdt!

http://spcl.inf.ethz.ch/Research/Parallel_Programming/foMPI

PART 1: SCALABLE WINDOW CREATION

Traditional windows



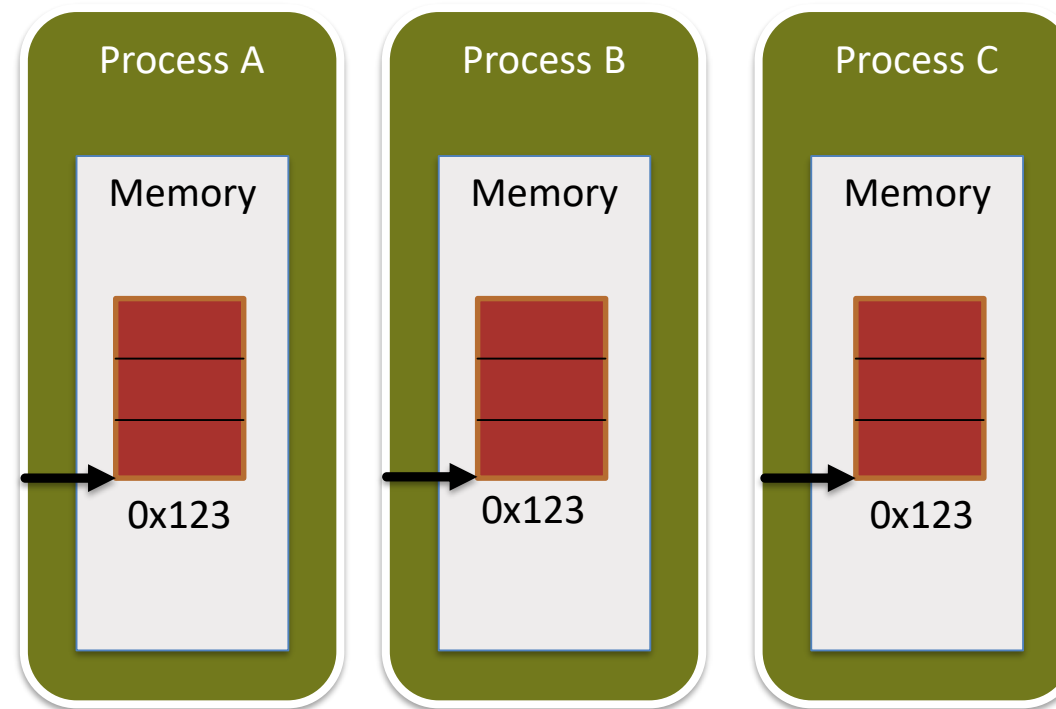
p = total number
of processes

backwards compatible
(MPI-2)

Time bound: $\mathcal{O}(p)$
Memory bound: $\mathcal{O}(p)$

PART 1: SCALABLE WINDOW CREATION

Allocated windows



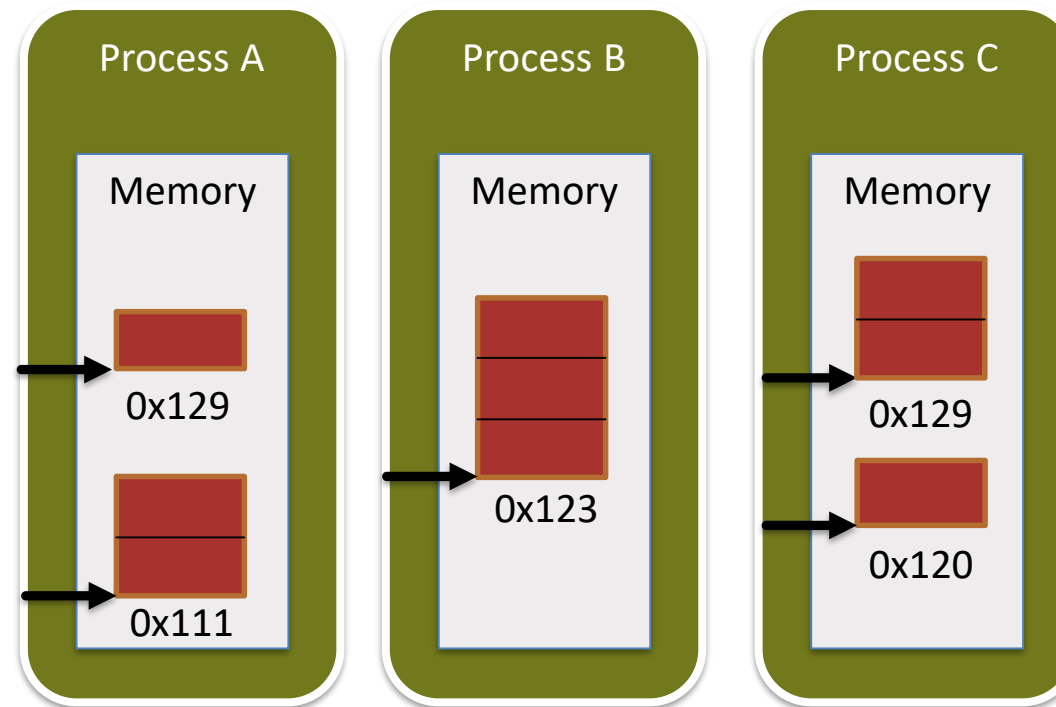
p = total number
of processes

Allows MPI
to allocate memory

Time bound: $\mathcal{O}(\log p)$ (*whp*)
Memory bound: $\mathcal{O}(1)$

PART 1: SCALABLE WINDOW CREATION

Dynamic windows



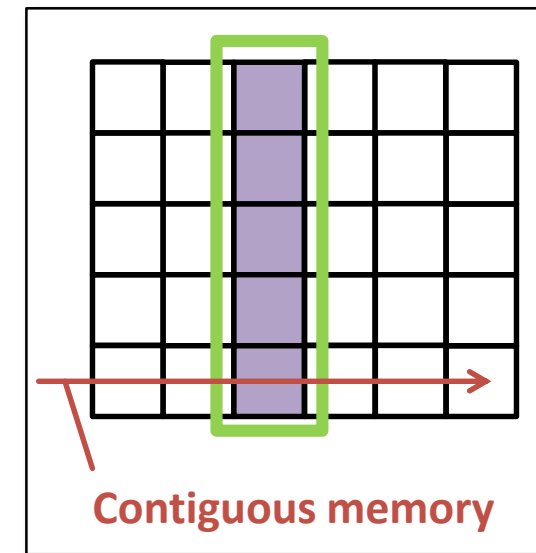
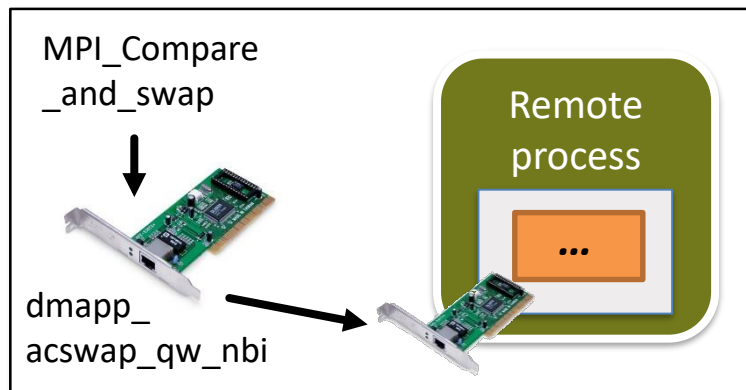
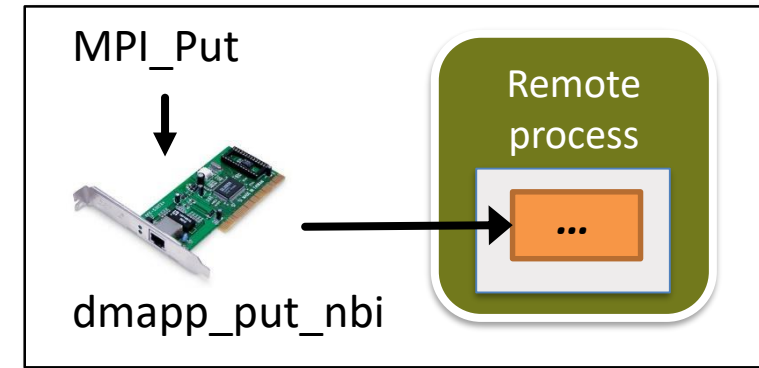
p = total number
of processes

Local attach/detach
Most flexible

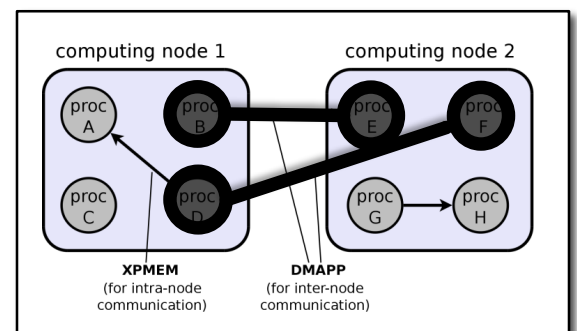
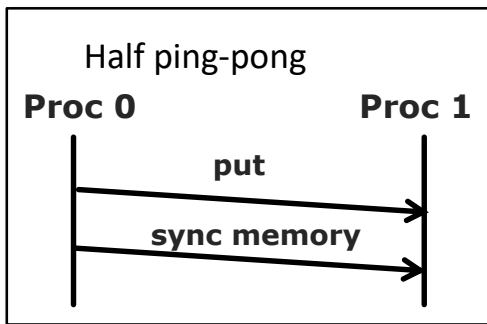
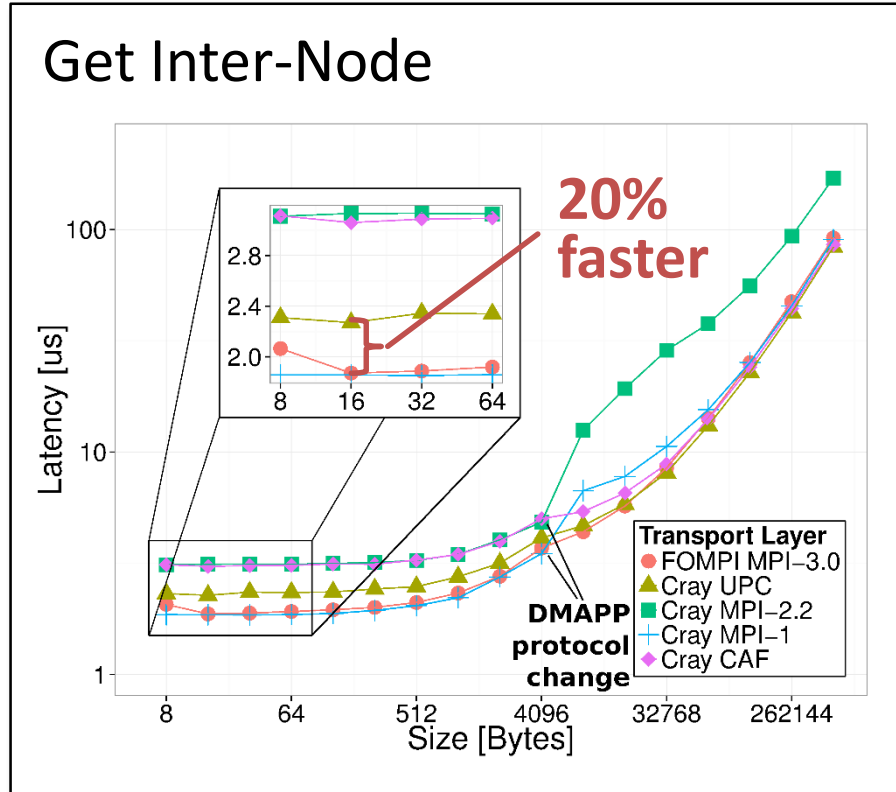
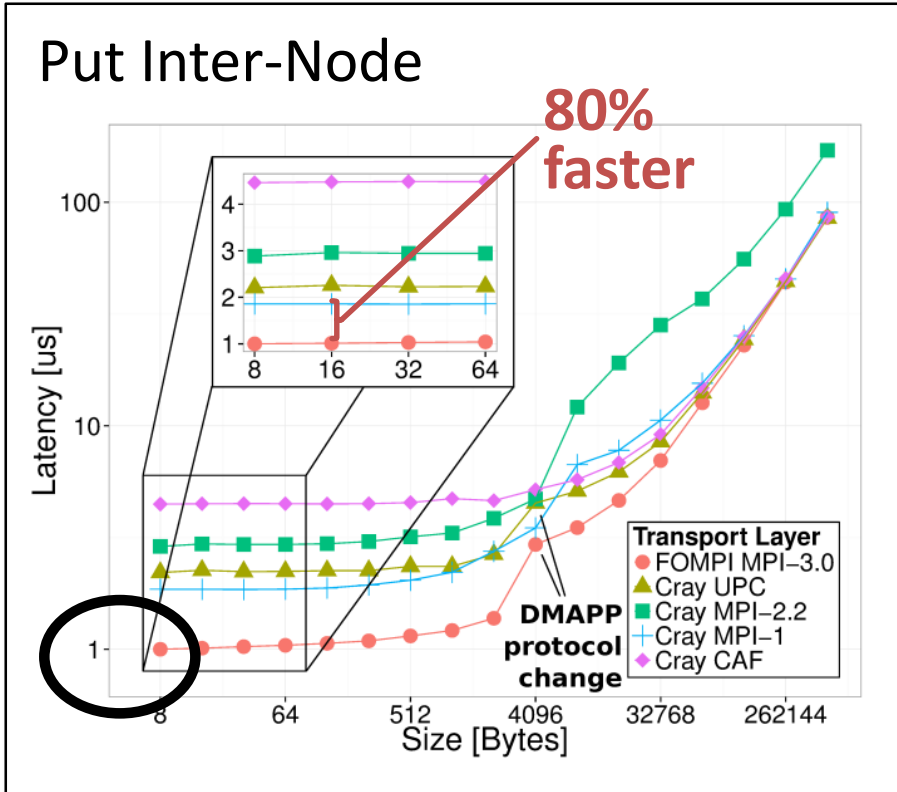
Time bound: $\mathcal{O}(p)$
Memory bound: $\mathcal{O}(p)$

PART 2: COMMUNICATION

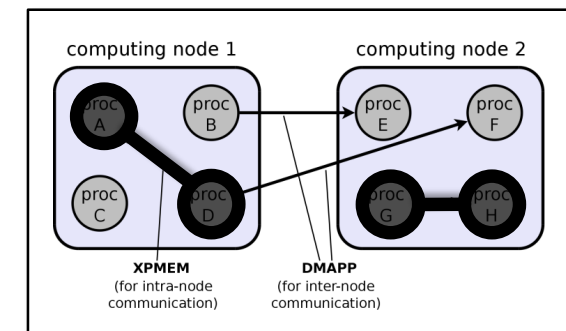
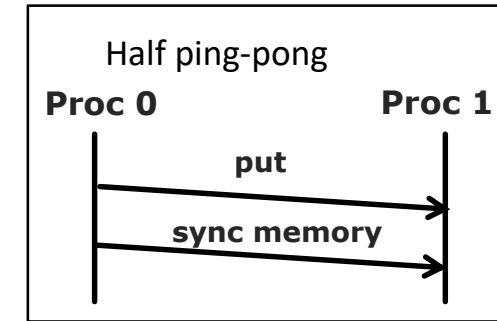
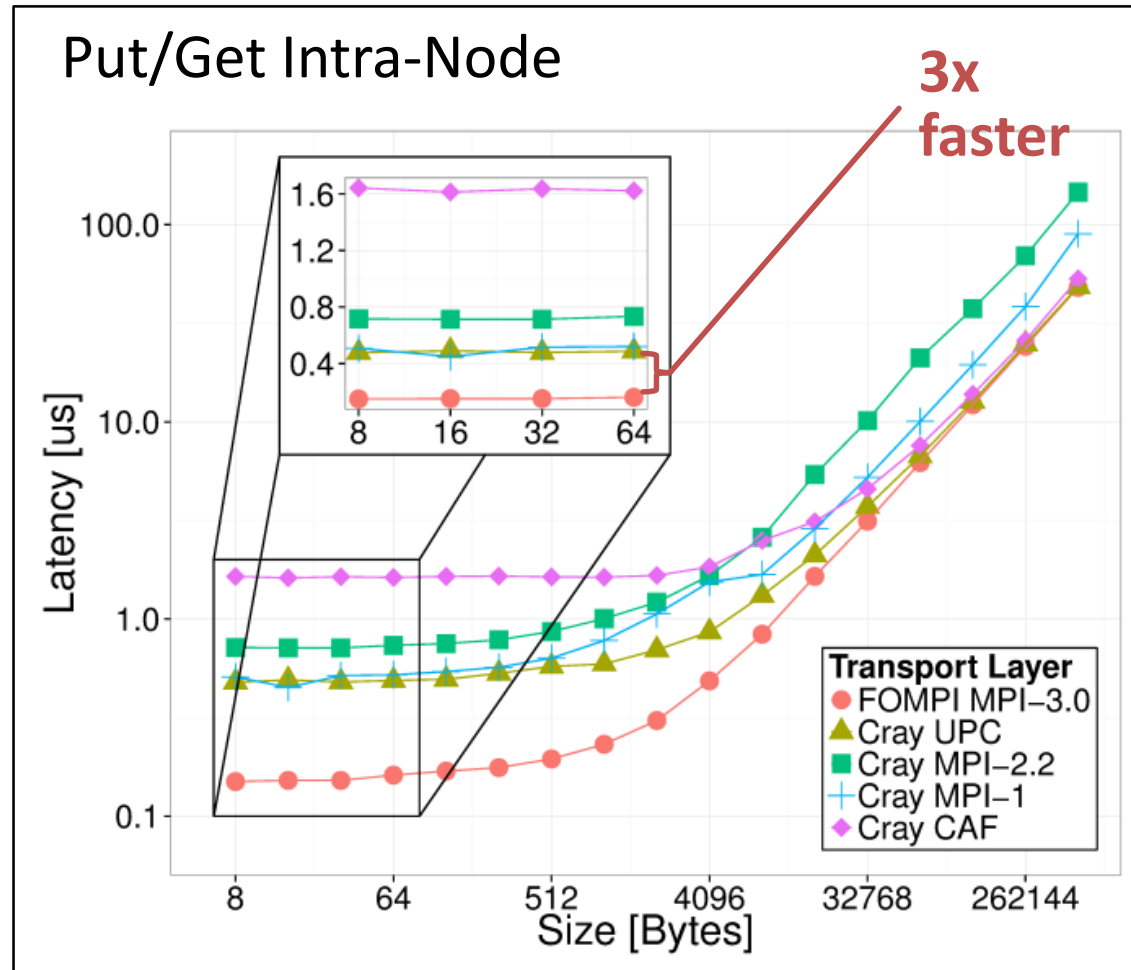
- Put and Get:
 - Direct DMAPP put and get operations or local (blocking) memcpy (XPMEM)
- Accumulate:
 - DMAPP atomic operations for 64 bit types
 - ...or fall back to remote locking protocol
- MPI datatype handling with MPITypes library [1]
 - Fast path for contiguous data transfers of common intrinsic datatypes (e.g., MPI_DOUBLE)



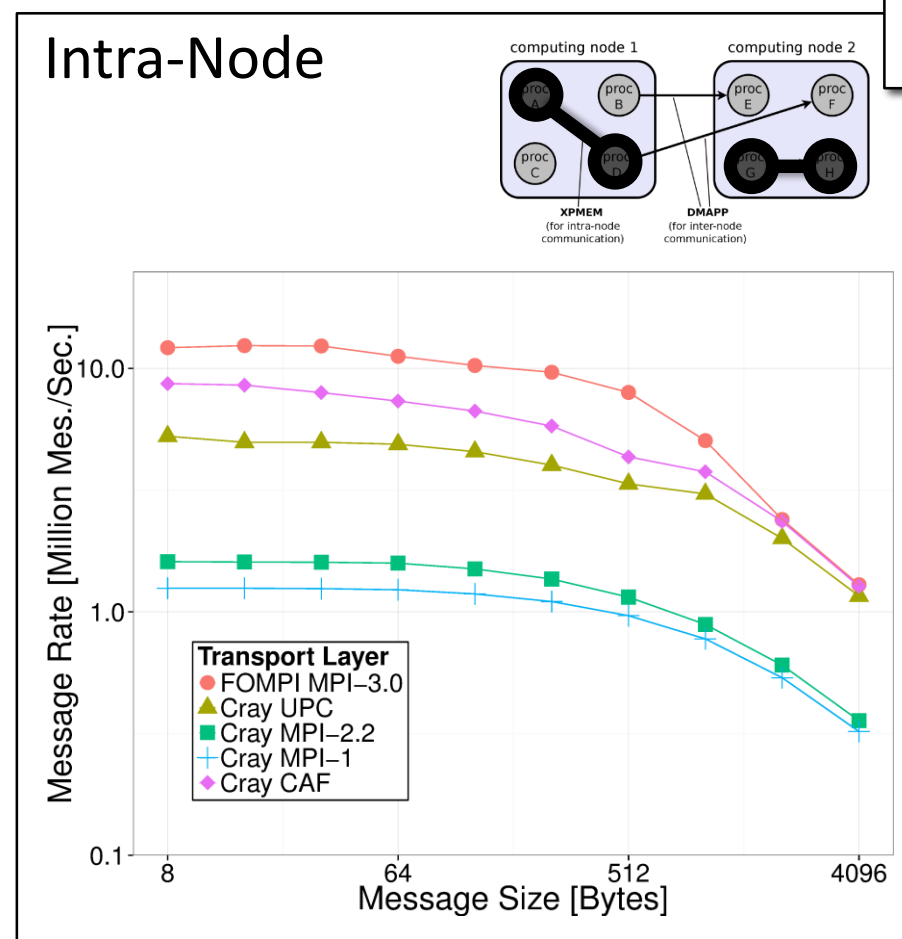
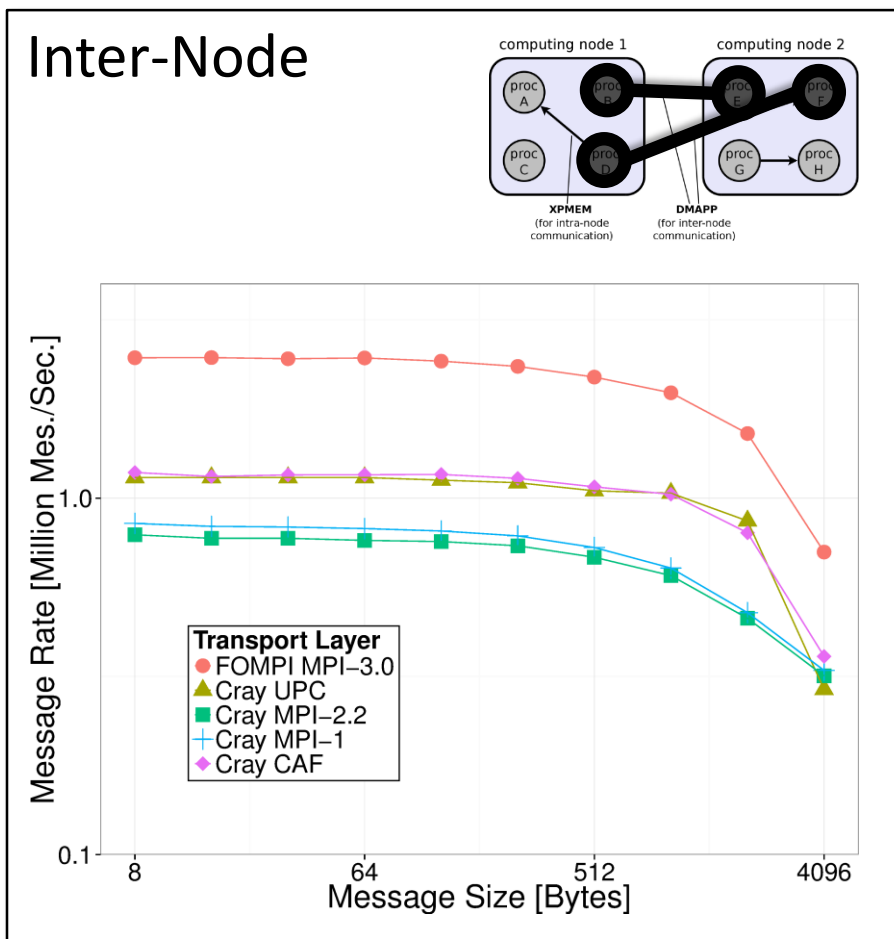
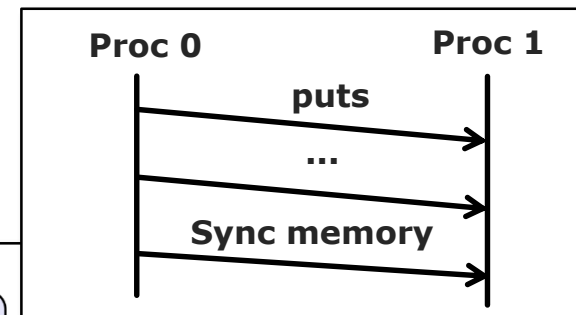
PERFORMANCE INTER-NODE: LATENCY



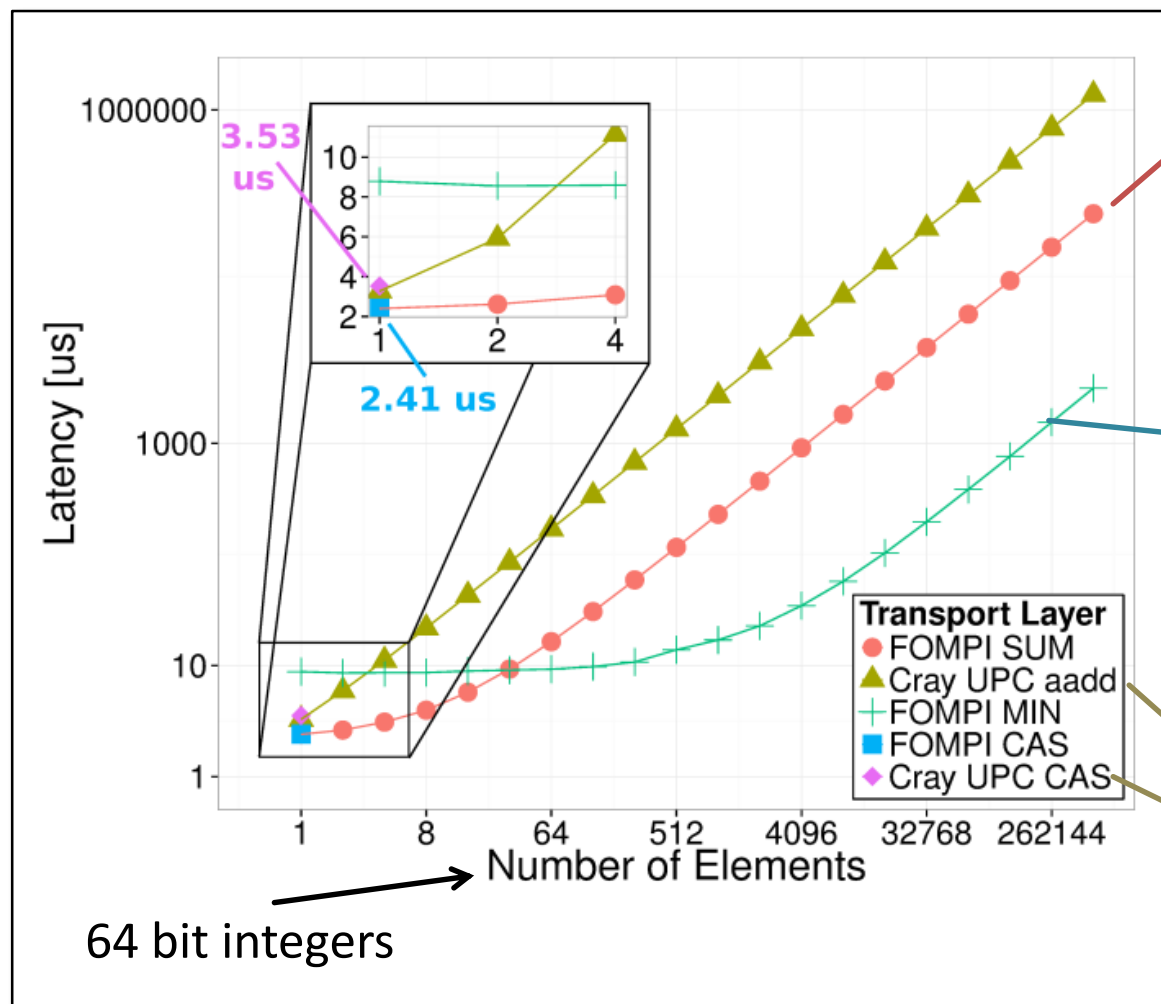
PERFORMANCE INTRA-NODE: LATENCY



PERFORMANCE: MESSAGE RATE



PERFORMANCE: ATOMICS

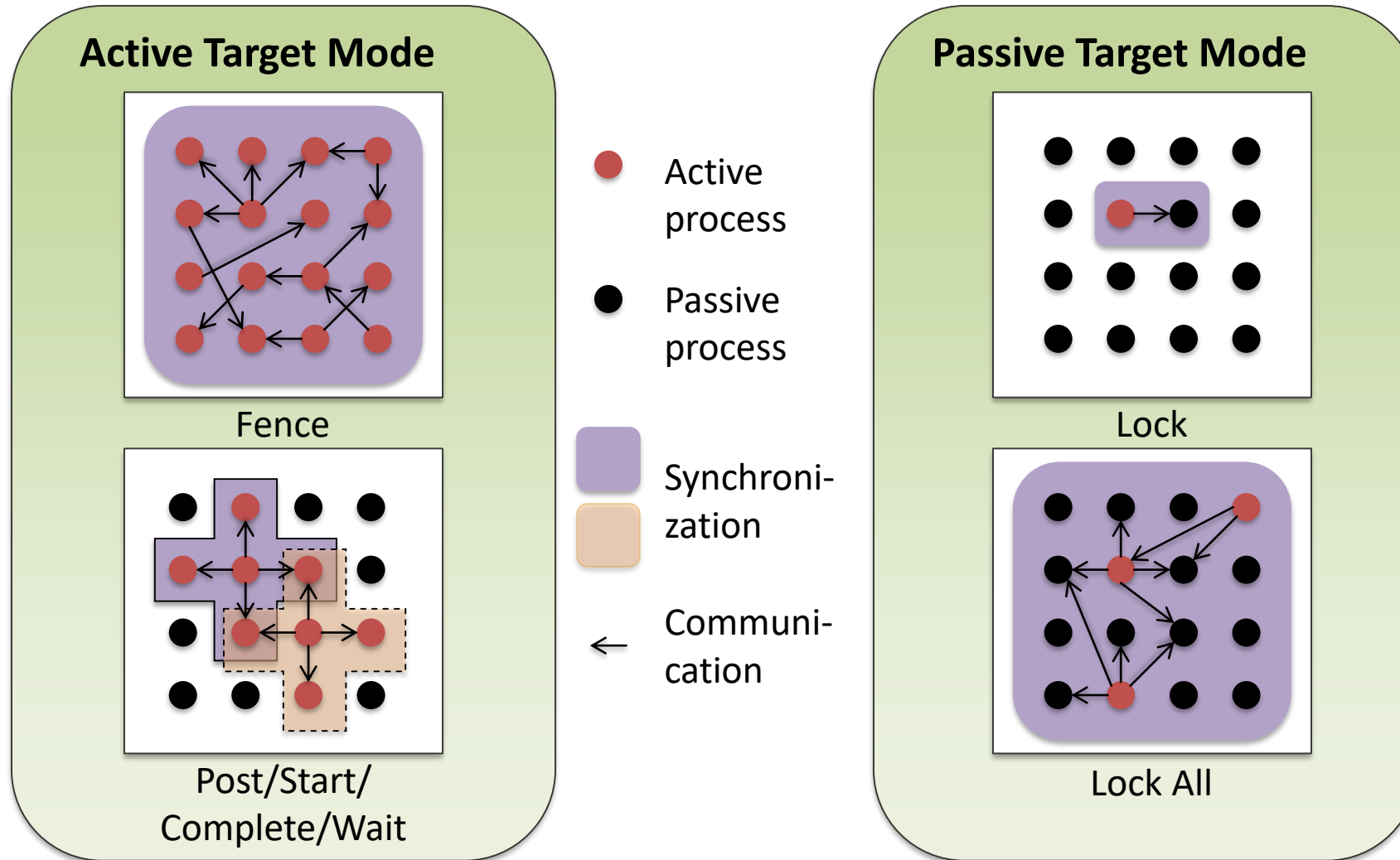


**hardware-
accelerated
protocol:**
lower latency

**fall back
protocol:**
*higher
bandwidth*

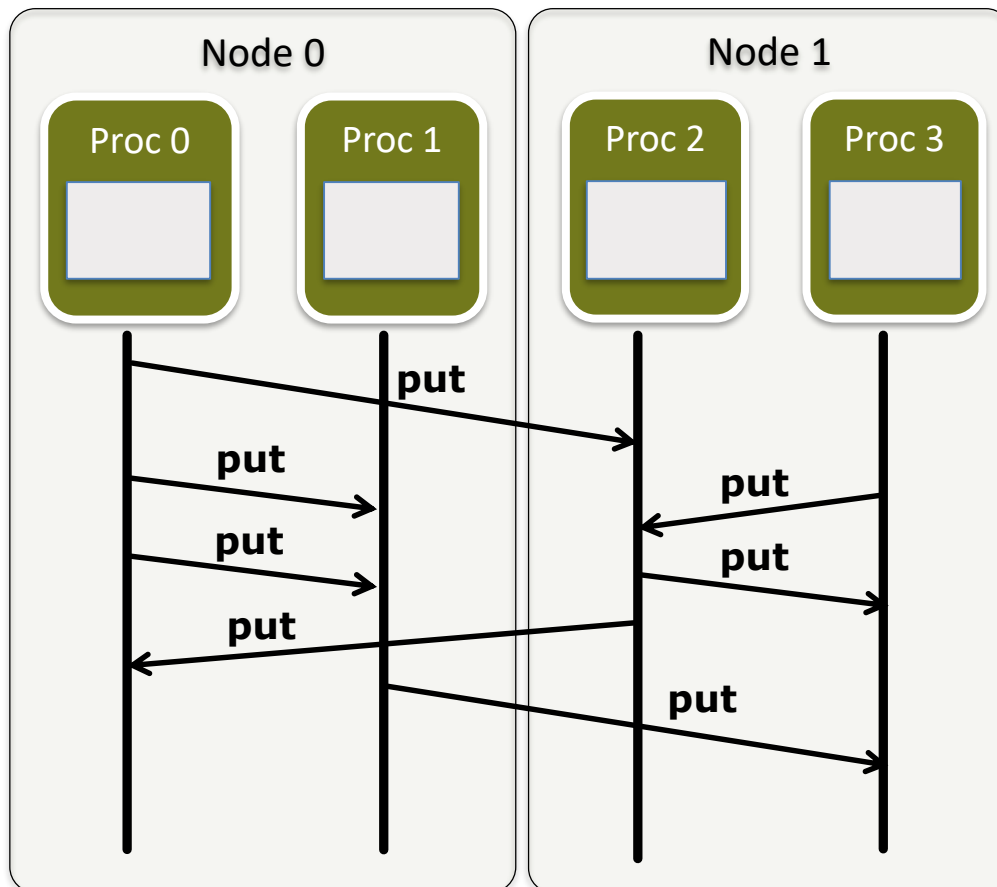
proprietary

PART 3: SYNCHRONIZATION



SCALABLE FENCE IMPLEMENTATION

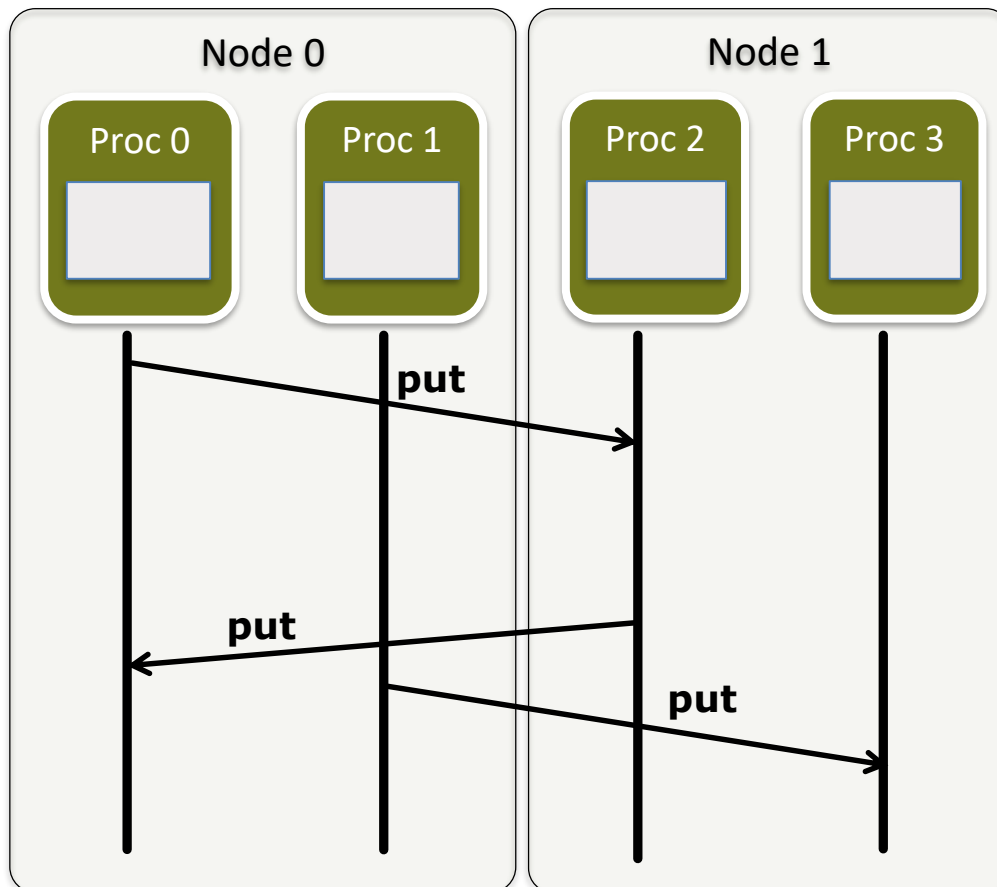
- Collective call
- Completes all outstanding memory operations



```
int MPI_Win_fence(...) {
    asm( mfence );
    dmapp_gsync_wait();
    MPI_Barrier(...);
    return MPI_SUCCESS;
}
```

SCALABLE FENCE IMPLEMENTATION

- Collective call
- Completes all outstanding memory operations

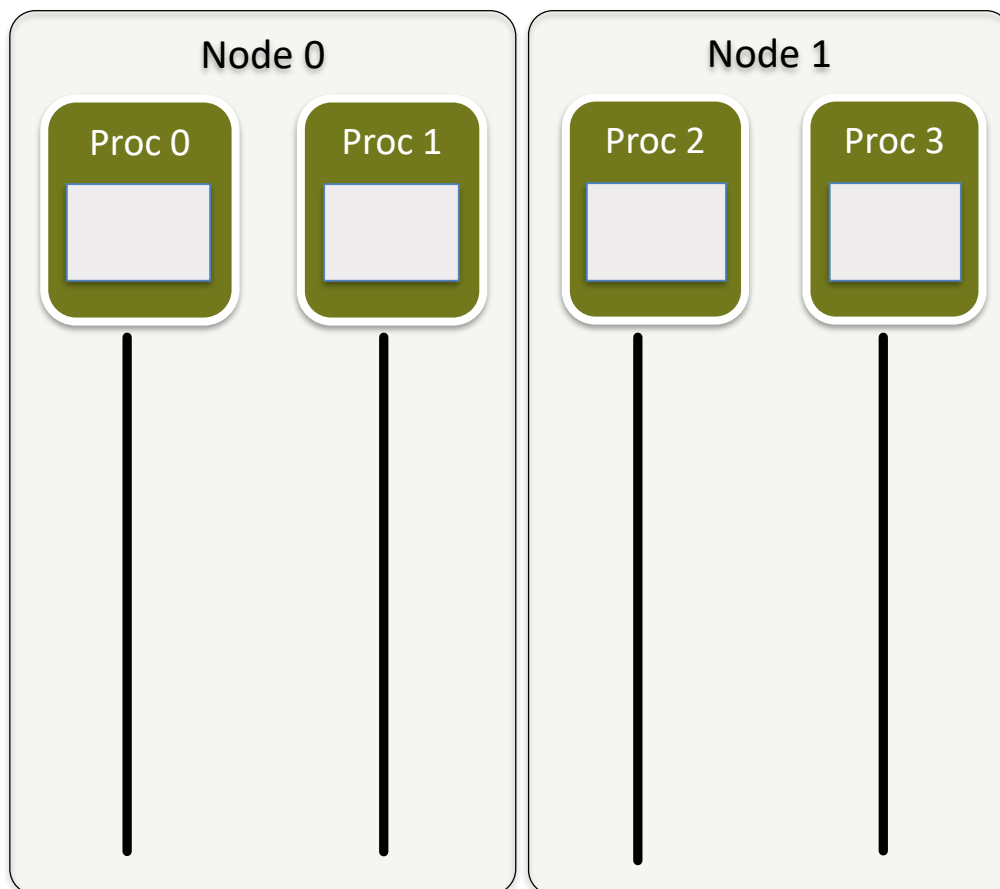


```
int MPI_win_fence(...) {
    asm( mfence );
    dmapp_gsync_wait();
    MPI_Barrier(...);
    return MPI_SUCCESS;
}
```

Local completion
(XPMEM)

SCALABLE FENCE IMPLEMENTATION

- Collective call
- Completes all outstanding memory operations

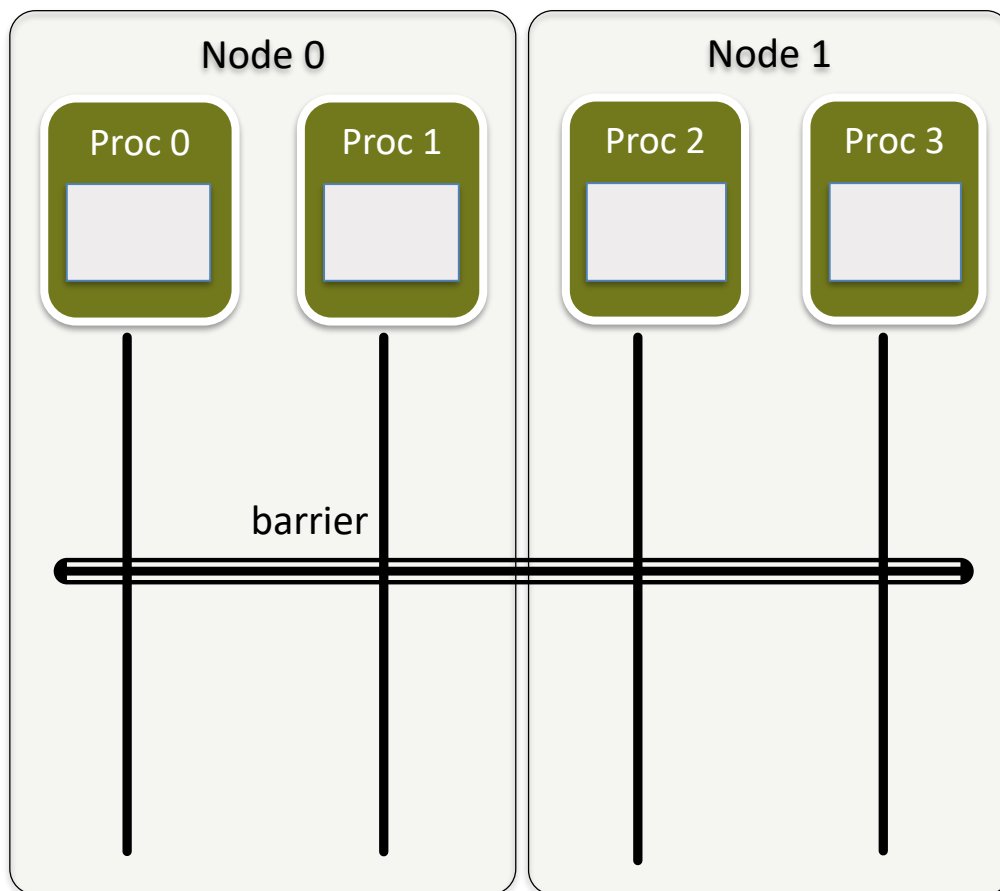


```
int MPI_win_fence(...) {
    asm( mfence );
    dmapp_gsync_wait();
    MPI_Barrier(...);
    return MPI_SUCCESS;
}
```

Local completion
(DMAPP)

SCALABLE FENCE IMPLEMENTATION

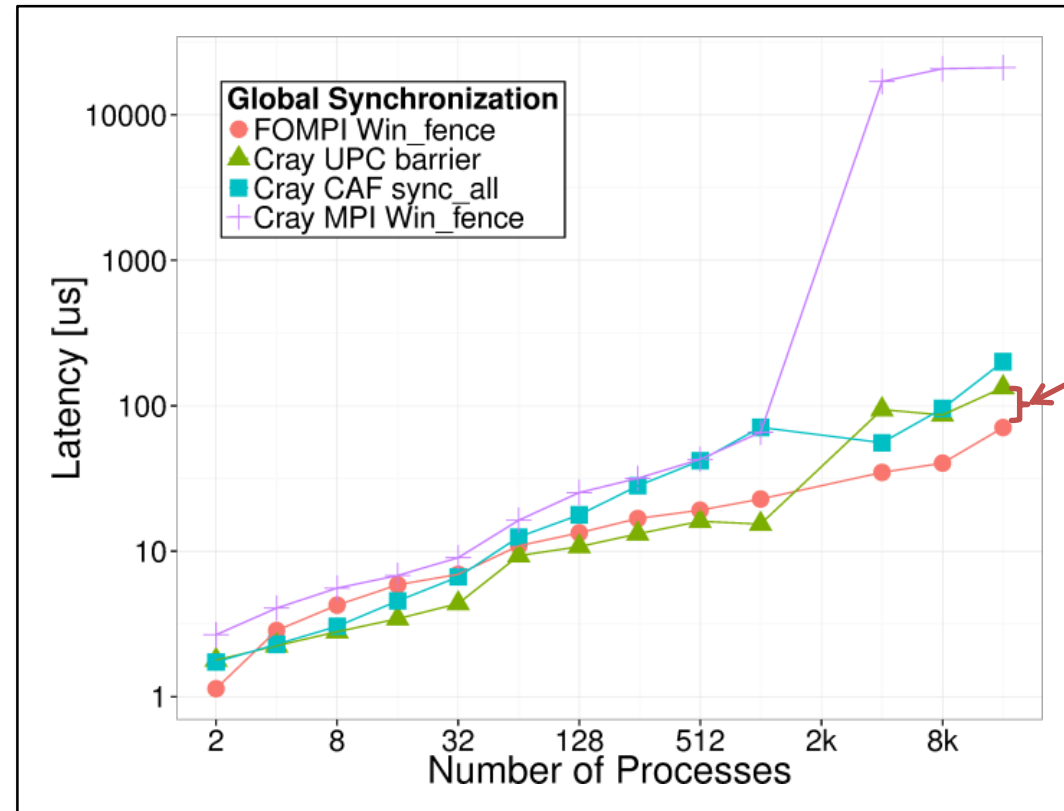
- Collective call
- Completes all outstanding memory operations



```
int MPI_Win_fence(...) {
    asm( mfence );
    dmapp_gsync_wait();
    MPI_Barrier(...);
    return MPI_SUCCESS;
}
```

Global
completion

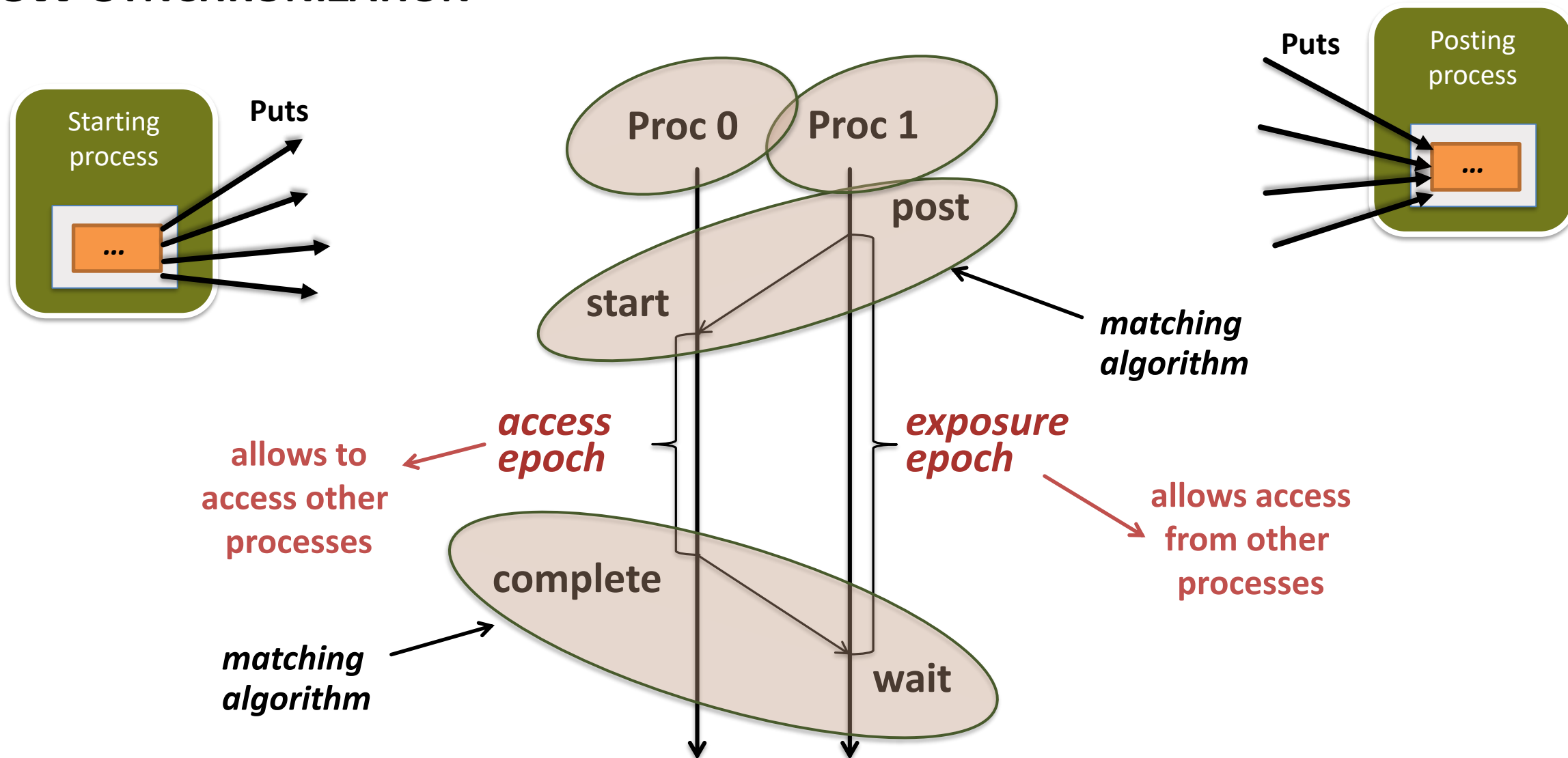
SCALABLE FENCE PERFORMANCE



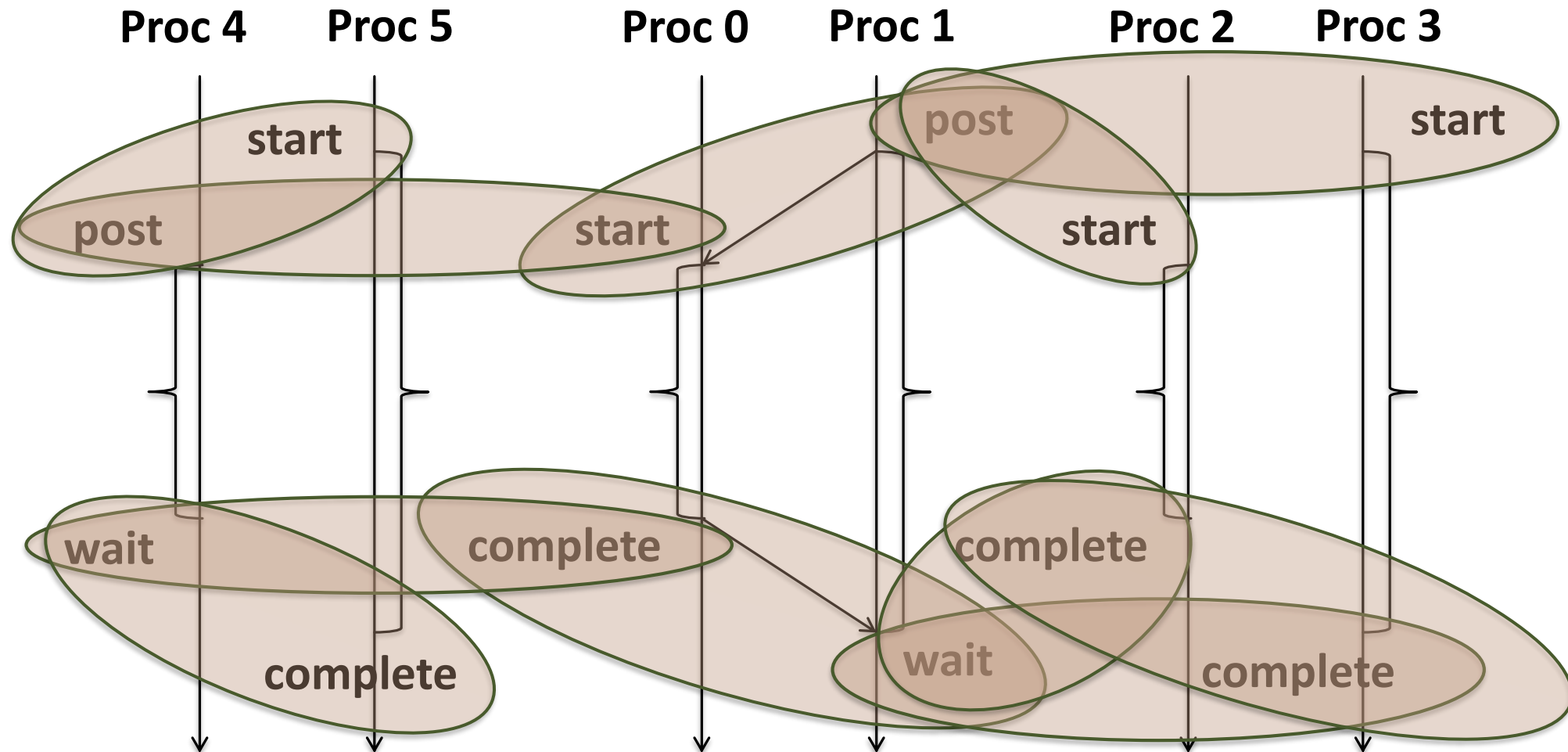
90%
faster

Time bound	$\mathcal{O}(\log p)$
Memory bound	$\mathcal{O}(1)$

PSCW SYNCHRONIZATION

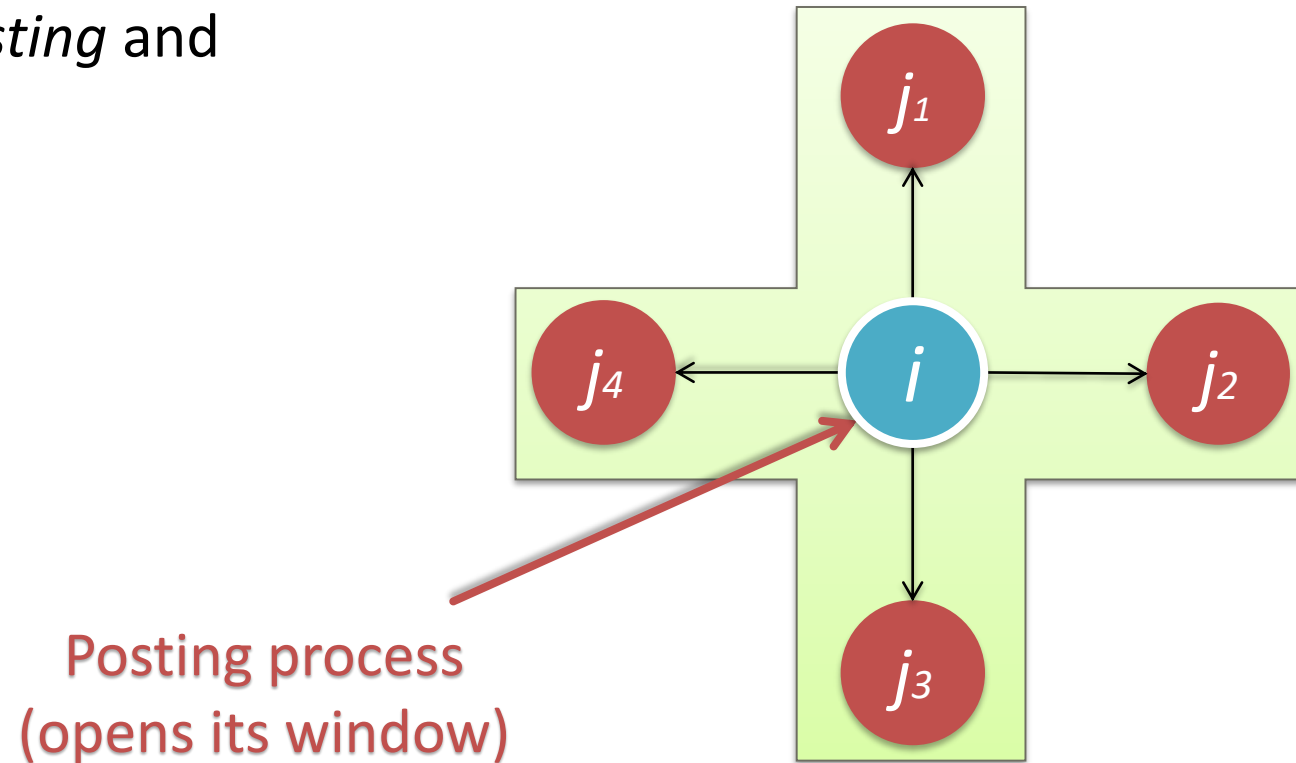


PSCW SYNCHRONIZATION



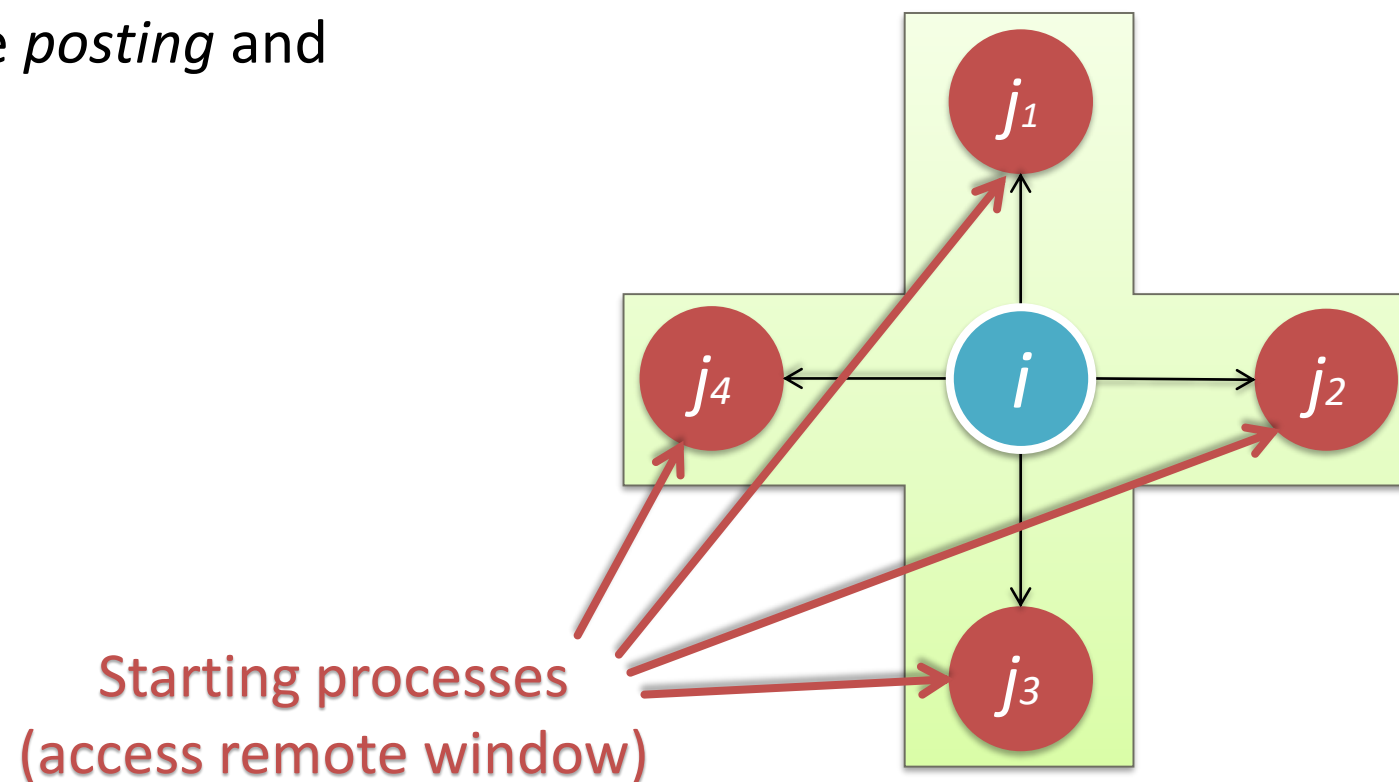
PSCW SCALABLE POST/START MATCHING

- In general, there can be n *posting* and m *starting* processes
- In this example there is one *posting* and 4 *starting* processes



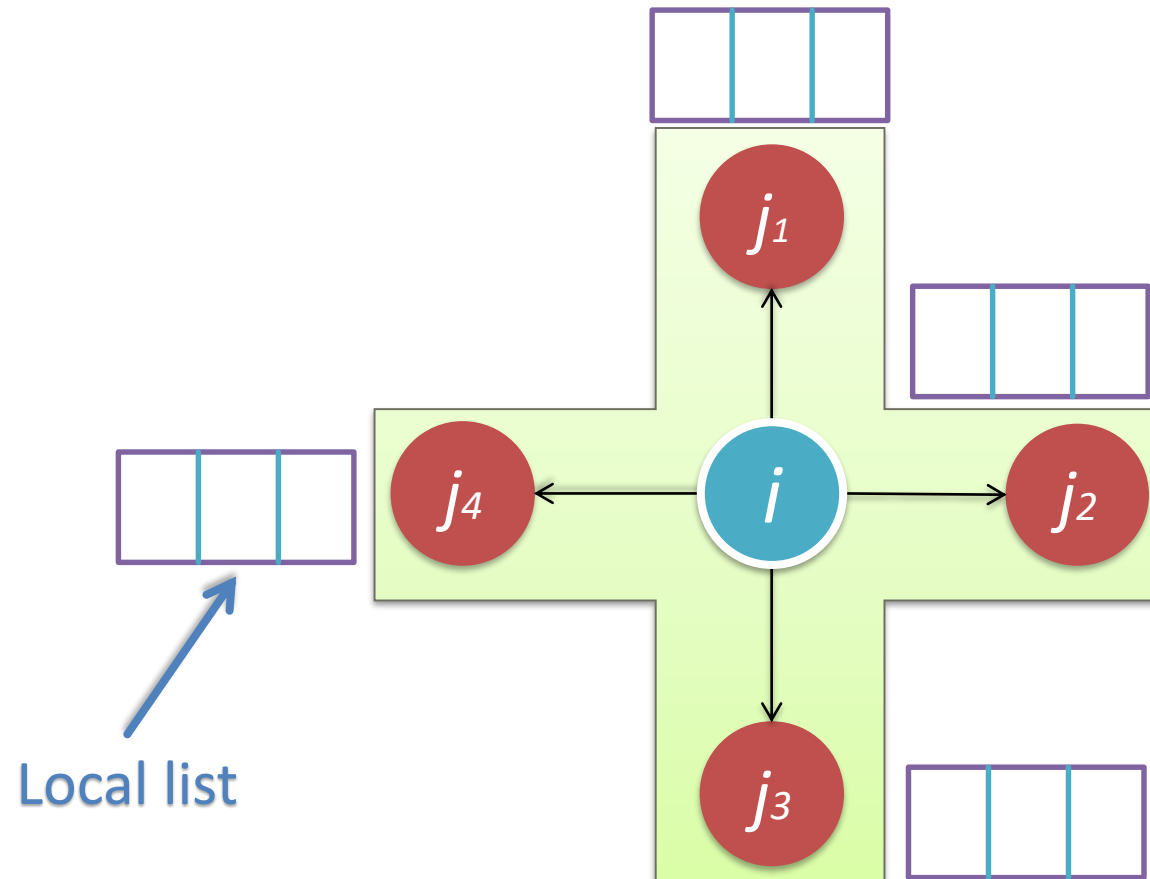
PSCW SCALABLE POST/START MATCHING

- In general, there can be n *posting* and m *starting* processes
- In this example there is one *posting* and 4 *starting* processes



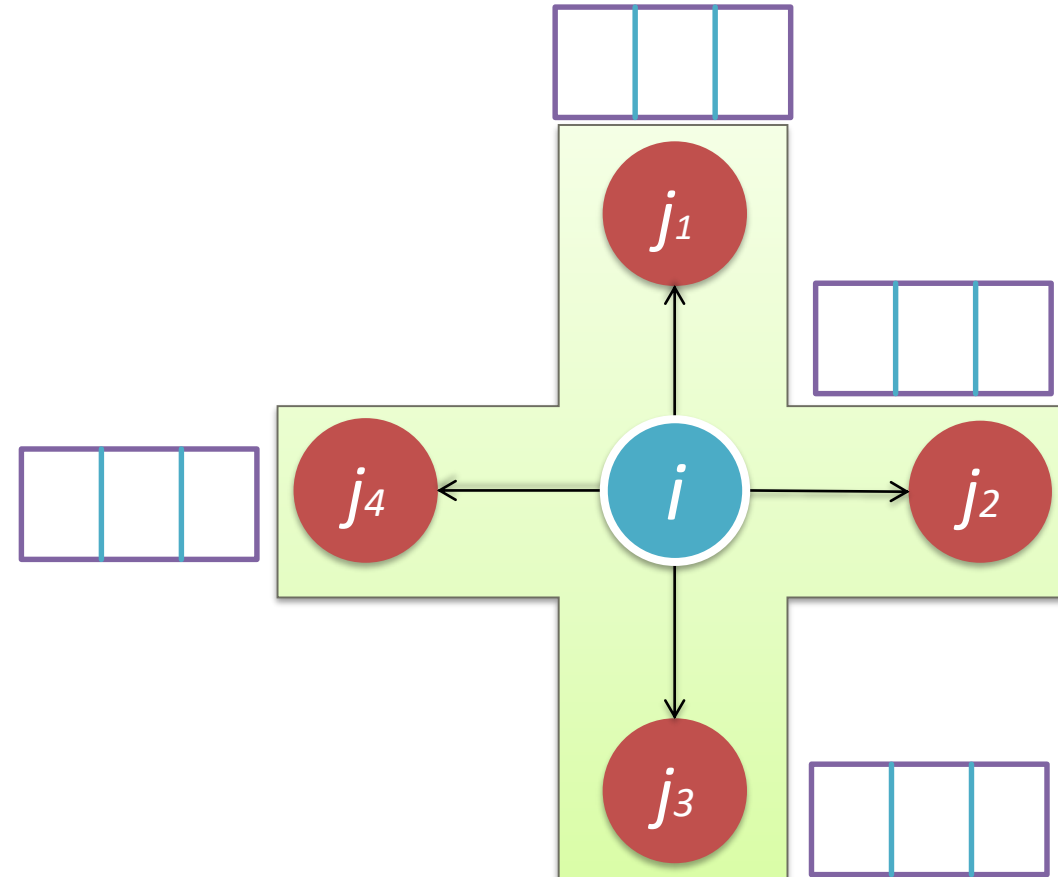
PSCW SCALABLE POST/START MATCHING

- Each starting process has a local list



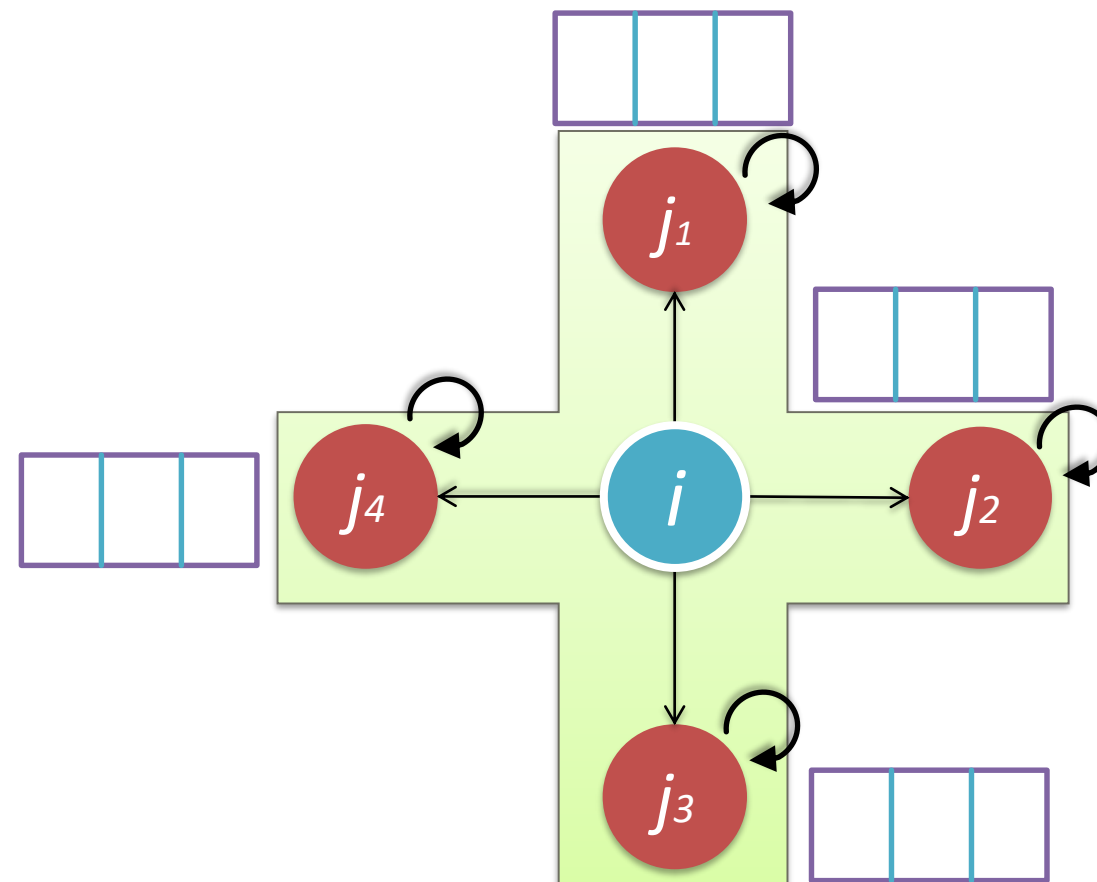
PSCW SCALABLE POST/START MATCHING

- *Posting* process i adds its rank i to a list at each *starting* process j_1, \dots, j_4
- Each *starting* process j waits until the rank of the *posting* process i is present in its local list



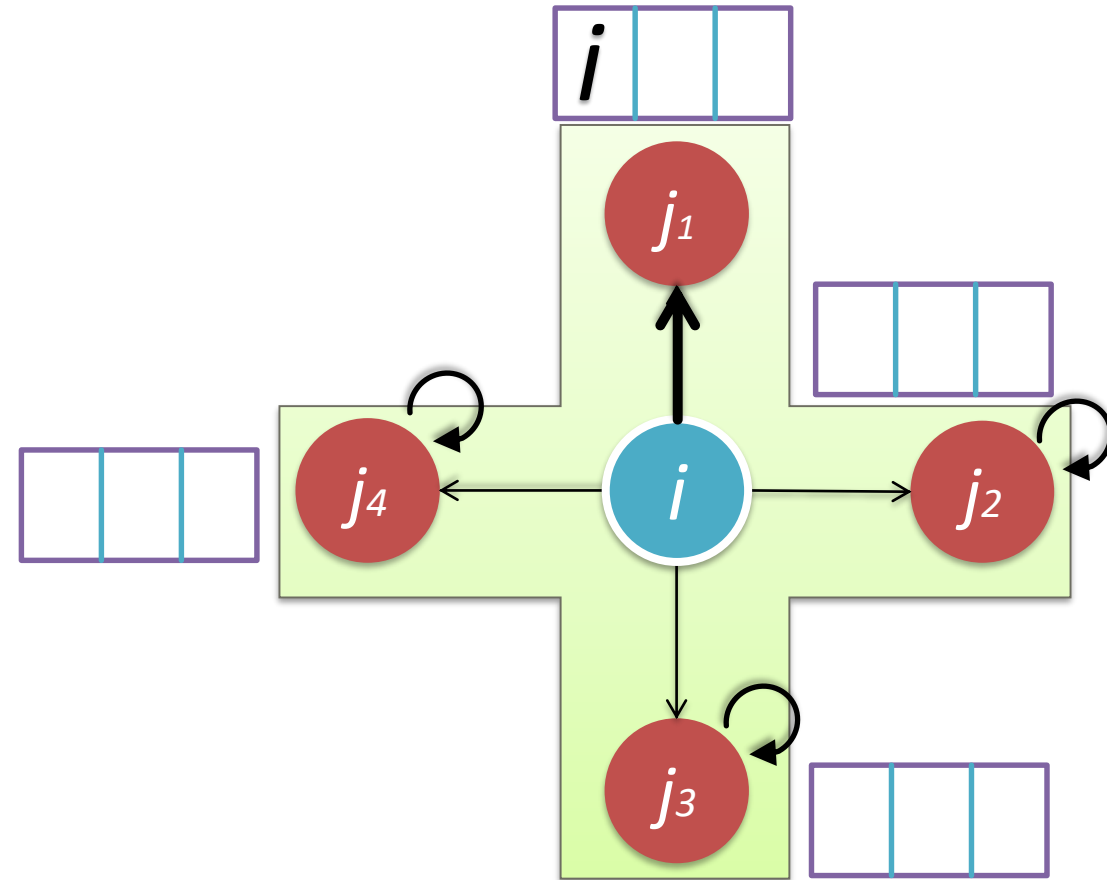
PSCW SCALABLE POST/START MATCHING

- *Posting* process i adds its rank i to a list at each *starting* process j_1, \dots, j_4
- Each *starting* process j waits until the rank of the *posting* process i is present in its local list



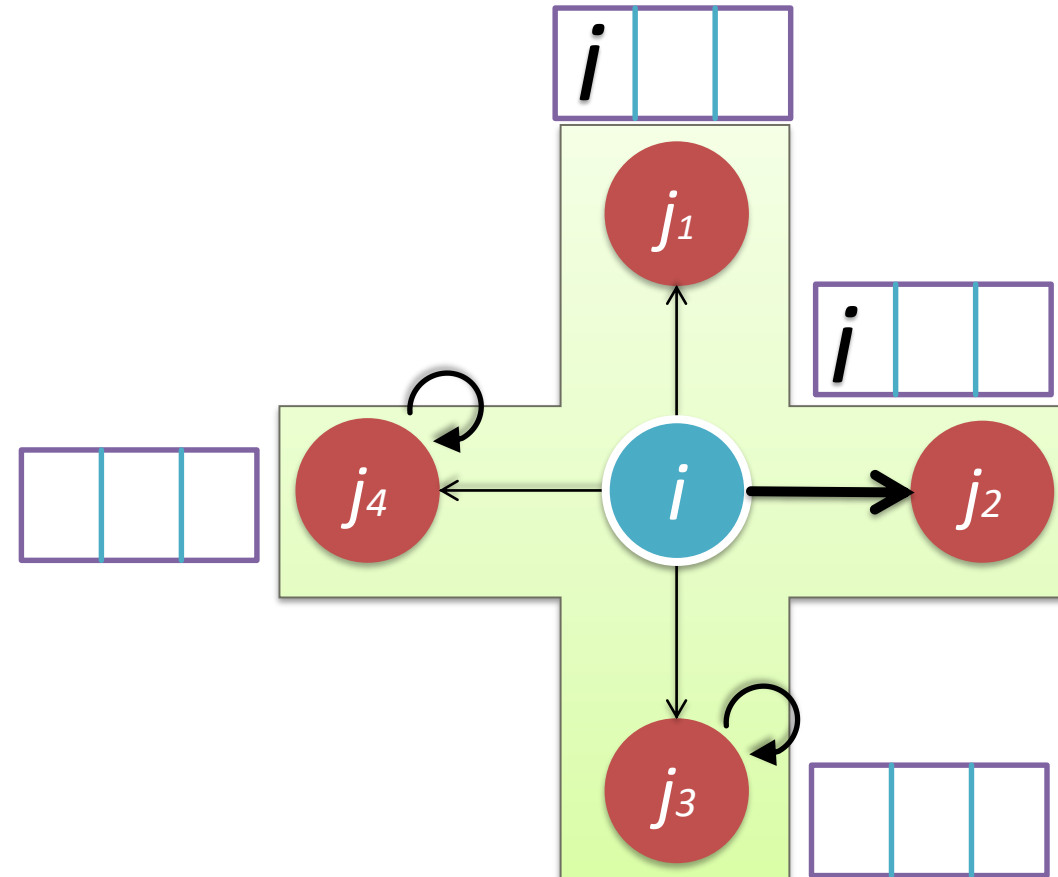
PSCW SCALABLE POST/START MATCHING

- *Posting* process i adds its rank i to a list at each *starting* process j_1, \dots, j_4
- Each *starting* process j waits until the rank of the *posting* process i is present in its local list



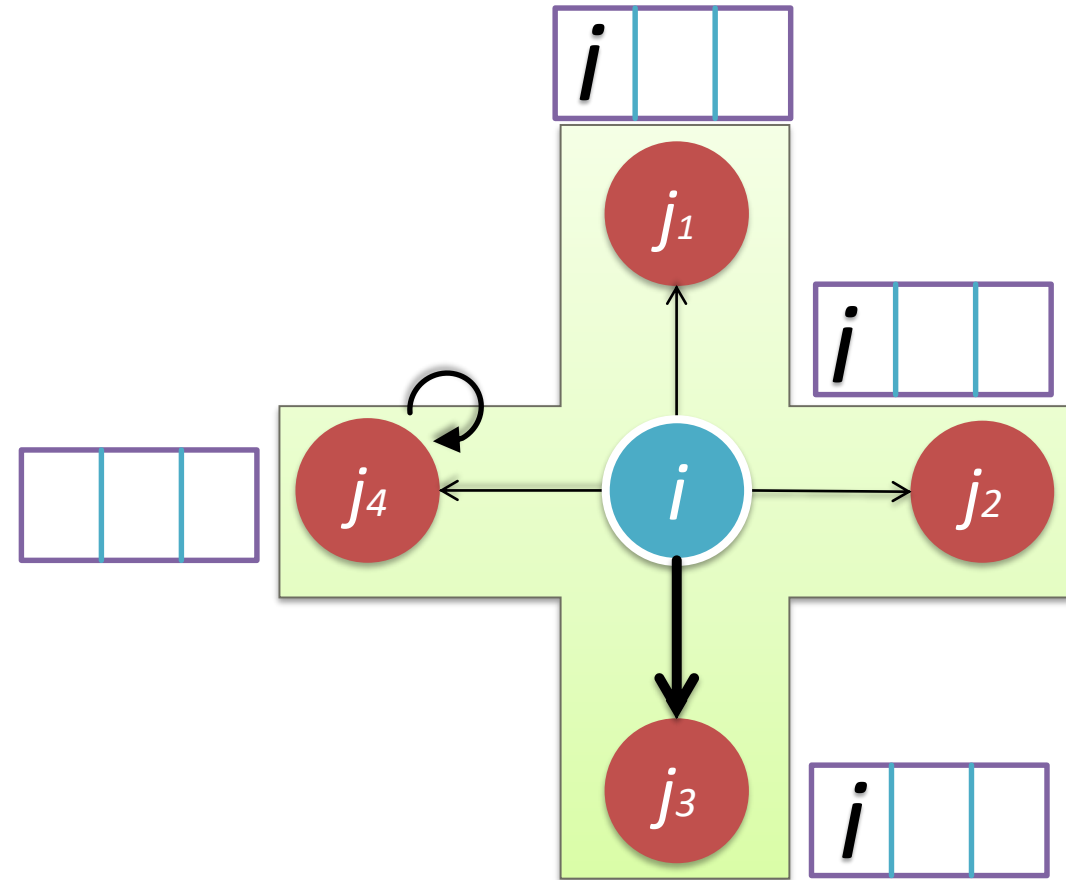
PSCW SCALABLE POST/START MATCHING

- *Posting* process i adds its rank i to a list at each *starting* process j_1, \dots, j_4
- Each *starting* process j waits until the rank of the *posting* process i is present in its local list



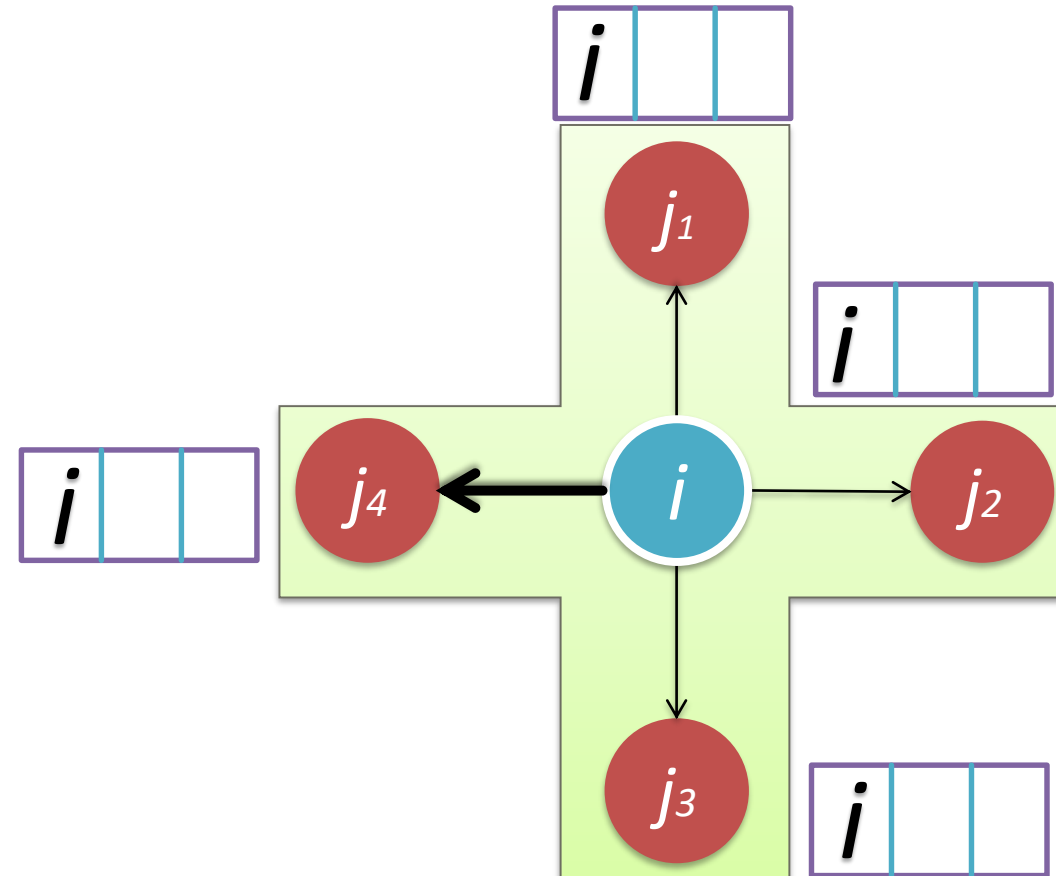
PSCW SCALABLE POST/START MATCHING

- *Posting* process i adds its rank i to a list at each *starting* process j_1, \dots, j_4
- Each *starting* process j waits until the rank of the *posting* process i is present in its local list



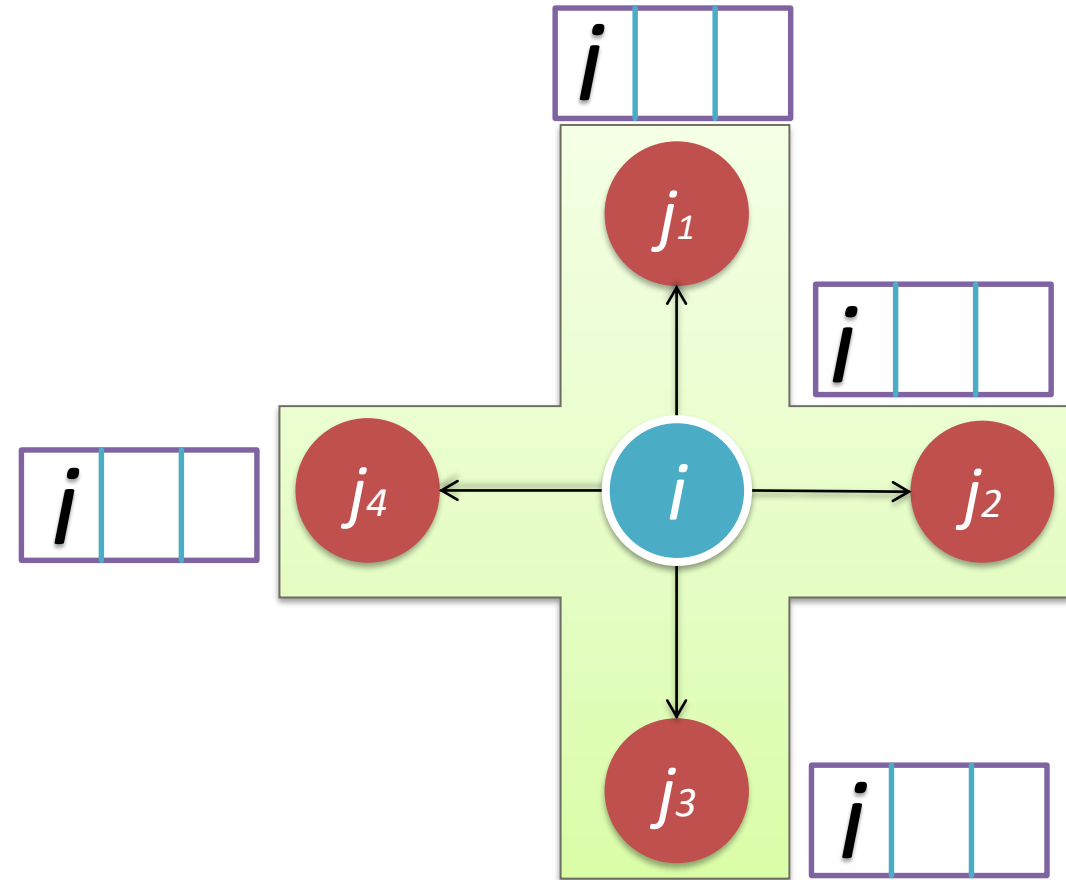
PSCW SCALABLE POST/START MATCHING

- *Posting* process i adds its rank i to a list at each *starting* process j_1, \dots, j_4
- Each *starting* process j waits until the rank of the *posting* process i is present in its local list



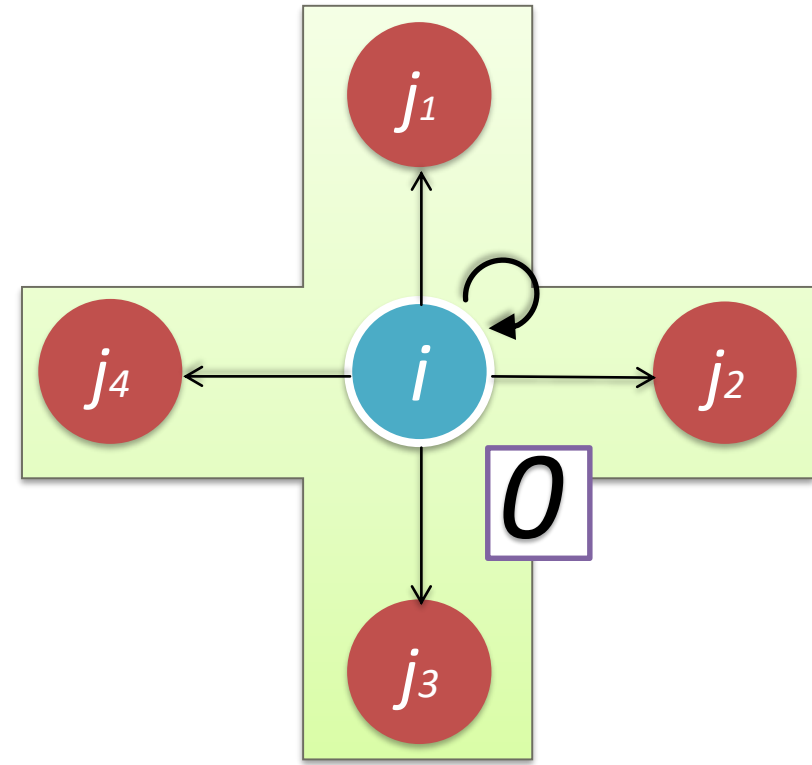
PSCW SCALABLE POST/START MATCHING

- *Posting* process i adds its rank i to a list at each *starting* process j_1, \dots, j_4
- Each *starting* process j waits until the rank of the *posting* process i is present in its local list



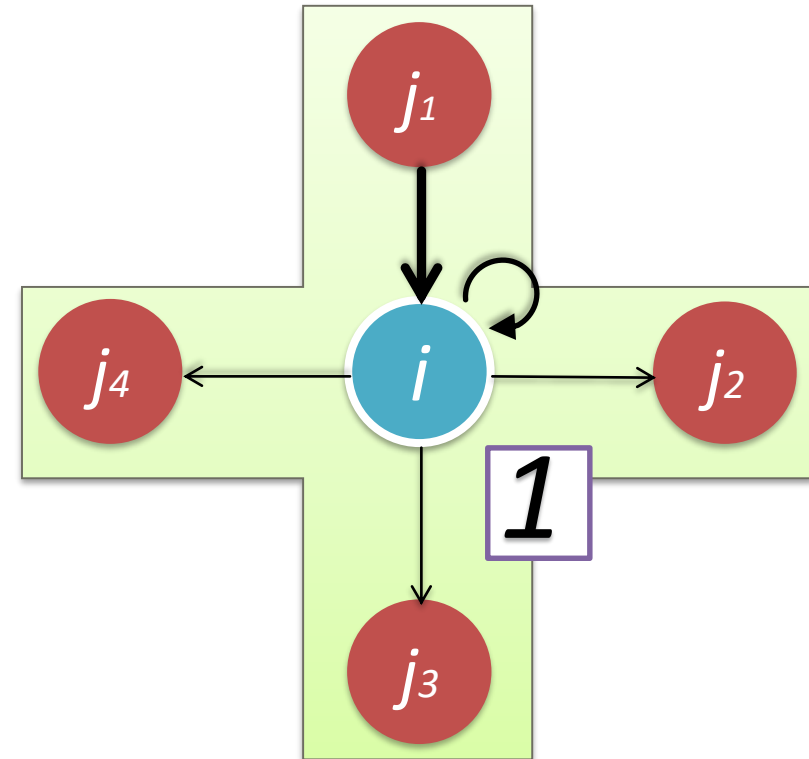
PSCW SCALABLE COMPLETE/WAIT MATCHING

- Each *completing* process increments a counter stored at the *posting* process



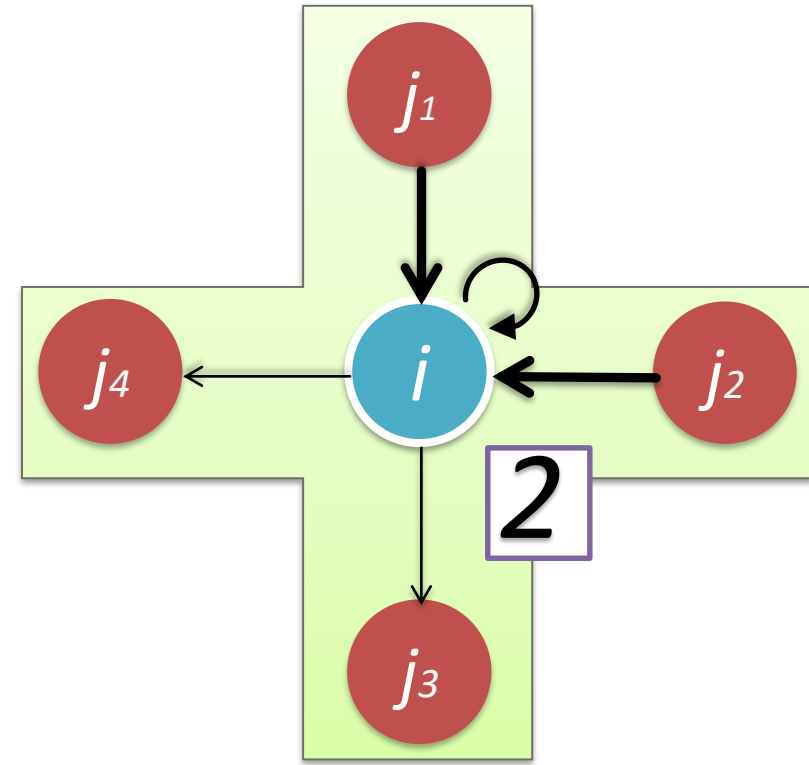
PSCW SCALABLE COMPLETE/WAIT MATCHING

- Each *completing* process increments a counter stored at the *posting* process



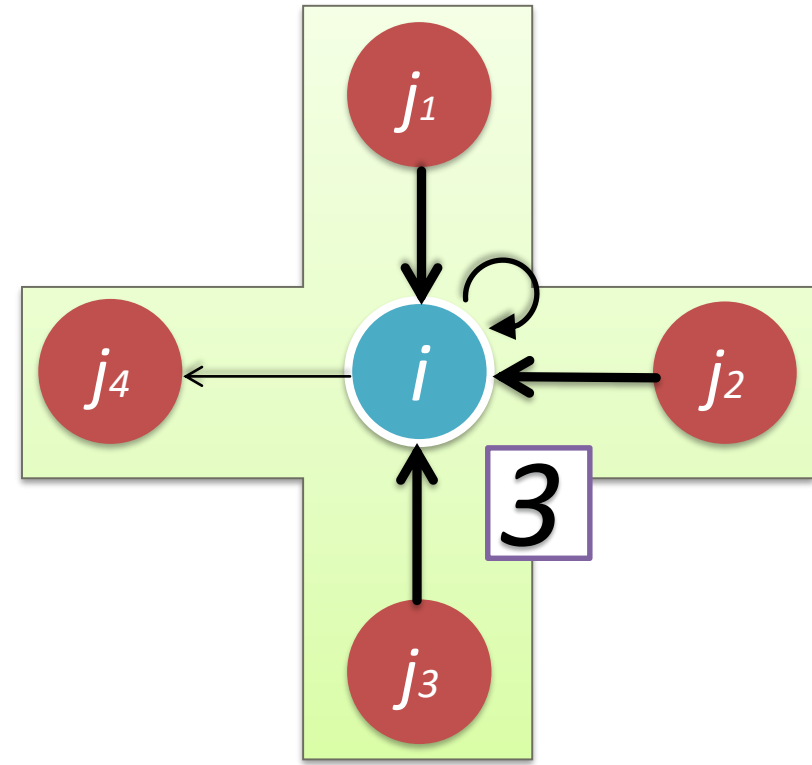
PSCW SCALABLE COMPLETE/WAIT MATCHING

- Each *completing* process increments a counter stored at the *posting* process



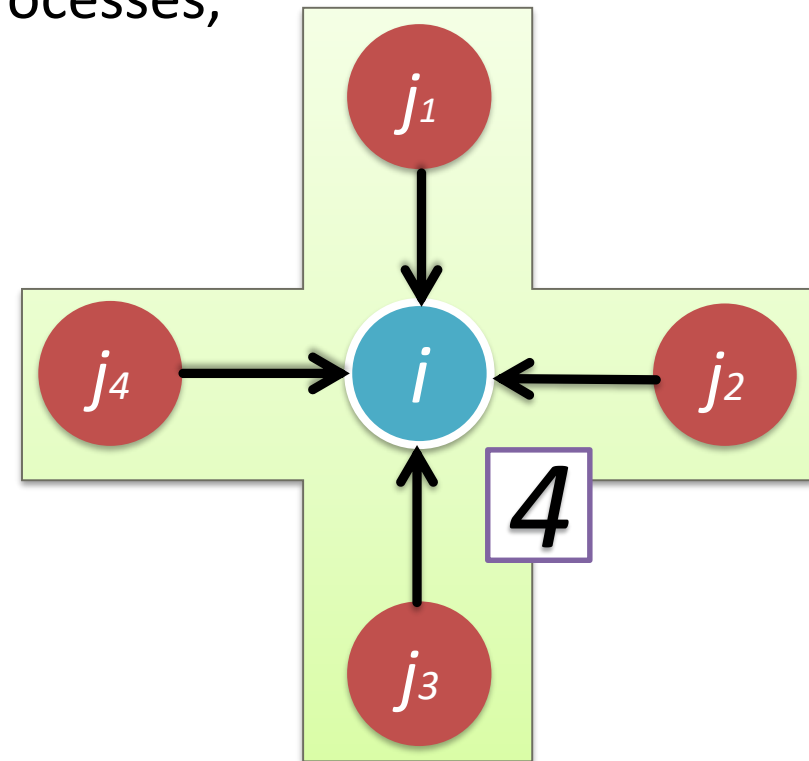
PSCW SCALABLE COMPLETE/WAIT MATCHING

- Each *completing* process increments a counter stored at the *posting* process



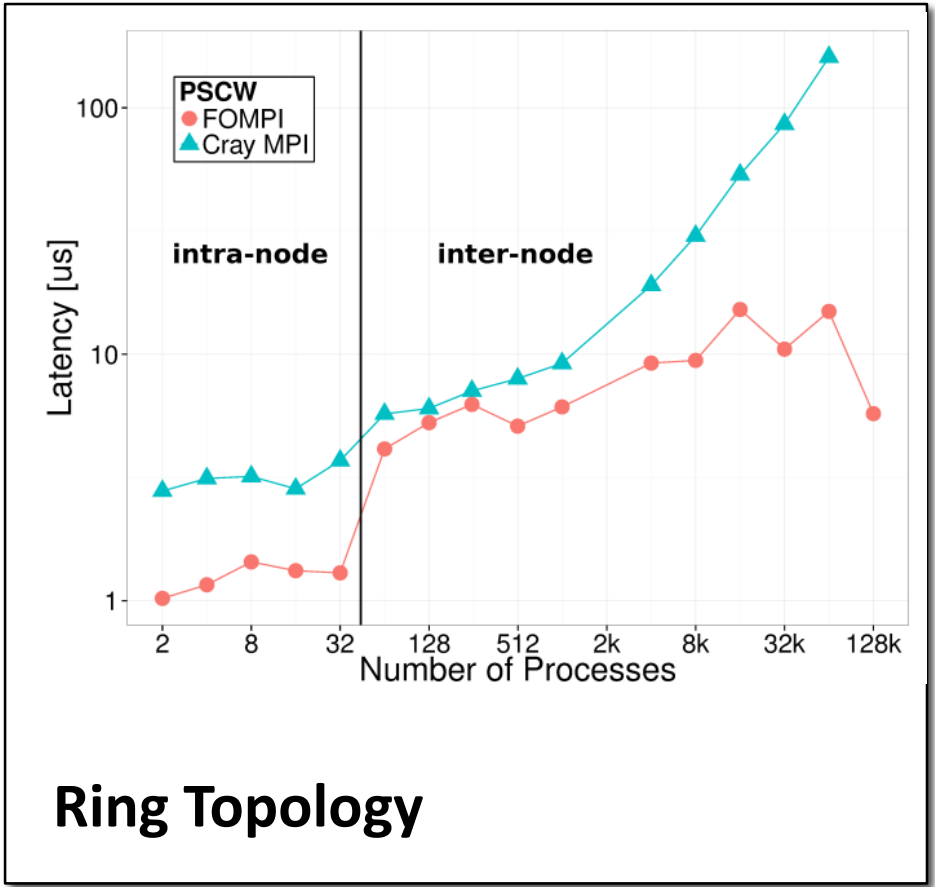
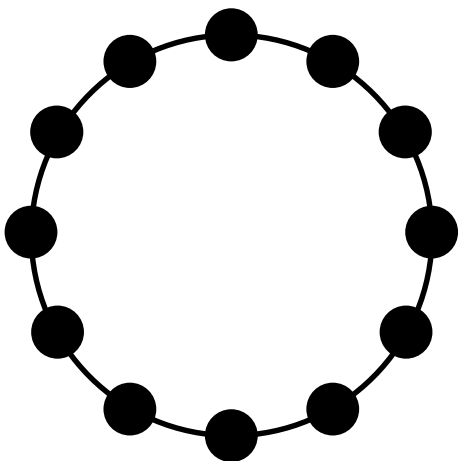
PSCW SCALABLE COMPLETE/WAIT MATCHING

- Each *completing* process increments a counter stored at the *posting* process
- When the counter is equal to the number of *starting* processes, the *posting* process returns from wait



PSCW PERFORMANCE

Time bound
$\mathcal{P}_{start} = \mathcal{P}_{wait} = \mathcal{O}(1)$ $\mathcal{P}_{post} = \mathcal{P}_{complete} = \mathcal{O}(\log p)$
Memory bound
$\mathcal{O}(\log p)$ (for scalable programs)



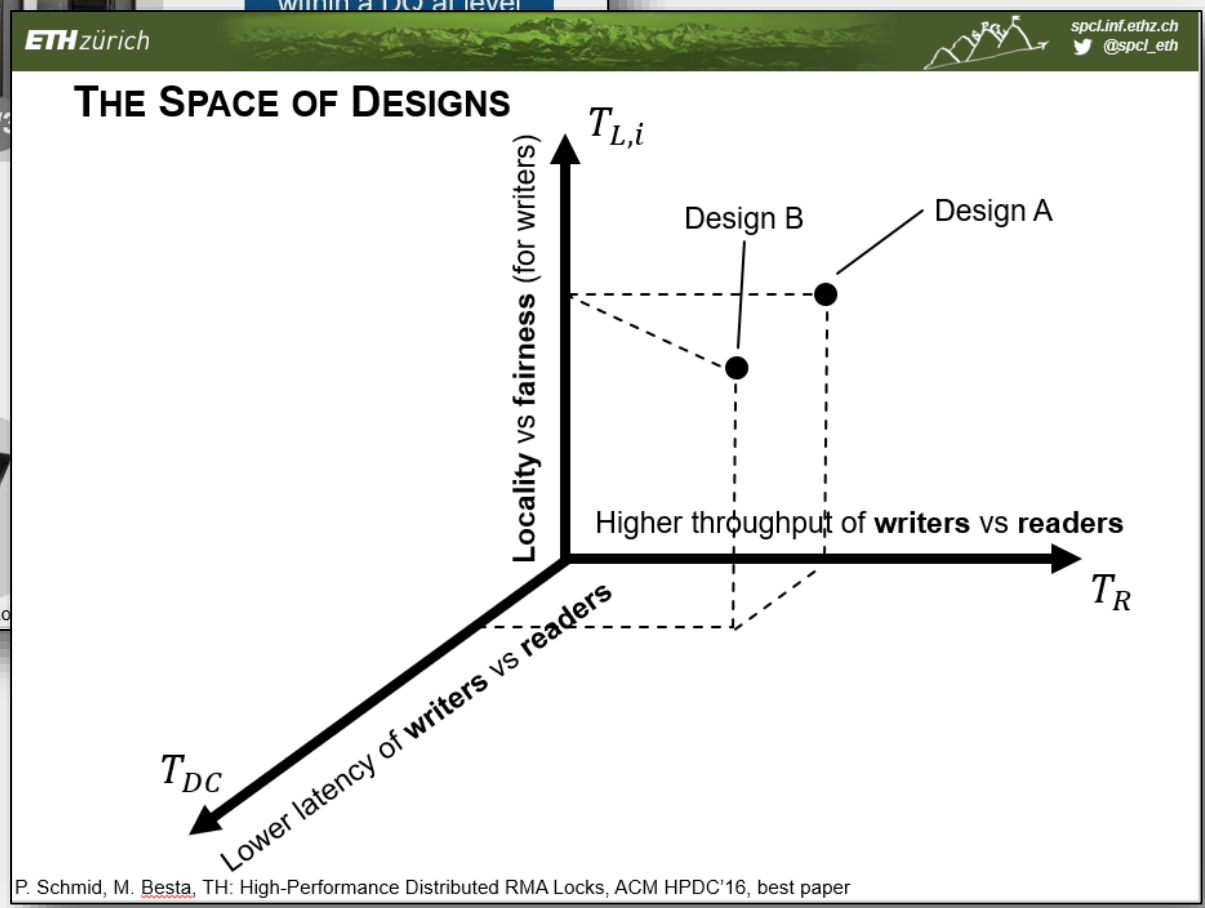
SCALABLE LOCK SYNCHRONIZATION

DISTRIBUTED MCS QUEUES (DQs)
Throughput vs Fairness

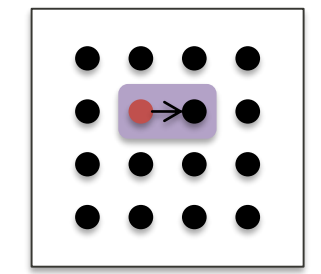
Each DQ: The maximum number of lock passings within a DQ at level

Larger $T_{L,i}$: more throughput at level i .
Smaller $T_{L,i}$: more fairness at level i .

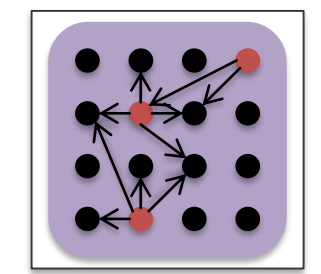
P. Schmid, M. Besta, TH: High-Performance Distributed RMA Locks



- Active process
- Passive process



Lock/Unlock
(shared/exclusive)

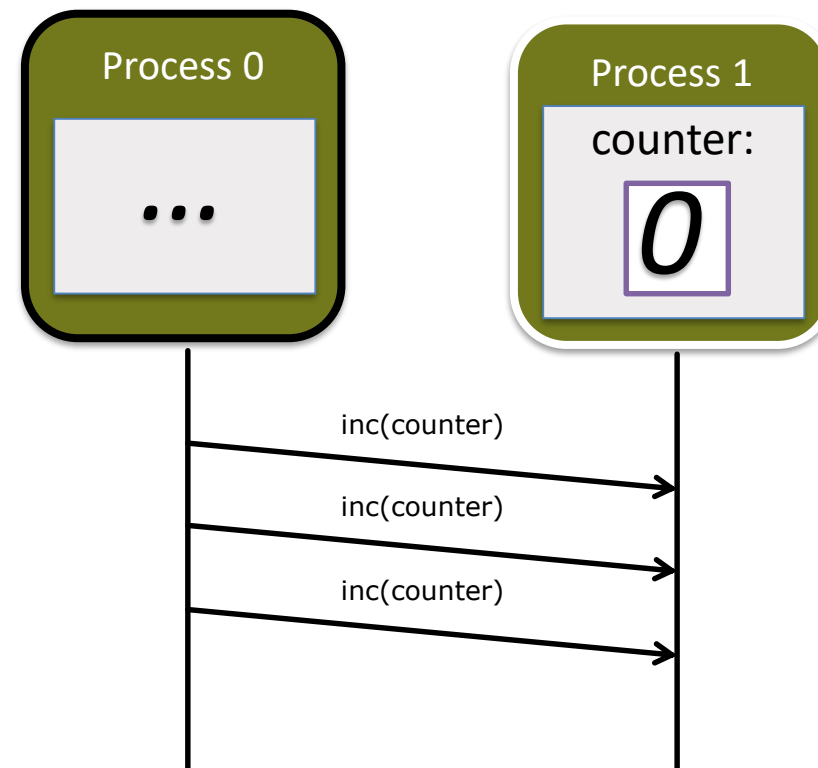


Lock All
(always shared)

FLUSH SYNCHRONIZATION

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

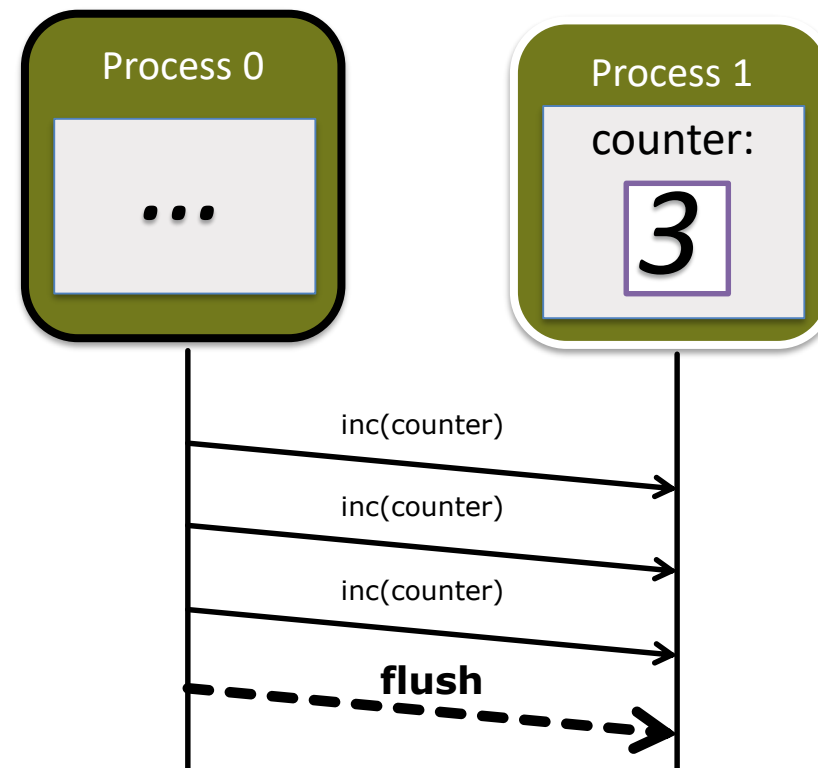
- Guarantees remote completion
- Issues a remote bulk synchronization and an x86 mfence
- One of the most performance critical functions, we add only 78 x86 CPU instructions to the critical path



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

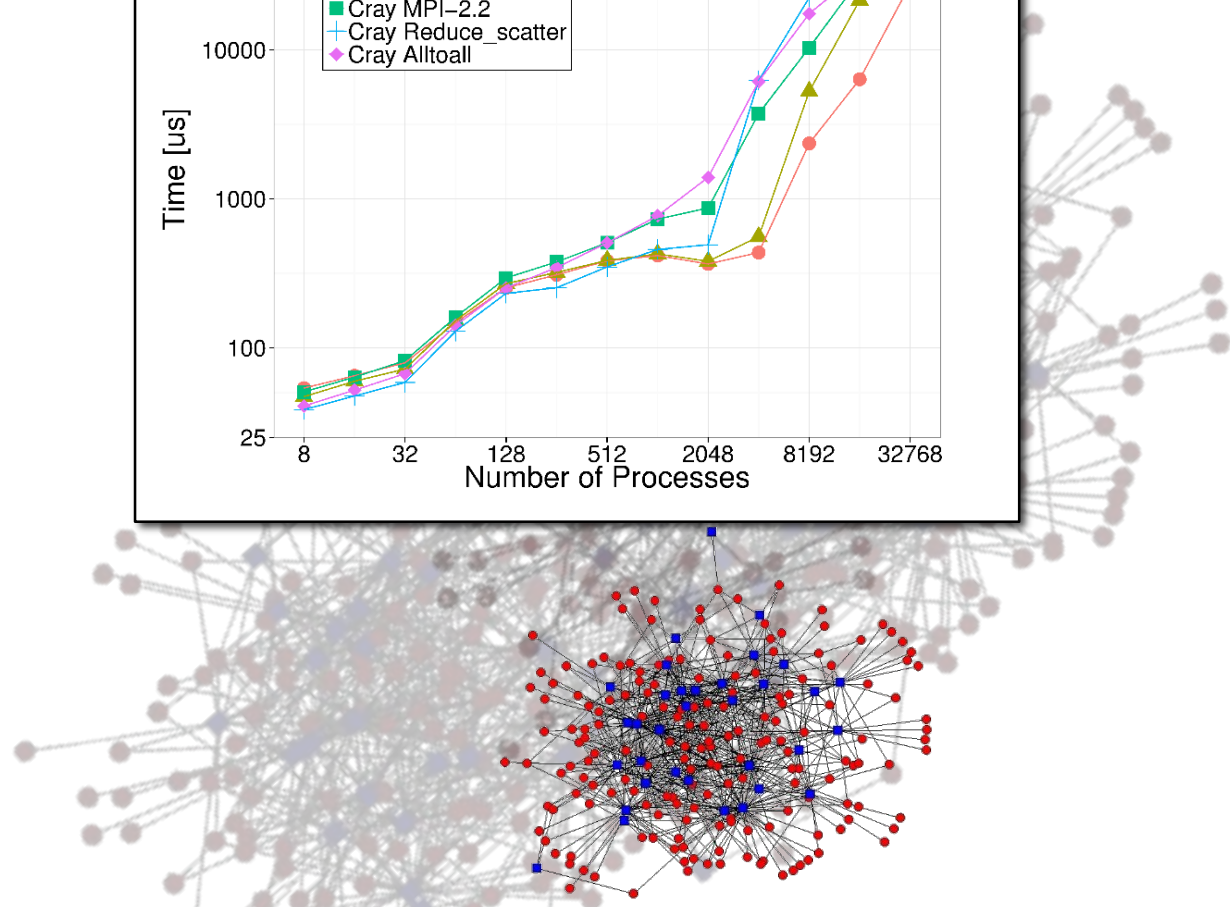
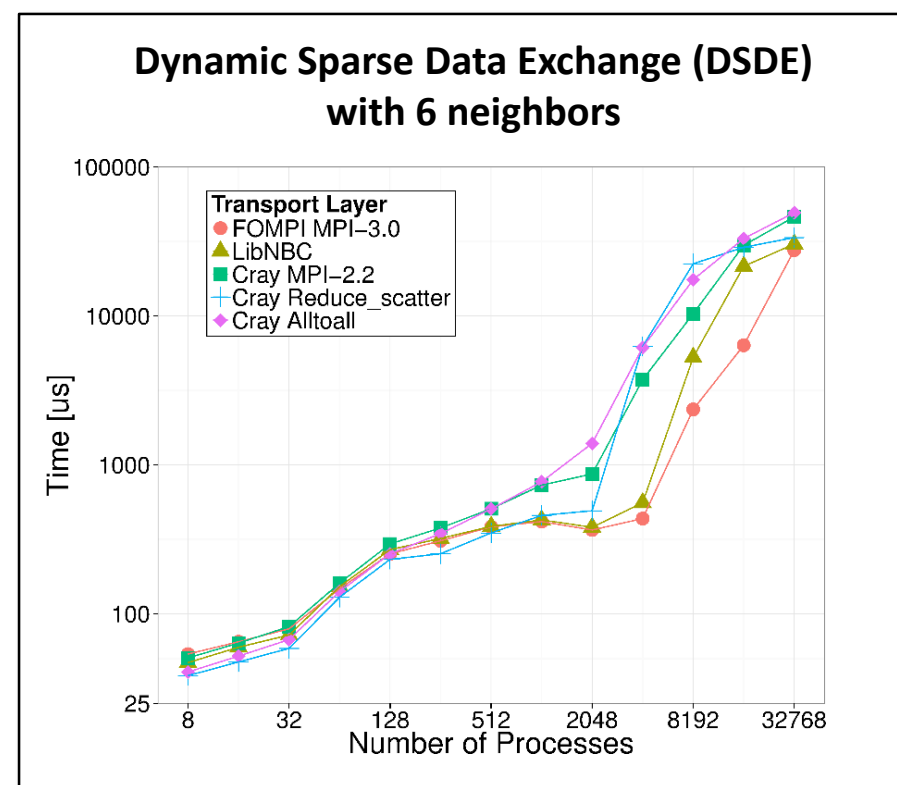
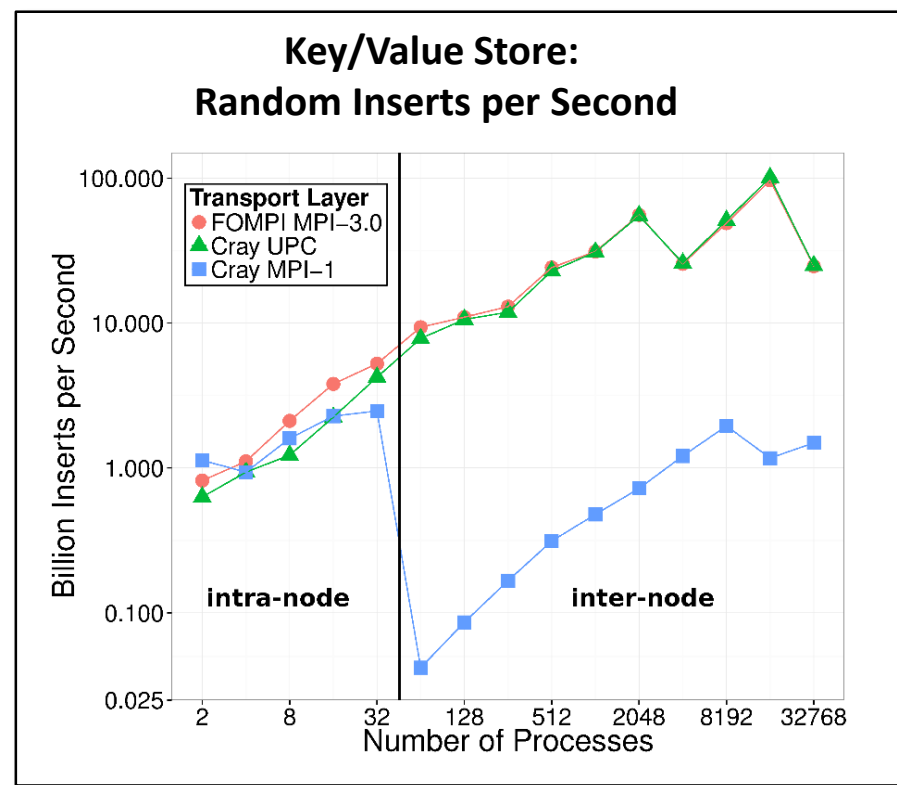


PERFORMANCE

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

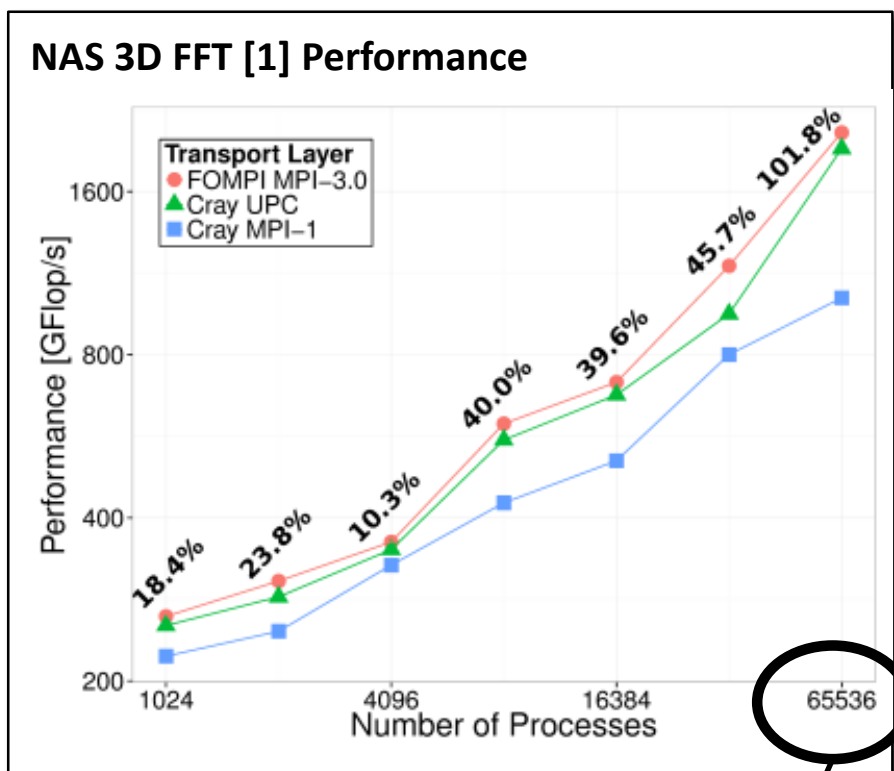
BLUE WATERS
SUSTAINED PETASCALE COMPUTING

PERFORMANCE: MOTIF APPLICATIONS

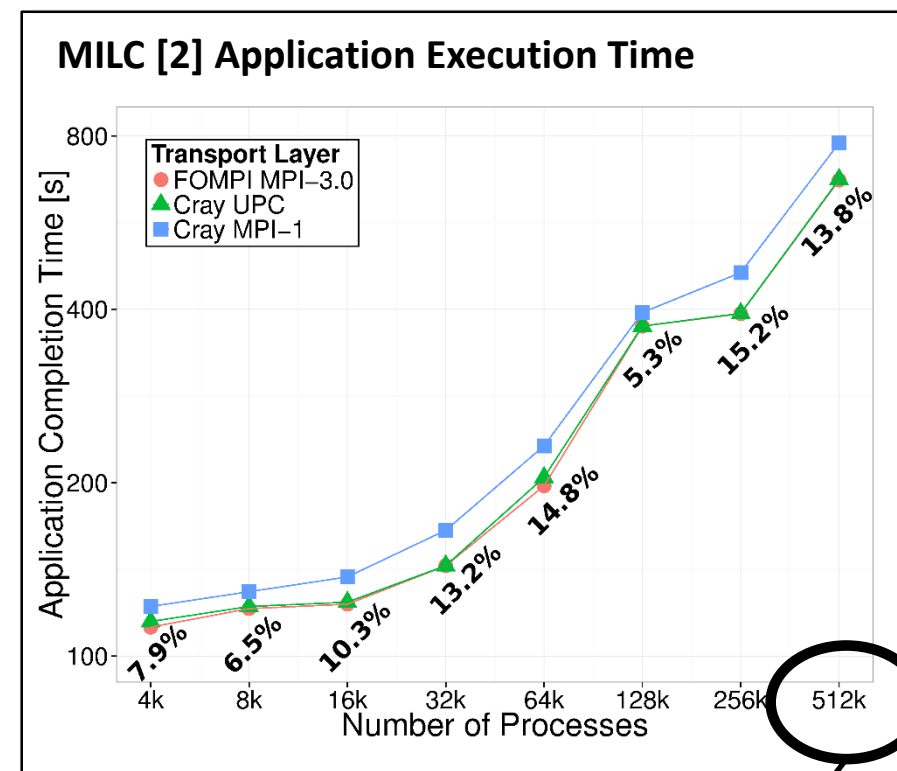


PERFORMANCE: APPLICATIONS

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



scale
to 65k procs

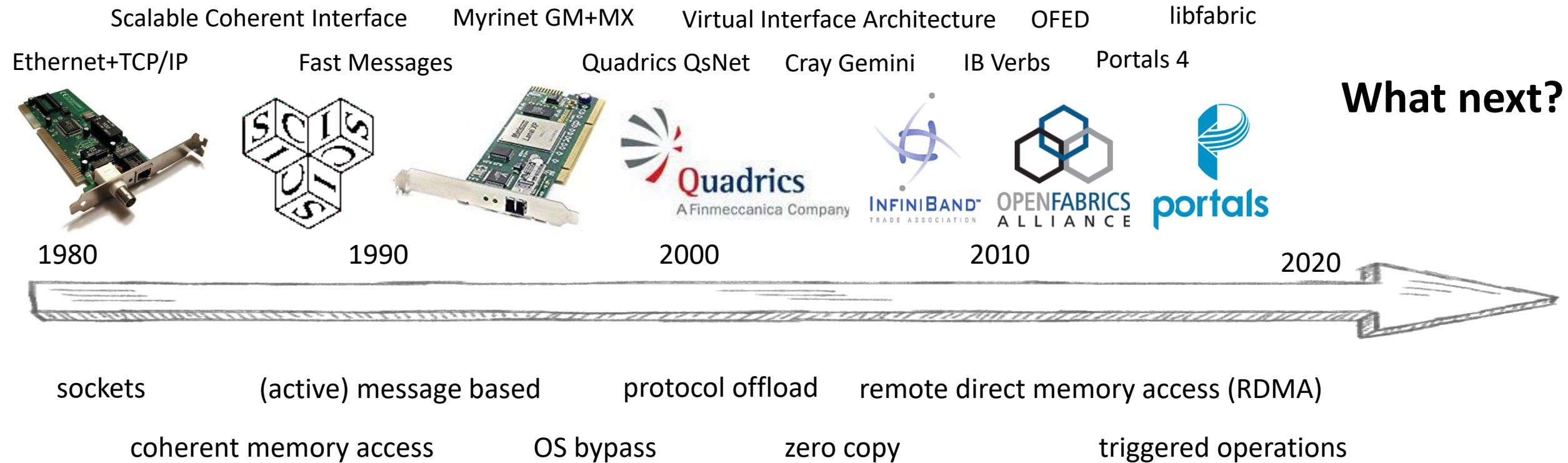


scale
to 512k procs

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

[2] Shan et al. Accelerating applications at scale using one-sided communication. PGAS'12

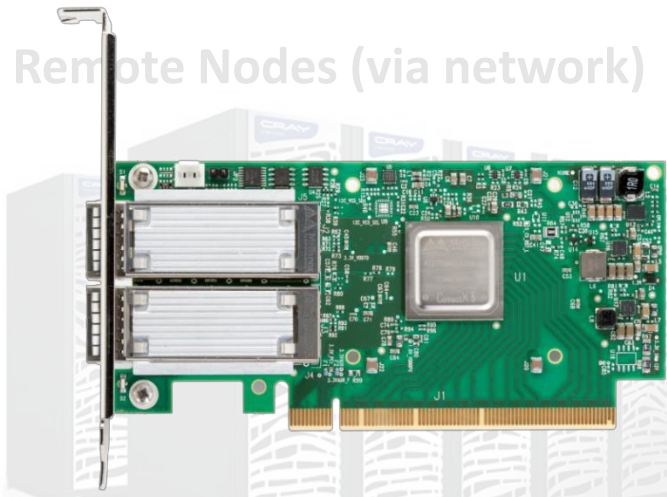
The Development of High-Performance Networking Interfaces



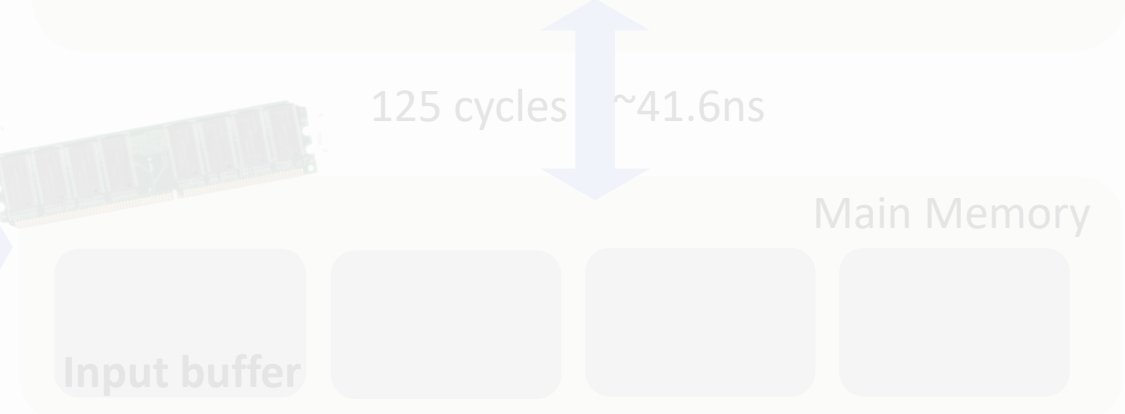
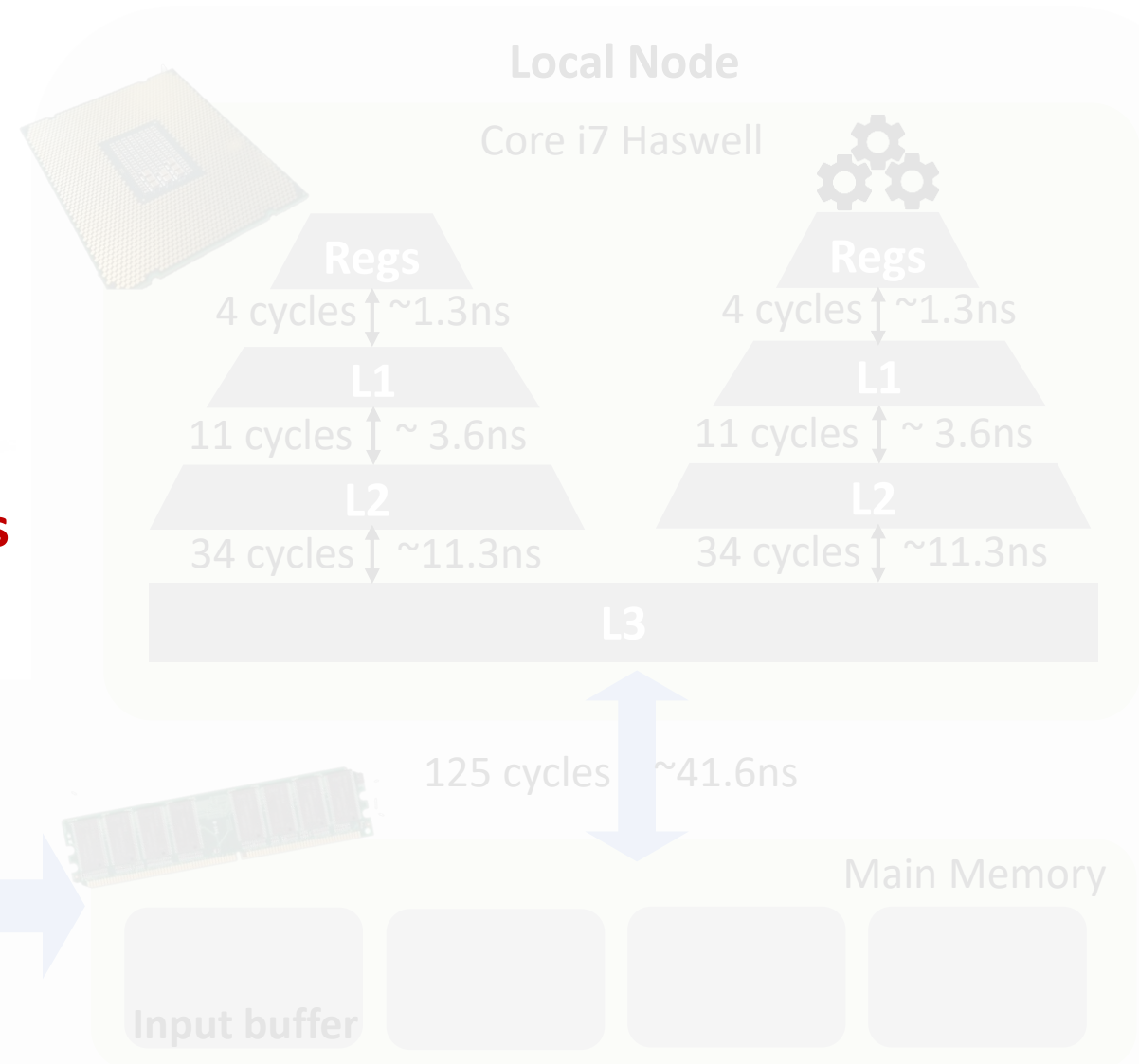
What next?

Fundamental development: CPU core speed cannot keep up with the growing network bandwidths! Already today, cannot process packets at line rate!

Data Processing in modern RDMA networks

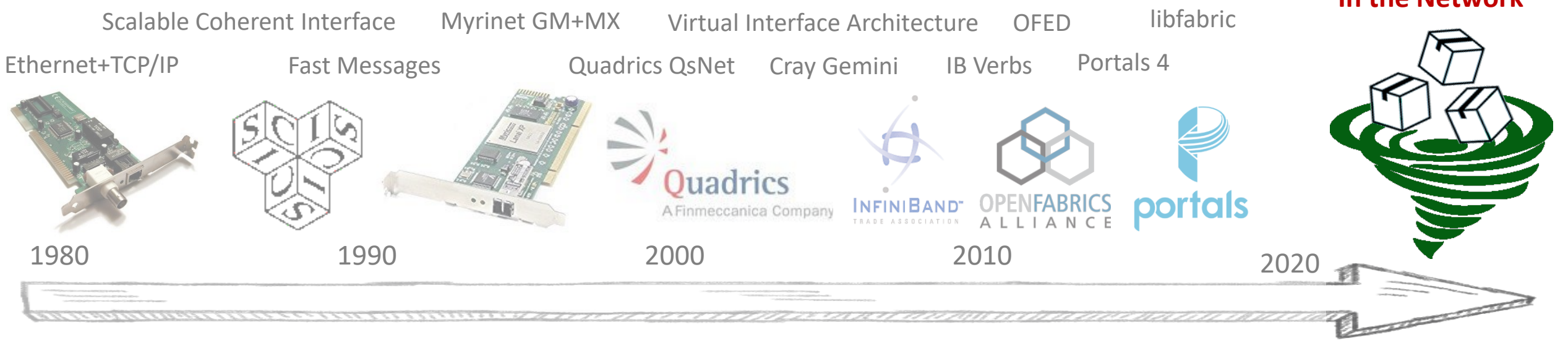


Mellanox Connect-X5: 1 packet/5ns
Tomorrow (400G): 1 packet/1.2ns



The future of High-Performance Networking Interfaces

SPIN
Streaming Processing
In the Network



sockets (active) message based protocol offload remote direct memory access (RDMA) fully programmable NIC acceleration

coherent memory access OS bypass zero copy triggered operations

Established Principles for Compute Acceleration

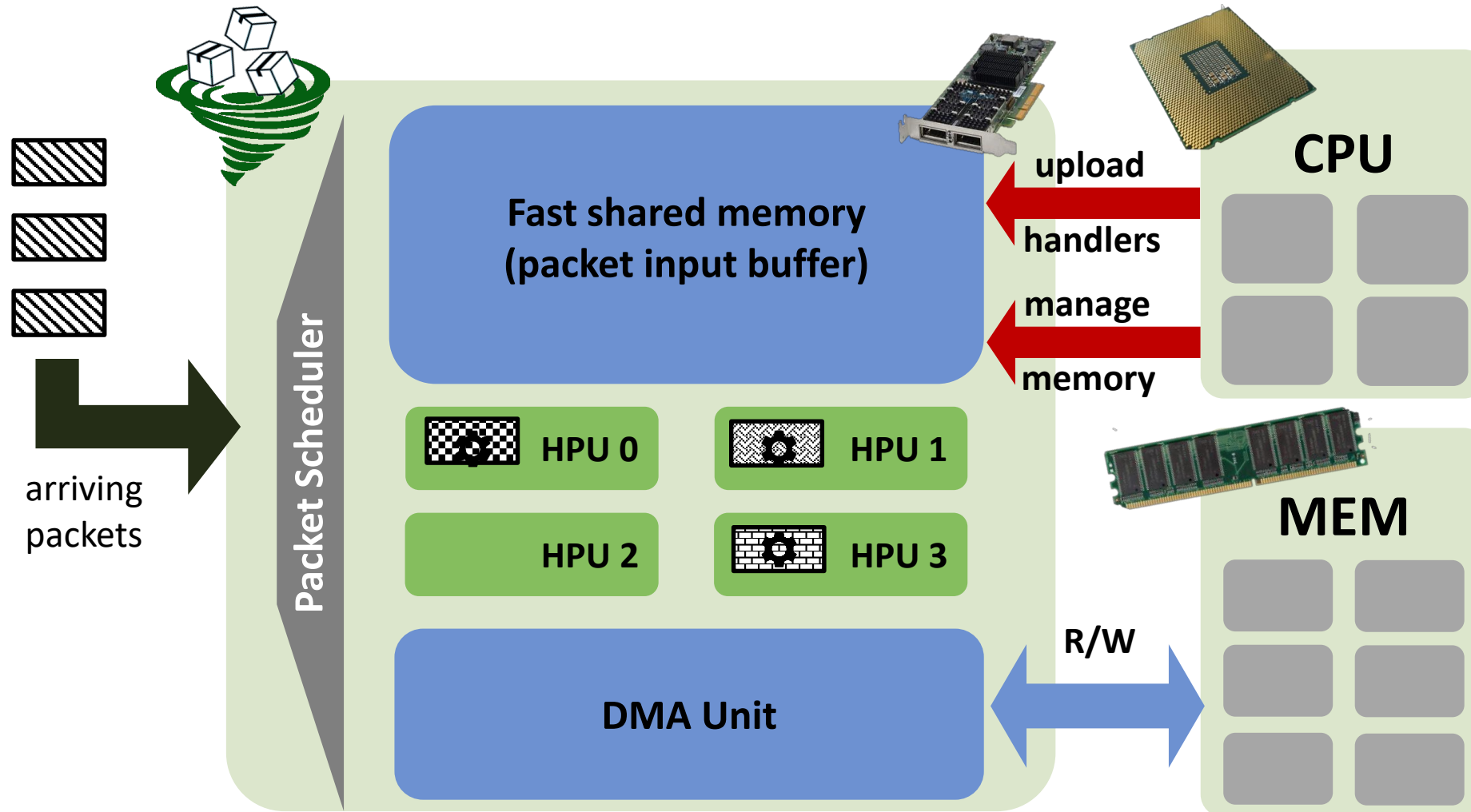
Specialization Programmability Libraries
Ease-of-use Portability Efficiency



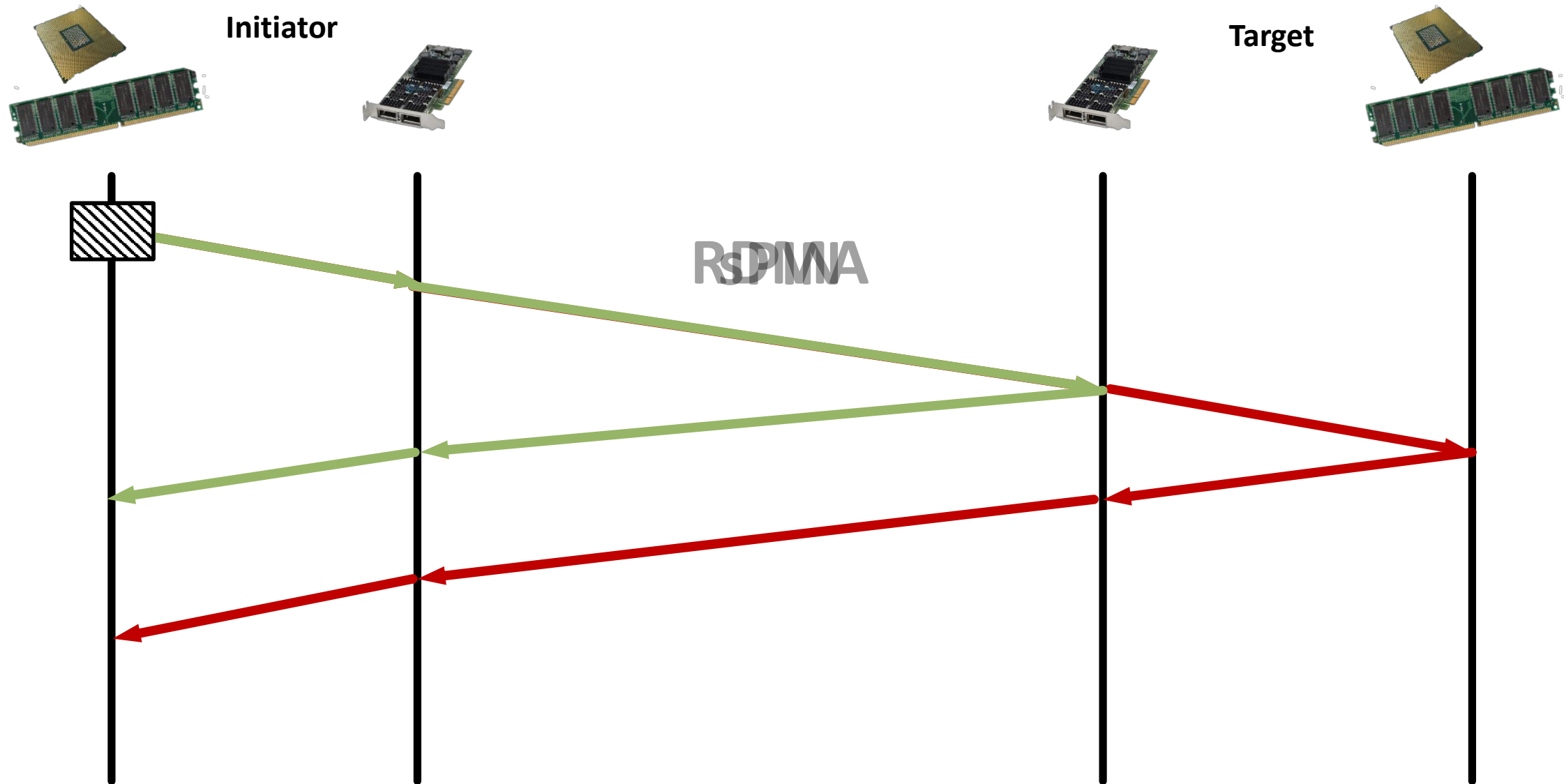
TOP 500[®] SUPERCOMPUTER SITES June 2017

95 / top-100 systems use RDMA
>285 / top-500 systems use RDMA

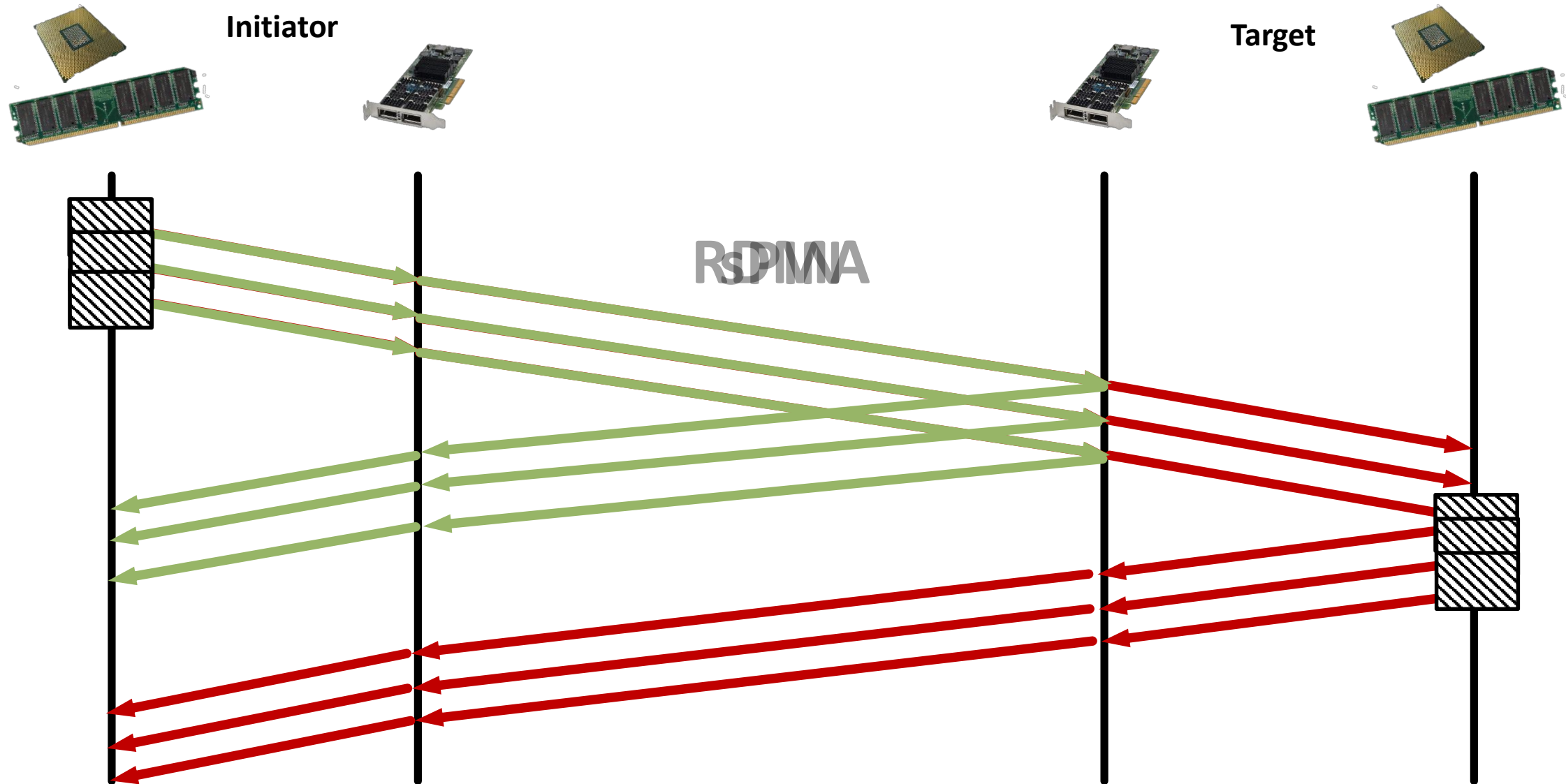
sPIN NIC - Abstract Machine Model



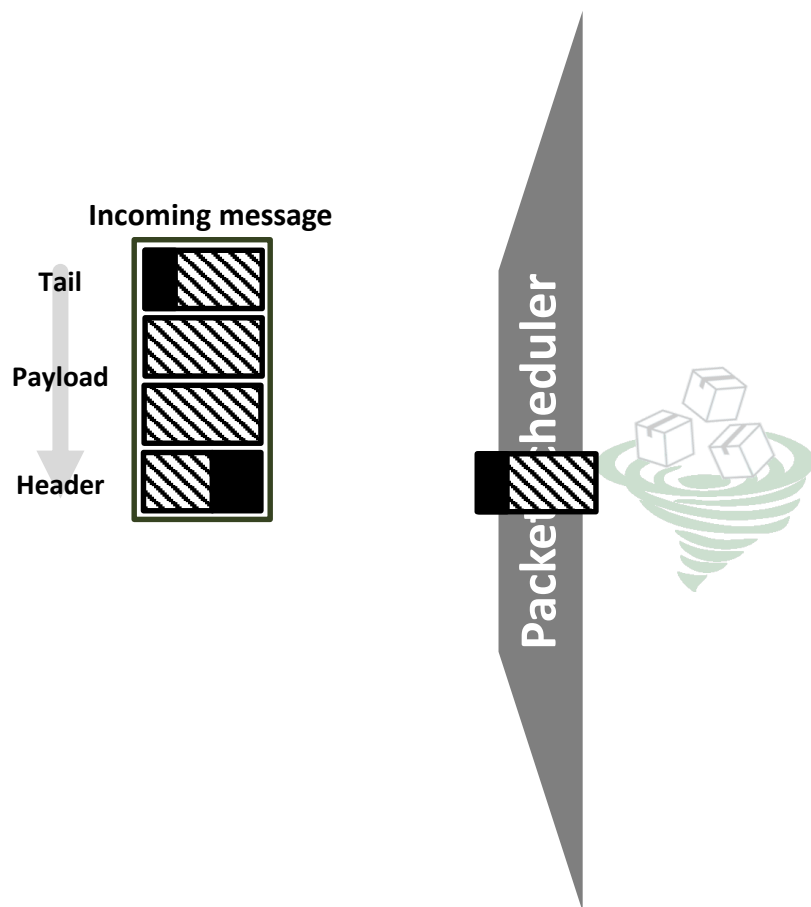
RDMA vs. sPIN in action: Simple Ping Pong



RDMA vs. sPIN in action: Streaming Ping Pong



sPIN – Programming Interface



Header handler

```
__handler int pp_header_handler(const ptl_header_t h, void *state) {
    pingpong_info_t *i = state;
    i->source = h.source_id;
    return PROCESS_DATA; // execute payload handler to put from device
}
```

Payload handler

```
__handler int pp_payload_handler(const ptl_payload_t p, void *state) {
    pingpong_info_t *i = state;
    PtlHandlerPutFromDevice(p.base, p.length, 1, 0, i->source, 10, 0, NULL, 0);
    return SUCCESS;
}
```

Completion handler

```
__handler int pp_completion_handler(int dropped_bytes,
                                     bool flow_control_triggered, void *state) {
    return SUCCESS;
}
```

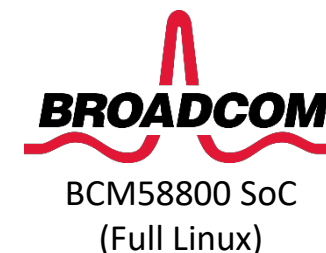
```
connect(peer, /* ... */, &pp_header_handler, &pp_payload_handler, &pp_completion_handler);
```

Possible sPIN implementations



- **sPIN is a programming abstraction, similar to CUDA or OpenCL combined with OFED or Portals 4**
 - It enables a large variety of NIC implementations!
 - For example, massively multithreaded HPUs
Including warp-like scheduling strategies
- **Main goal: sPIN must not obstruct line-rate**
 - Programmer must limit processing time per packet
Little's Law: 500 instructions per handler, 2.5 GHz, IPC=1, 1 Tb/s → 25 kiB memory
 - Relies on fast shared memory (processing in packet buffers)
Scratchpad or registers
 - Quick (single-cycle) handler invocation on packet arrival
Pre-initialized memory & context
- **Can be implemented in most RDMA NICs with a firmware update**
 - Or in software in programmable (Smart) NICs

at 400G, process more than
833 million messages/s



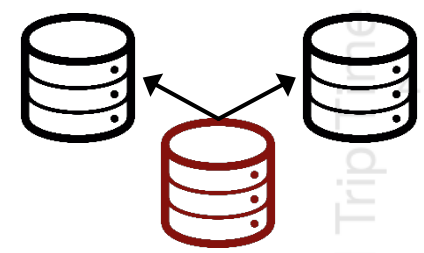
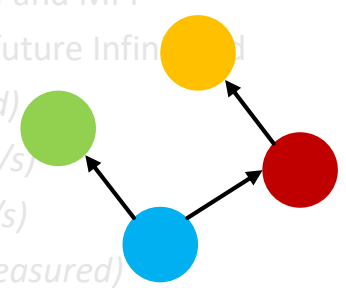
Simulating a sPIN NIC – Ping Pong

- LogGOPSim v2 [1]: combine LogGOPSim (packet-level network) with gem5 (cycle accurate CPU simulation)

Network (LogGOPSim):

- Supports Portals 4 and MPI
- Parametrized for future InfiniBand
- $\sigma=65ns$ (measured)
- $g=6.7ns$ (150 MM/s)
- $G=2.5ps$ (400 Gib/s)
- Switch $L=50ns$ (measured)
- Wire $L=33.4ns$ (100 Gb/s)

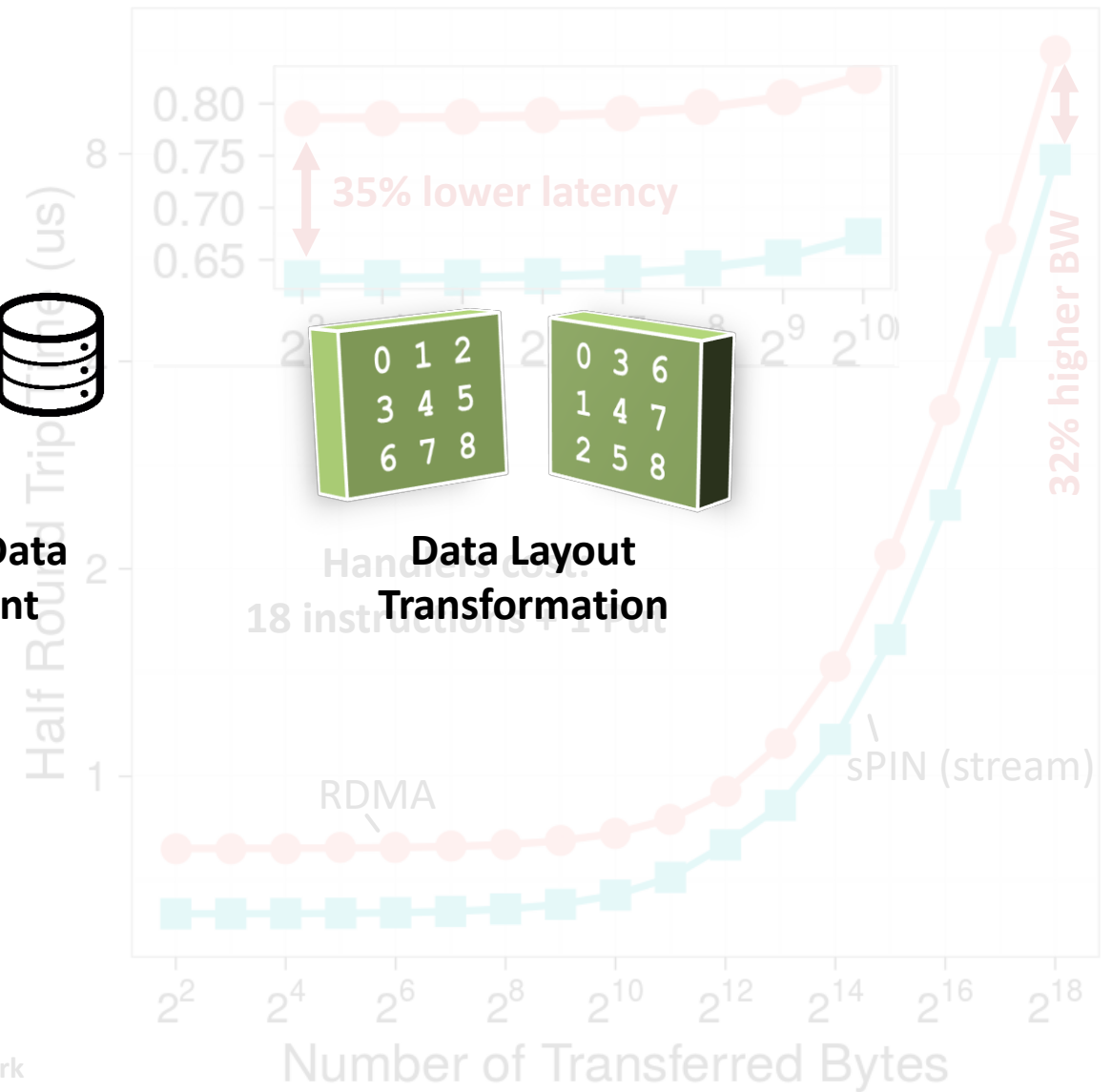
Network Group Communication



Distributed Data Management

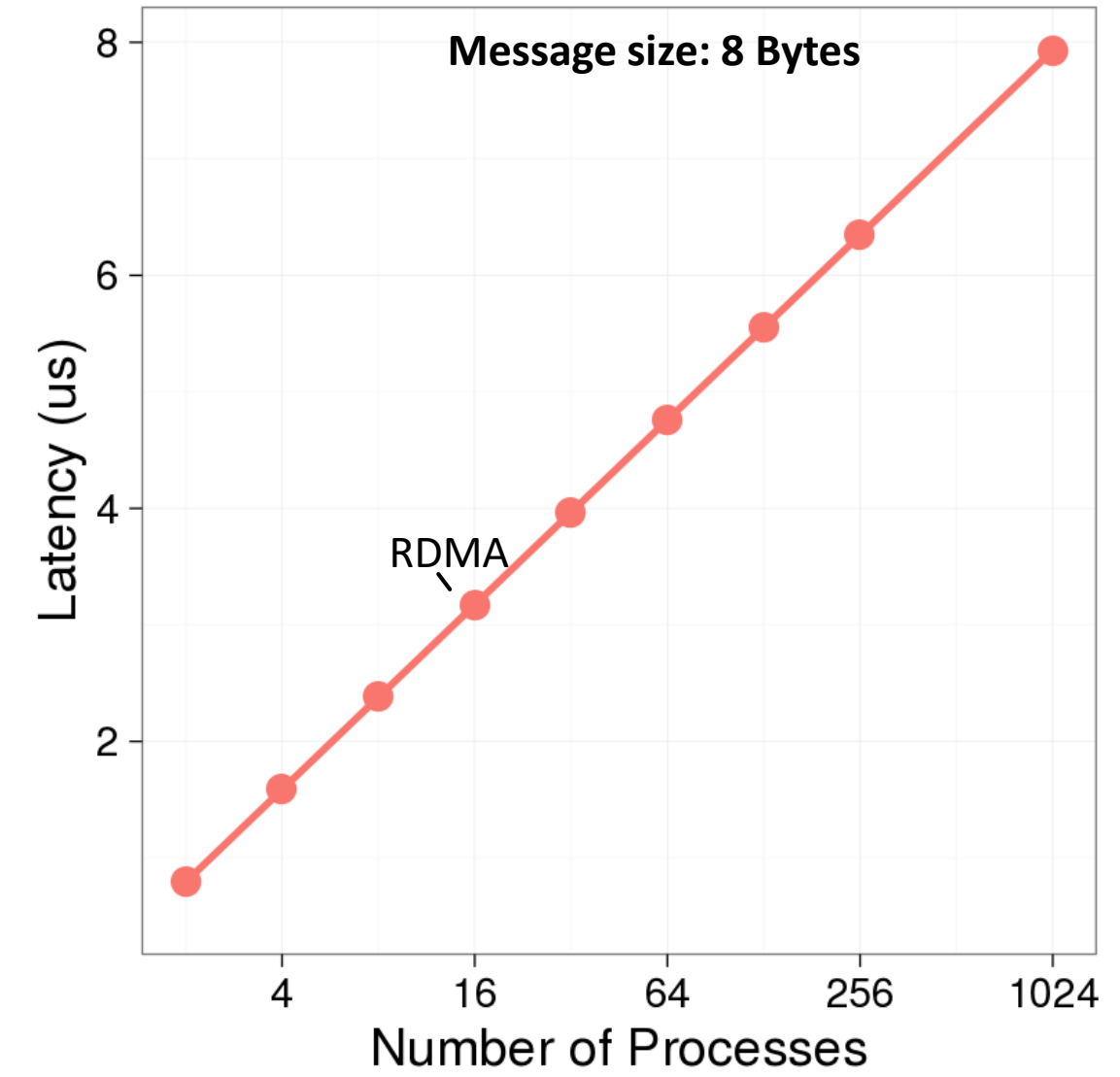
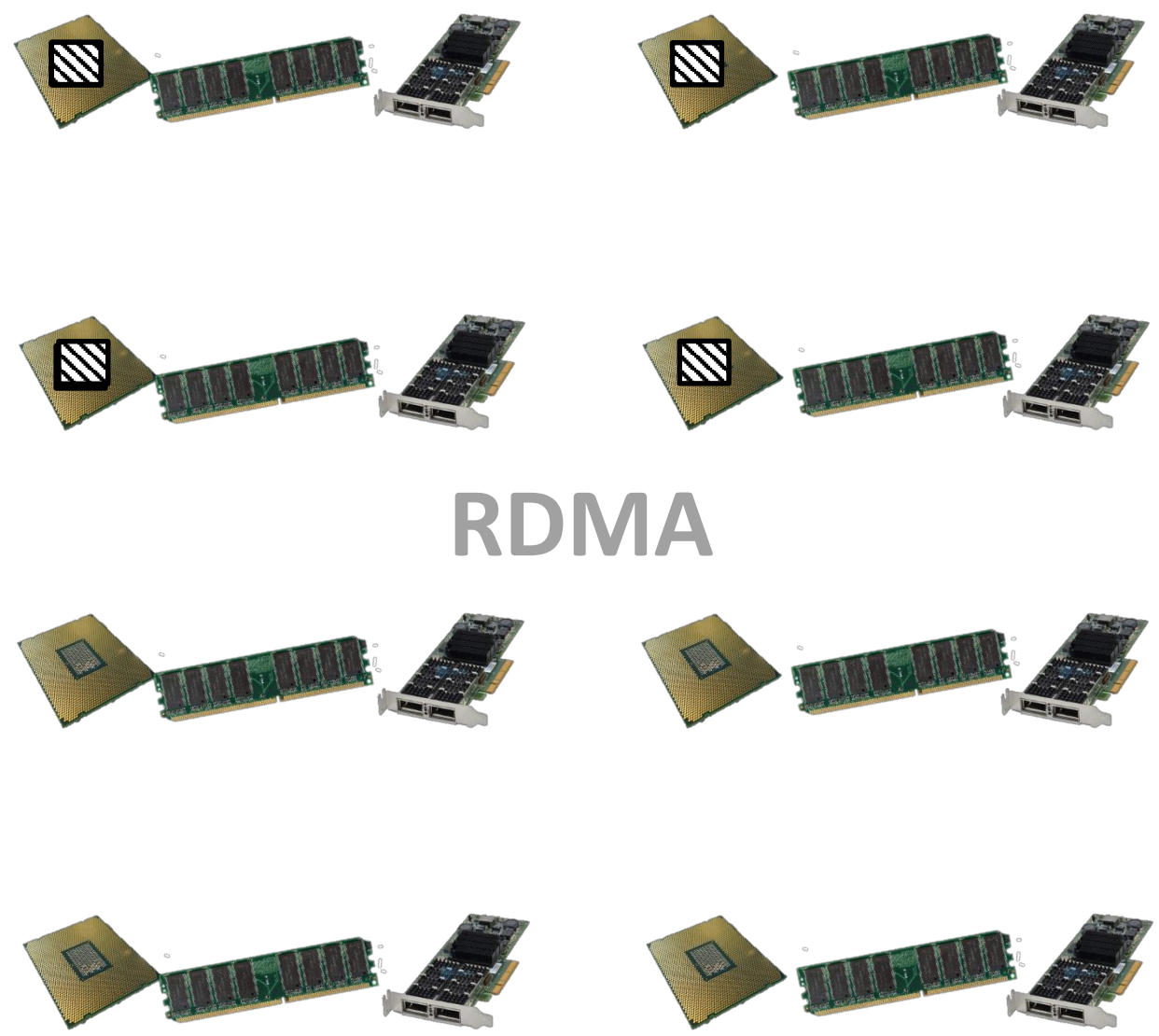
NIC HPU

- 2.5 GHz ARM Cortex A15 OOO
- ARMv8-A 32 bit ISA
- Single-cycle access SRAM (no DRAM)
- Header matching $m=30ns$, per packet 2ns
- In parallel with g!*

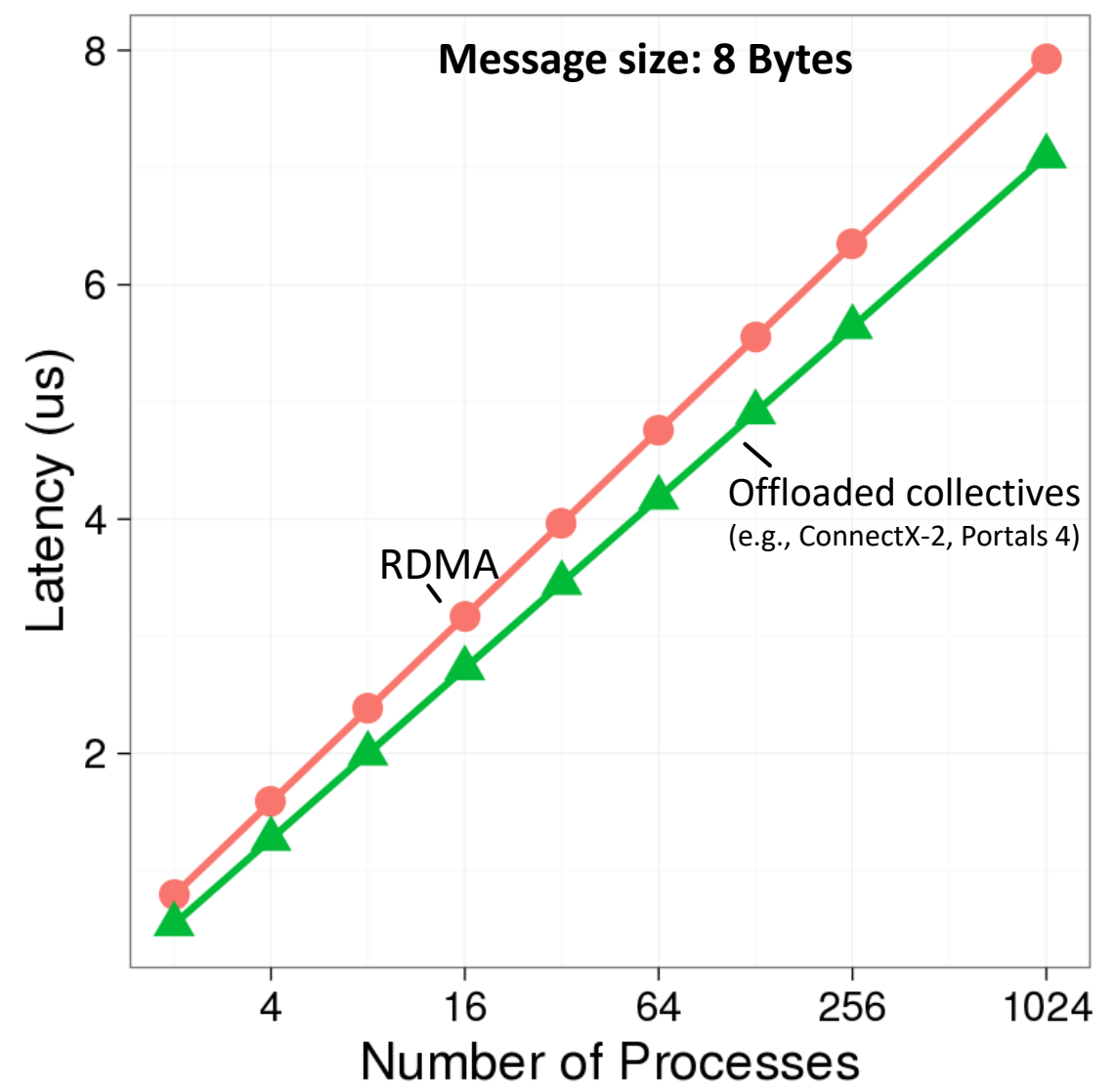


[1] S. Di Girolamo, K. Taranov, T. Schneider, E. Stalder, T. Hoefler, LogGOPSim+gem5: Simulating Network Offload Engines Over Packet-Switched Networks. Presented at ExaMPI'17

Use Case 1: Broadcast acceleration



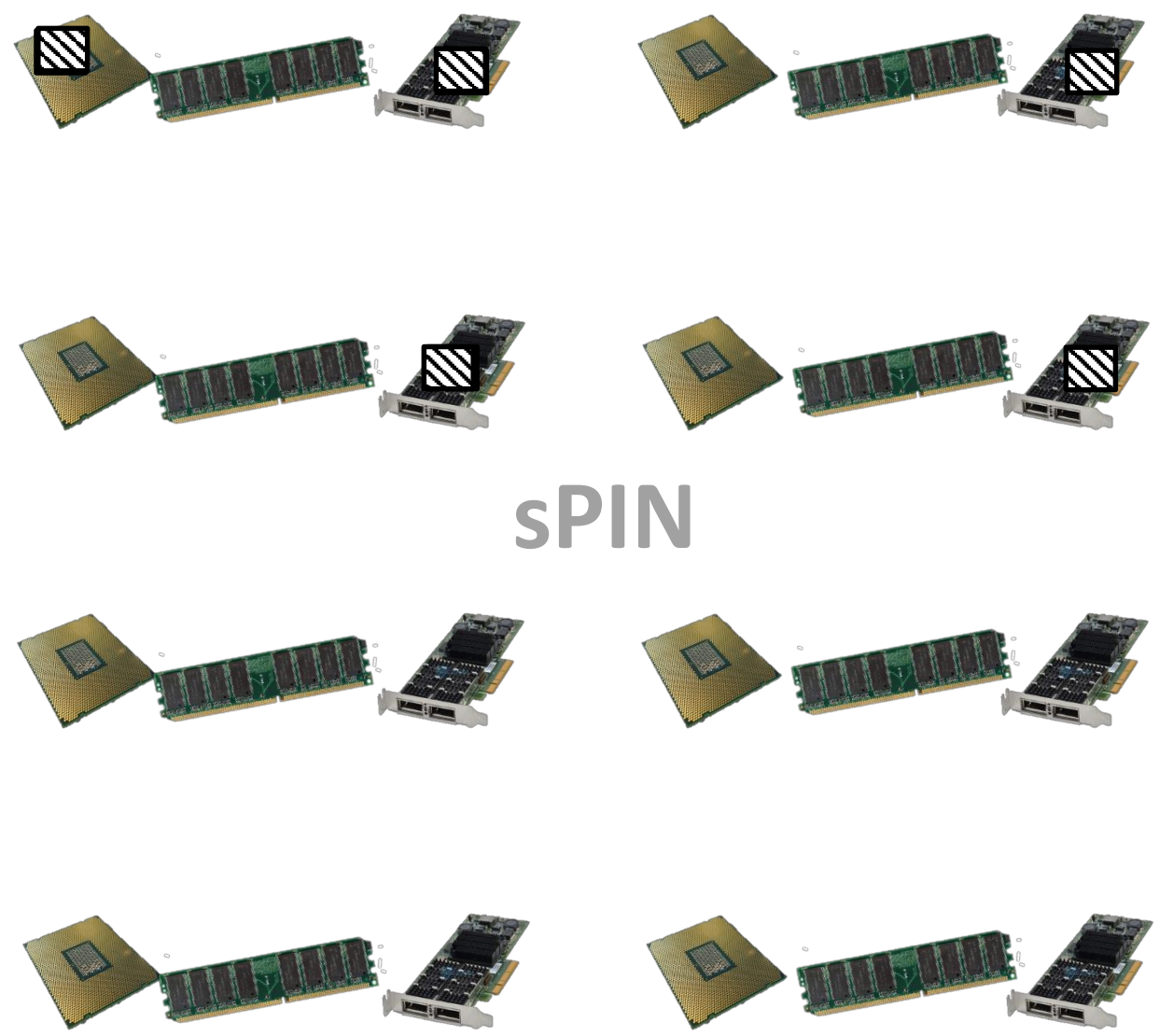
Use Case 1: Broadcast acceleration



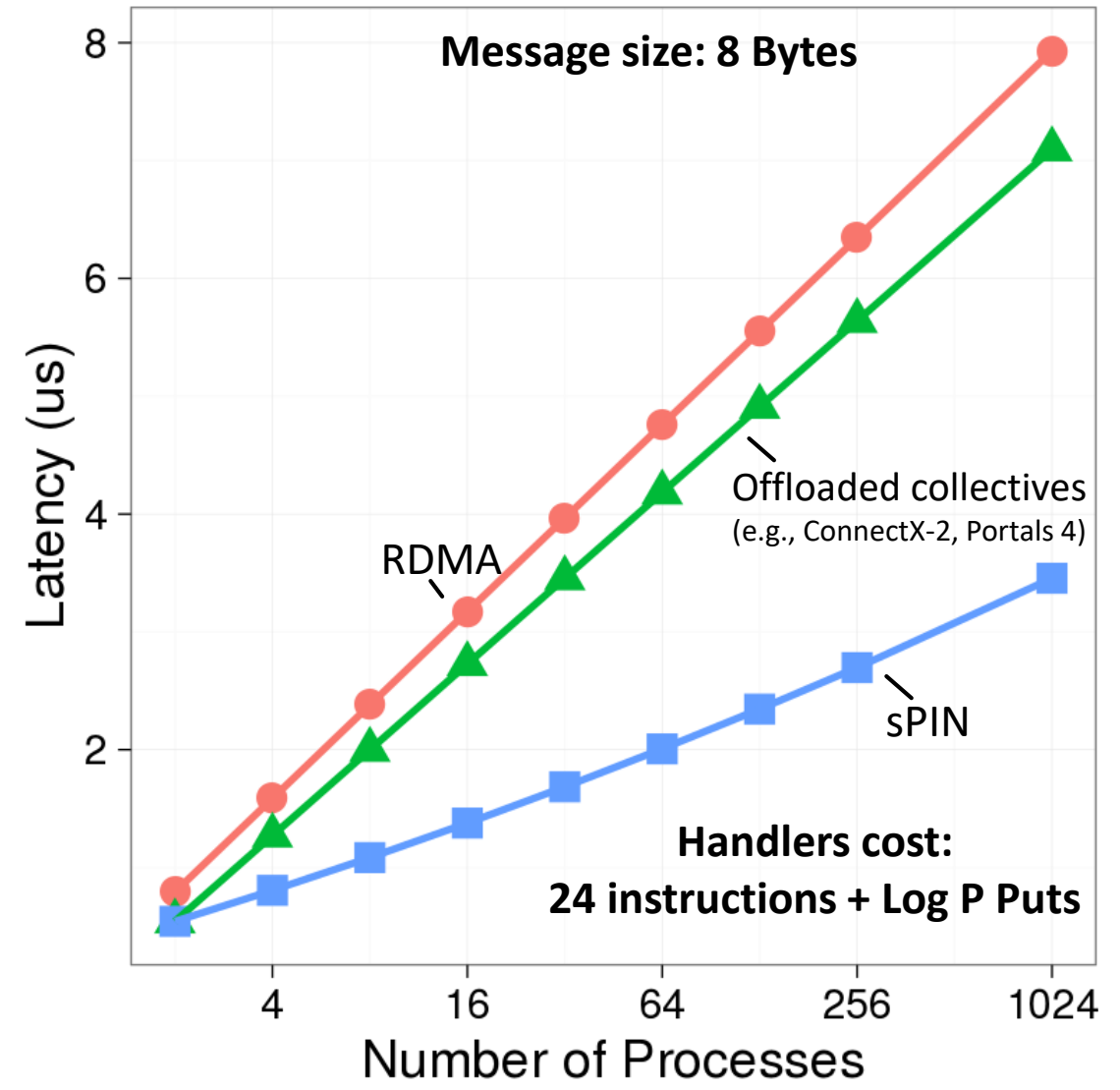
Underwood, K.D., et al., Enabling flexible collective communication offload with triggered operations. *HOTI'11*

Liu, J., et al., High performance RDMA-based MPI implementation over InfiniBand. *International Journal of Parallel Programming* 2004

Use Case 1: Broadcast acceleration



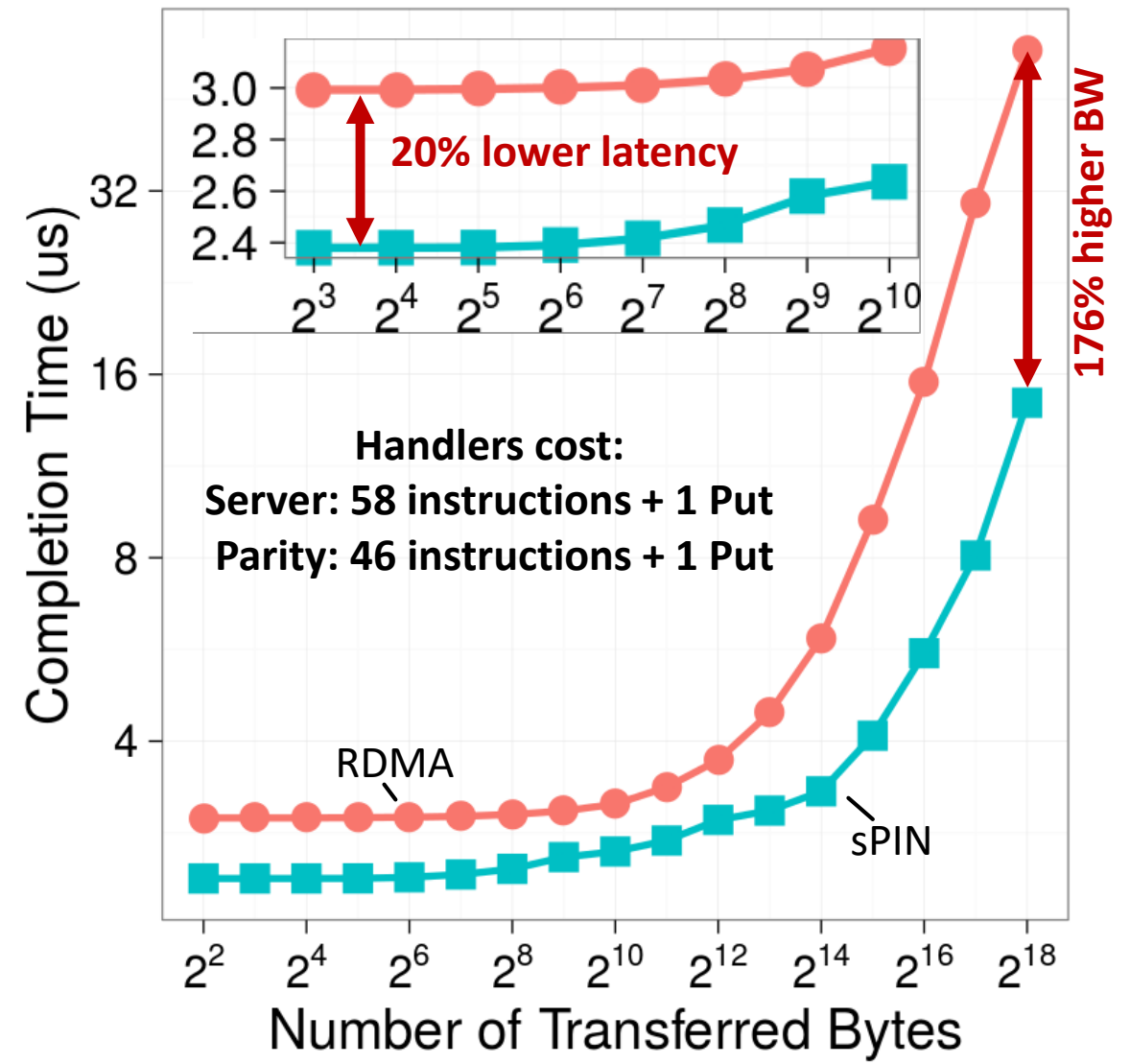
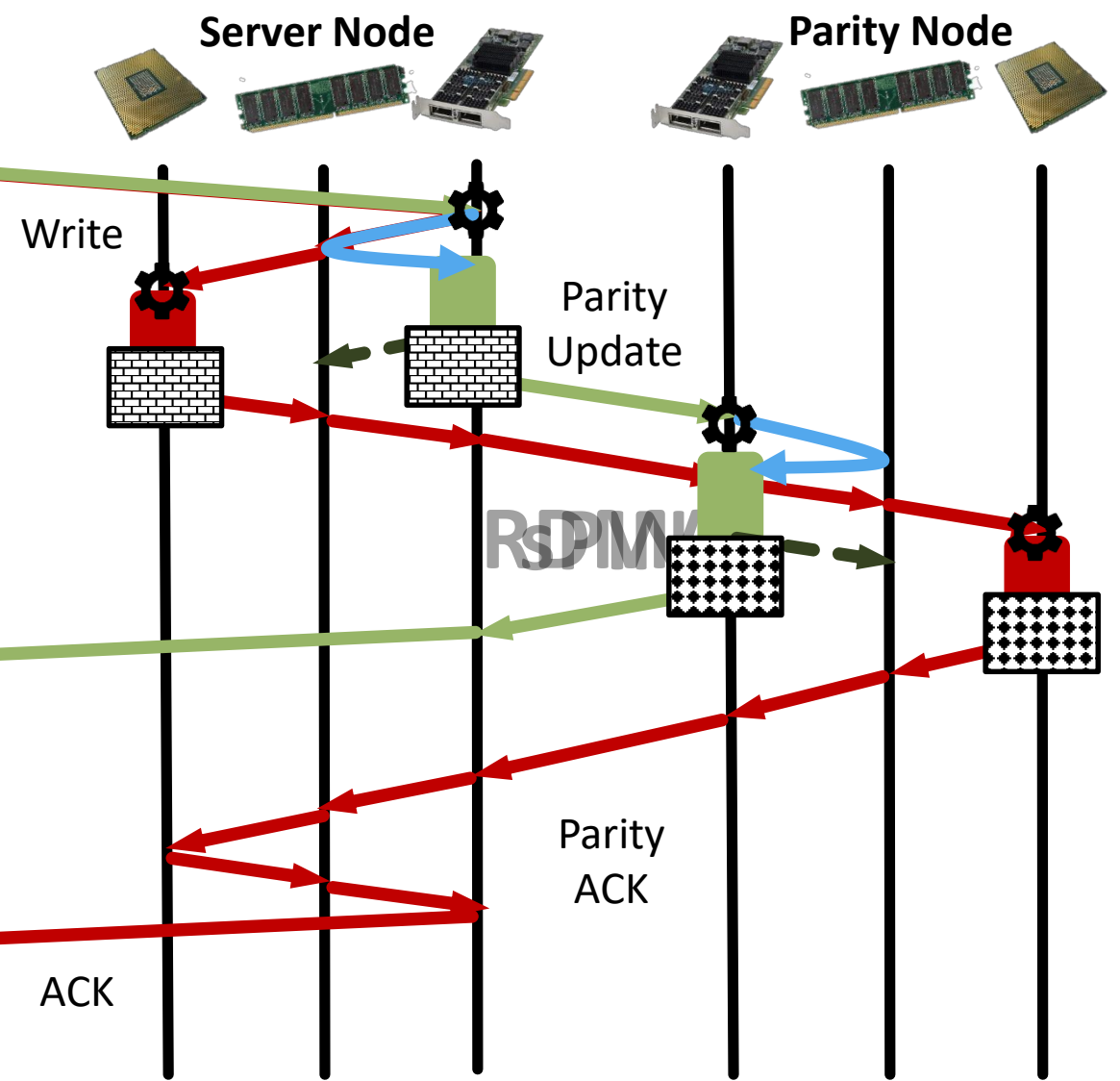
sPIN



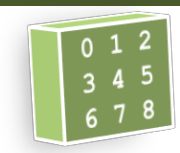
Underwood, K.D., et al., Enabling flexible collective communication offload with triggered operations. *HOTI'11*

Liu, J., et al., High performance RDMA-based MPI implementation over InfiniBand. *International Journal of Parallel Programming* 2004

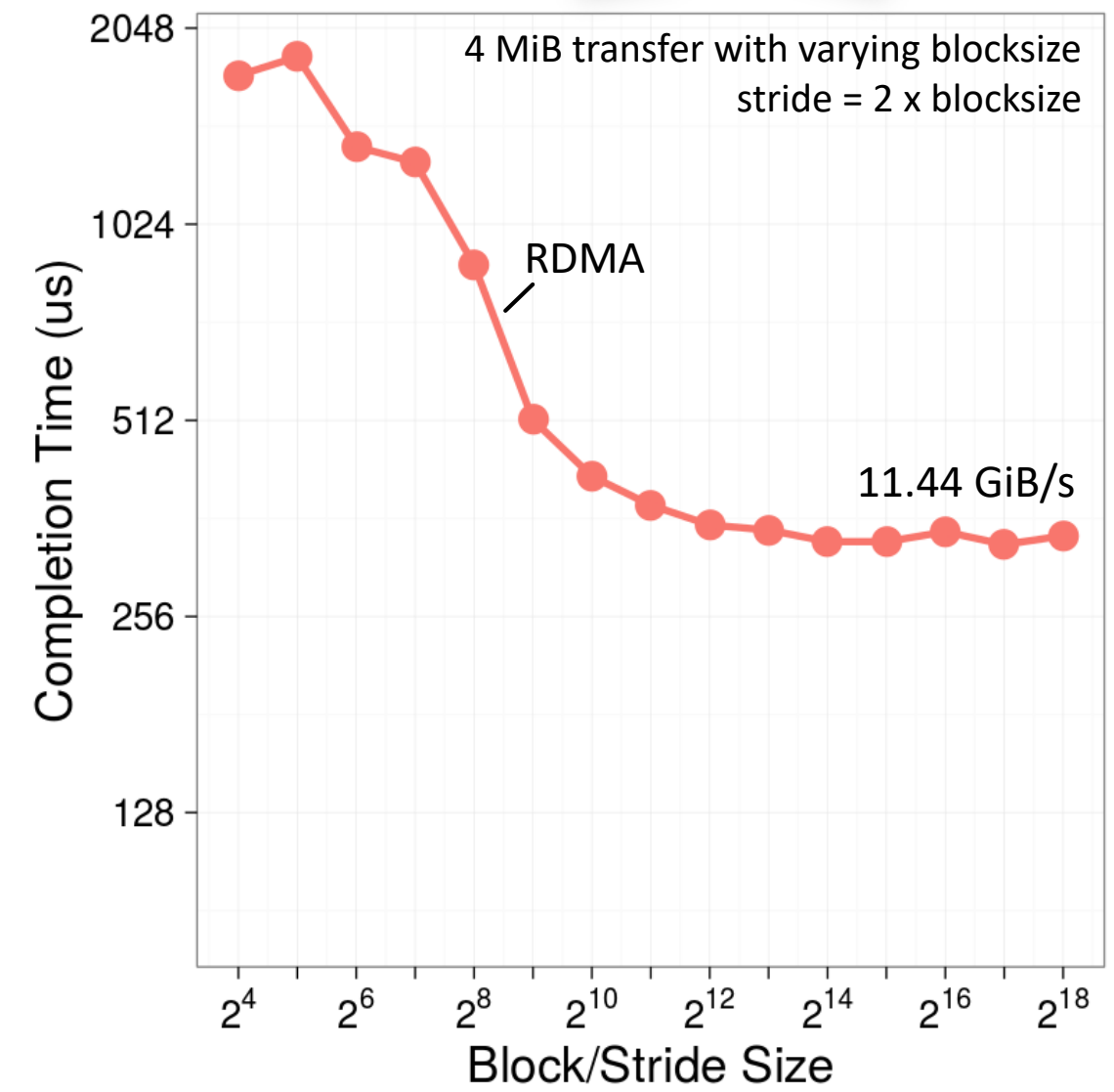
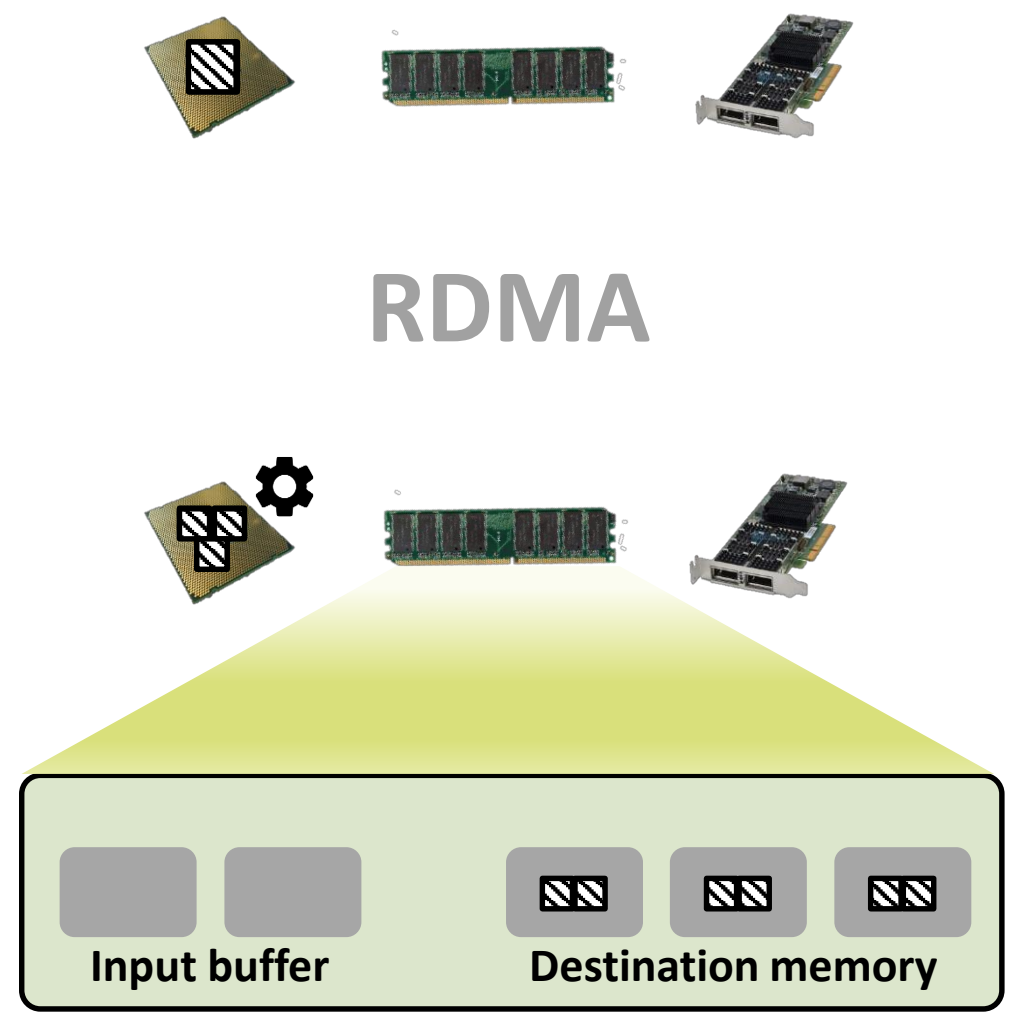
Use Case 2: RAID acceleration



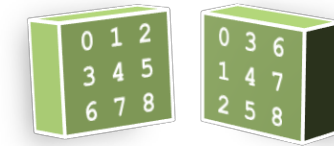
Use Case 3: MPI Datatypes acceleration



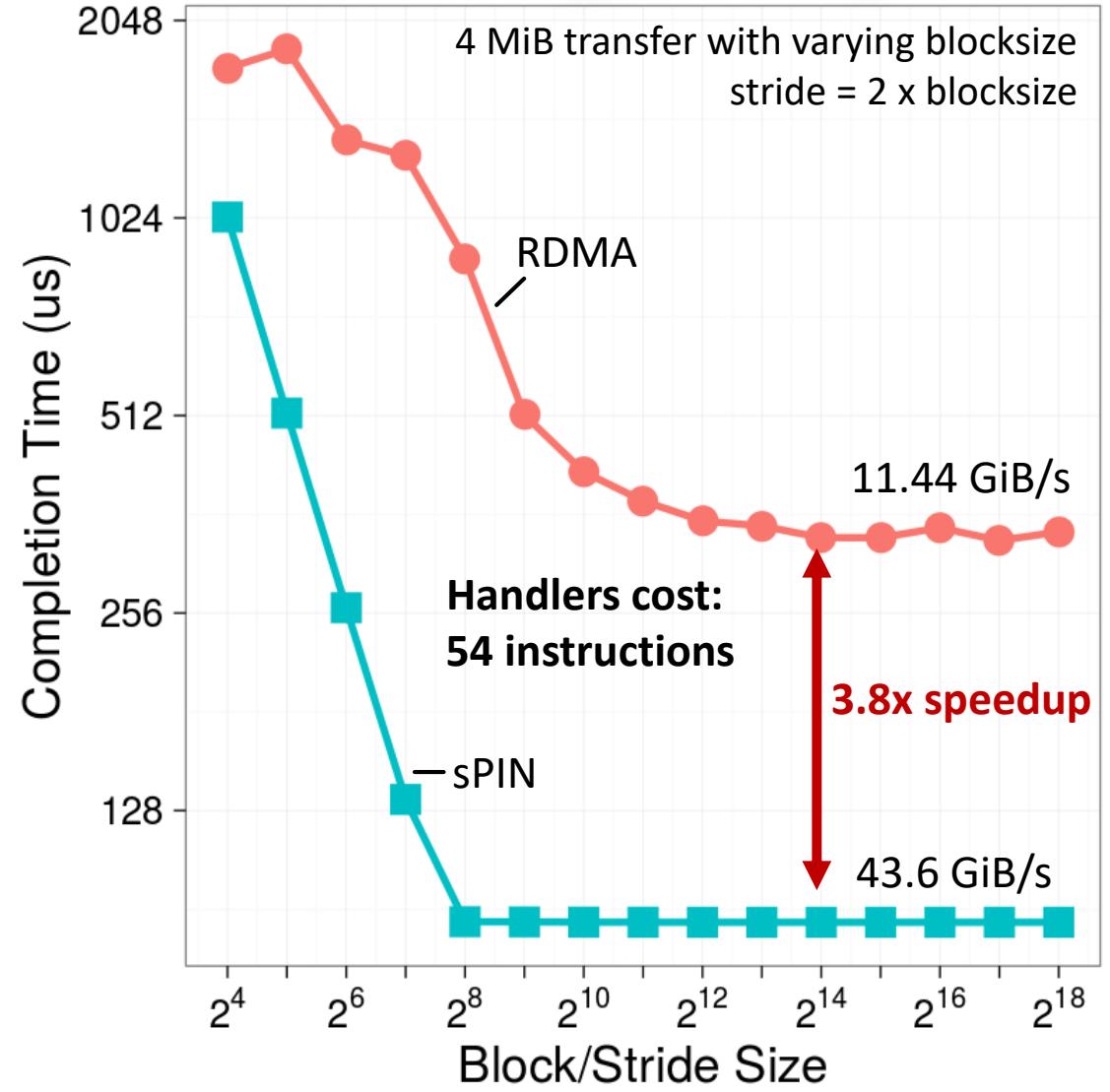
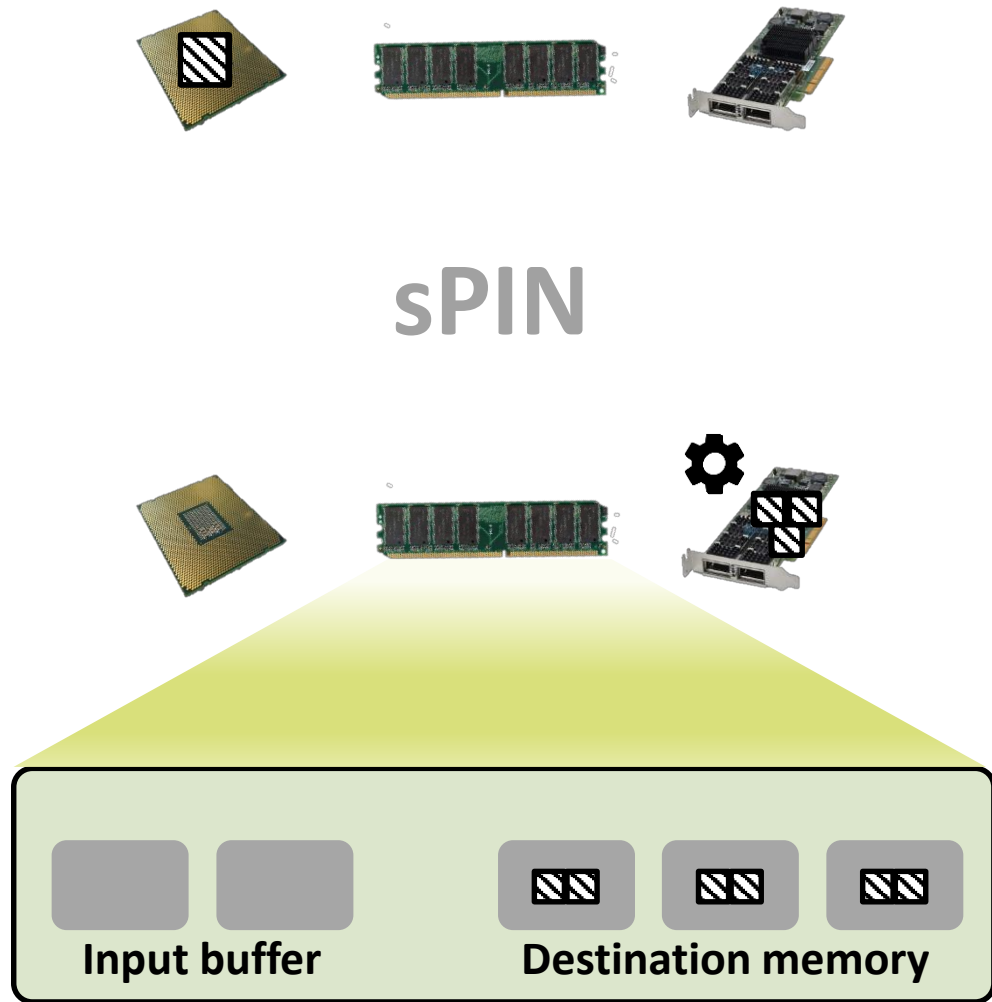
Data Layout Transformation



Use Case 3: MPI Datatypes acceleration



Data Layout Transformation




Further results and use-cases

Further results and use-cases

SPCL ETH zürich


Use Case 4: MPI Rendezvous Protocol



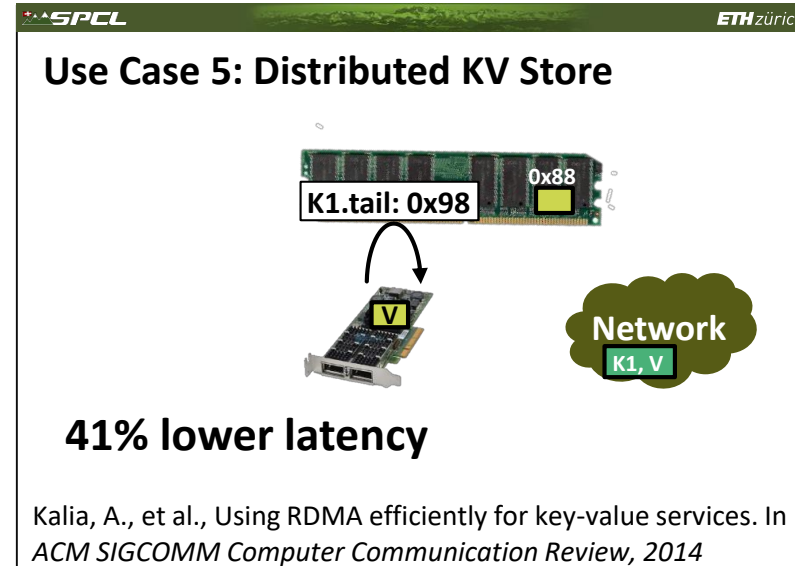
program	p	msgs	ovhd	ovhd	red
MILC	64	5.7M	5.5%	1.9%	65%
POP	64	772M	3.1%	2.4%	22%
coMD	72	5.3M	6.1%	2.4%	60%
coMD	360	28.1M	6.5%	2.8%	58%
Cloverleaf	72	2.7M	5.2%	2.4%	53%
Cloverleaf	360	15.3M	5.6%	3.2%	42%

Further results and use-cases

Use Case 4: MPI Rendezvous Protocol




program	p	msgs	ovhd	ovhd	red
MILC	64	5.7M	5.5%	1.9%	65%
POP	64	772M	3.1%	2.4%	22%
coMD	72	5.3M	6.1%	2.4%	60%
coMD	360	28.1M	6.5%	2.8%	58%
Cloverleaf	72	2.7M	5.2%	2.4%	53%
Cloverleaf	360	15.3M	5.6%	3.2%	42%




Further results and use-cases

Use Case 4: MPI Rendezvous Protocol



program	p	msgs	ovhd	ovhd	red
MILC	64	5.7M	5.5%	1.9%	65%
POP	64	772M	3.1%	2.4%	22%
coMD	72	5.3M	6.1%	2.4%	60%
coMD	360	28.1M	6.5%	2.8%	58%
Cloverleaf	72	2.7M	5.2%	2.4%	53%
Cloverleaf	360	15.3M	5.6%	3.2%	42%

Use Case 5: Distributed KV Store

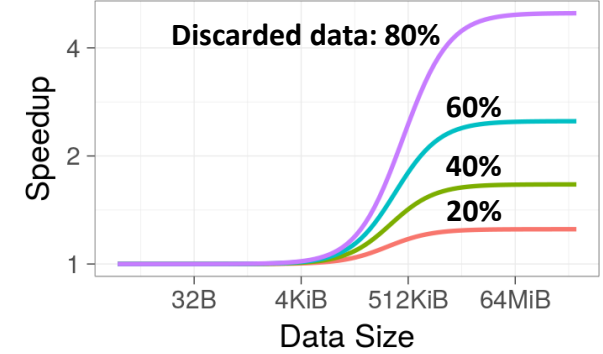


41% lower latency

Network

Kalia, A., et al., Using RDMA efficiently for key-value services. In *ACM SIGCOMM Computer Communication Review*, 2014

Use Case 6: Conditional Read



Discarded data: 80%

60%

40%

20%

Speedup

Data Size

Barthels, C., et al., Designing Databases for Future High-Performance Networks. *IEEE Data Eng. Bulletin*, 2017

Further results and use-cases

Use Case 4: MPI Rendezvous Protocol

program	p	msgs	ovhd	ovhd	red
MILC	64	5.7M	5.5%	1.9%	65%
POP	64	772M	3.1%	2.4%	22%
coMD	72	5.3M	6.1%	2.4%	60%
coMD	360	28.1M	6.5%	2.8%	58%
Cloverleaf	72	2.7M	5.2%	2.4%	53%
Cloverleaf	360	15.3M	5.6%	3.2%	42%

Use Case 5: Distributed KV Store

41% lower latency

Kalia, A., et al., Using RDMA efficiently for key-value services. In *ACM SIGCOMM Computer Communication Review*, 2014

Use Case 6: Conditional Read

Discarded data: 80%

Speedup

Data Size

Barthels, C., et al., Designing Databases for Future High-Performance Networks. *IEEE Data Eng. Bulletin*, 2017

Use Case 7: Distributed Transactions

data pkts

log pkts

Dragojević, A, et al., No compromises: distributed transactions with consistency, availability, and performance. *SOSP'15*

Further results and use-cases

Use Case 4: MPI Rendezvous Protocol

program	p	msgs	ovhd	ovhd	red
MILC	64	5.7M	5.5%	1.9%	65%
POP	64	772M	3.1%	2.4%	22%
coMD	72	5.3M	6.1%	2.4%	60%
coMD	360	28.1M	6.5%	2.8%	58%
Cloverleaf	72	2.7M	5.2%	2.4%	53%
Cloverleaf	360	15.3M	5.6%	3.2%	42%

Use Case 5: Distributed KV Store

41% lower latency

Kalia, A., et al., Using RDMA efficiently for key-value services. In *ACM SIGCOMM Computer Communication Review*, 2014

Use Case 6: Conditional Read

Discarded data: 80%

Speedup

Data Size

Barthels, C., et al., Designing Databases for Future High-Performance Networks. *IEEE Data Eng. Bulletin*, 2017

Use Case 7: Distributed Transactions

Dragojević, A, et al., No compromises: distributed transactions with consistency, availability, and performance. *SOSP'15*

Use Case 8: FT Broadcast

Bosilca, G., et al., Failure Detection and Propagation in HPC systems. *SC'16*

Further results and use-cases

Use Case 4: MPI Rendezvous Protocol

program	p	msgs	ovhd	ovhd	red
MILC	64	5.7M	5.5%	1.9%	65%
POP	64	772M	3.1%	2.4%	22%
coMD	72	5.3M	6.1%	2.4%	60%
coMD	360	28.1M	6.5%	2.8%	58%
Cloverleaf	72	2.7M	5.2%	2.4%	53%
Cloverleaf	360	15.3M	5.6%	3.2%	42%

Use Case 5: Distributed KV Store

The Next 700 sPIN use-cases

... just think about sPIN graph kernels ...

41% lower latency

Kalia, A., et al., Using sPIN for distributed key-value services. In ACM SIGCOMM Conference on Data Communication Systems, 2017.

Use Case 6: Conditional Read

Barthels, C., et al., Designing Databases for Future High-Performance Networks. *IEEE Data Eng. Bulletin*, 2017

Use Case 7: Distributed Transactions

Dragojević, A, et al., No compromises: distributed transactions with consistency, availability, and performance. SOSP'15

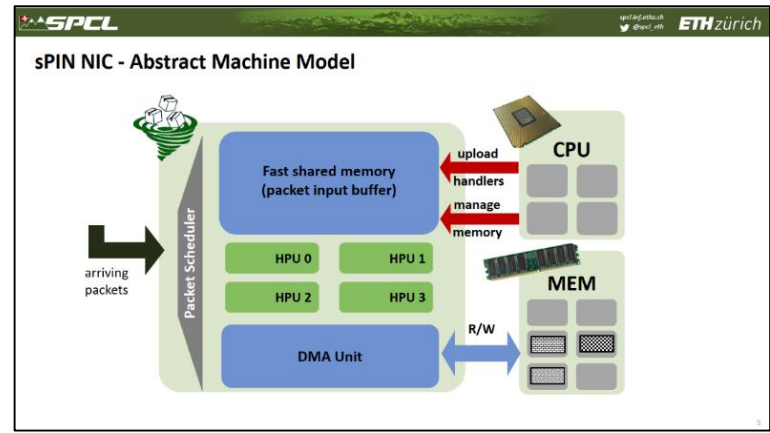
Use Case 8: Failure Detection and Propagation in HPC systems

Bosilca, G., et al., Failure Detection and Propagation in HPC systems. SC'16

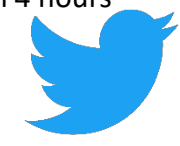
Use Case 9: Distributed Consensus

István, Z., et al., Consensus in a Box: Inexpensive Coordination in Hardware. NSDI'16

sPIN Streaming Processing in the Network for Network Acceleration



3,800+ reads in less than 4 hours



sPIN

beyond RDMA

```

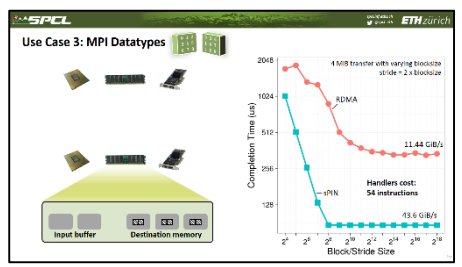
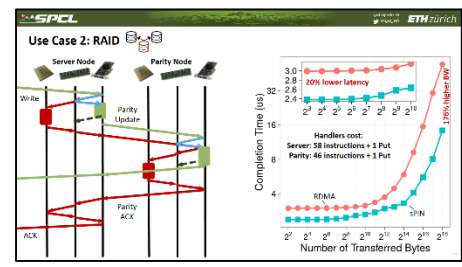
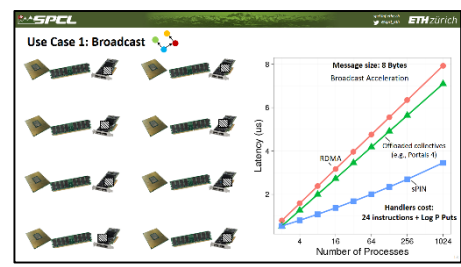
sPIN - Programming Interface

Header handler
_handler int pp_header_handler(const pti_header_t h, void *state) {
    pingpong_info_t *i = state;
    i->source = h.source_id;
    return PROCESS_DATA; // execute payload handler to put from device
}

Payload handler
_handler int pp_payload_handler(const pti_payload_t b, void *state) {
    pingpong_info_t *i = state;
    PtiHandlerPutFromDevice(p.base, p.length, 1, 0, i->source, 10, 0, NULL, 0);
    return SUCCESS;
}

Completion handler
_handler int pp_completion_handler(int dropped_bytes,
    bool flow_control_triggered, void *state) {
    return SUCCESS;
}

connect(peer, /* ... */, &pp_header_handler, &pp_payload_handler, &pp_completion_handler);
    
```



Possible sPIN implementations

- sPIN is a programming abstraction, similar to CUDA or OpenCL combined with OFED or Portals 4
- It enables a large variety of NIC implementations!
- For example, massively multithreaded HPUs
 - Including warp-like scheduling strategies
 - at 400G, process more than 833 million messages/s
- Main goal: sPIN must not obstruct line-rate
 - Programmer must limit processing time per packet
 - Little's Law: 500 instructions per handler, 2.5 GHz, IPC=1, 1 Tb/s → 25 kiB memory
 - Relies on fast shared memory (processing in packet buffers)
 - Scratchpad or registers
 - Quick (single-cycle) handler invocation on packet arrival
 - Pre-initialized memory & context
- Can be implemented in most RDMA NICs with a firmware update
 - Or in software in programmable (Smart) NICs
- Two implementation modes: integrated and discrete
 - Integrated on the same SoC (e.g., through CC mechanism)
 - Discrete is connected via bus interface (e.g., PCIe)

Logos: Mellanox Technologies (Innova Flex (Kintex FPGA)), Microsoft Catapult (Virtex FPGA), Broadcom BCM58800 SoC (Full Linux)

MPI – some new research challenges!

- More complex network interfaces

- sPIN – just discussed in detail
- How to use it in real applications?



EuroMPI'19

September 11-13 2019

Zurich, Switzerland

<https://eurompi19.inf.ethz.ch>

Submit papers by April 15th!

- How to implement sPIN

- Collaboration with Luca Benini, RISC-V SoC
- Full FPGA/simulation environment available as open-source

- New application use-cases – “Deep Learning is HPC [1]”

- Different communication requirements and opportunities for optimization
Mainly sparsity [2] and asynchrony!
- Need to be a bit careful though when entering a new field



T. Hoefler: “Twelve ways to fool the masses when reporting performance of deep learning workloads”
(my humorous guide to floptimize deep learning)

[1]: Ben-Nun, Hoefler: “Demystifying Parallel and Distributed Deep Learning: An In-Depth Concurrency Analysis, arXiv:1802.09941

[2]: Renggli et al.: “SparCML: High-Performance Sparse Communication for Machine Learning”, arXiv:1802.08021



EuroMPI'19

September 11-13 2019
Zurich, Switzerland
<https://eurompi19.inf.ethz.ch>

Important dates:

- Submission server opens:** January 14th, 2019
- Full paper submission:** April 15th, 2019 (AOE)
- Notification:** July 1st, 2019
- Camera-ready:** August 5th, 2019