



Efficient Implementation of MPI-2 Passive One-Sided Communication on InfiniBand Clusters

W. Jiang, J. Liu, H. Jin, D.K. Panda
D. Buntinas, R. Thakur, W. Gropp

The Ohio State University
Argonne National Laboratory

Outline

- ◆ Background
- ◆ Problem Statement
- ◆ Design and Implementation
- ◆ Experimental Results
- ◆ Conclusion



n i n

ag Pa ing n a



a u

- On Sided Communication
- Dynamic Process Management
- Parallel I/O support



P H

n a i n

P

L

One sided Communication

◆ a ing a i n i n u i d
a

◆ a a a i n a i d d i n i a
a i n

◆ u i u n i a i n a n

◆ u d i a n i n a i n a
a i a i n

One i e uni ti n

- **Terminology**

- Origin
- Target
- Window

- **Communication**

- MPI_Put
- MPI_Get
- MPI_Accumulate

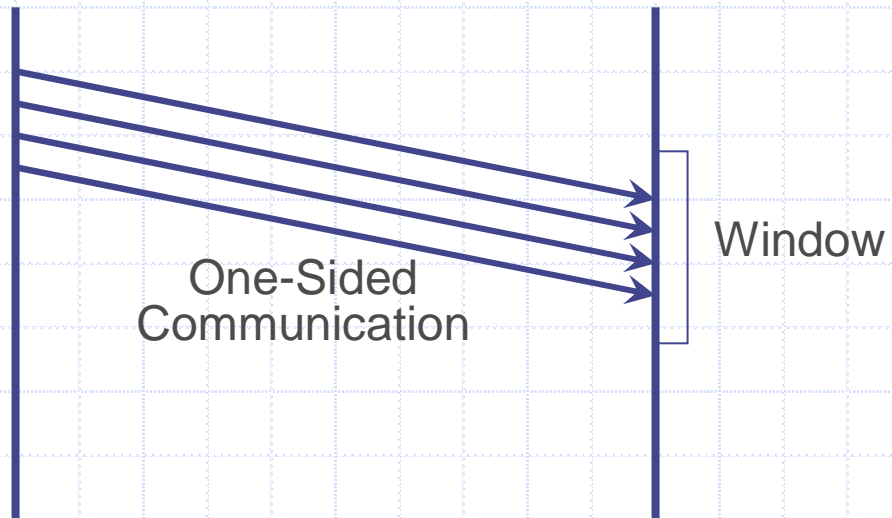
- **Synchronization**

- Active target
- Passive target

- *MPI-2 one-sided communication decouples synchronization from communication functions*

Origin

Target



Initiation
unit

One side



Initiation

- The target is actively involved in synchronization
- Previously implemented active mode one sided communication using RDMA based communication on InfiniBand

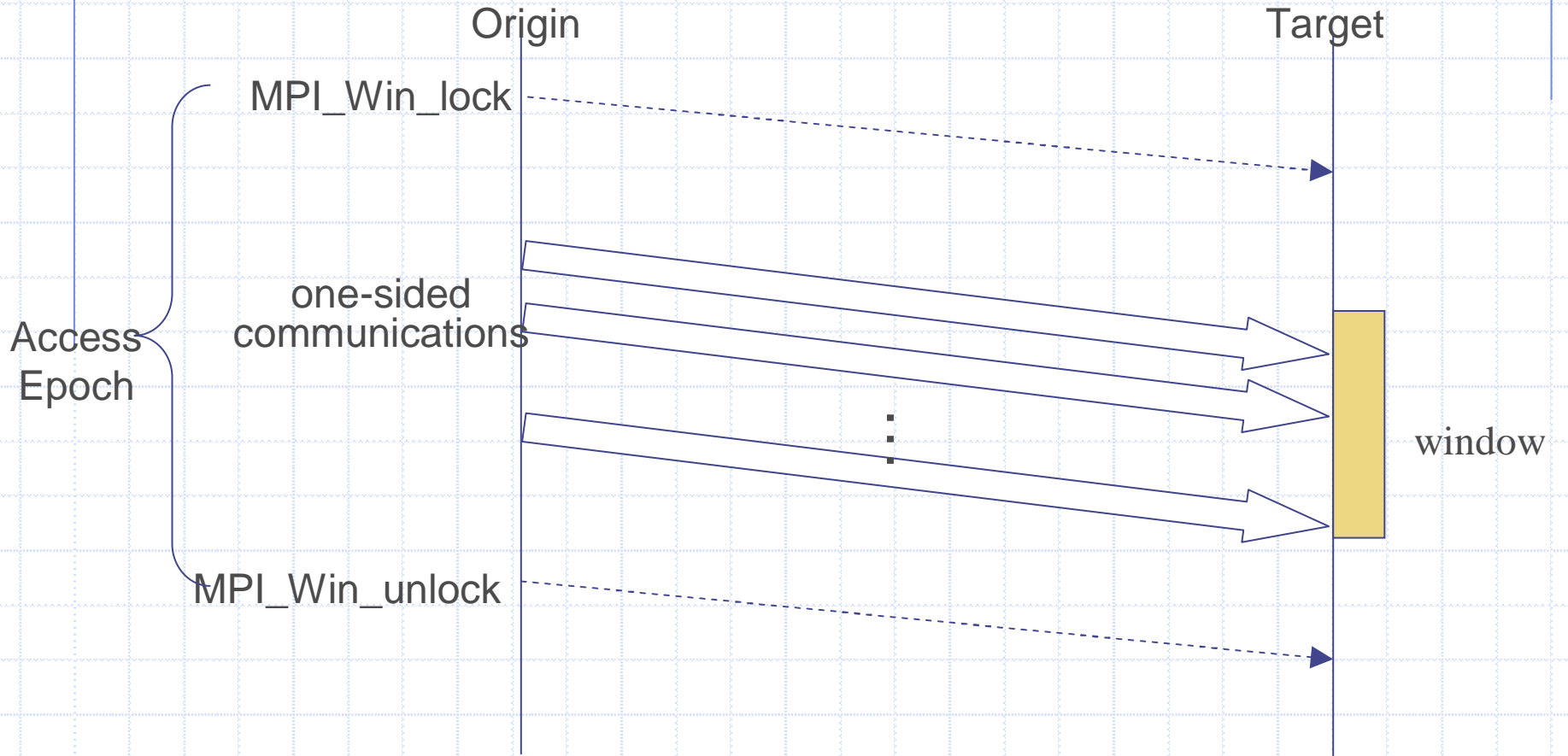
◆ *High Performance ~~MP~~ One-Sided Communication over InfiniBand. Weihang Jiang, Jiuxing Liu, Hyun-Wook Jin, Dhableswar K. Panda and William Gropp and Rajeev Thakur . CCGRID[2004]*



Passive Initiation

- The target is not explicitly involved in synchronization
- Main focus of our work here

One initiator unit



InfiniBand

- ◆ Performance and Scalability
 - ◆ Performance
 - 5 usec latency
 - 10Gbps throughput
 - ◆ Scalability
 - Remote Direct Memory Access (RDMA)
 - ◆ RDMA write
 - ◆ RDMA read
 - Atomic operations
 - ◆ Compare and Swap
 - ◆ Fetch and Add

Outline

- ◆ Background
- ◆ Problem Statement
- ◆ Design and Implementation
- ◆ Evaluation results
- ◆ Conclusion

o em tatement

◆ esi n e i ient assive syn hroni ation y
ta in a vanta e o n ini an 's a van e
eat res

◆ ore i erent a roa hes an strate ies

◆ va ate these a roa hes

Outline

- ◆ Bac round
- ◆ Problem statement
- ◆ Design and Implementation
- ◆ Experimental results
- ◆ Conclusion

esi n

- ◆ esi n C allen es
 - o Overhea or syn hroni ation
 - n e en ent ro ress
 - on rrent o ni ation

- ◆ o road approac es
 - Threa ase esi n
 - to i O eration ase esi n

e e e i n

◆ Variations of thread based design

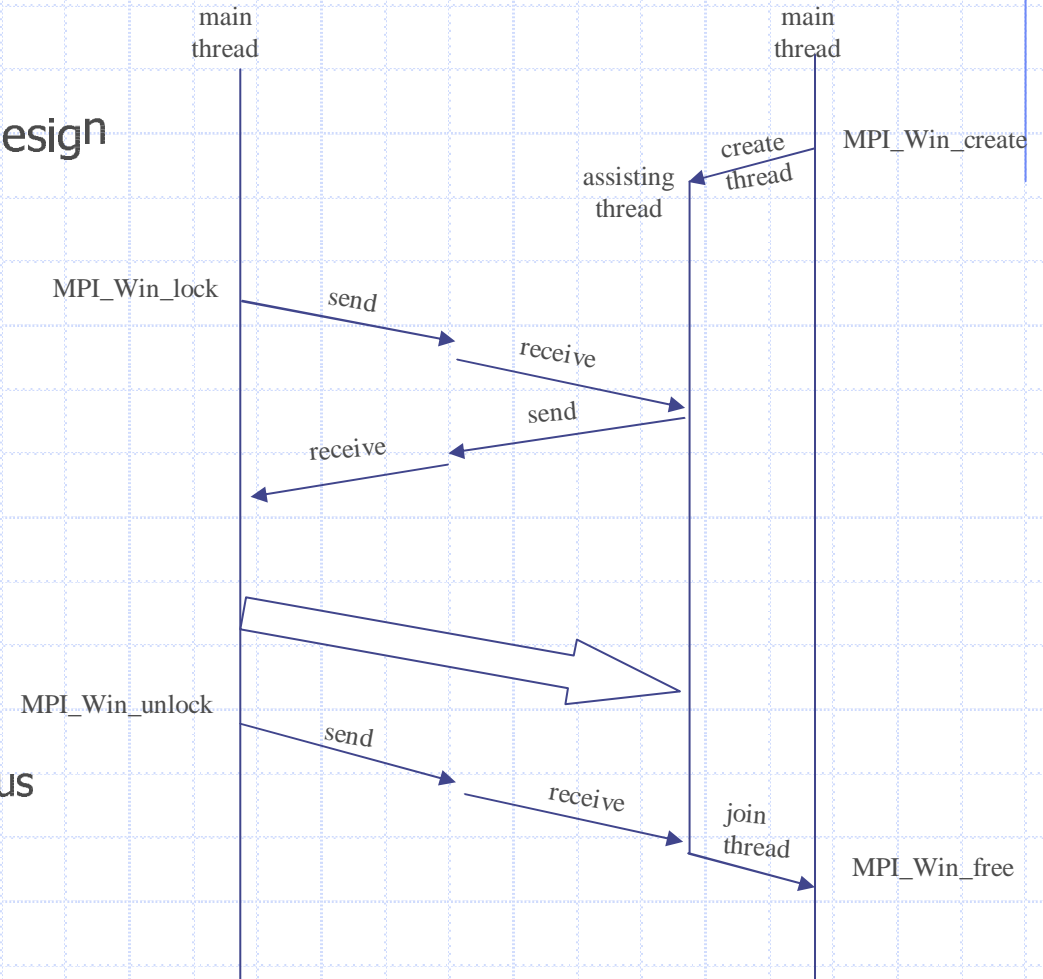
- Dedicated thread
- Multiple thread
- Event Driven thread

◆ Pros:

- Simple and portable

◆ Cons:

- CPU cycle consuming
- Can not handle simultaneous requests well



time contention

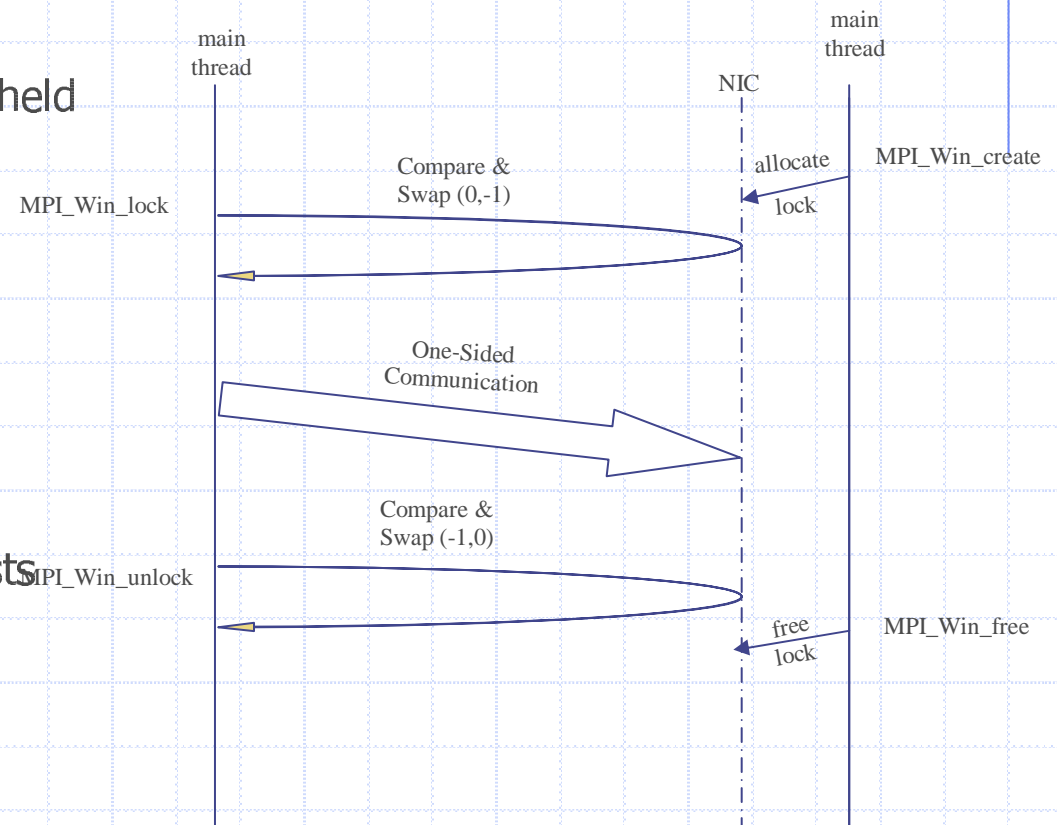
- ◆ ‘Test & Set’ based design
 - Flag used to indicate if a lock is held
 - Taking advantage of Infiniband atomic operations

◆ Pros:

- Simple
- Can handle simultaneous requests well (with hardware support)

◆ Cons:

- Large number of messages



Atomic Operation Based Design 2

J.M. Mellor-Crummey and Michael L. Scott

- Proposed for shared-memory environment
- Implement MCS by taking advantage of Infiniband atomic operations

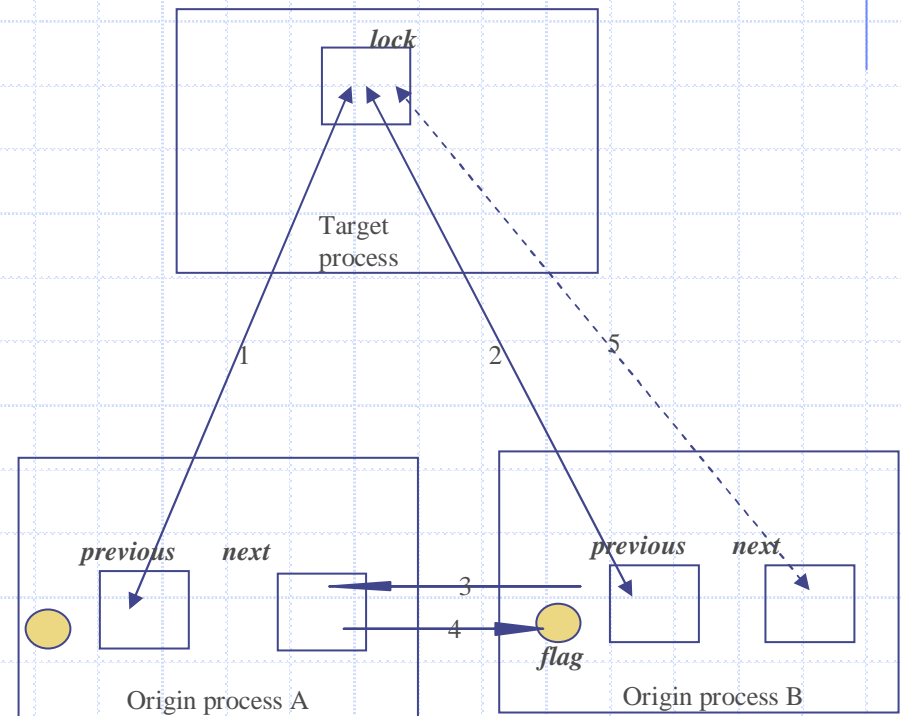
Pros : Scalable to large number of clients

Cons : Needs atomic swap operation

t i O e t i n e e i n

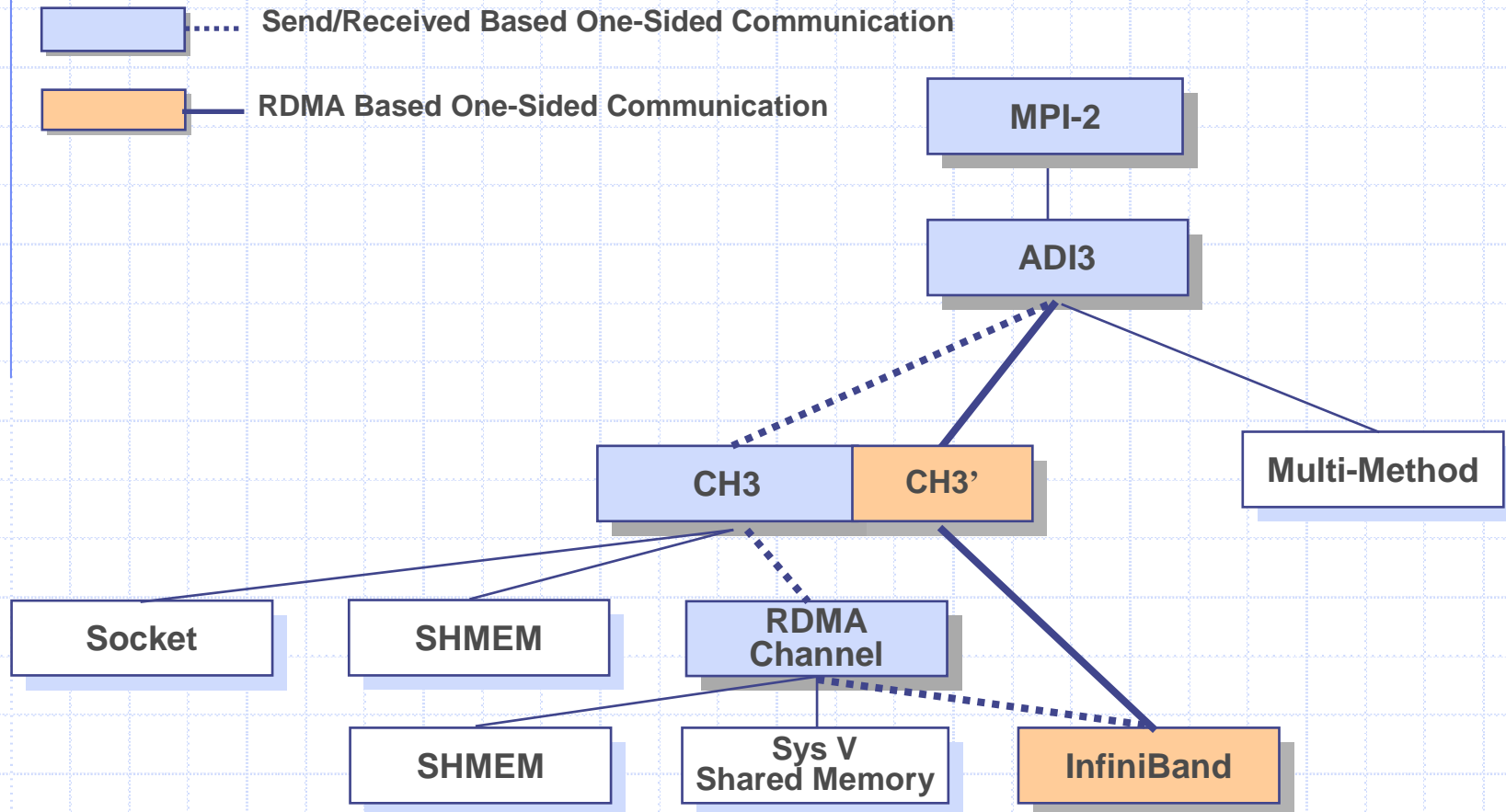
MCS

- ◆ Maintains a distributed queue
- ◆ Origin process maintains flag, previous and next
- ◆ Target process maintains lock
- ◆ Avoid spinning on the remote memory
- ◆ A scalable synchronization algorithm
 - Needs SWAP for efficient implementation



e

One sided design



Outline

- ◆ Background
- ◆ Motivation
- ◆ Design and Implementation
- ◆ Experimental Results
- ◆ Conclusion

Experimental Test bed



u i n d

- Dual Intel Xeon 2.40 GHz processors
- PCI-X 64-bit 133MHz interfaces
- 512K L2 cache and a 400 MHz front side bus



an n iniH

Dua P H



n ini a

ig

P n ini and i



Linu

d Ha

.

.

n ,

.

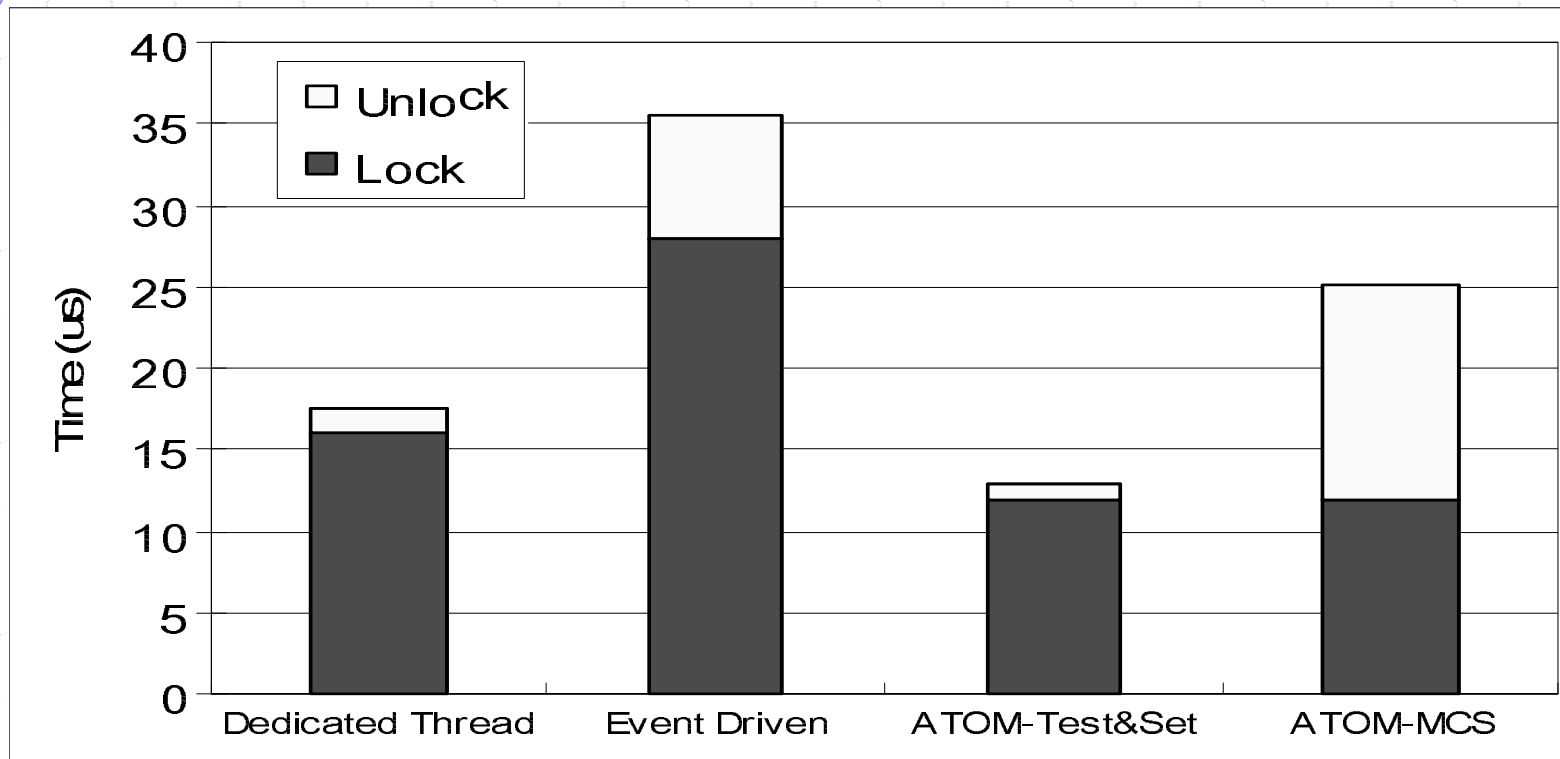
Experiments and tests

- ◆ Synchronization Overhead
- ◆ Synchronization delay
- ◆ semaphore efficiency
- ◆ concurrency
- ◆ Overhead in reference counting thread
- ◆ Semaphore

Synchronization Overhead

- ◆ One process calls only MPI-2 passive synchronization functions (MPI_Win_lock and MPI_Win_unlock) on a window on the other process for multiple iterations
- ◆ We then report the time taken for each iteration

Synchronization Overhead

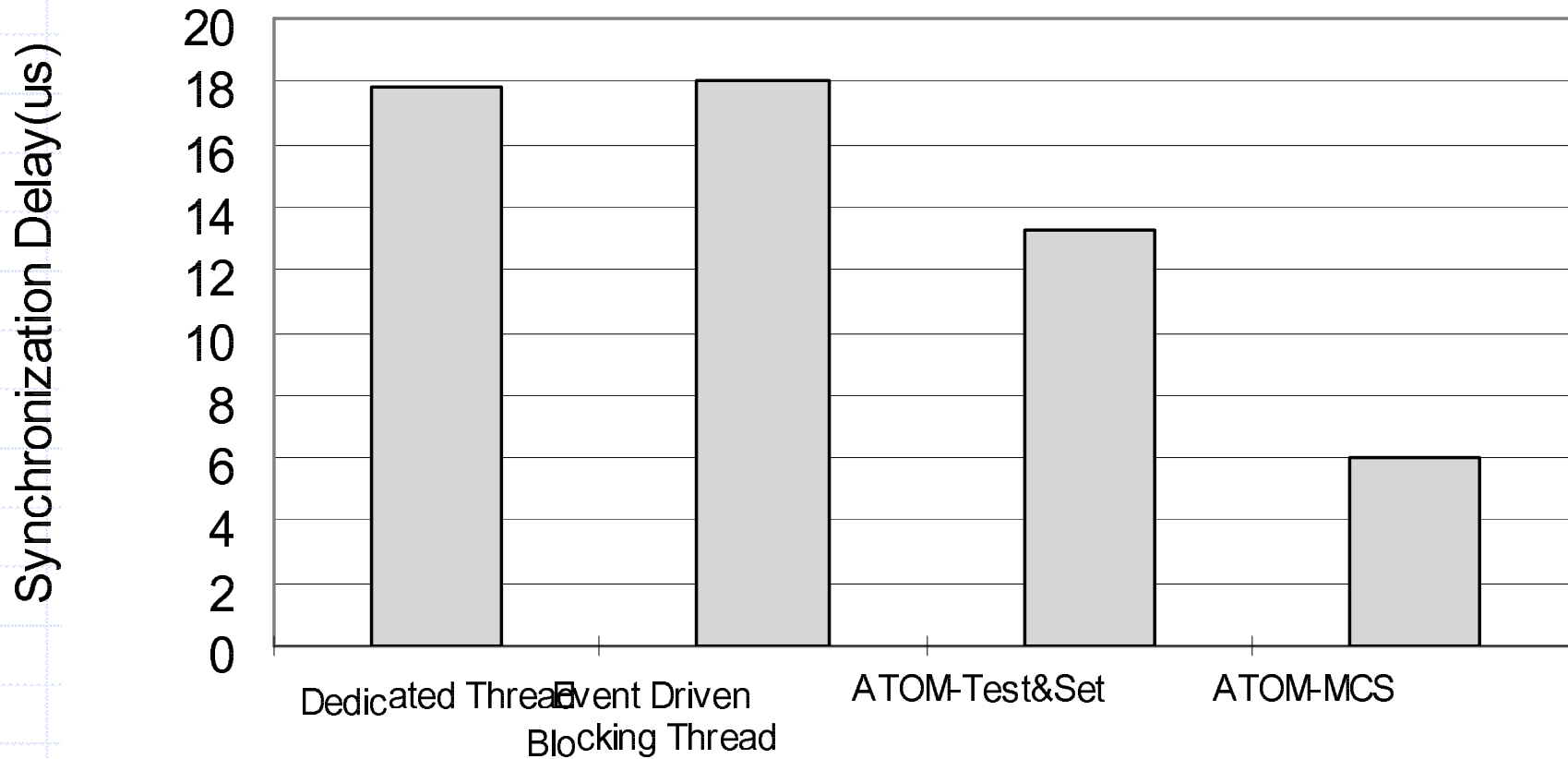


- Event driven brings heavy overhead
- “Test&Set” based design has best performance (12.8usecs)

Initialization

- ◆ Initialize the source and the target source and the source code for the first iteration
- ◆ Initialize the source code
- ◆ Record the delay to change the source code

Performance Comparison

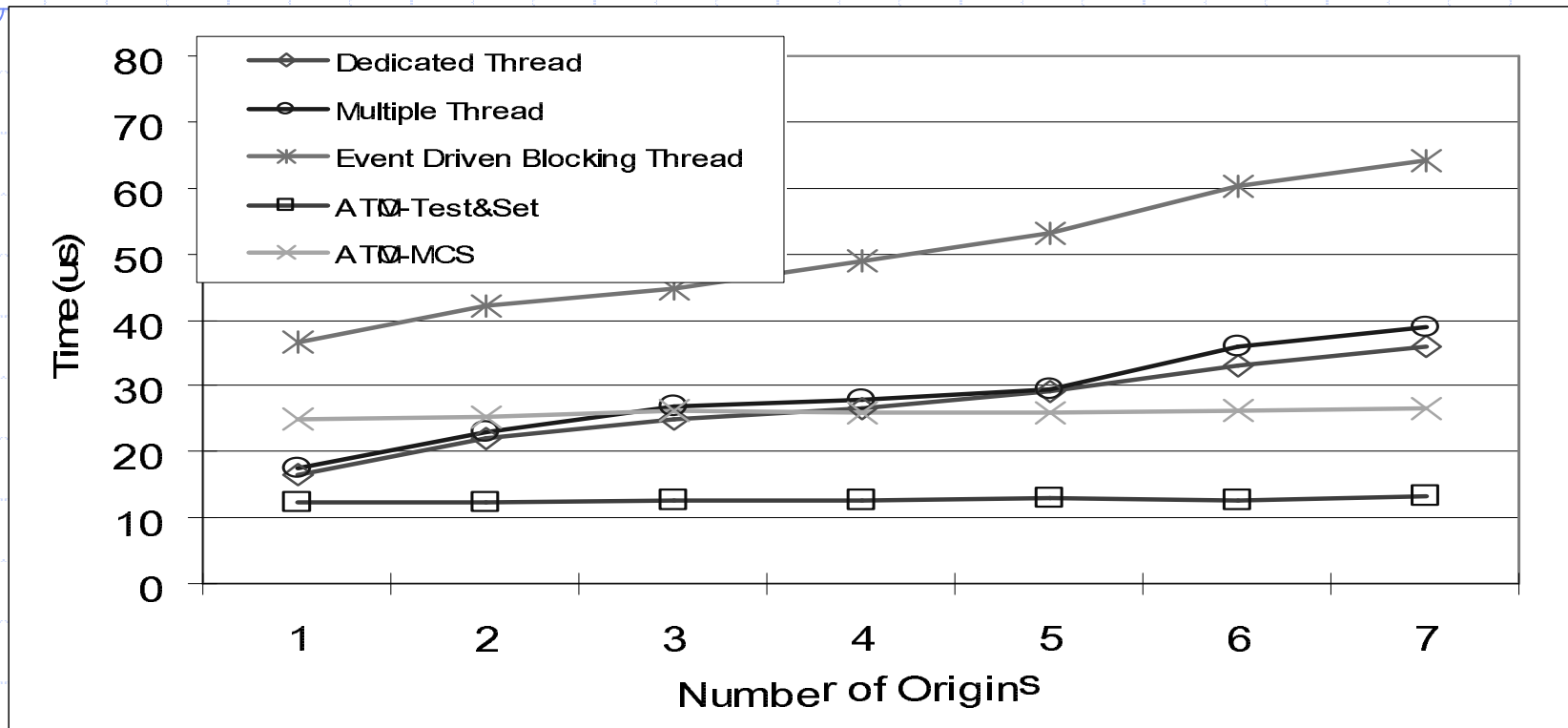


- ATOMIC-MCS shows the best synchronization delay
- The lock can be transferred to next process with a single message

Non-Blocking Test

- ◆ In the target process multiple windows are created and the number of windows is equal to the number of origin processes
- ◆ In each iteration each origin process calls only `MPI_Win_lock` and `MPI_Win_unlock` on the corresponding target window
- ◆ We then report the average time spent on each iteration

n rren

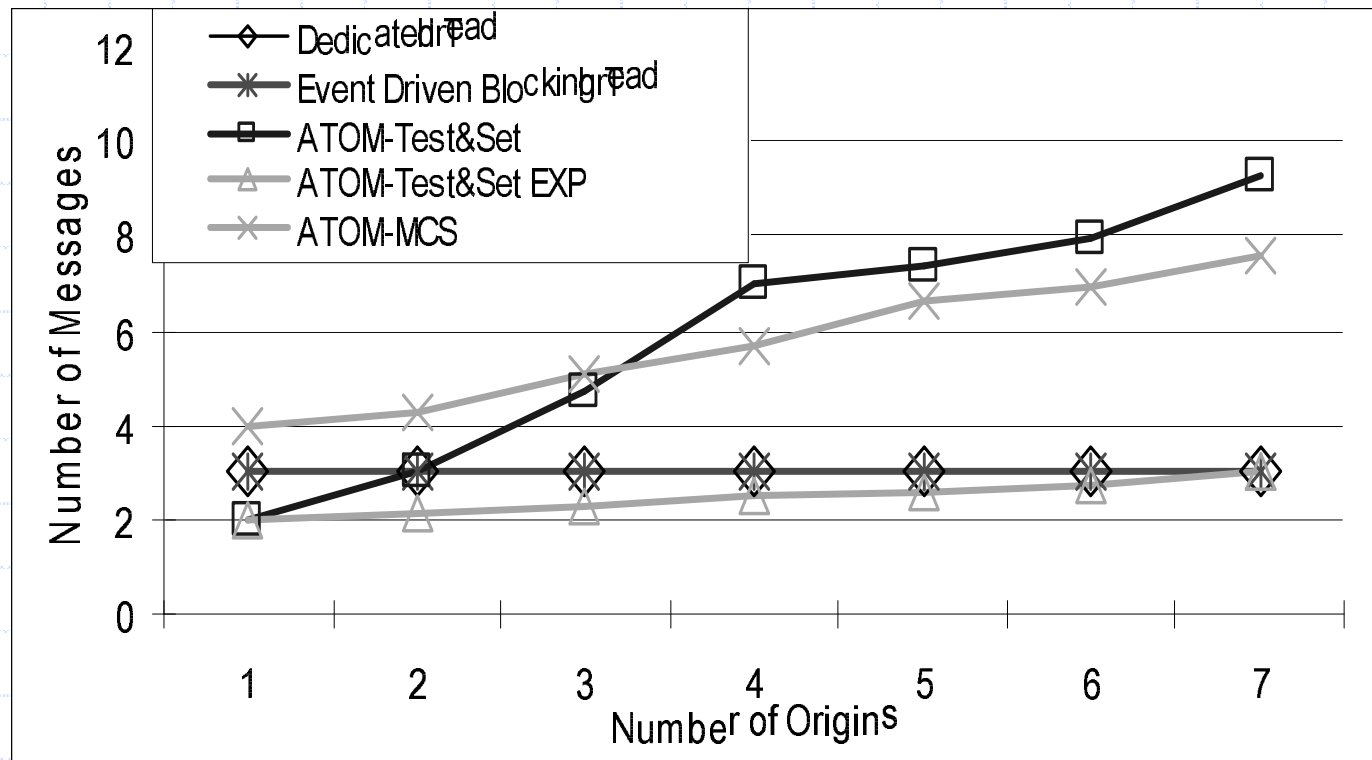


- The hardware support for atomic basedesign helps in handling concurrency with almost no increase in response time
- The response time of thread based designs increase with number of origins

essa e mplexit

- ◆ The test uses multiple origin processes and one target process and all the origin processes compete for the same lock on a target window
- ◆ We calculate the average number of messages exchanged between one origin process and one target process

Message Complexity

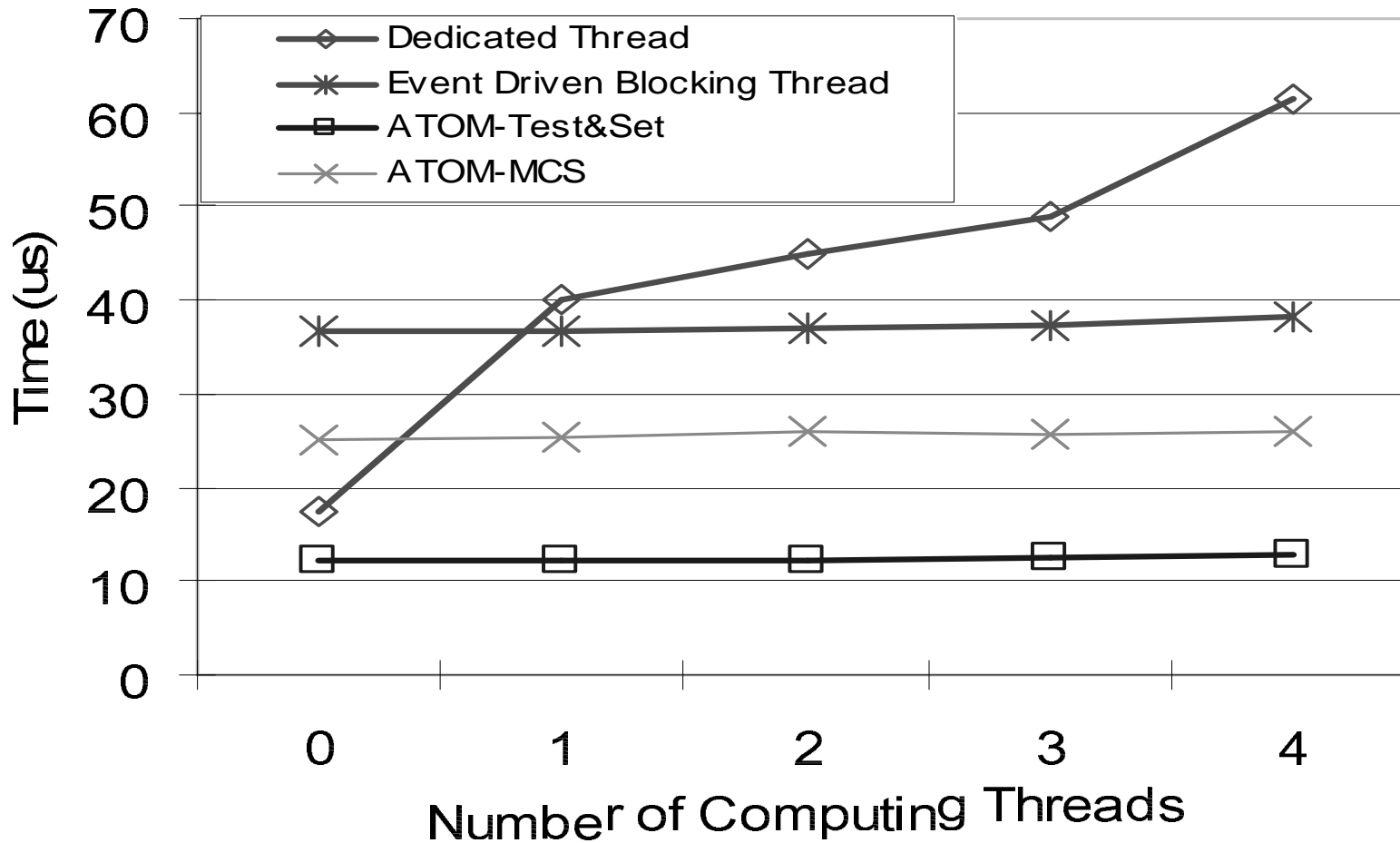


- Without atomic swap, MCS can not achieve its benefits
- With exponential back-off, the “Test&Set” design has low message complexity

mp tin T read

- ◆ we have one target process and one originator process
- ◆ the target process is a non-timed continuous thread
- ◆ the originator process can synchronize
function. These are the overhead in the
reference of the continuous thread

utine

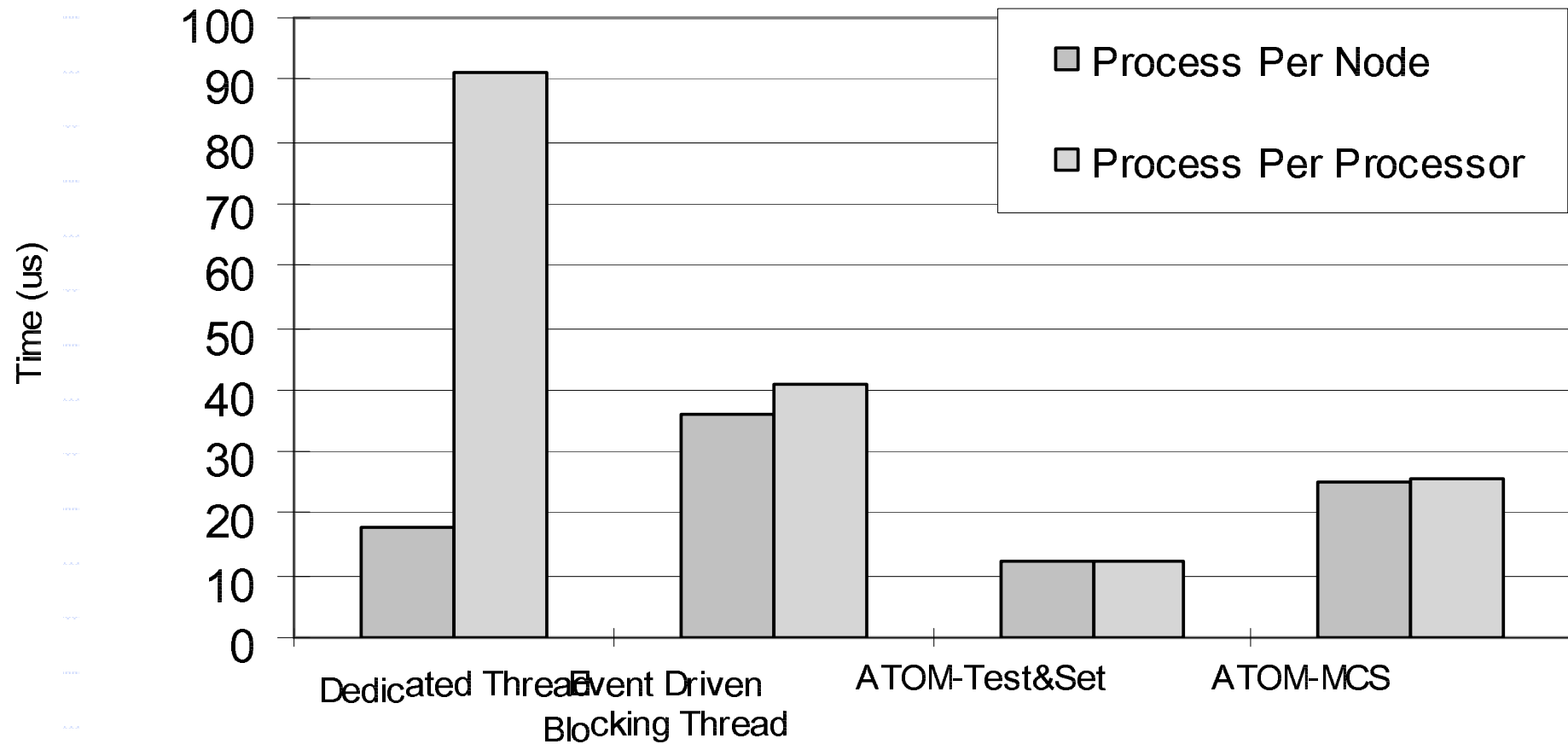


- Atomic based designs performance is not affected by computing threads
- The performance degrades for dedicated thread based design

S de

- ◆ Source processes and target processes
- ◆ Source processes call synchronization functions on target processes
- ◆ Parallel source in target processes in the different nodes
- ◆ SMP target in target processes in the same dual-CP nodes

v S ode



- The atomic based designs outperform the thread based design
- More suitable for SMP mode runs

Outline

- ◆ Background
- ◆ Motivation
- ◆ Design and Implementation
- ◆ Experimental Results
- ◆ Conclusion

Analysis

◆ Comparison of Design and Analysis

- Achieves lower overhead (~12.8 us)
- Reasonable Message Complexity (less than 3)
- Less CPU utilization

So t are i tri tion

◆ P H
High performance MP2 Implementation over InfiniBand based
on MPICH2

◆ u n i n P H . a

◆ d an gani a i n ing P
n ini and

◆ W i a P H . . in a

- This release will support active one-sided communication
- Passive one-sided communication will be incorporated in the next release

◆ n a . i . i a . du i
i a ind .

t r e r

◆ Data union identifier

- Packing/Unpacking
- InfiniBand specific datatype solutions
 - Multi RDMA
 - Gather/Scatter
 - Hybrid
- identifier



Thank You!

For more information, please visit the

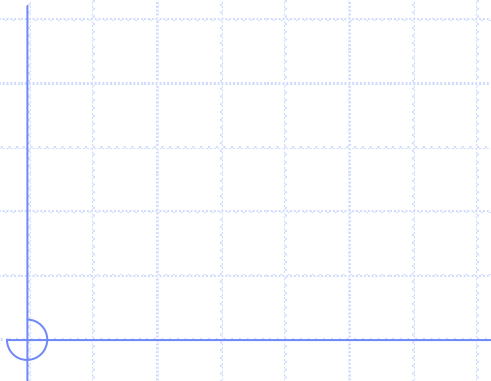
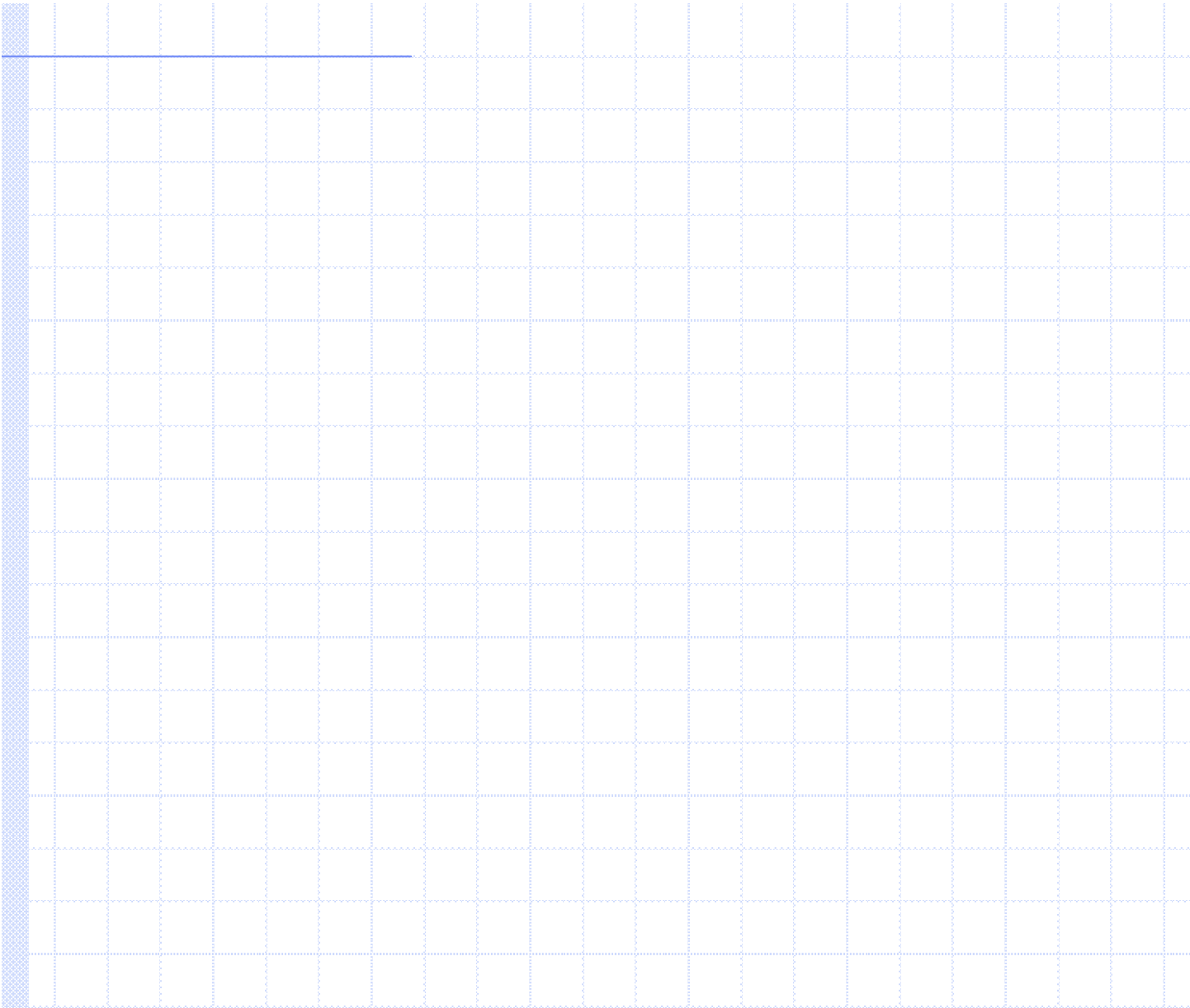
NBC

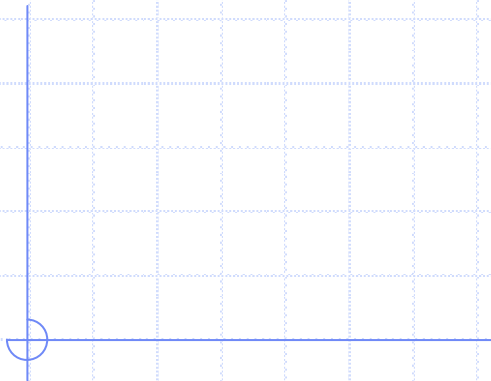
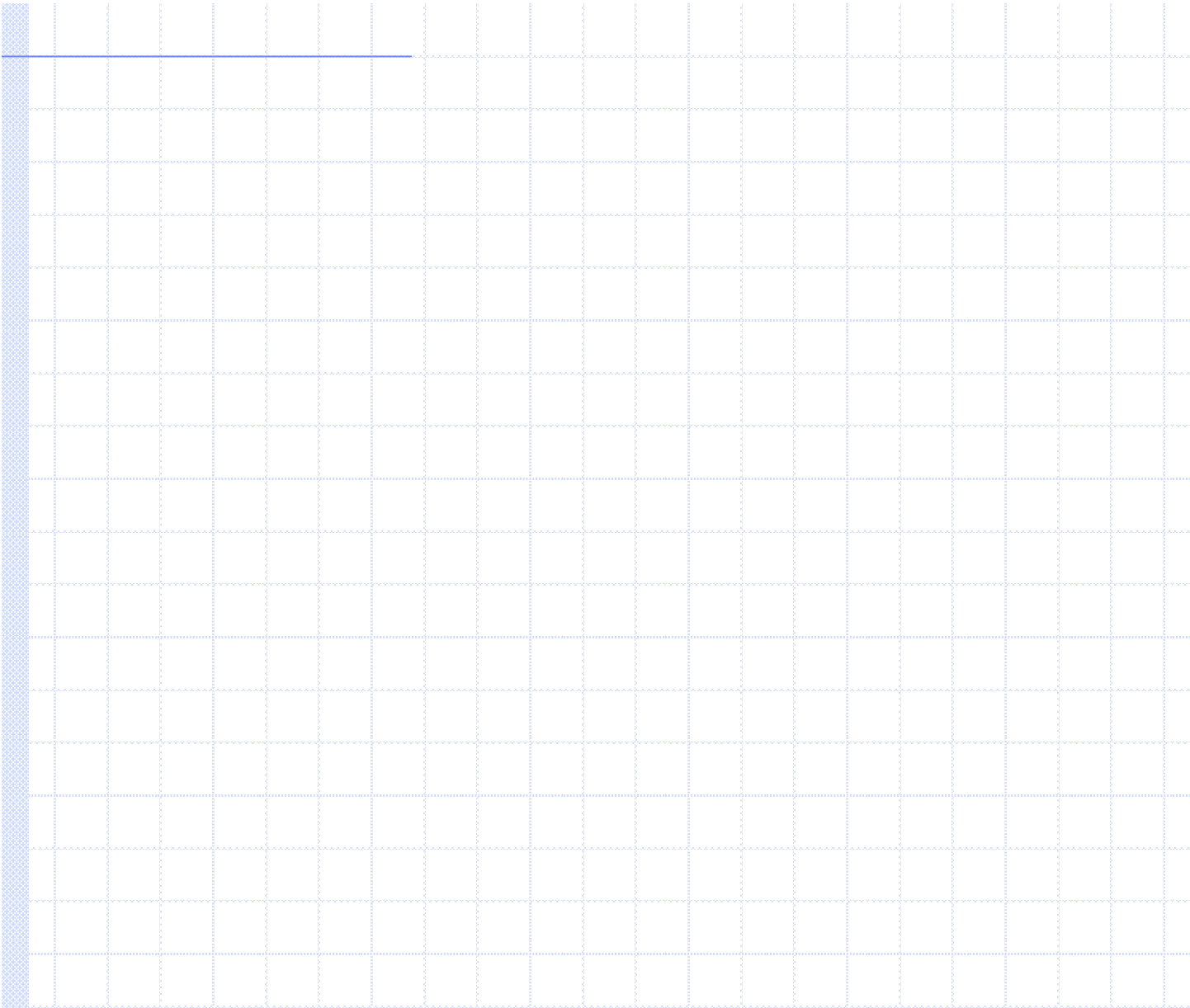
Home Page

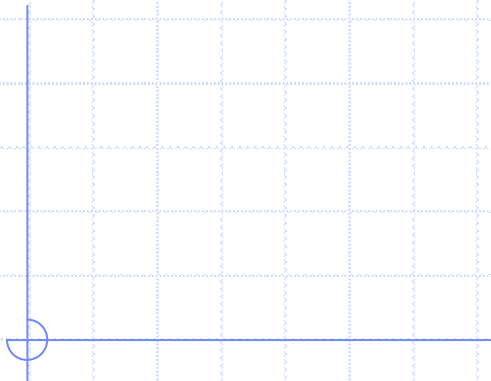
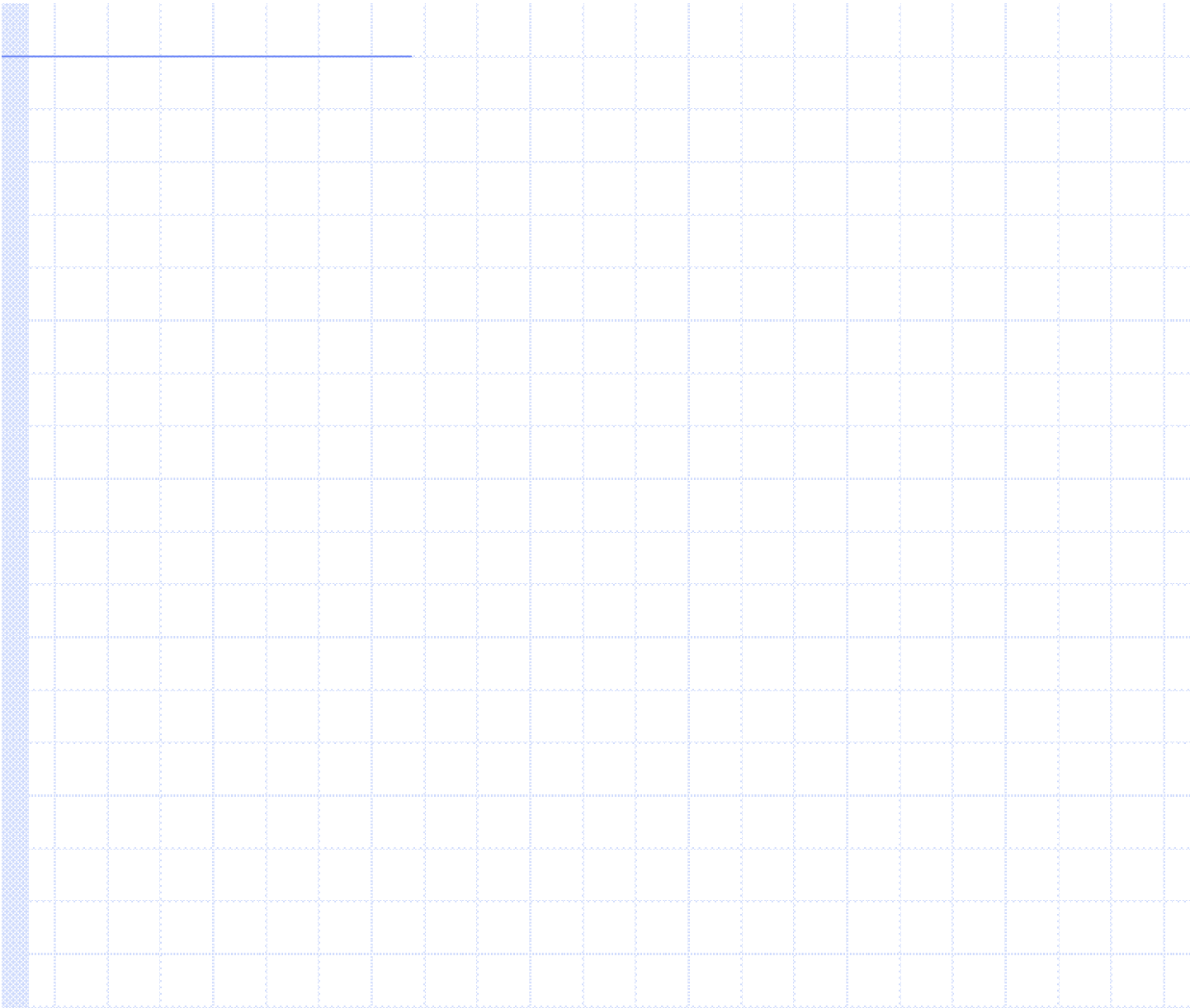
<http://nowlab.cis.ohio-state.edu>

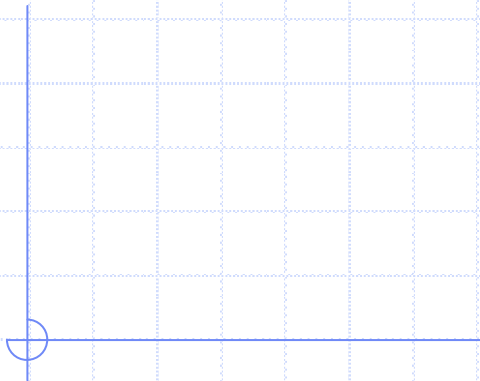
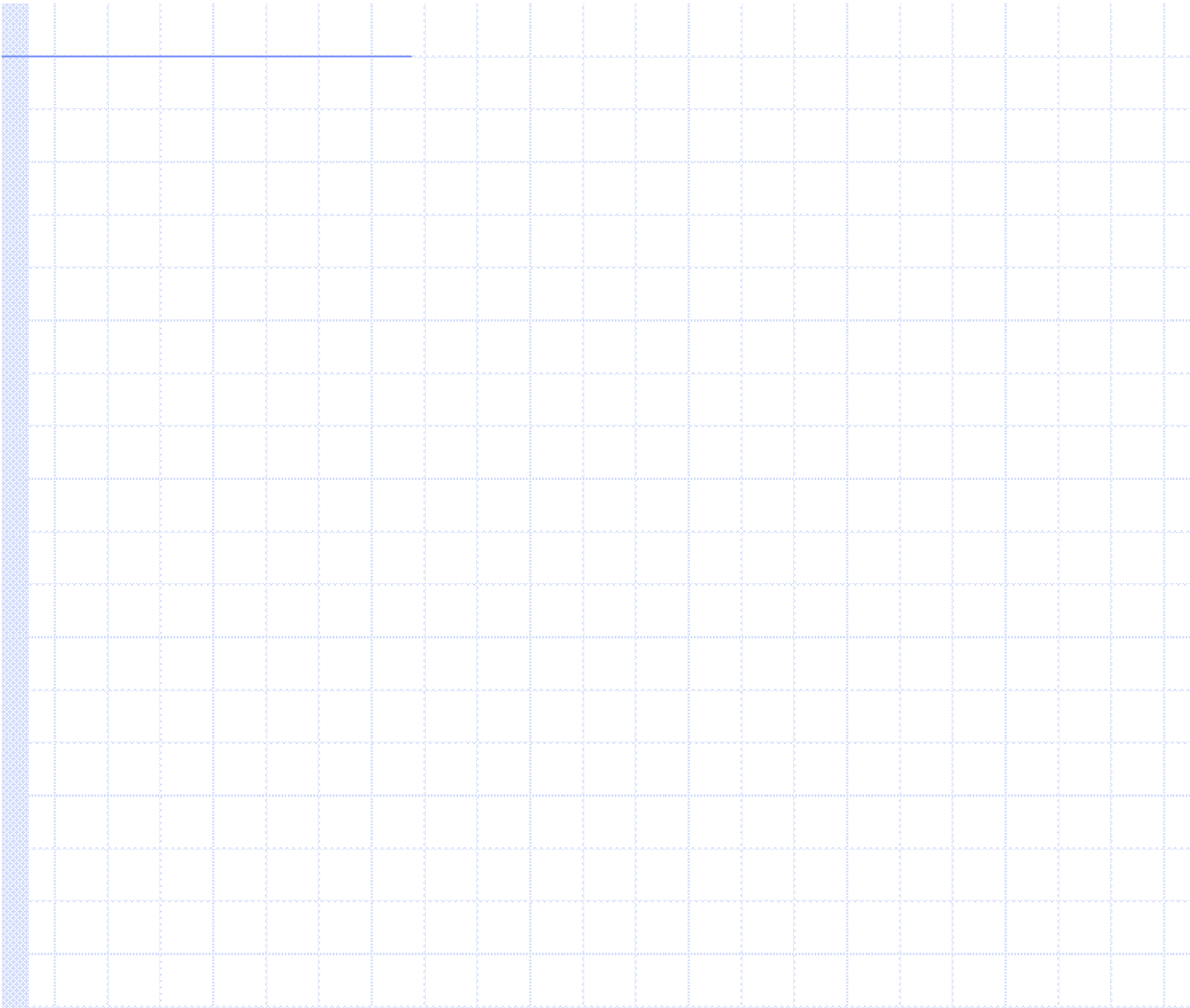
Network Based Computing Laboratory

The Ohio State University









re i s r

◆ i h er or ance One Sided

o nication over n ini and

eihan ian i in i y n oo in

ha a e ar anda and i ia ro and

a eev ha r rid

◆ M - ased implementation can ac ieve

■ etter er or ance or one ided

co nication and active ynchronization

