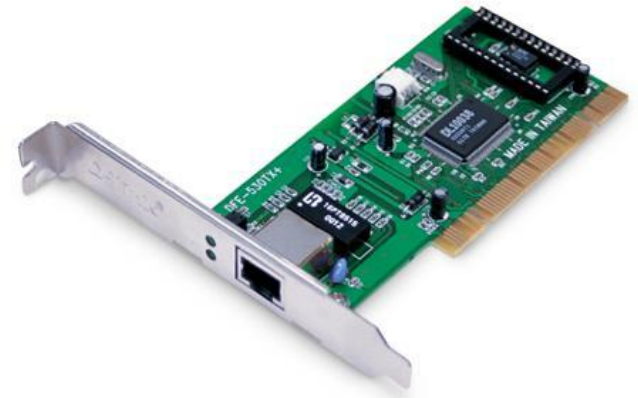# Enabling Highly-Scalable Remote Memory Access Programming with MPI-3 One Sided

ROBERT GERSTENBERGER, MACIEJ BESTA, TORSTEN HOEFLER

# MPI-3.0 RMA

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

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

# MPI-3.0 RMA

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







[1] http://www.mpi-forum.org/docs/mpi-3.0/mpi30-report.pdf
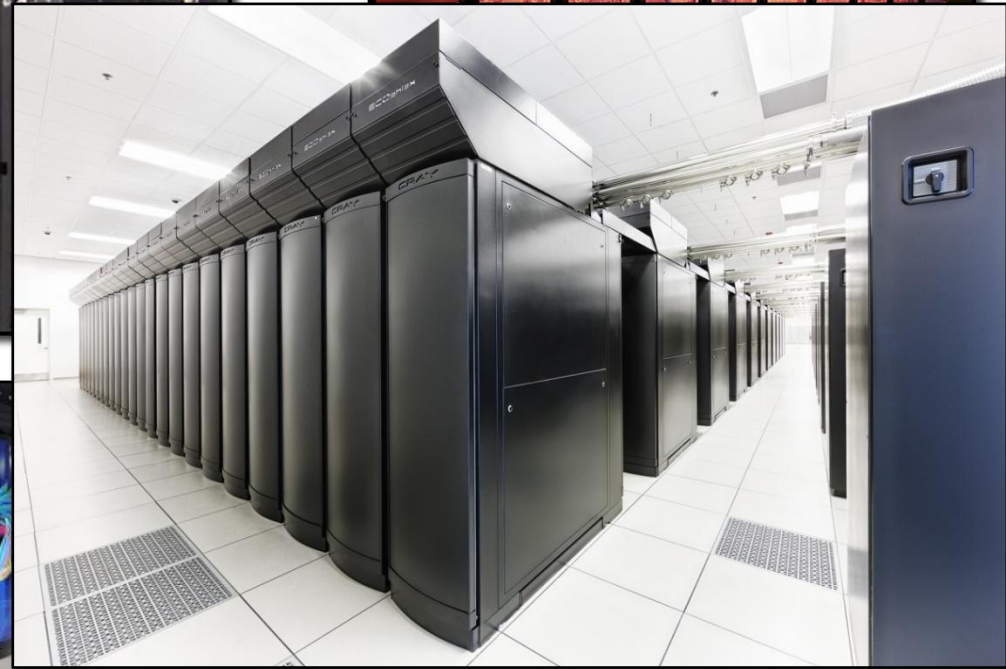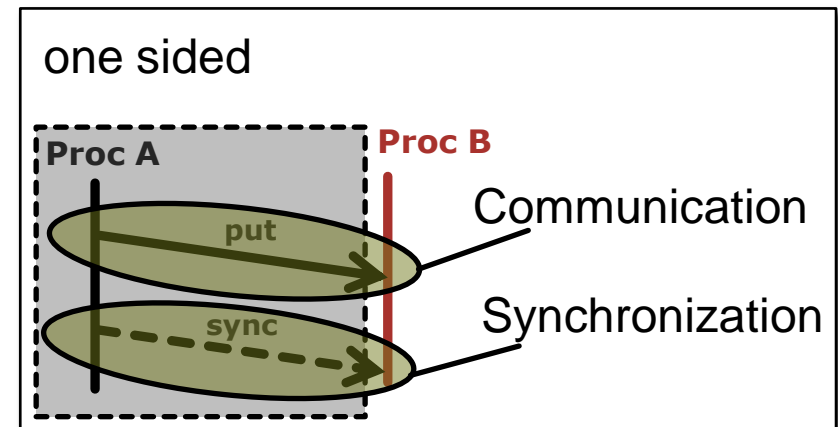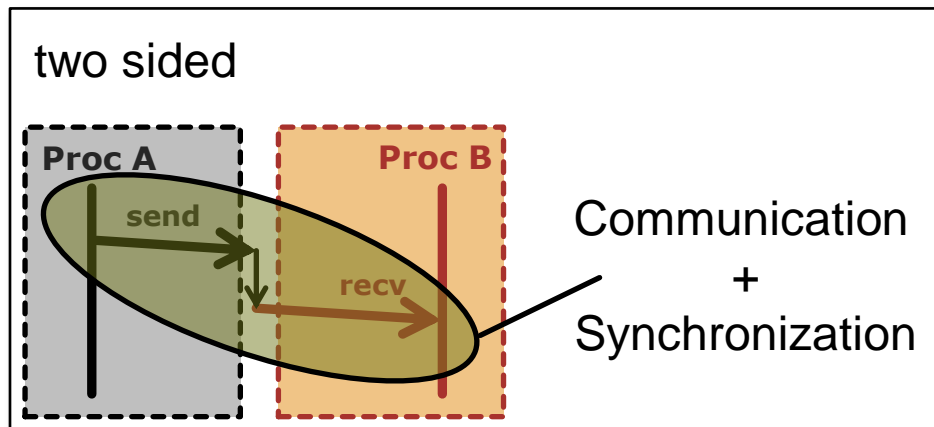
# MPI-3.0 RMA

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

# MPI-3.0 RMA

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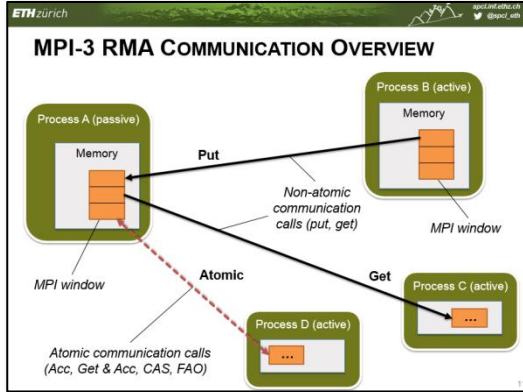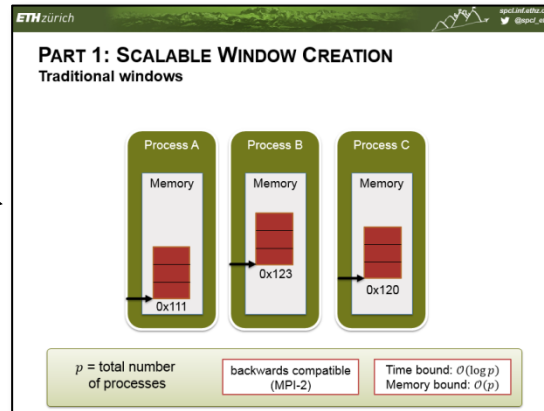    - Majority of HPC networks support RDMA
- Communication is „one sided" (no involvement of destination)
- RMA decouples communication & synchronization
    - Different from message passing



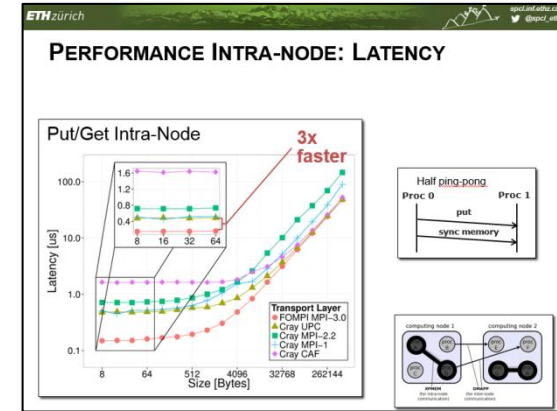[1] http://www.mpi-forum.org/docs/mpi-3.0/mpi30-report.pdf
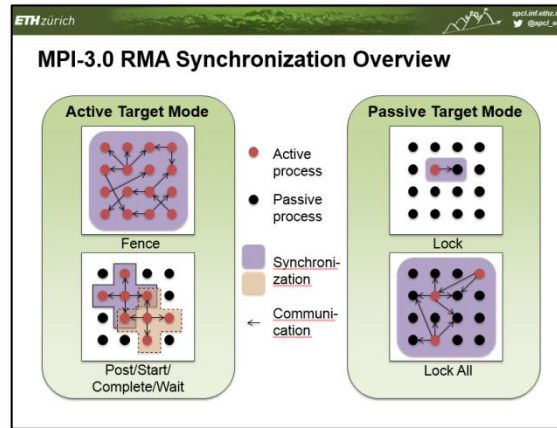
# PRESENTATION OVERVIEW



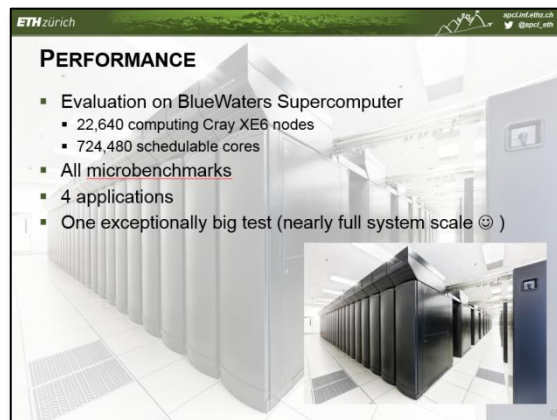1. Overview of three MPI-3 RMA concepts

2. MPI window creation

3. Communication

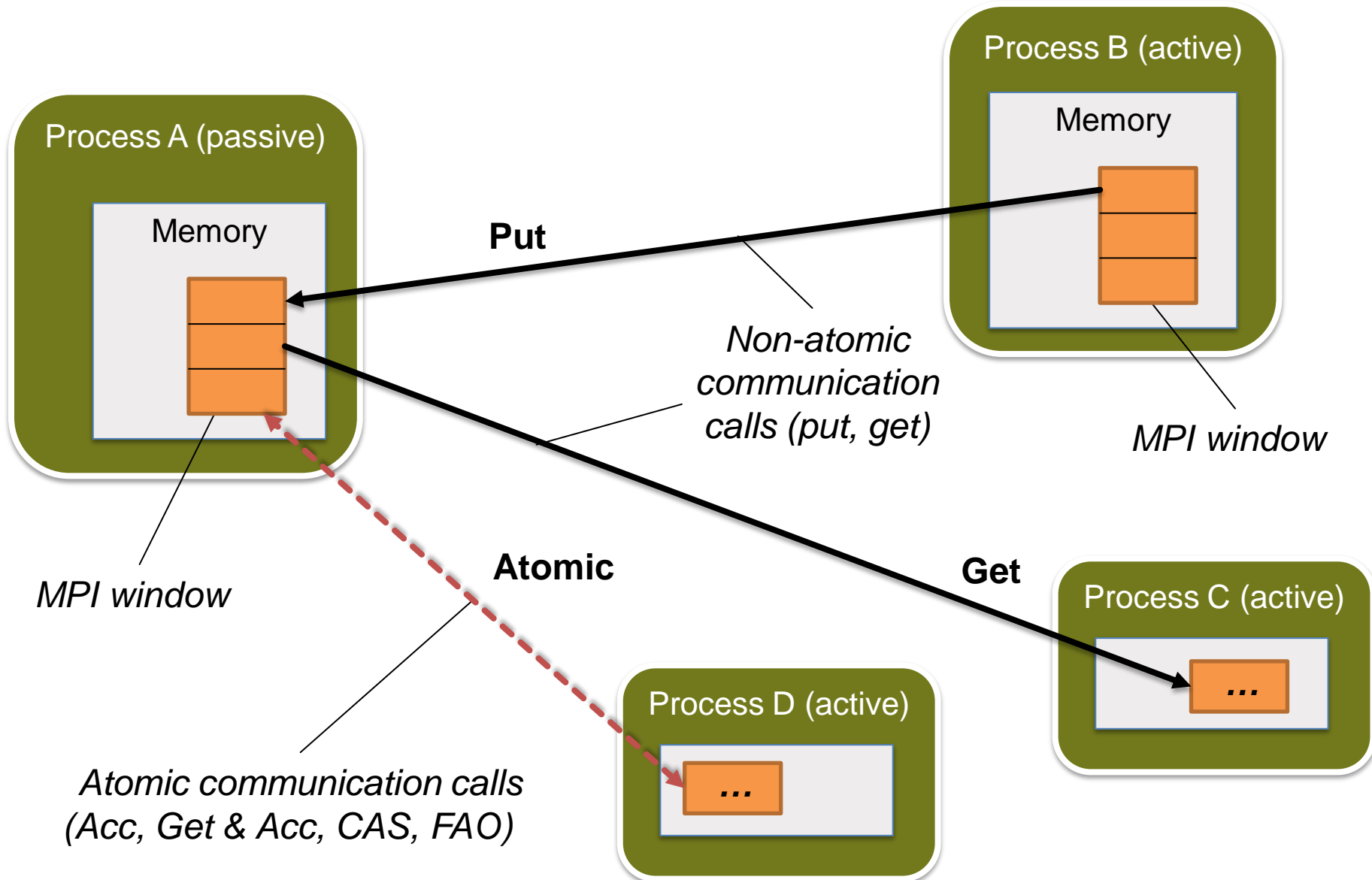4. Synchronization

5. Application evaluation

# MPI-3 RMA COMMUNICATION OVERVIEW



Process A (passive)

Memory

Put

Process B (active)

Memory

Non-atomic communication calls (put, get)

MPI window

MPI window

Atomic

Get

Process C (active)

...

MPI window

Atomic communication calls (Acc, Get & Acc, CAS, FAO)
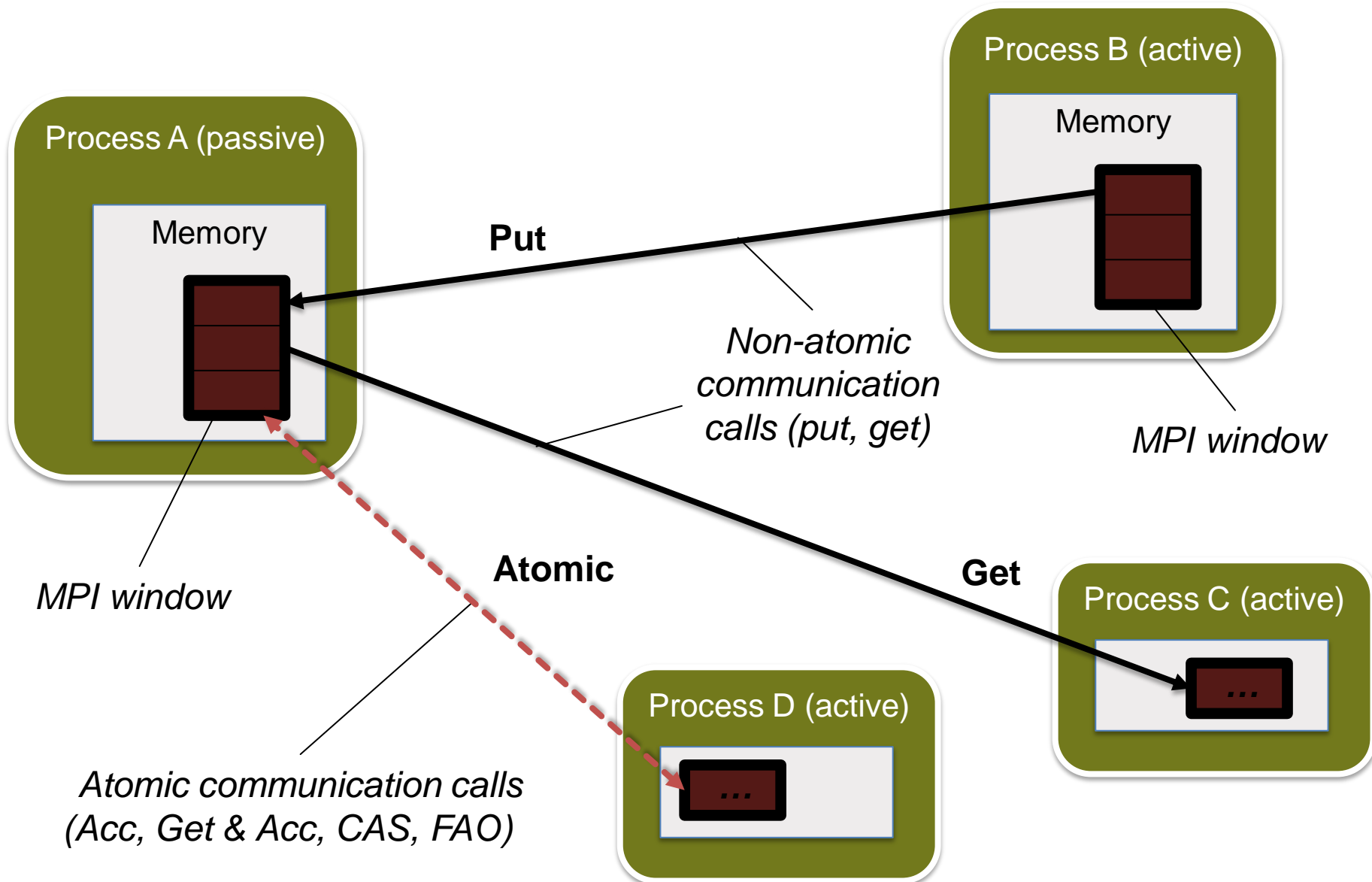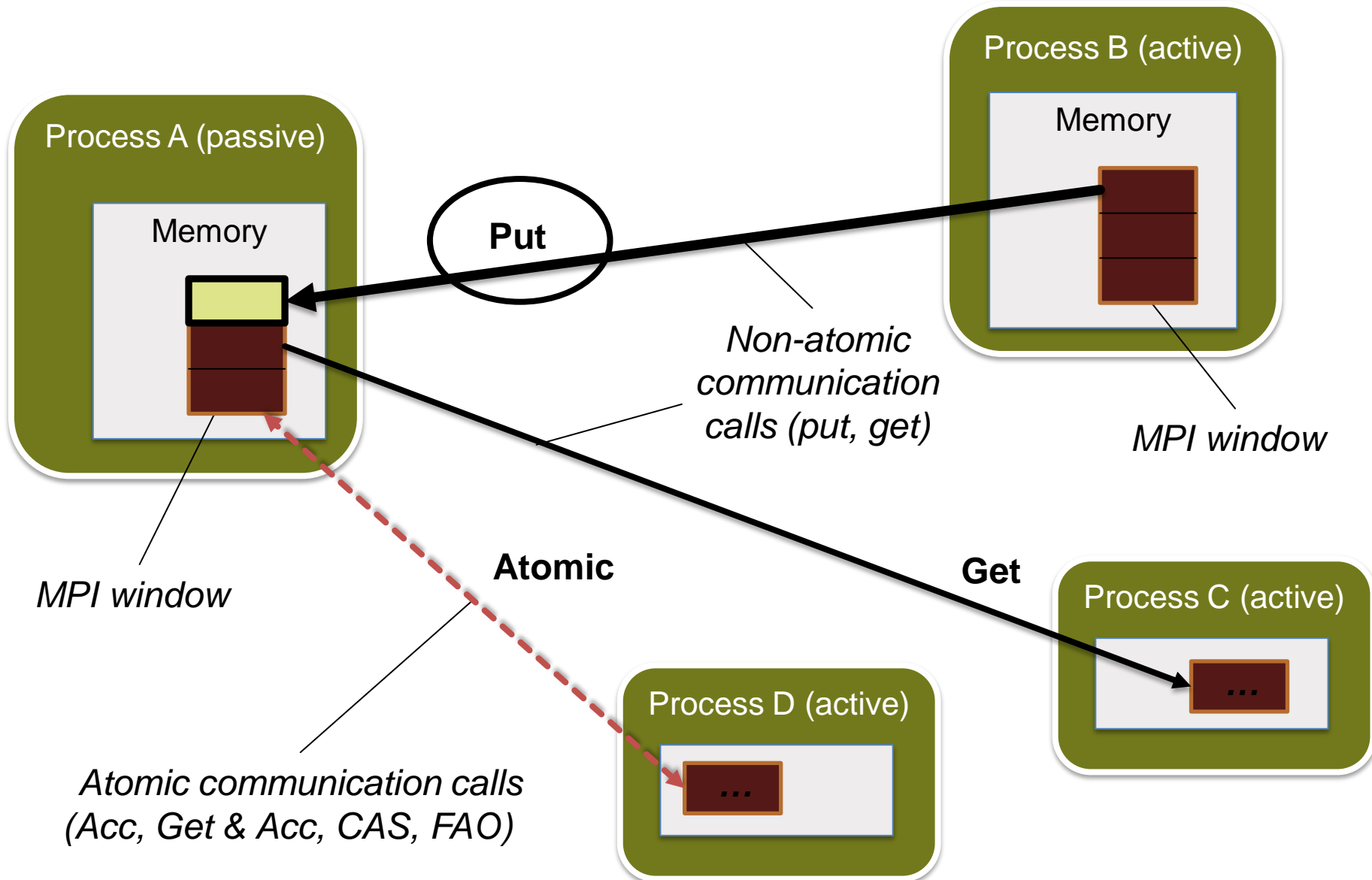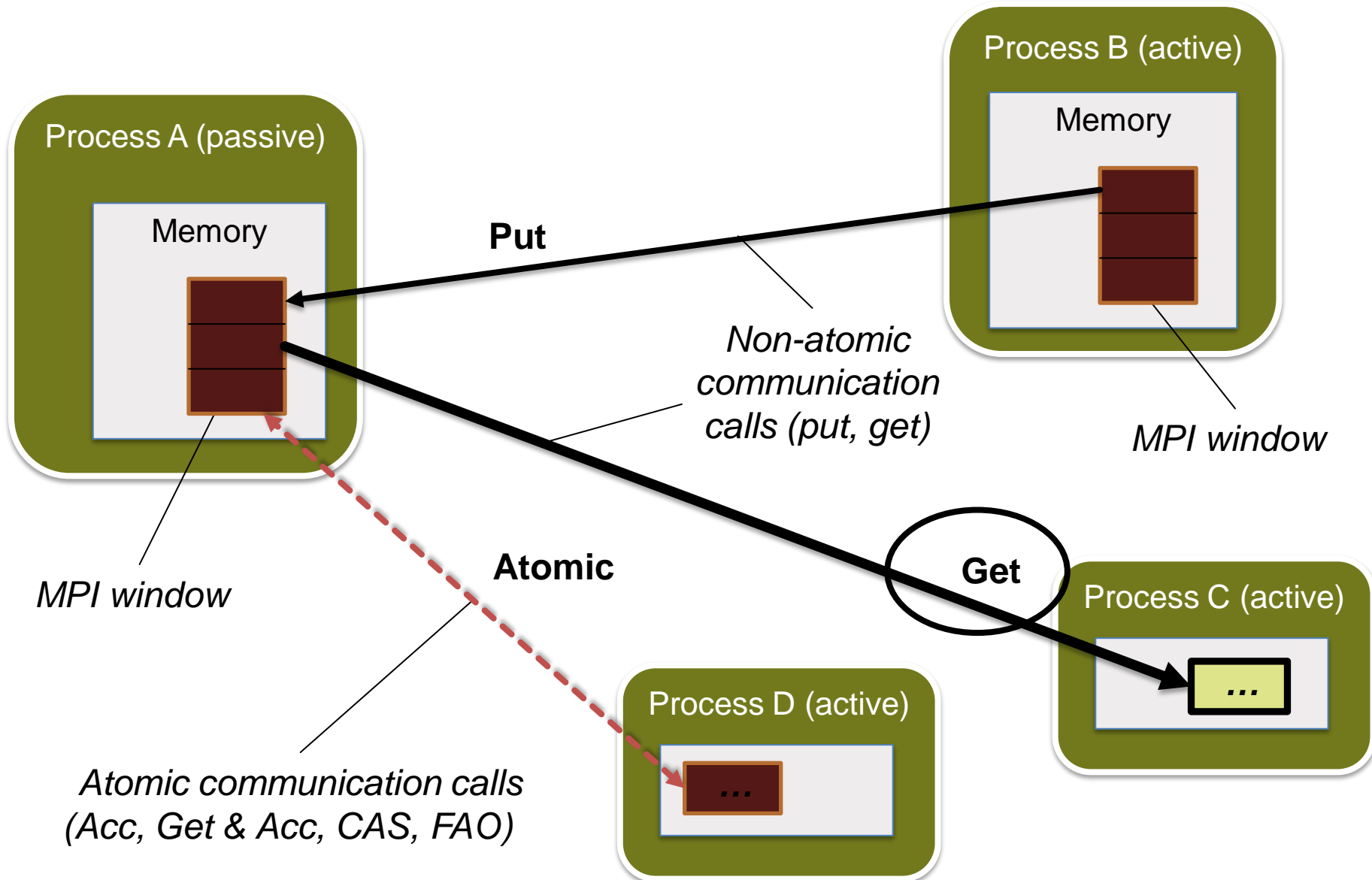
Process D (active)

...

7

# MPI-3 RMA COMMUNICATION OVERVIEW

# MPI-3 RMA COMMUNICATION OVERVIEW

# MPI-3 RMA COMMUNICATION OVERVIEW



Process B (active)

Memory

Process A (passive)

Memory

**Put**

*Non-atomic communication calls (put, get)*

*MPI window*

*MPI window*

**Atomic**

**Get**

Process C (active)

...

*MPI window*

Process D (active)

...

*Atomic communication calls (Acc, Get & Acc, CAS, FAO)*

10

# MPI-3 RMA COMMUNICATION OVERVIEW



Process A (passive)

Memory

Put

Process B (active)

Memory

*Non-atomic communication calls (put, get)*

*MPI window*

*MPI window*

**Atomic**

Get

Process C (active)

Process D (active)

...

...

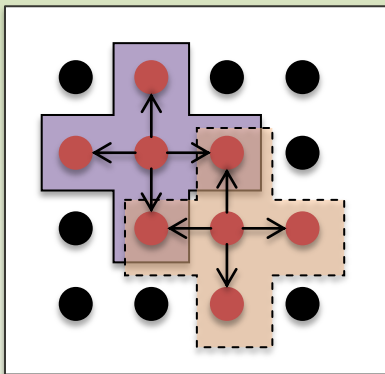*Atomic communication calls (Acc, Get & Acc, CAS, FAO)*

11

# MPI-3.0 RMA SYNCHRONIZATION OVERVIEW
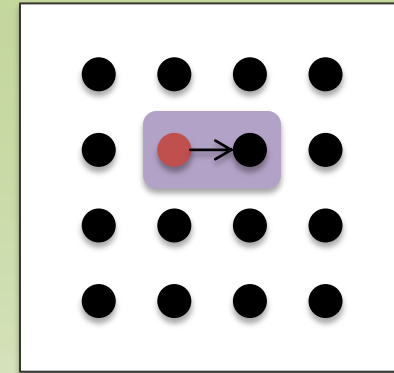


**Active Target Mode**

Fence
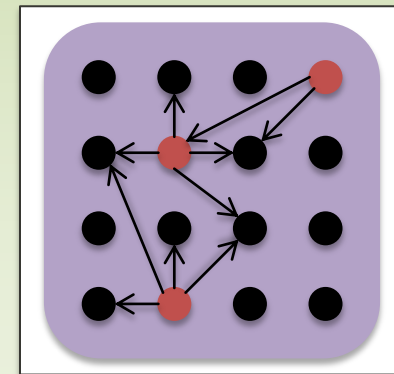
Post/Start/
Complete/Wait

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

● Active process

● Passive process
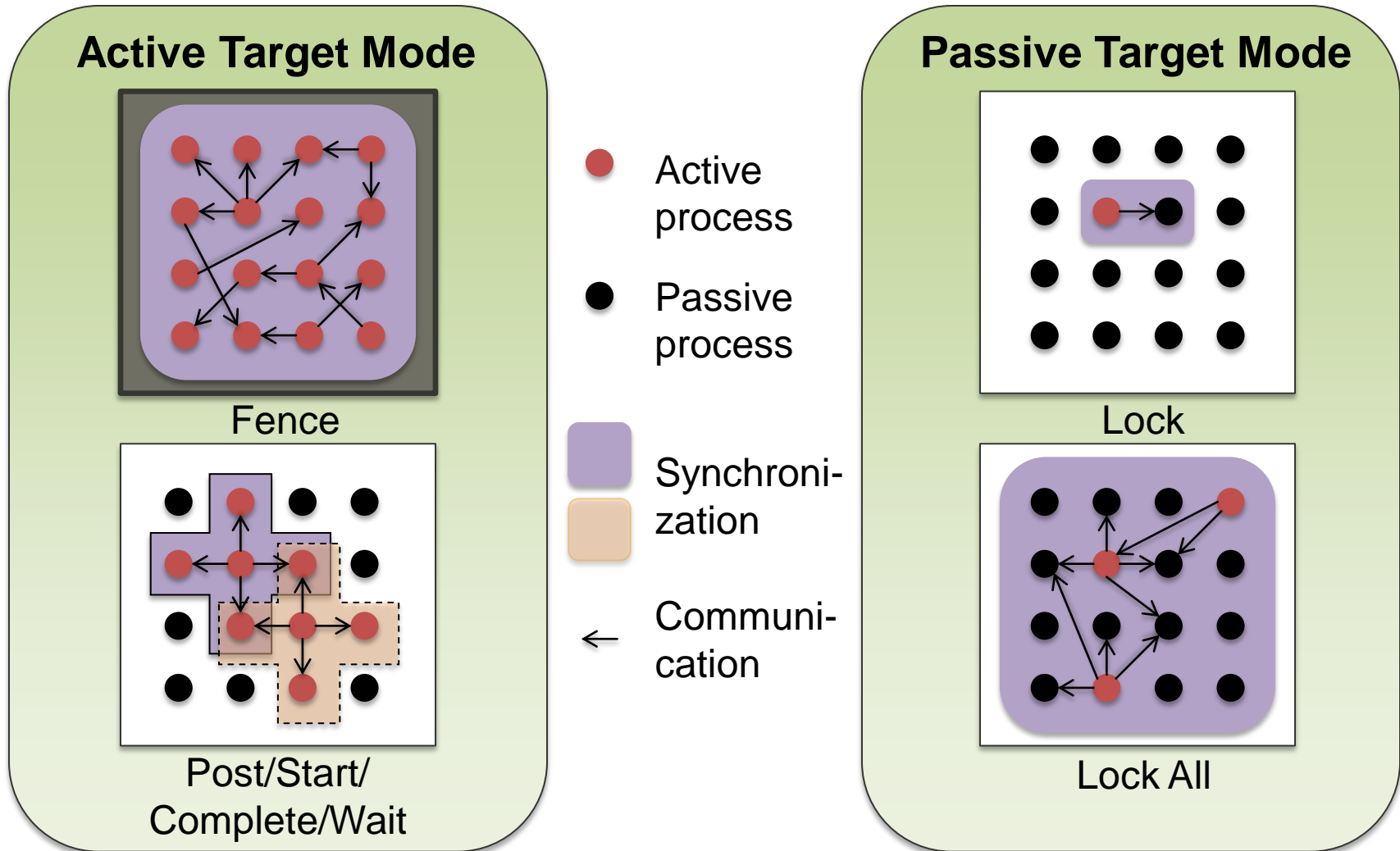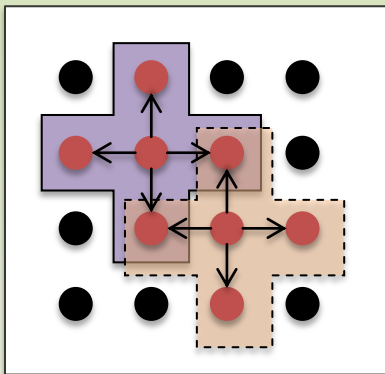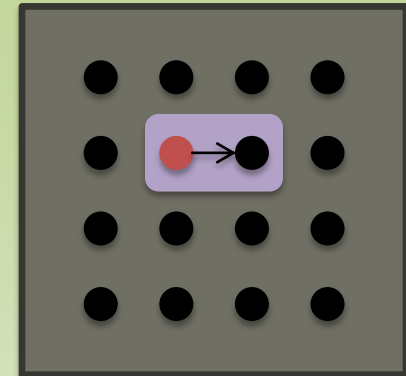
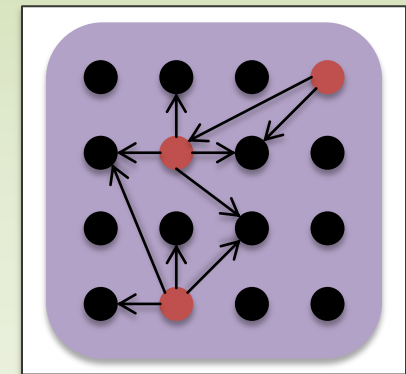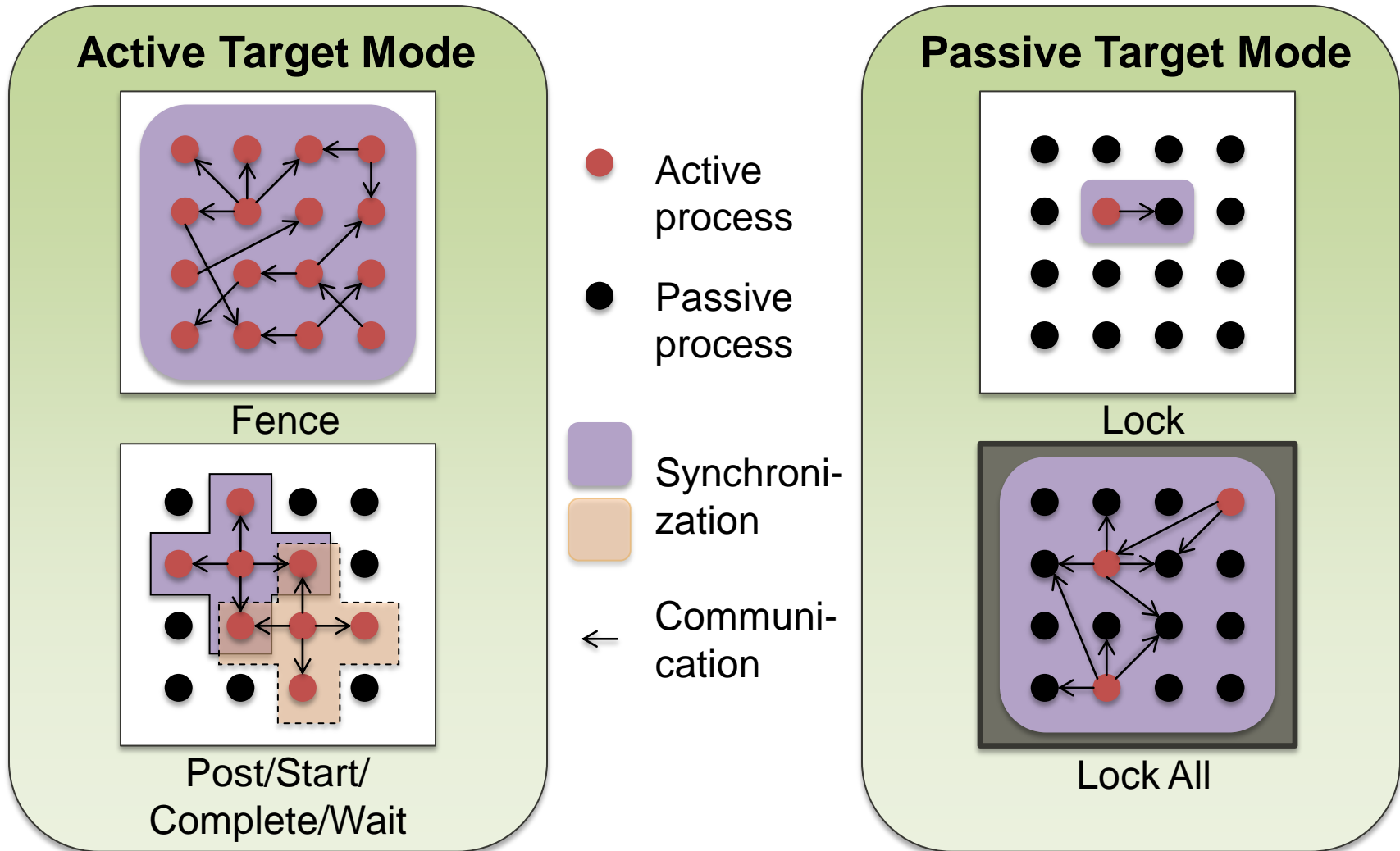■ Synchronization

■ Communication

**Passive Target Mode**
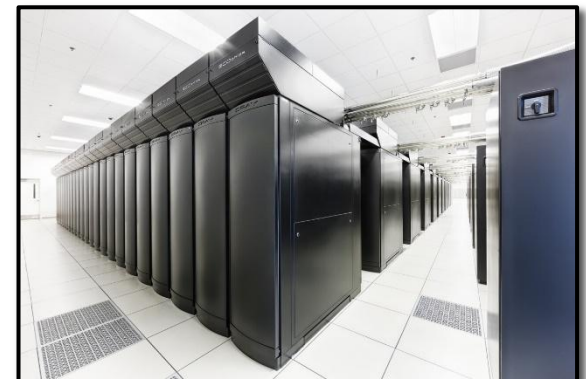
Lock

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

# 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

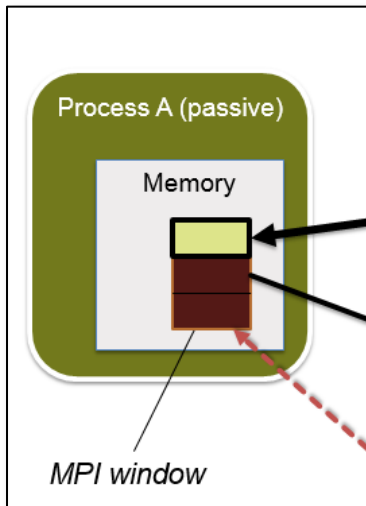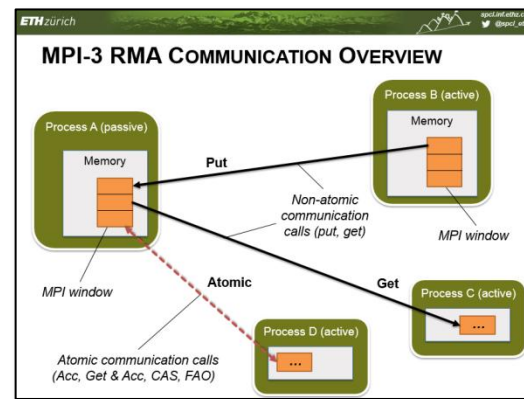**Passive Target Mode**

Lock

Lock All

# SCALABLE PROTOCOLS & REFERENCE IMPLEMENTATION

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

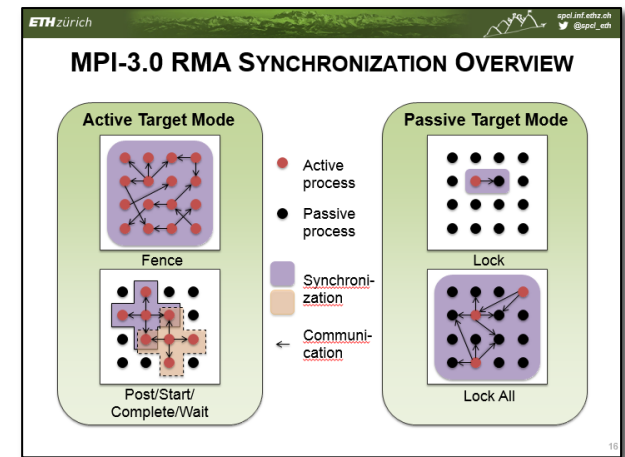# SCALABLE PROTOCOLS & REFERENCE IMPLEMENTATION

- 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



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



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

# PART 1: SCALABLE WINDOW CREATION
## Traditional windows



Process A — Memory — 0x111

Process B — Memory — 0x123

Process C — Memory — 0x120

$p$ = total number of processes

backwards compatible (MPI-2)

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

# PART 1: SCALABLE WINDOW CREATION
## Allocated windows



Process A — Memory — 0x123

Process B — Memory — 0x123

Process C — Memory — 0x123

$p$ = 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: $\mathcal{O}(\log p)$
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 64Bit types
  - ...or fall back to remote locking protocol
- MPI datatype handling with MPITypes library [1]
  - Fast path for contiguous data transfers of common intrinisic datatypes (e.g., MPI_DOUBLE)

MPI_Put

dmapp_put_nbi

Remote process

...

MPI_Compare _and_swap

dmapp_ acswap_qw_nbi

Remote process

...

**Contiguous memory**

[1] Ross, Latham, Gropp, Lusk, Thakur. Processing MPI datatypes outside MPI. EuroMPI/PVM'09

# PERFORMANCE INTER-NODE: LATENCY

# PERFORMANCE INTRA-NODE: LATENCY



Put/Get Intra-Node

**3x faster**

Half ping-pong

# PERFORMANCE: OVERLAP



Inter-Node Overlap in %

Useful for, e.g., scientific codes:

AWM-Olsen seismic

3D FFT

MILC

# PERFORMANCE: MESSAGE RATE



## Inter-Node



## Intra-Node

# PERFORMANCE: ATOMICS



hardware-
accelerated
protocol:
*lower latency*

fall-back
protocol:
*higher
bandwidth*

proprietary

64Bit integers

# PART 3: SYNCHRONIZATION



**Active Target Mode**

Fence

Post/Start/
Complete/Wait

● Active process

● Passive process

Synchroni-zation

Communi-cation

**Passive Target Mode**

Lock

Lock All

# 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;
}
```

Remote 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



**Puts**

Posting process

Starting process

**Puts**

...

...

Proc 0   Proc 1

post

start

*matching algorithm*

allows to access other processes

*access epoch*

*exposure epoch*

allows access from other processes

complete

*matching algorithm*

wait

# 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



Posting process
(opens its window)

# 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



Starting processes
(access remote window)

# PSCW SCALABLE POST/START MATCHING

- Each starting process has a local list



Local list

# PSCW SCALABLE POST/START MATCHING

- *Posting* process *i* adds its rank *i* to a list at each *starting* process $j_1, \ldots, j_4$

- 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 $j_1, \ldots, j_4$

- 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 $j_1, \ldots, j_4$

- 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 $j_1, \ldots, j_4$

- 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 $j_1, \ldots, j_4$

- 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 $j_1, \ldots, j_4$

- Each *starting* process *j* waits until the rank of the *posting* process *i* is present in its local list
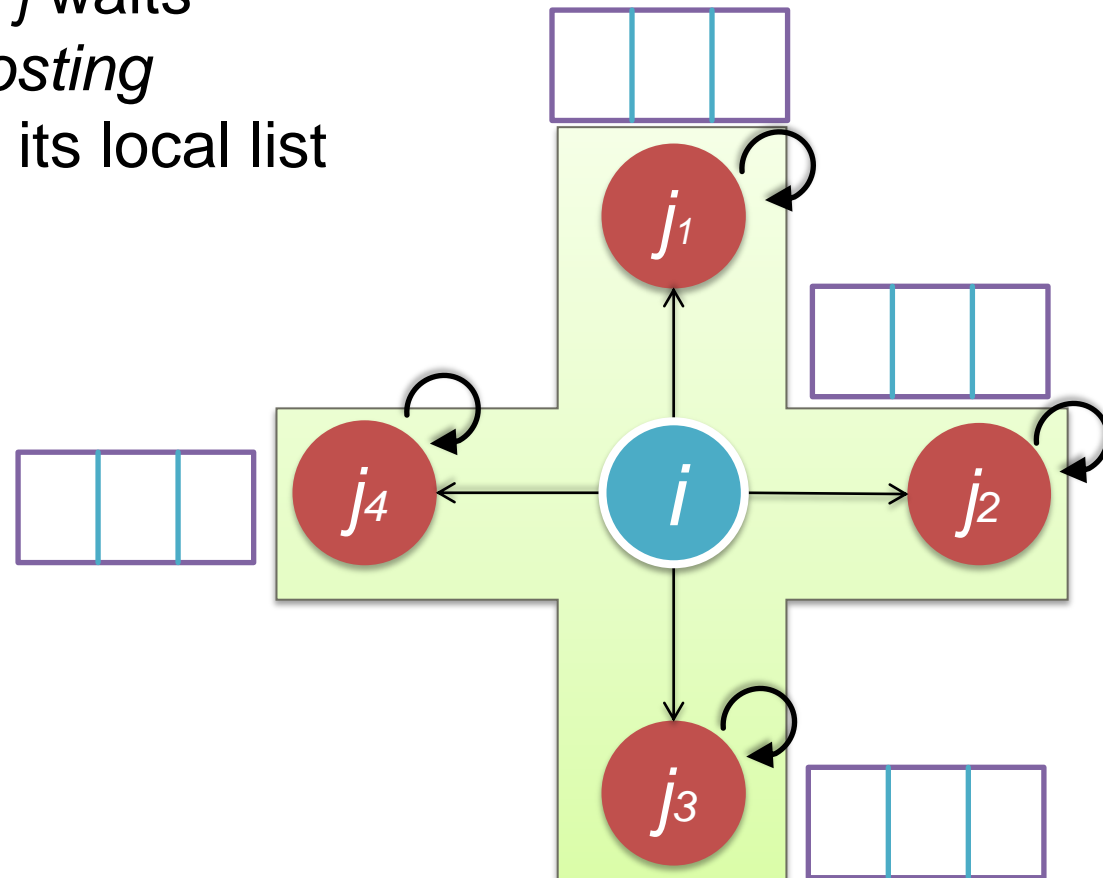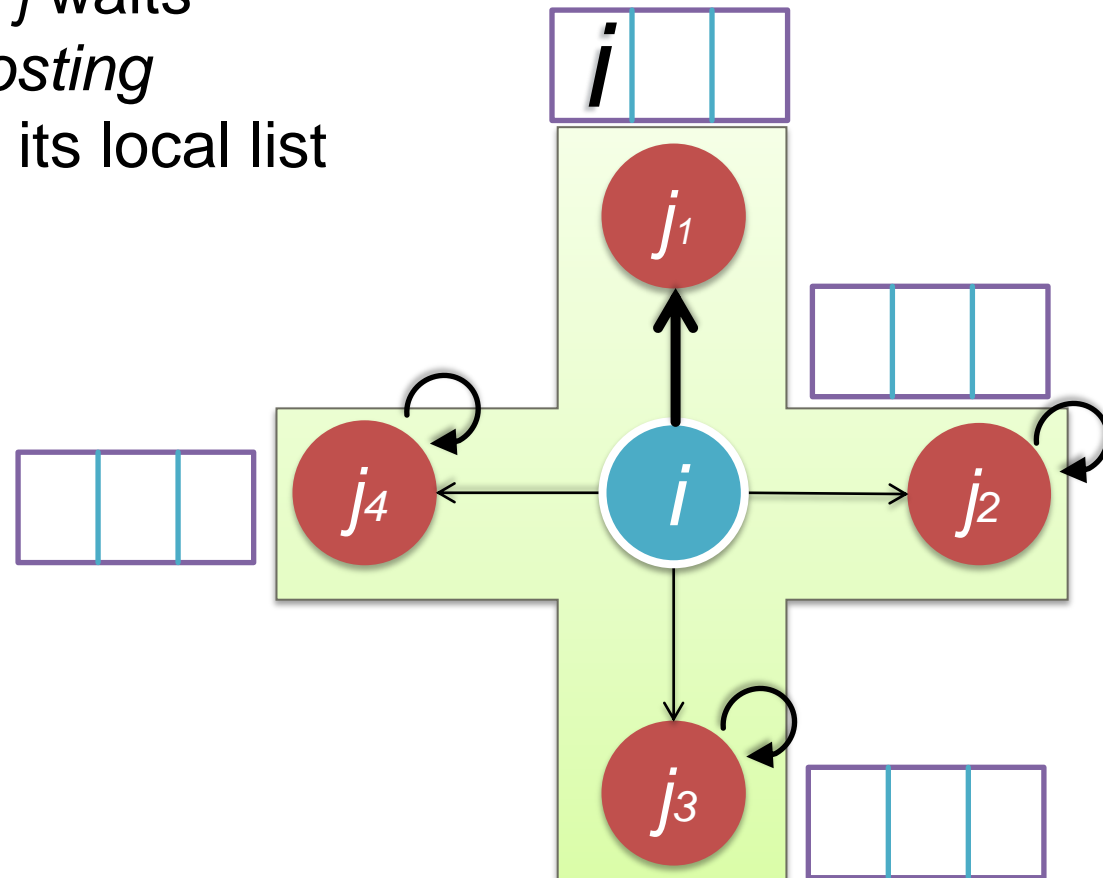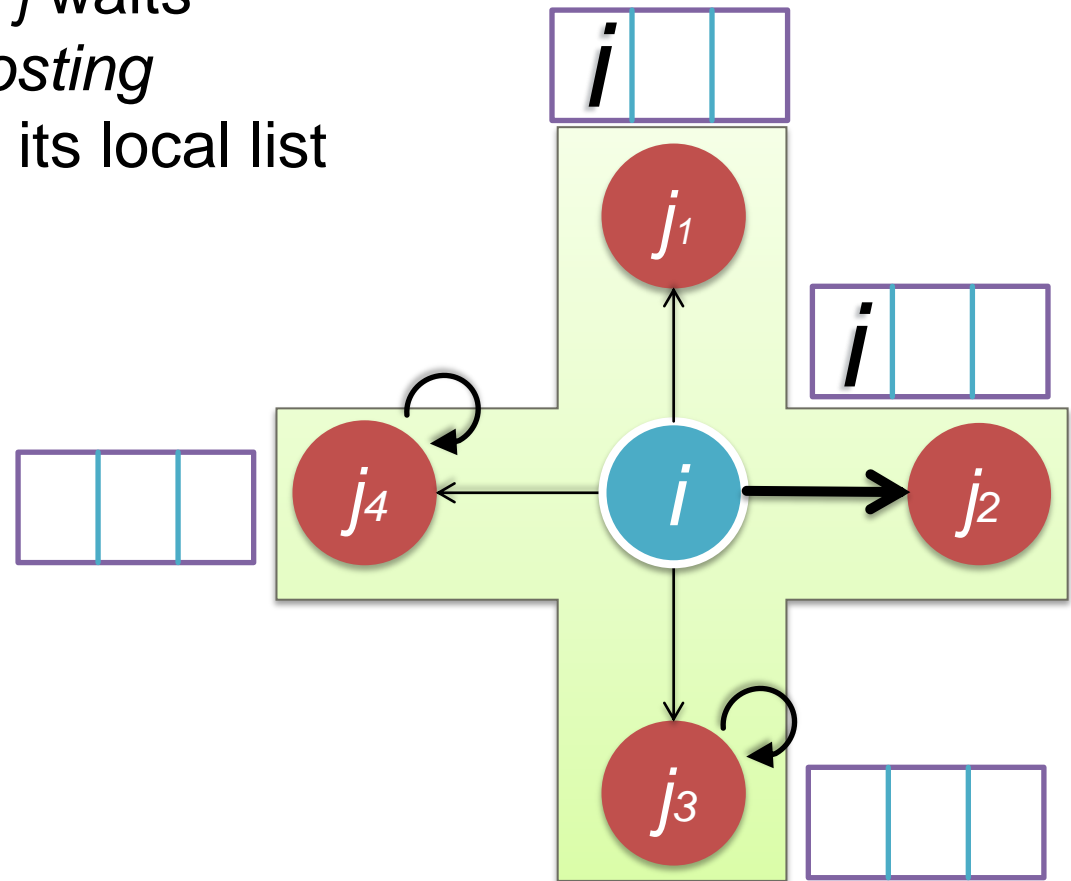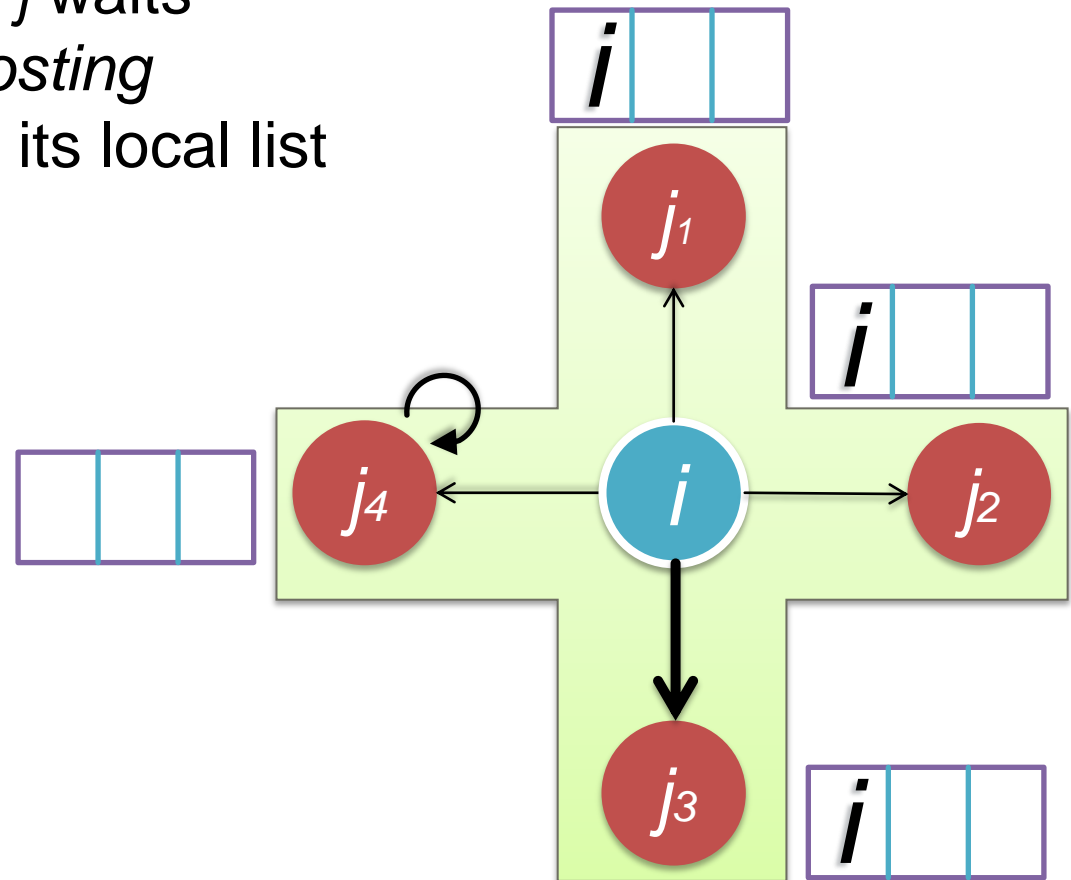
# PSCW SCALABLE POST/START MATCHING

- *Posting* process *i* adds its rank *i* to a list at each *starting* process $j_1, \ldots, j_4$

- Each *starting* process *j* waits until the rank of the *posting* process *i* is present in its local list
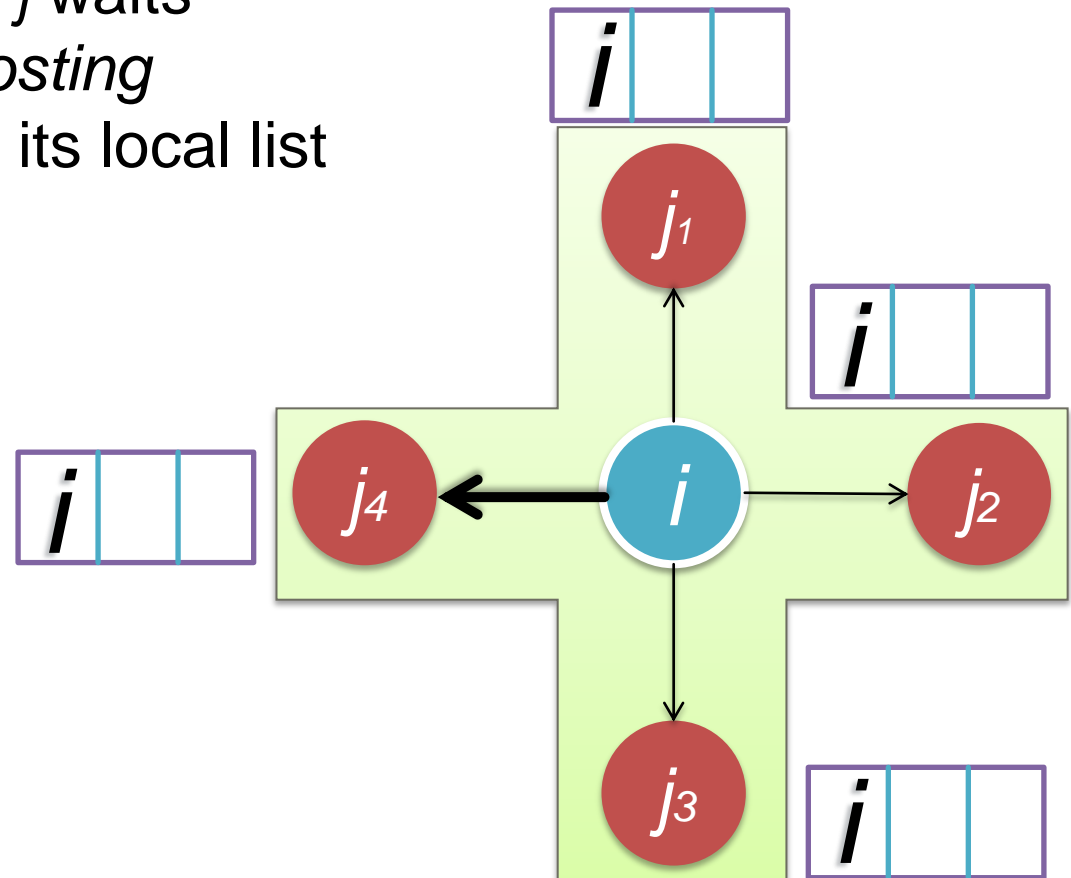
# PSCW SCALABLE COMPLETE/WAIT MATCHING

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

# PSCW Scalable Complete/Wait Matching

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

# PSCW Scalable Complete/Wait Matching

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

# PSCW Scalable Complete/Wait Matching

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

# PSCW SCALABLE COMPLETE/WAIT MATCHING

- Each *starting* 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
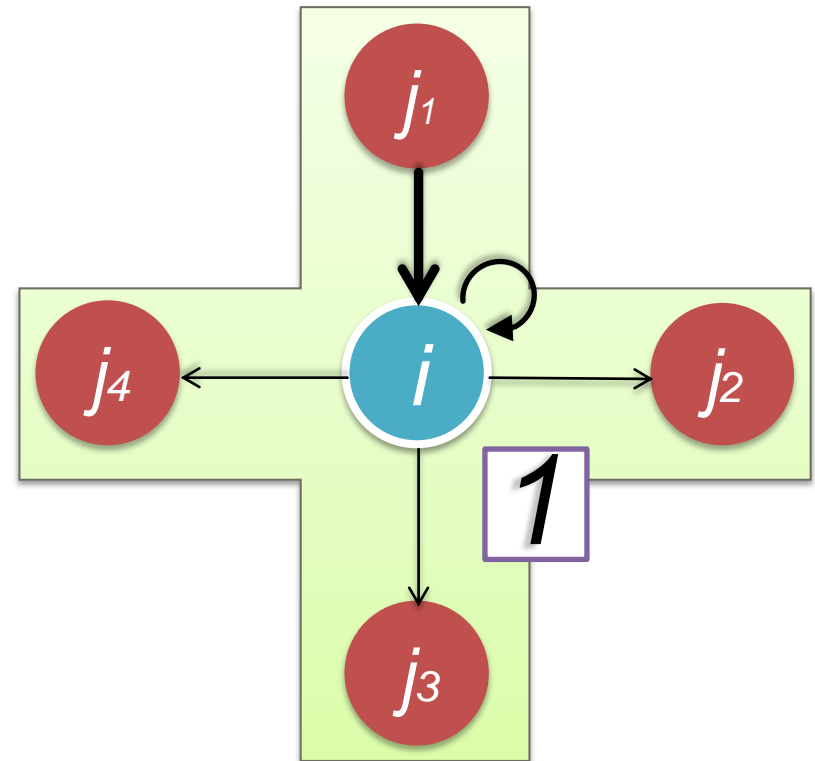
# PSCW PERFORMANCE

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



**Ring Topology**

# SCALABLE LOCK SYNCHRONIZATION

Lock/Unlock
(shared/exclusive)

Lock All
(always shared)

● Active process

● Passive process

Two-level lock hierarchy:

Master Process



| | Process 0 | Process 1 | Process P-1 |
|---|---|---|---|

Shared Counter

Exclusive Counter

Memory

global: | 000 | 000 |

local: | 00000 | 0 |

local: | 00000 | 0 |

local: | 00000 | 0 |

Shared Counter    Exclusive Bit

Shared Counter    Exclusive Bit

Shared Counter    Exclusive Bit

# EXCLUSIVE LOCAL LOCK: TWO PHASES

- PHASE 1: increment the global exclusive counter
  (Invariant 1: no global shared lock held concurrently)

Proc 2 wants to lock Proc 1 exclusively



Process 1

00000 | 0

Process 0

000 | 000

00000 | 0

Process 2

00000 | 0

58

# EXCLUSIVE LOCAL LOCK: TWO PHASES

- **PHASE 1: increment the global exclusive counter**
  (Invariant 1: no global shared lock held concurrently)



Proc 2 wants to lock Proc 1 exclusively

Process 1

Process 0

Process 2

000  001

00000  0

00000  0

00000  0

fetch-add

MPI_Win_lock( EXCL, 1 )

000  000

# EXCLUSIVE LOCAL LOCK: TWO PHASES

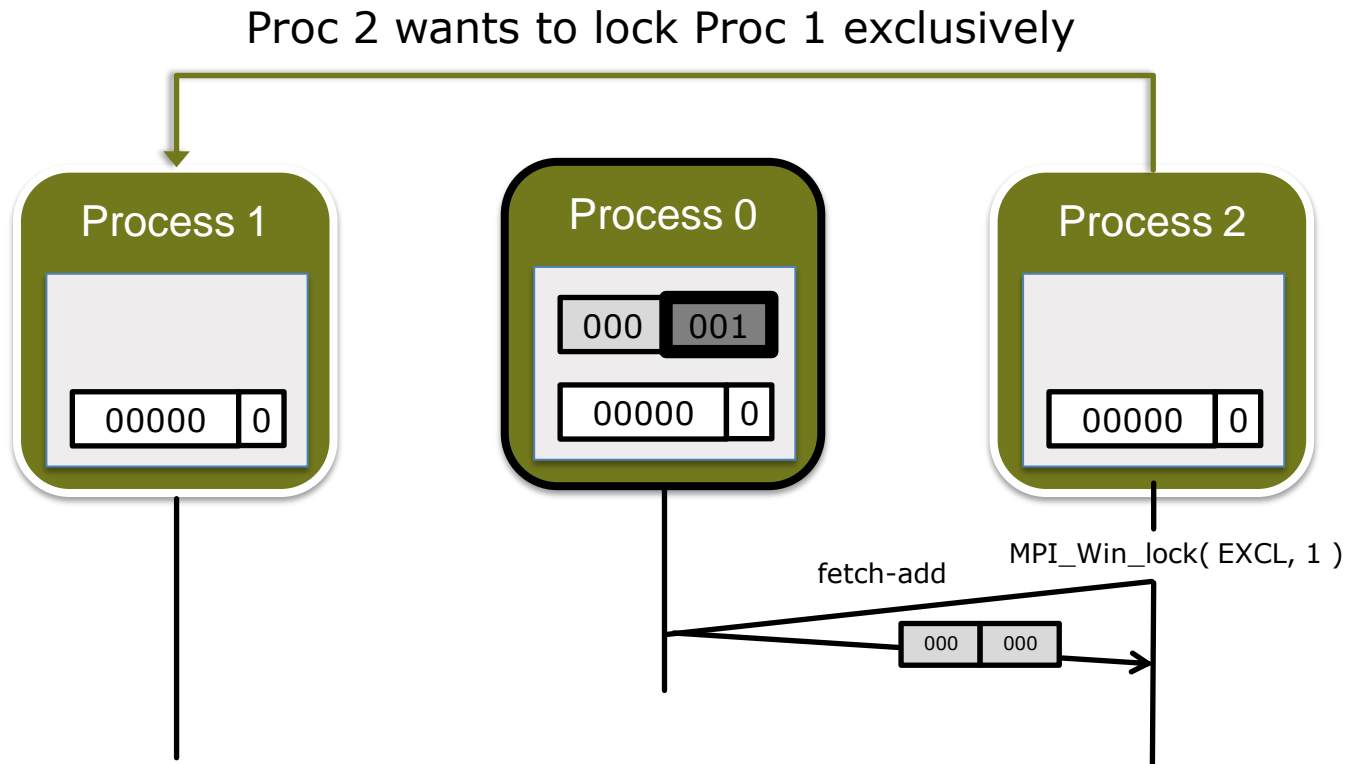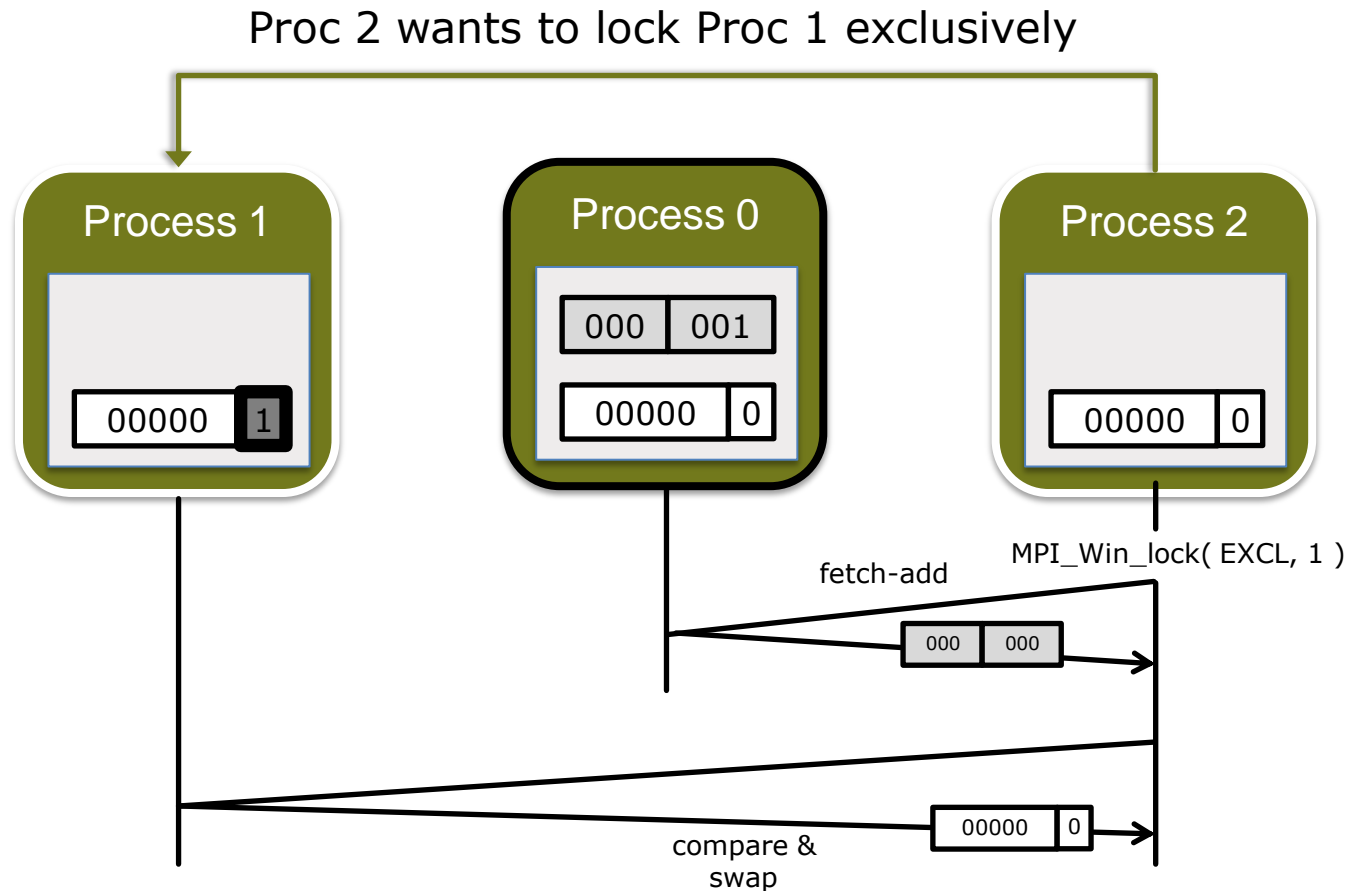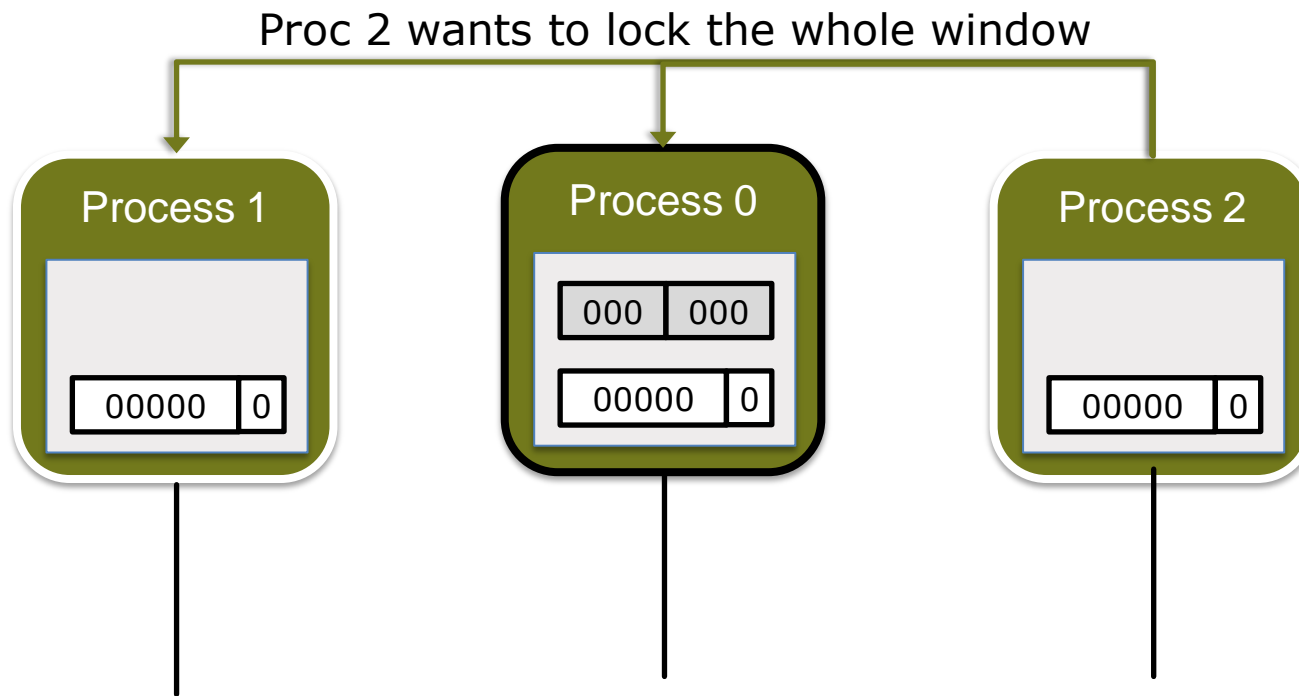- PHASE 1: increment the global exclusive counter
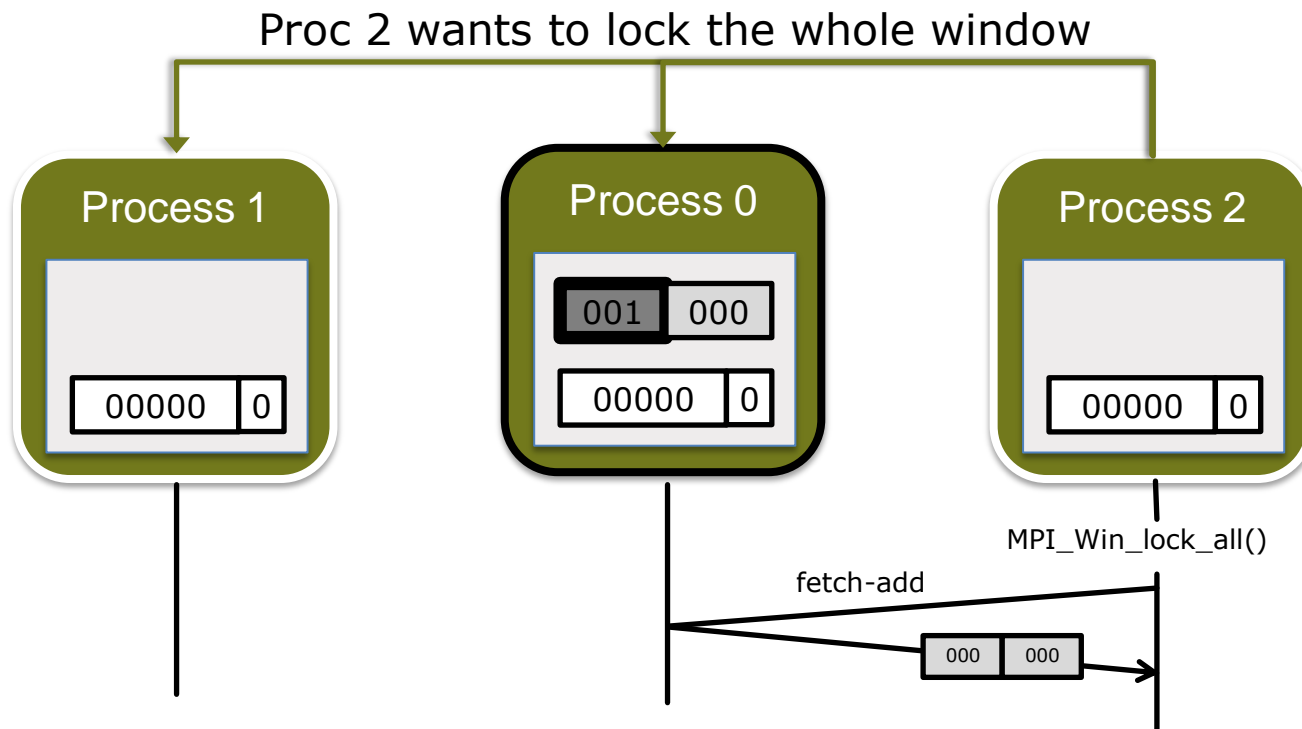  (Invariant 2: no local shared/exclusive lock held concurrently)



Proc 2 wants to lock Proc 1 exclusively

# SHARED GLOBAL LOCK: ONE PHASE

- Increment global shared counter
  (Invariant: no local exclusive lock is held concurrently)

Proc 2 wants to lock the whole window

# SHARED GLOBAL LOCK: ONE PHASE

- Increment global shared counter
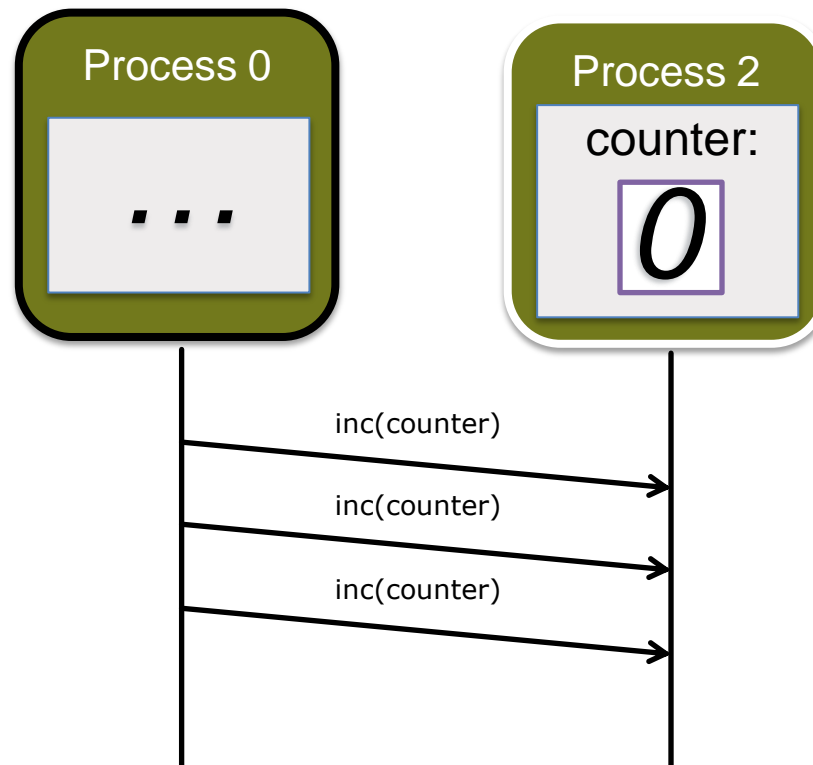  (Invariant: no local exclusive lock is held concurrently)



Proc 2 wants to lock the whole window

Process 1

Process 0

001   000

00000   0

Process 2

00000   0

00000   0

MPI_Win_lock_all()

fetch-add

000   000

- Constant number of operations for $p$ processes ☺

# FLUSH SYNCHRONIZATION

| Time bound | $\mathcal{O}(1)$ |
|---|---|
| Memory bound | $\mathcal{O}(1)$ |

- 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

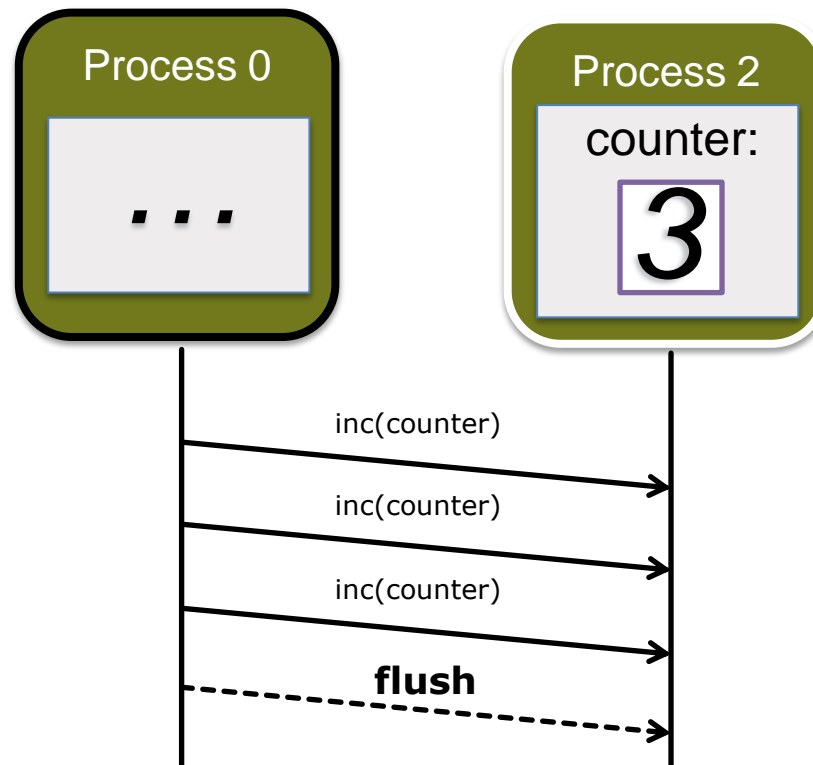| Time bound | $\mathcal{O}(1)$ |
|---|---|
| Memory bound | $\mathcal{O}(1)$ |

- 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
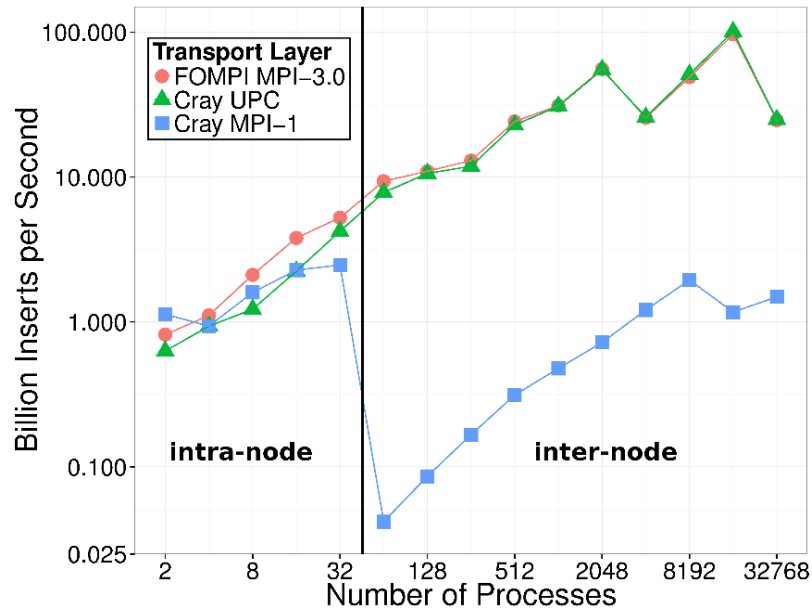
# PERFORMANCE

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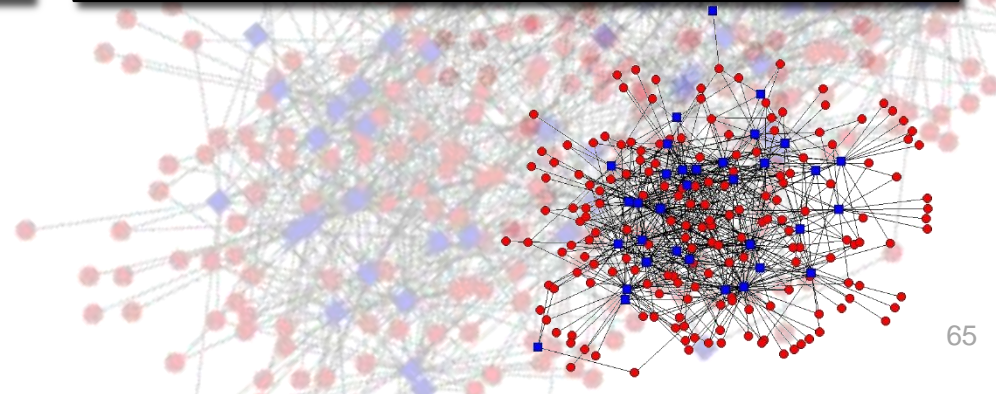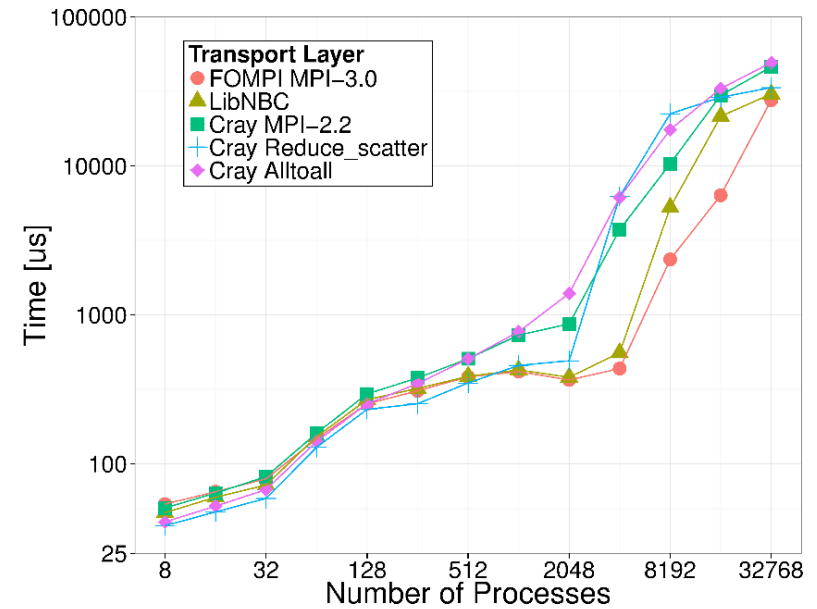- 4 applications
- One nearly full-scale run ☺

# PERFORMANCE: MOTIF APPLICATIONS
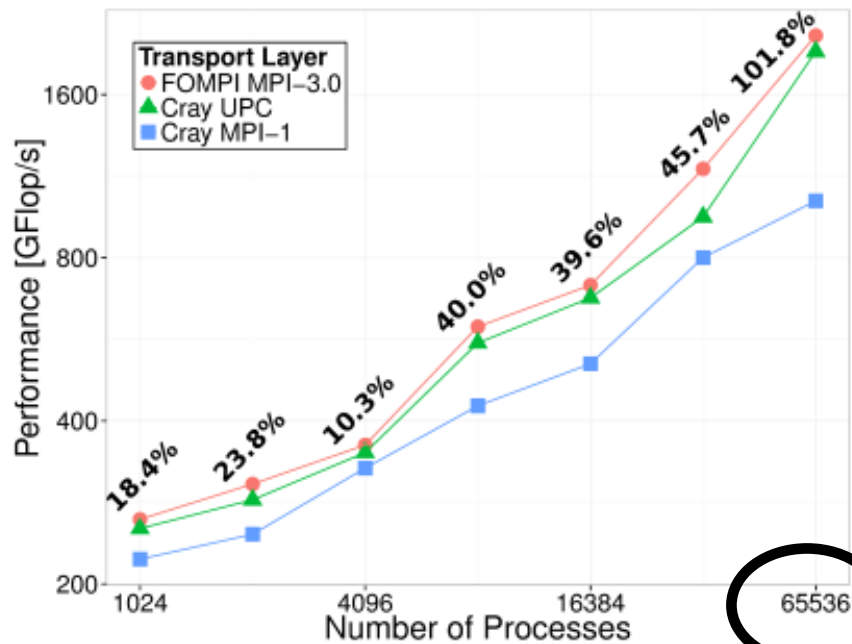


Key/Value Store: Random Inserts per Second

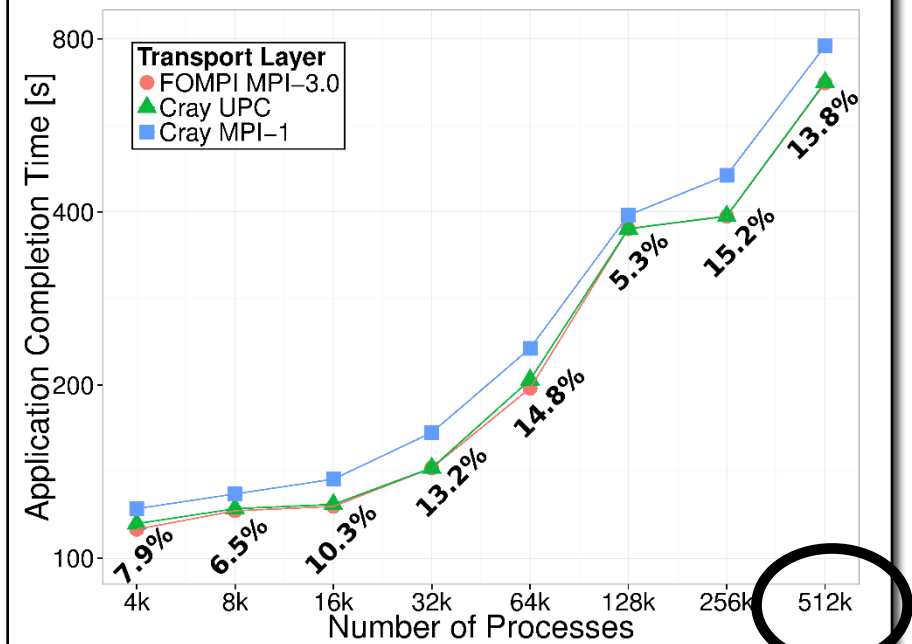Dynamic Sparse Data Exchange (DSDE) with 6 neighbors

# PERFORMANCE: APPLICATIONS

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



NAS 3D FFT [1] Performance

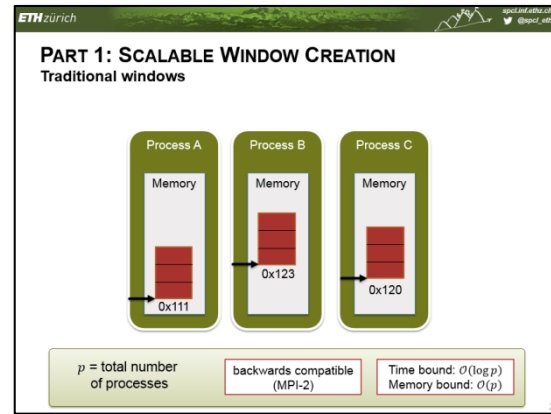MILC [2] Application Execution Time

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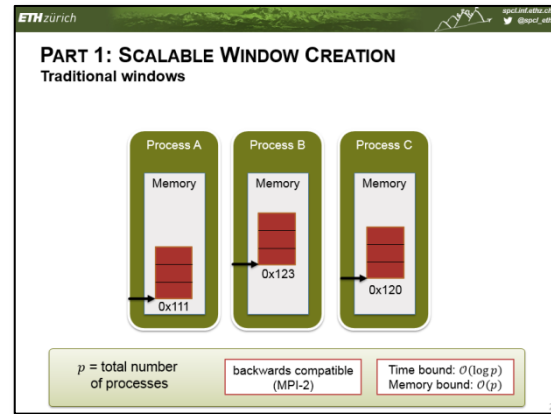[2] Shan et al. Accelerating applications at scale using one-sided communication. PGAS'12

# CONCLUSIONS & SUMMARY

# CONCLUSIONS & SUMMARY
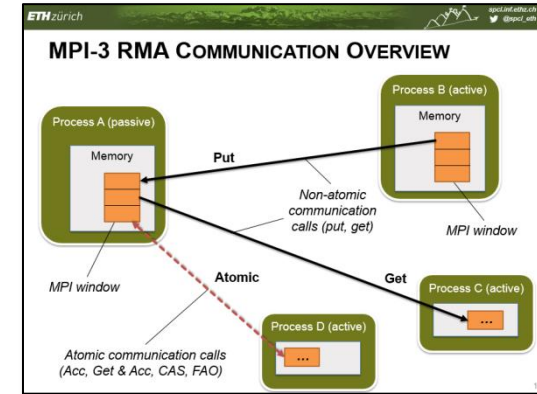


1. MPI window creation routines

# CONCLUSIONS & SUMMARY



1. MPI window creation routines



2. Non-atomic & atomic communication

# CONCLUSIONS & SUMMARY



3. Fence / PSCW



1. MPI window creation routines



2. Non-atomic & atomic communication

# CONCLUSIONS & SUMMARY



3. Fence / PSCW

1. MPI window ... & atomic ... communication routines

4. Locks

# CONCLUSIONS & SUMMARY



3. ... & atomic ...cation

4. Locks

5. foMPI reference implementation

# CONCLUSIONS & SUMMARY



3.

## PERFORMANCE: APPLICATIONS

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

**NAS 3D FFT [1] Performance**

**MILC [2] Application Execution Time**

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

**PSCW Scalable Post/Start Matching**

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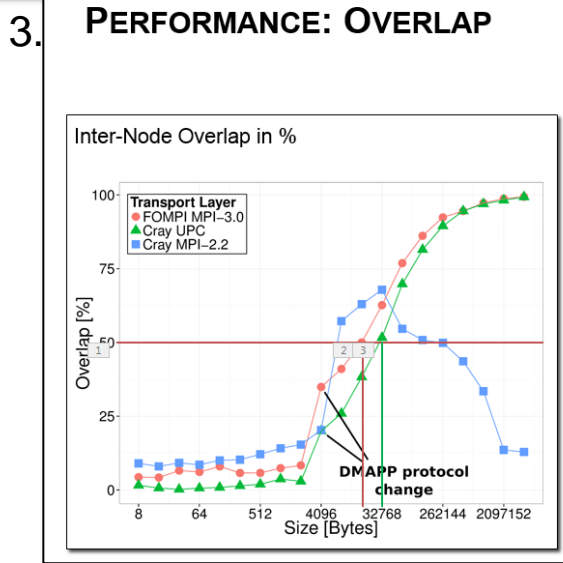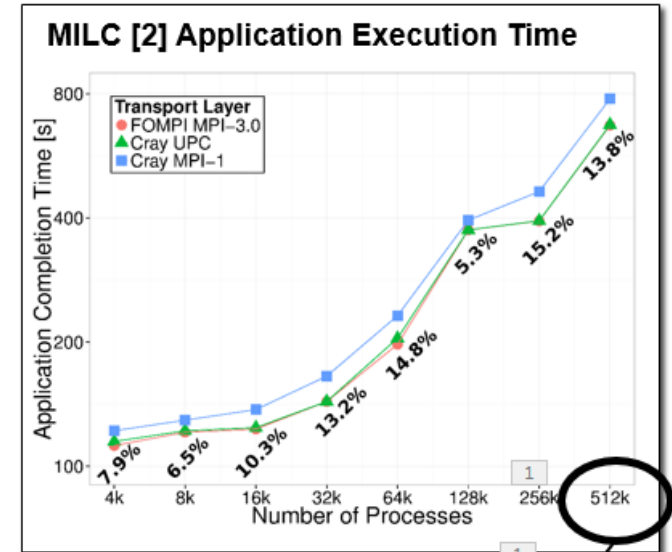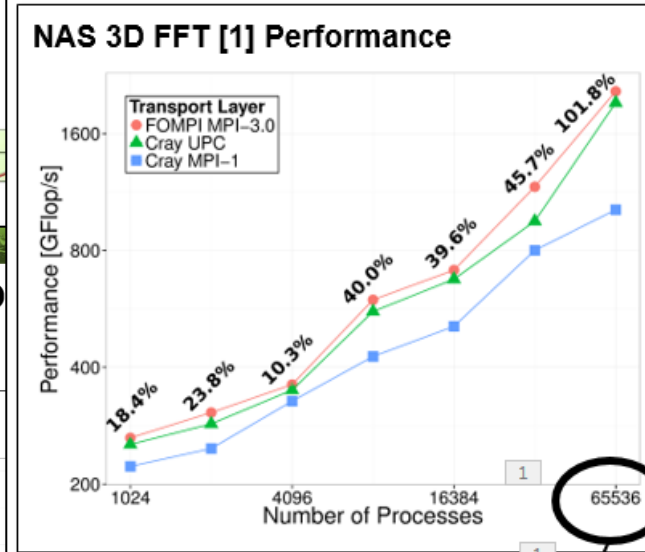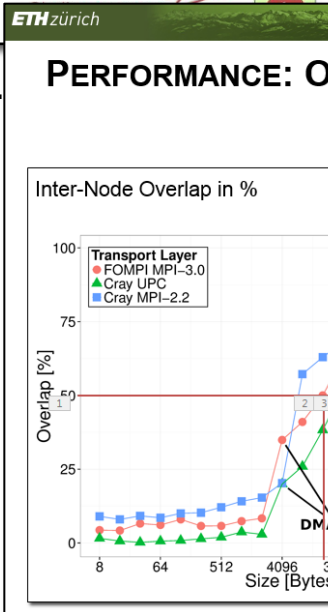- In this example there is one *posting* and 4 *starting* processes

**PERFORMANCE: O**

Inter-Node Overlap in %

5. foMPI reference implementation

6. Application implementation & evaluation

# CONCLUSIONS & SUMMARY



## PERFORMANCE MODELLING

### Performance functions for synchronization protocols

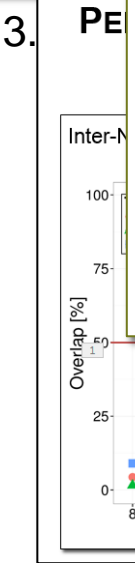| Fence | $\mathcal{P}_{fence} = 2.9\mu s \cdot \log_2(p)$ |
|---|---|
| PSCW | $\mathcal{P}_{start} = 0.7\mu s, \mathcal{P}_{wait} = 1.8\mu s$ $\mathcal{P}_{post} = \mathcal{P}_{complete} = 350ns \cdot k$ |
| Locks | $\mathcal{P}_{lock,excl} = 5.4\mu s$ $\mathcal{P}_{lock,shrd} = \mathcal{P}_{lock\_all} = 2.7\mu s$ $\mathcal{P}_{unlock} = \mathcal{P}_{unlock\_all} = 0.4\mu s$ $\mathcal{P}_{flush} = 76ns$ $\mathcal{P}_{sync} = 17ns$ |

### Performance functions for communication protocols

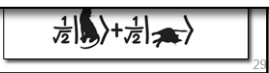| Put/get | $\mathcal{P}_{put} = 0.16ns \cdot s + 1\mu s$ $\mathcal{P}_{get} = 0.17ns \cdot s + 1.9\mu s$ |
|---|---|
| Atomics | $\mathcal{P}_{acc,sum} = 28ns \cdot s + 2.4\mu s$ $\mathcal{P}_{acc,min} = 0.8ns \cdot s + 7.3\mu s$ |

[1] Nishtala et al. Scaling communication-intensive applications on BlueGene/P using one-sided communication and overlap. IPDPS'09
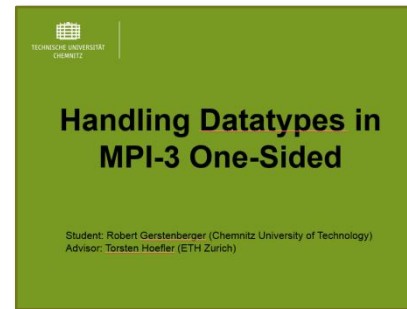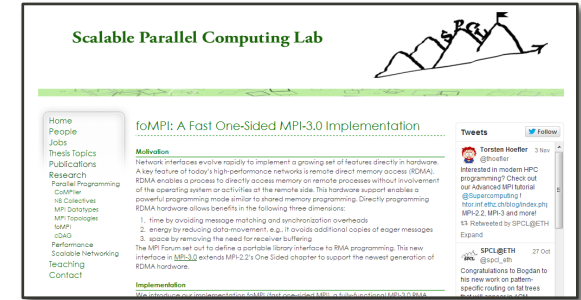[2] Shan et al. Accelerating applications at scale using one-sided communication. PGAS'12

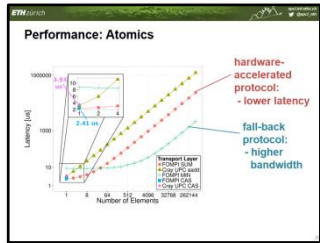5. foMPI reference implementation

6. Application implementation & evaluation

74

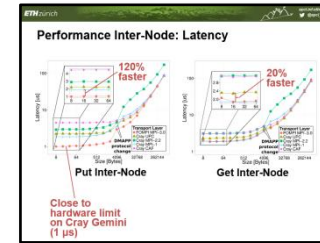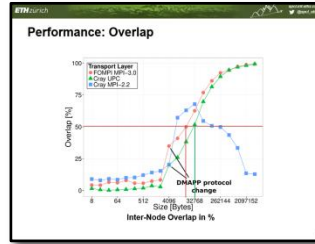**ETH**_zürich_

# ACKNOWLEDGMENTS

- Try foMPI yourself:
http://spcl.inf.ethz.ch/Research/
Parallel_Programming/foMPI

- Ongoing work on
DDT (see ACM Students
Research Competition Poster)

- Thanks to: Timo Schneider,
Greg Bauer, Bill Kramer,
Duncan Roweth, Nick Wright,
Paul Hargrove (and the whole UPC team)
and the MPI Forum RMA WG …

… and the institutions:

# Thank you
# for your attention