

TORSTEN HOEFLER

# Active RDMA - new tricks for an old dog

with M. Besta, R. Belli, S. di Girolamo @ SPCL

presented at Salishan Meeting, Gleneden Beach, OR, USA, April 2016





TORSTEN HOEFLER

# Active RDMA - new tricks for an old dog

with M. Besta, R. Belli, S. di Girolamo @ SPCL

presented at Salishan Meeting, Salishan, OR, USA, April 2016

**Alternative (better) title: Beyond RDMA**

### Remote Operations

- put, get, atomics



### Remote Matching

- partial control at target



Lossy Networks  
Ethernet

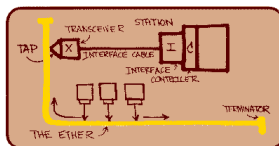
1980's

Lossless Networks  
RDMA

2000's

Full Device Programs  
Offload

2020's



## Remote Operations

- put, get, atomics



## Remote Matching

- partial control at target



## Remote Synchronization

[IPDPS'15]

- Extend RMA semantics
- Fully one-sided (in HW)
- Synchronization

## Remote Transactions

[HPDC'15]

- Similar to HTM
- Extend across nodes
- Think active messages?

## Remote Invocation

[ICS'15]

- Utilizes IOMMUs
- Control transfer
- Active memory

# RDMA IN CASE YOU WANT TO LEARN MORE ABOUT RMA

- PGAS and RMA are produced by:
  - PGAS as language extension
  - RMA as library (integrated)
- Offer abstraction for:
  - Data placement, read, write
  - Target has very little control

SCIENTIFIC AND ENGINEERING COMPUTATION SERIES

## Using Advanced MPI

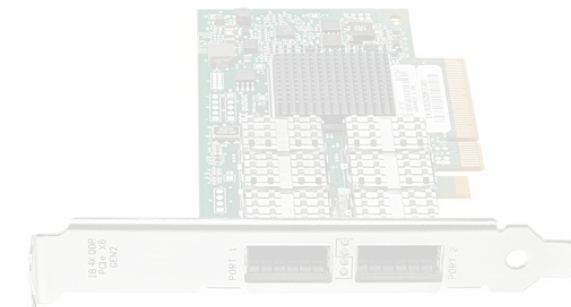
Modern Features of the Message-Passing Interface

William Gropp  
Torsten Hoefler  
Rajeev Thakur  
Ewing Lusk



**How to implement producer/consumer in passive mode?**

- RDMA is a hardware mechanism:
  - Often accessible through
- Specific to a (set of) hardware:
  - Offers varying levels of functionality
  - Most common: read, write
  - Address-space management
  - Common denominator is

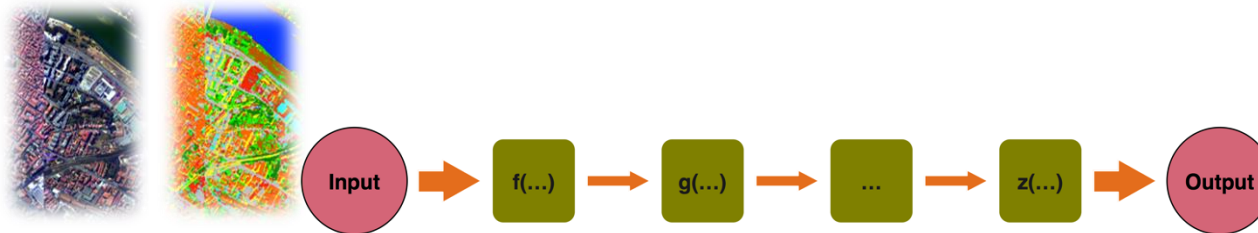
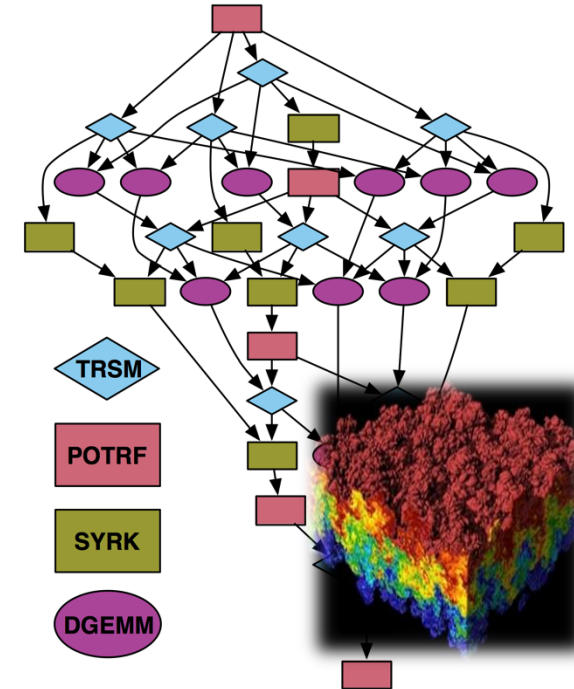
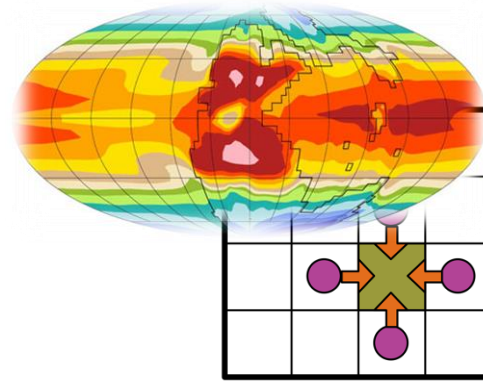
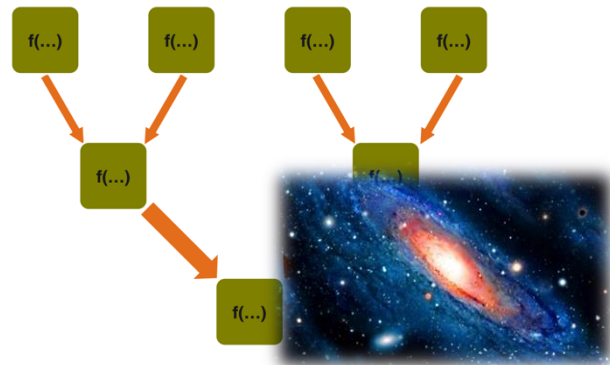




# PRODUCER-CONSUMER RELATIONS

- Most important communication idiom

- Some examples:



- Perfectly supported by MPI-1 Message Passing

- But how does this actually work over RDMA?

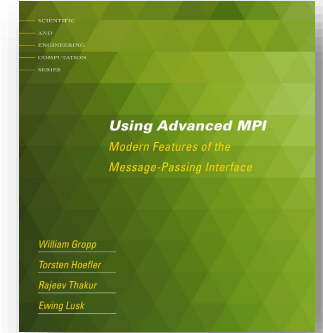
# Remote Synchronization

# ONE SIDED – PUT + SYNCHRONIZATION

Producer

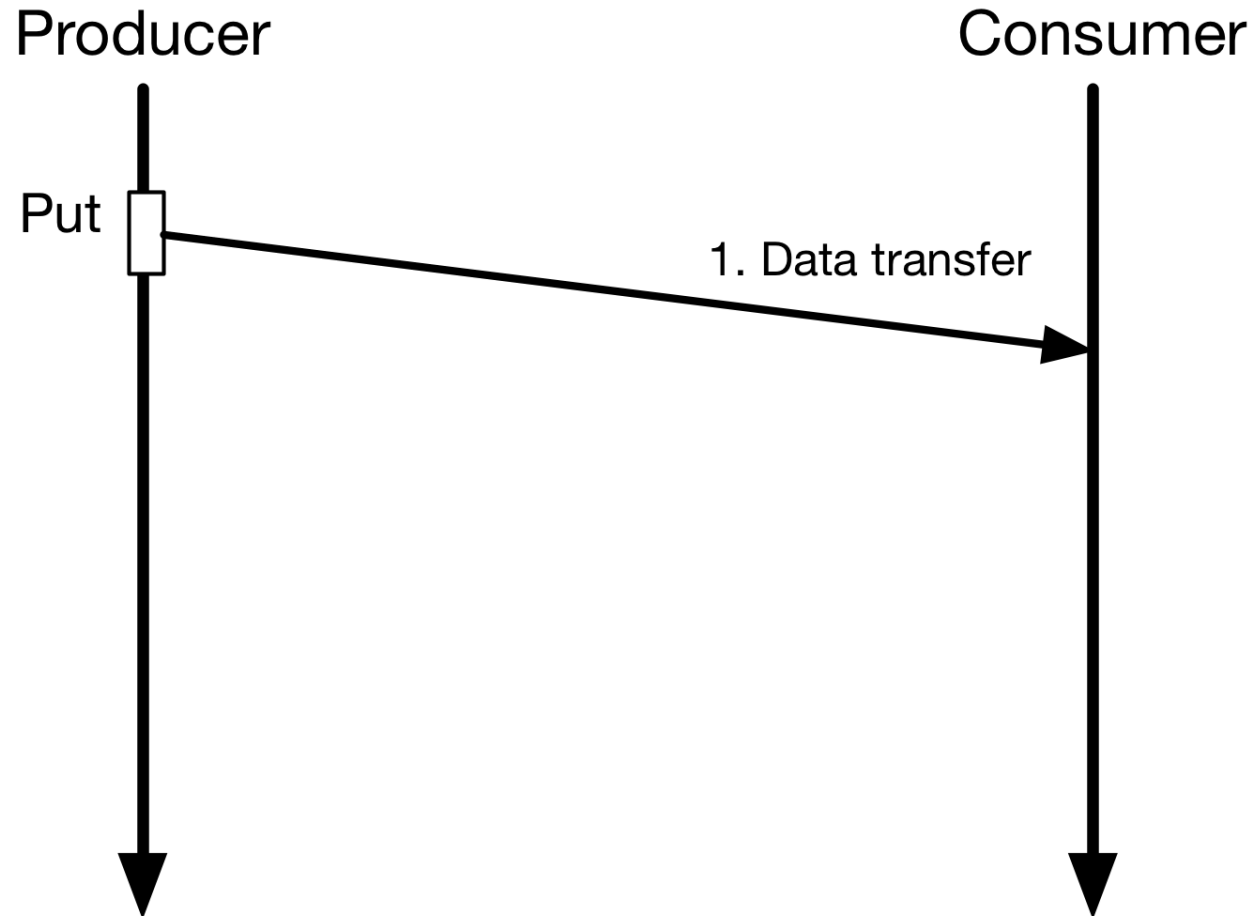


Consumer

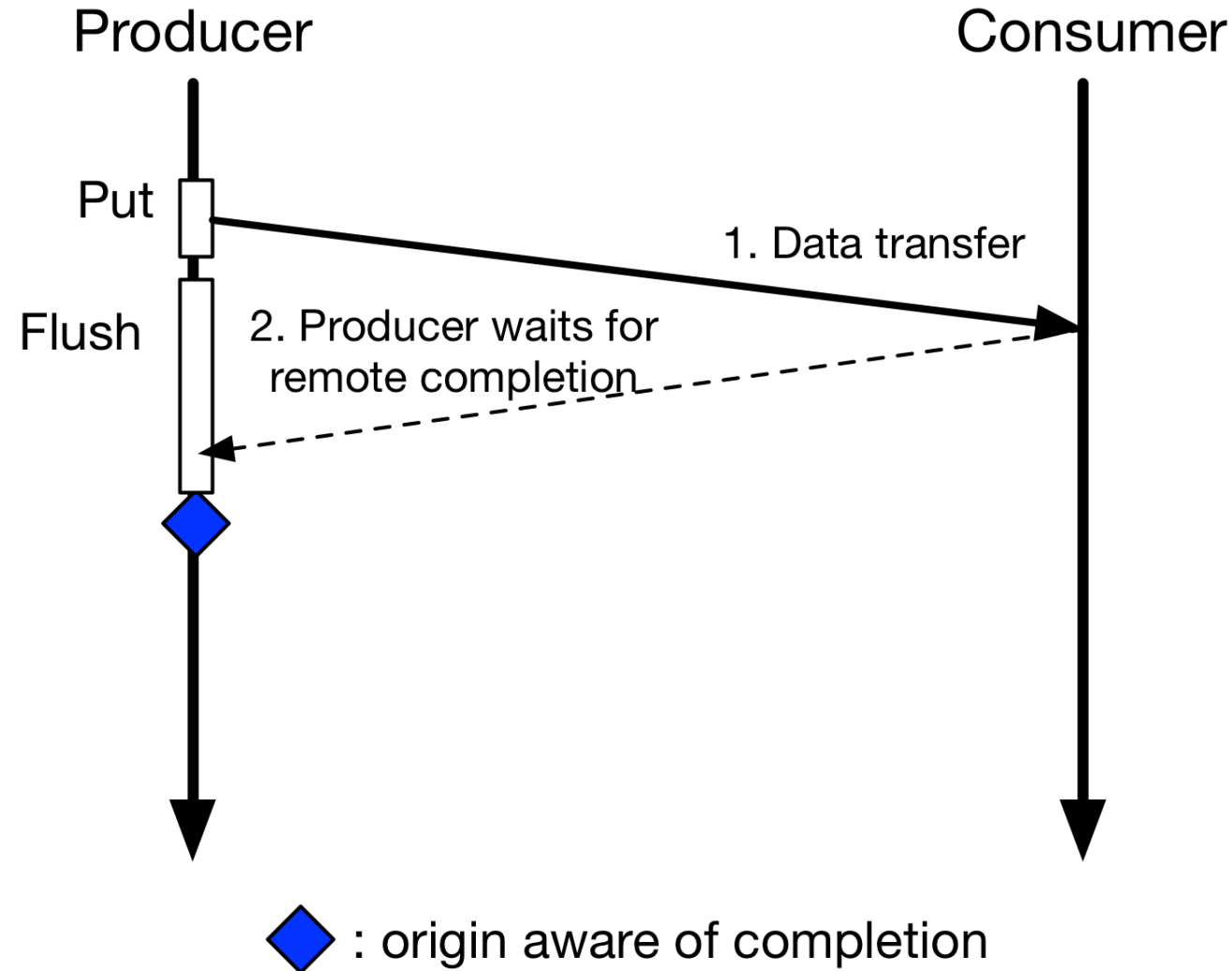




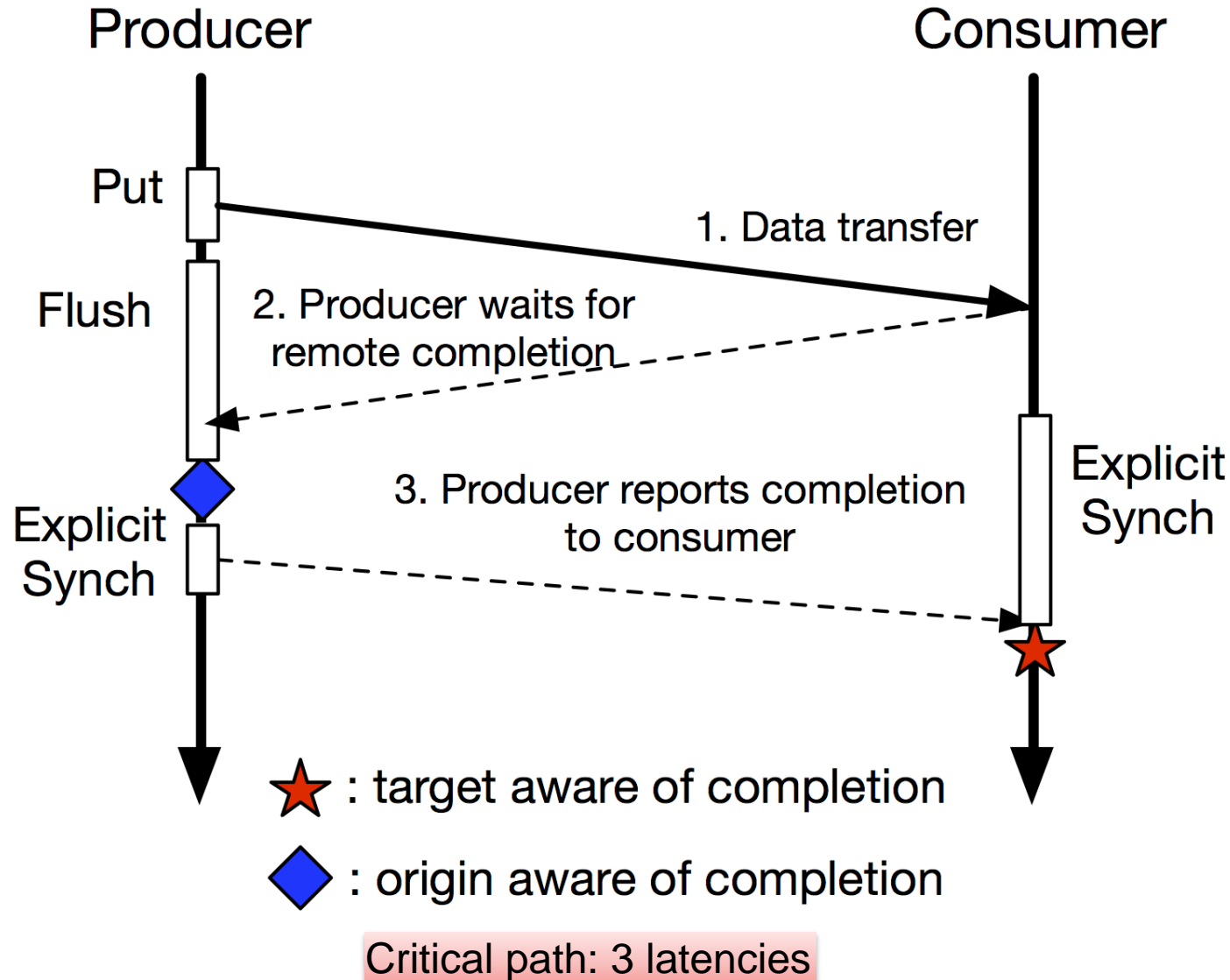
# ONE SIDED – PUT + SYNCHRONIZATION



# ONE SIDED – PUT + SYNCHRONIZATION

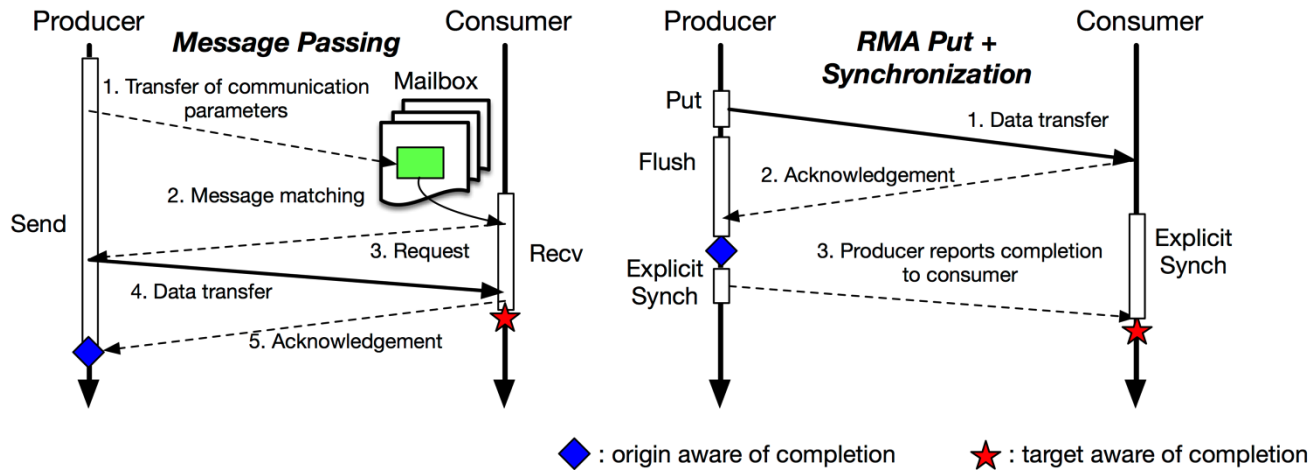


# ONE SIDED – PUT + SYNCHRONIZATION





# COMPARING APPROACHES

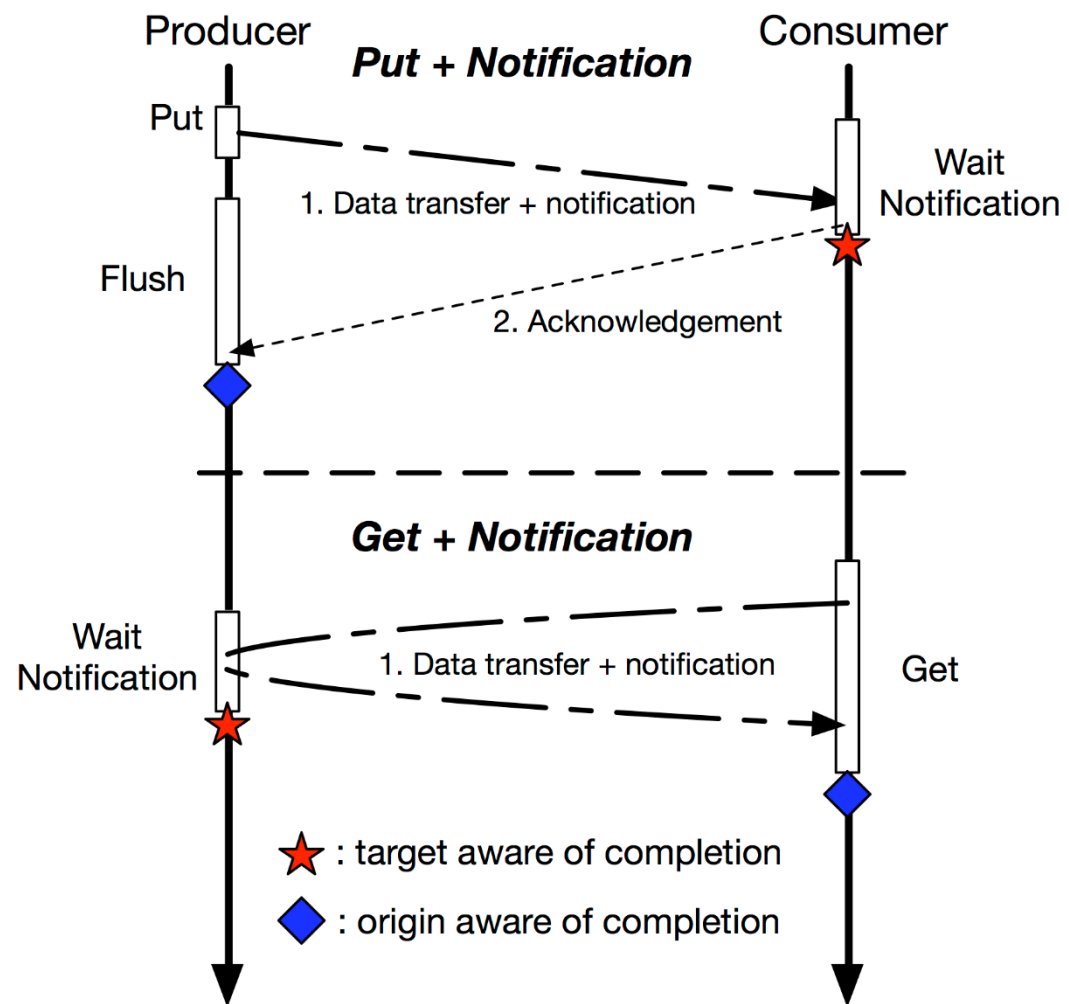


**Message Passing**  
1 latency + copy /  
3 latencies

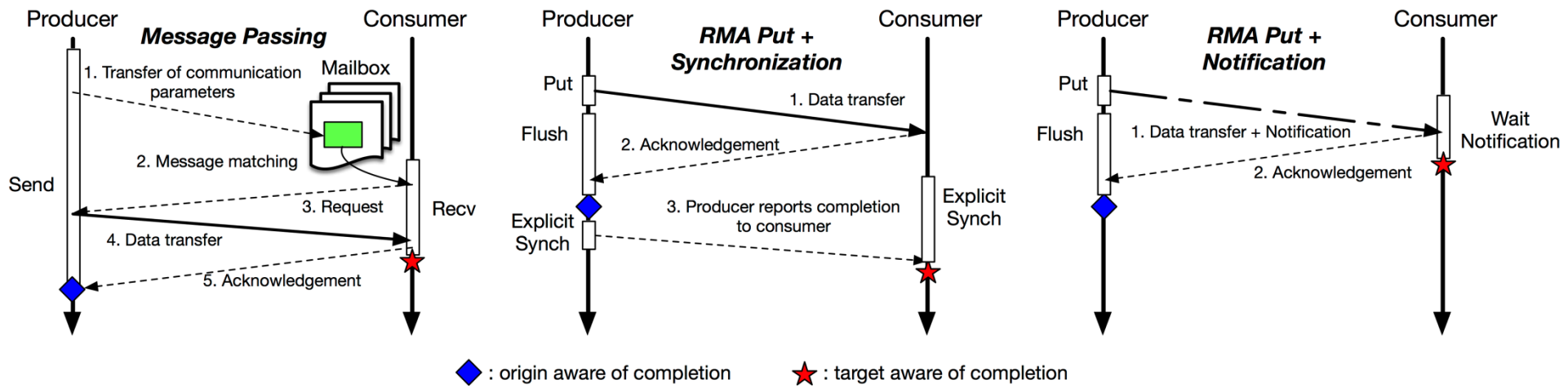
**One Sided**  
3 latencies

# IDEA: RMA NOTIFICATIONS

- First seen in Split-C (1992)
- Combine communication and synchronization using RDMA
- RDMA networks can provide various notifications
  - Flags
  - Counters
  - Event Queues



# COMPARING APPROACHES



**Message Passing**  
1 latency + copy /  
3 latencies

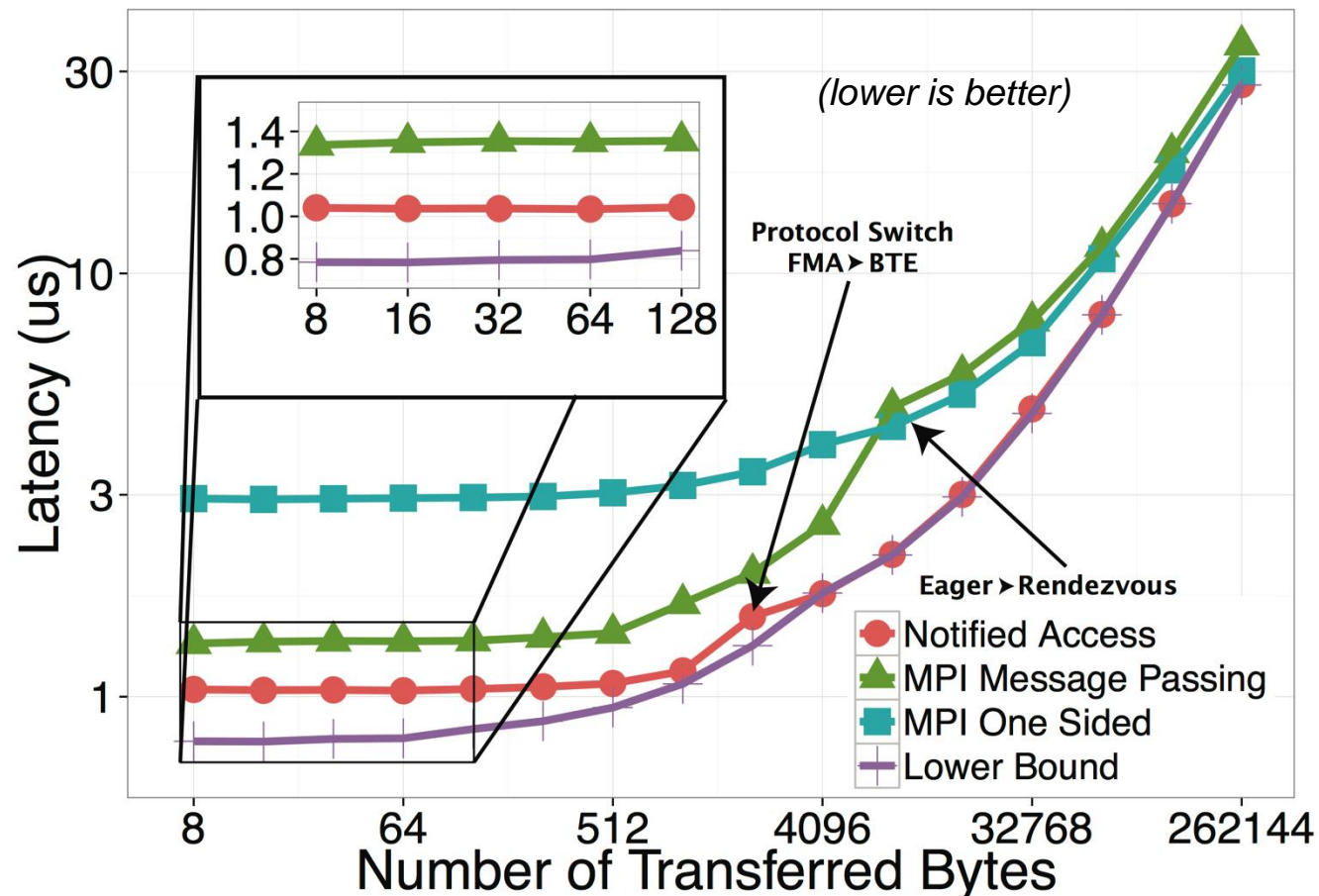
**One Sided**  
3 latencies

**Notified Access**  
1 latency



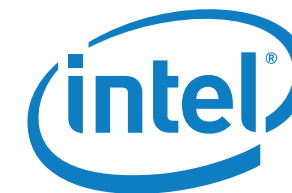
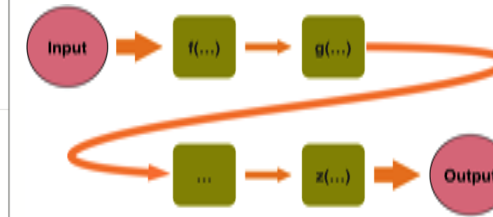
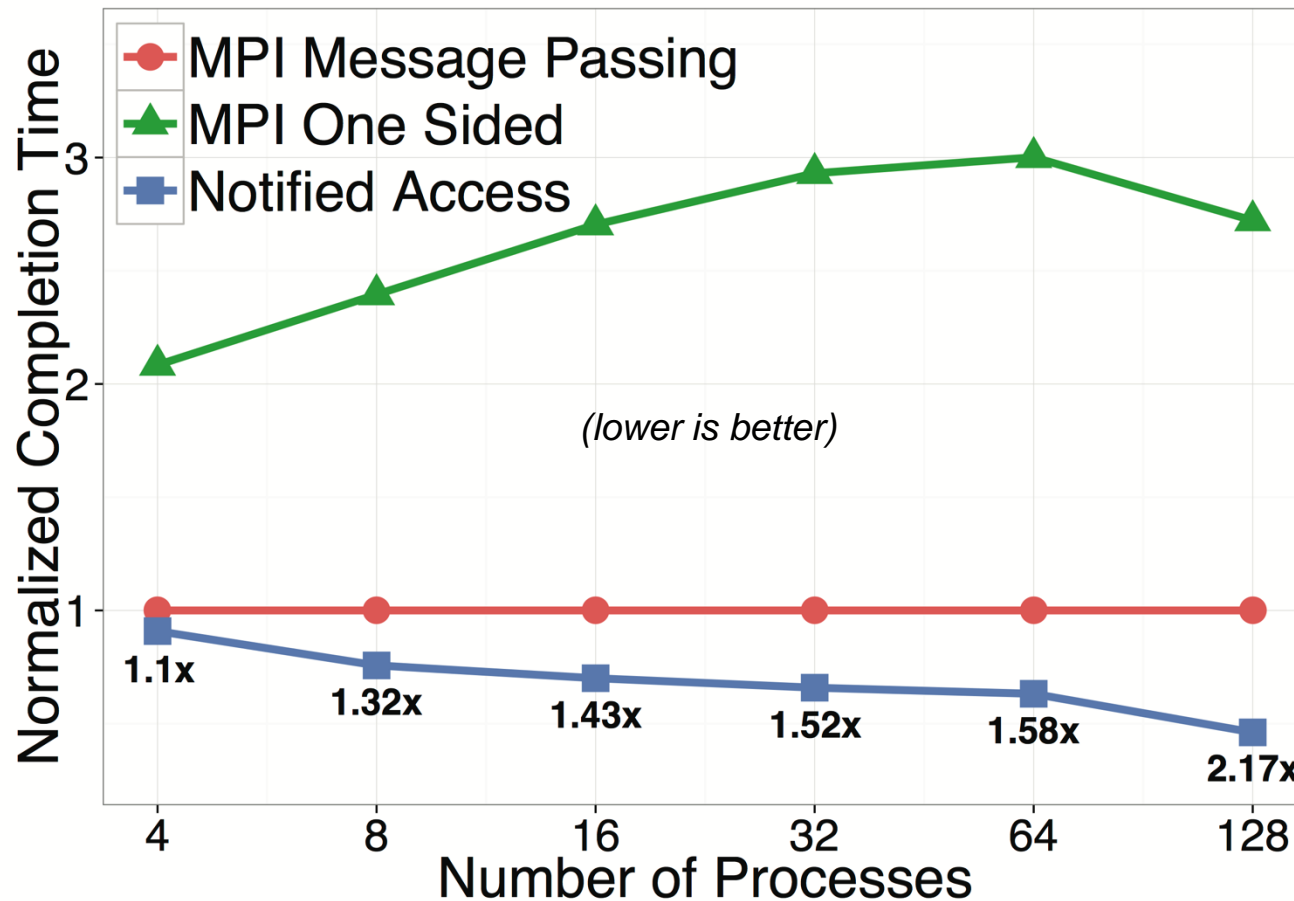
# PING PONG PERFORMANCE (INTER-NODE)

- 1000 repetitions, each timed separately, RDTSC timer
- 95% confidence interval always within 1% of median



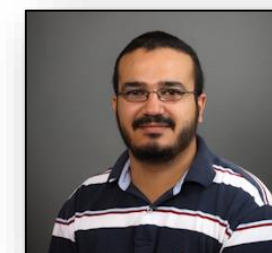
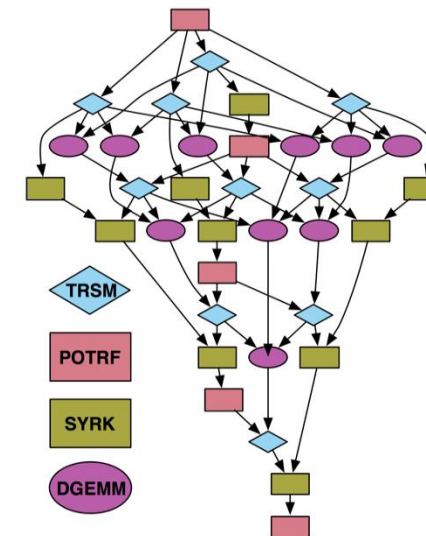
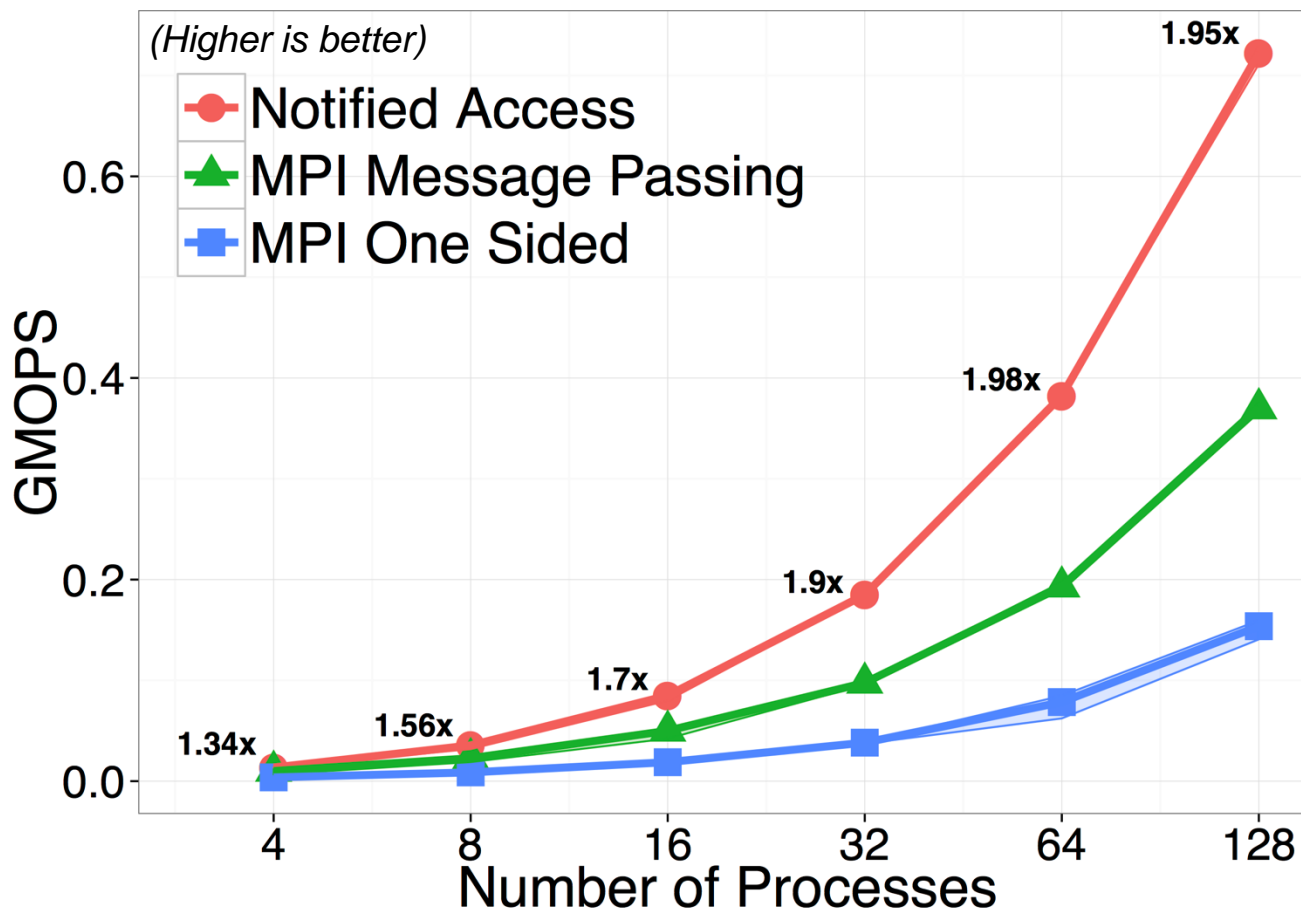
# PIPELINE – ONE-TO-ONE SYNCHRONIZATION

- 1000 repetitions, each timed separately, RDTSC timer
- 95% confidence interval always within 1% of median



# CHOLESKY – MANY-TO-MANY SYNCHRONIZATION

- 1000 repetitions, each timed separately, RDTSC timer
- 95% confidence interval always within 10% of median

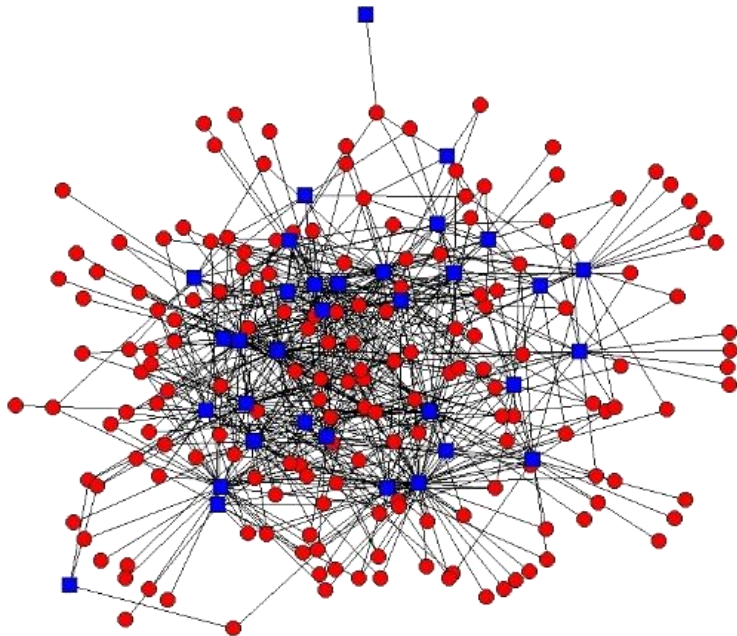




# (Remote) Transactions

# LARGE-SCALE IRREGULAR GRAPH PROCESSING

- Becoming more important [1]
  - Machine learning
  - Computational science
  - Social network analysis

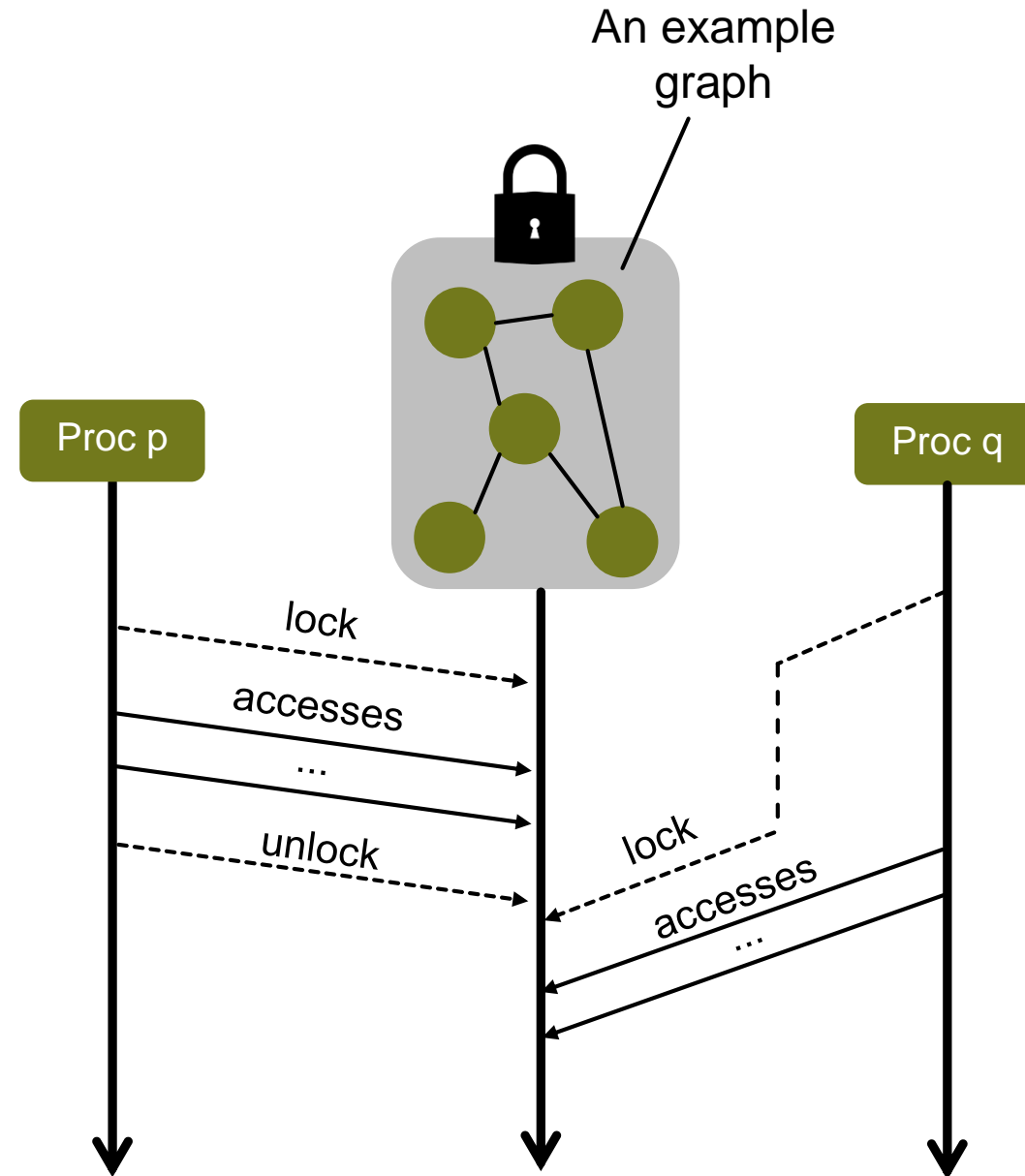


$$\frac{1}{\sqrt{2}} |\text{cat}\rangle + \frac{1}{\sqrt{2}} |\text{dog}\rangle$$

# SYNCHRONIZATION MECHANISMS

## COARSE LOCKS

- ✓ Simple protocols
- ✗ Serialization
- ✗ Detrimental performance



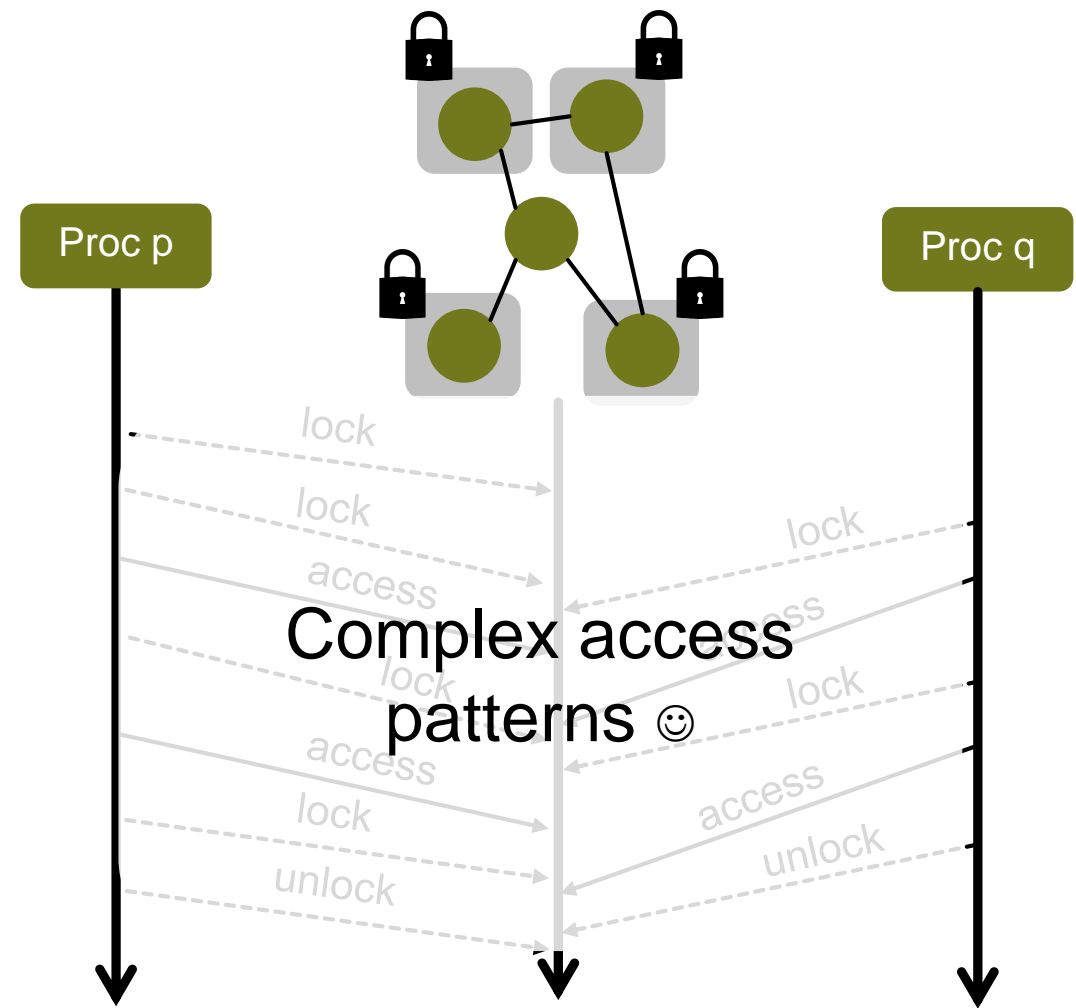
# SYNCHRONIZATION MECHANISMS

## FINE LOCKS

✓ Higher performance possible

✗ Complex protocols

✗ Risk of deadlocks





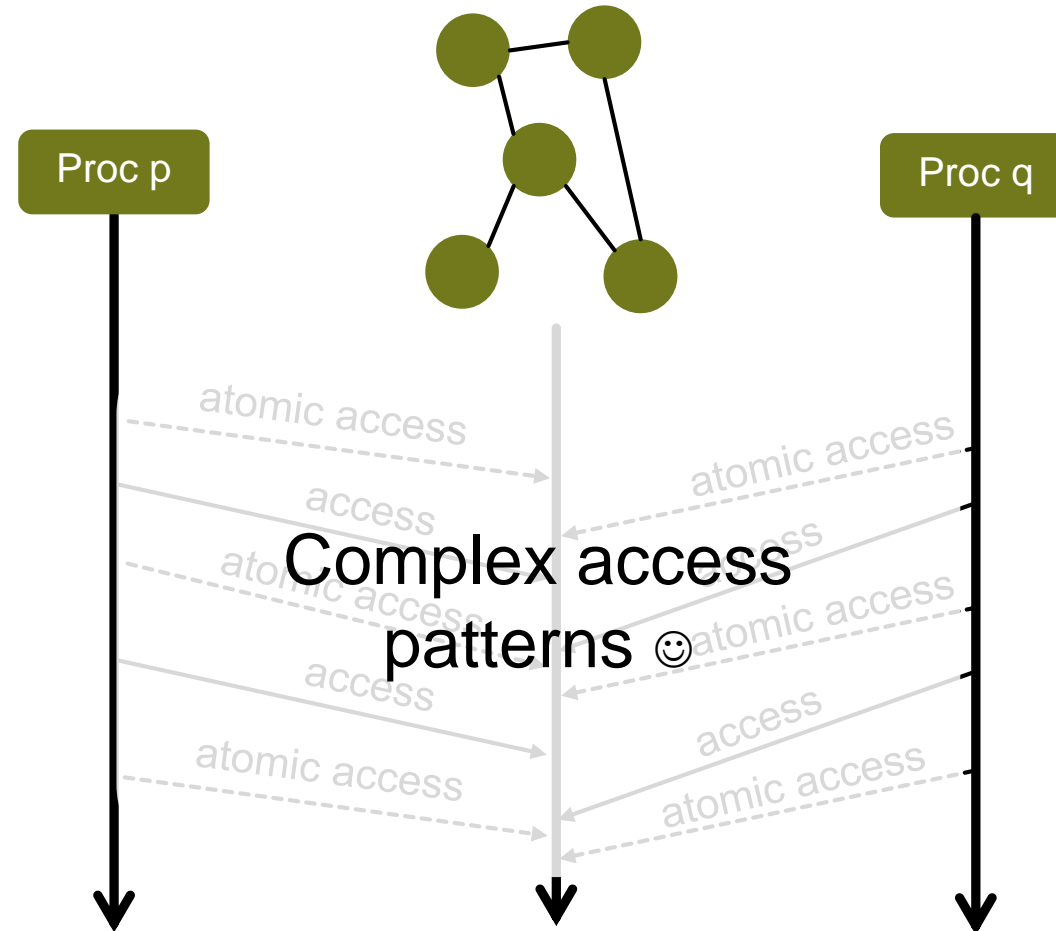
# SYNCHRONIZATION MECHANISMS

## ATOMIC OPERATIONS

✓ High performance (may be challenging to get)

✗ Complex protocols

✗ Subtle issues (ABA, ...)



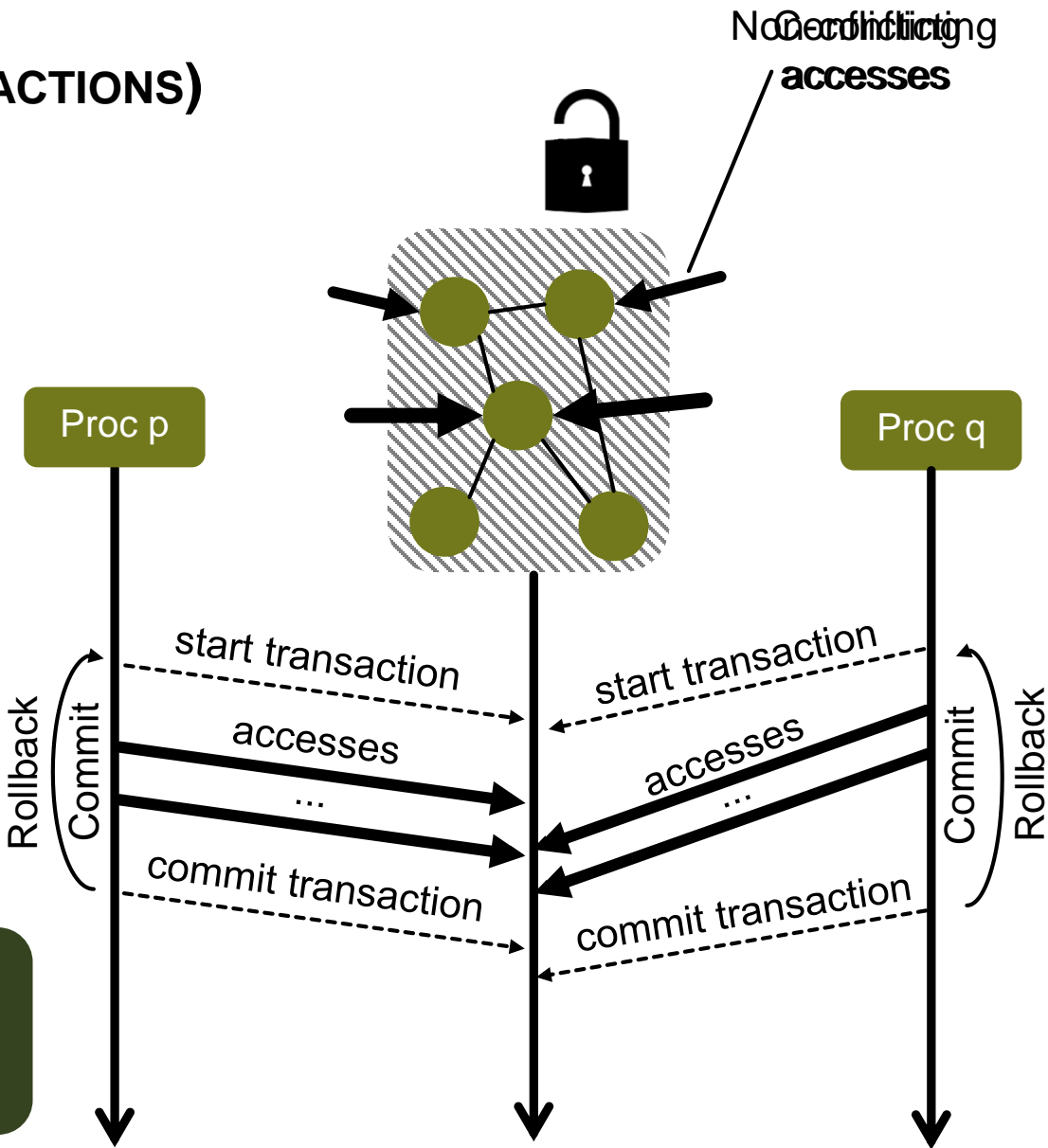
# SYNCHRONIZATION MECHANISMS

## TRANSACTIONAL MEMORY (CF. DB TRANSACTIONS)

Conflicts solved with rollbacks and/or serialization.

**✗** Software overheads

**✓** Simple protocols



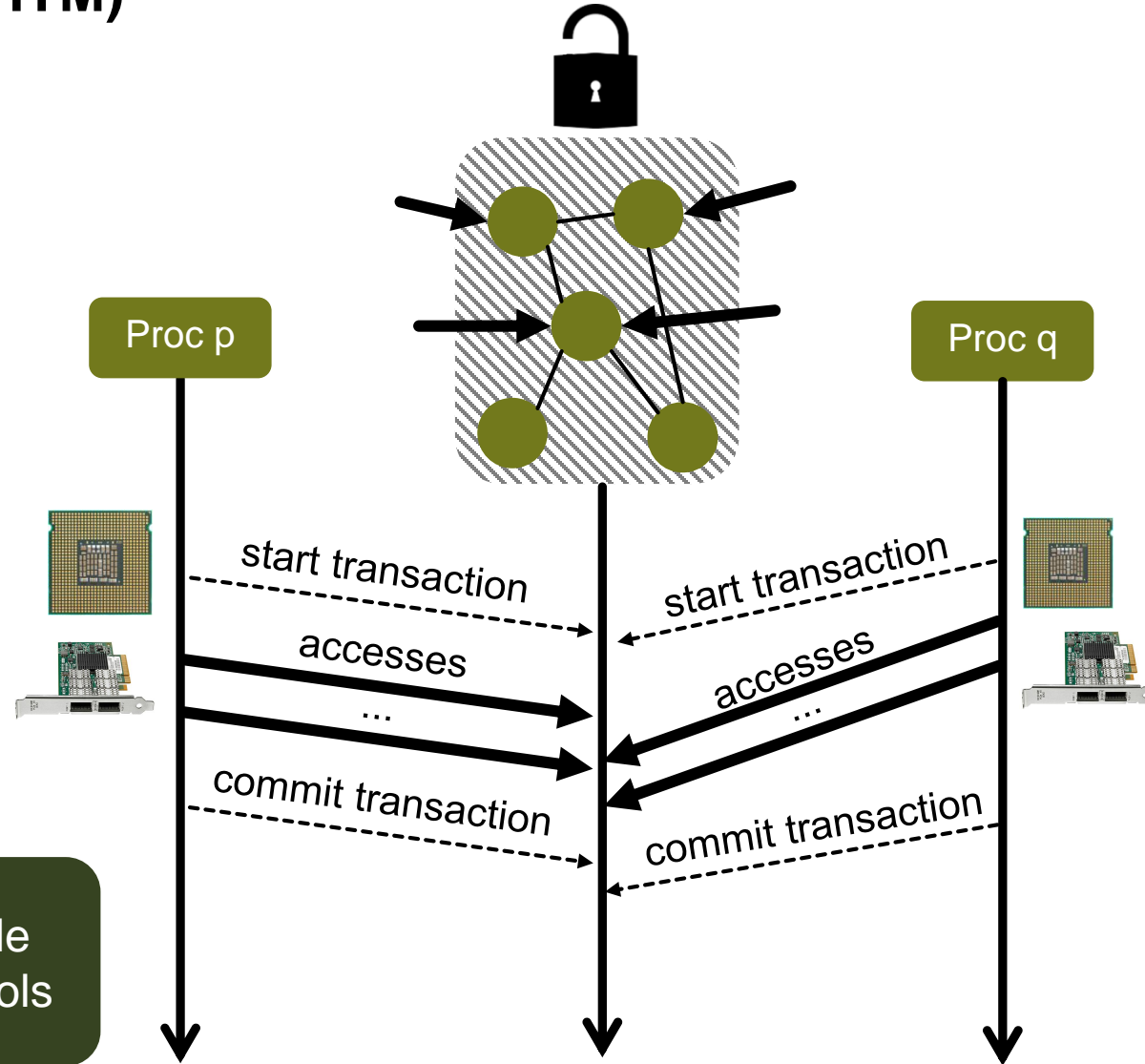
# SYNCHRONIZATION MECHANISMS

## HARDWARE TRANSACTIONAL MEMORY (HTM)

Conflicts solved with rollbacks and/or HW serialization.

? High performance? For graphs?

✓ Simple protocols

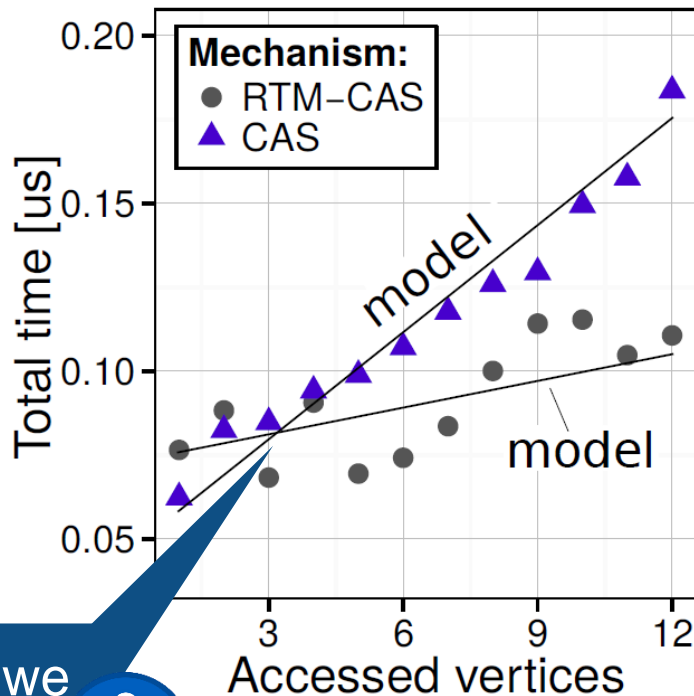


# PERFORMANCE MODEL

## ATOMICS VS TRANSACTIONS

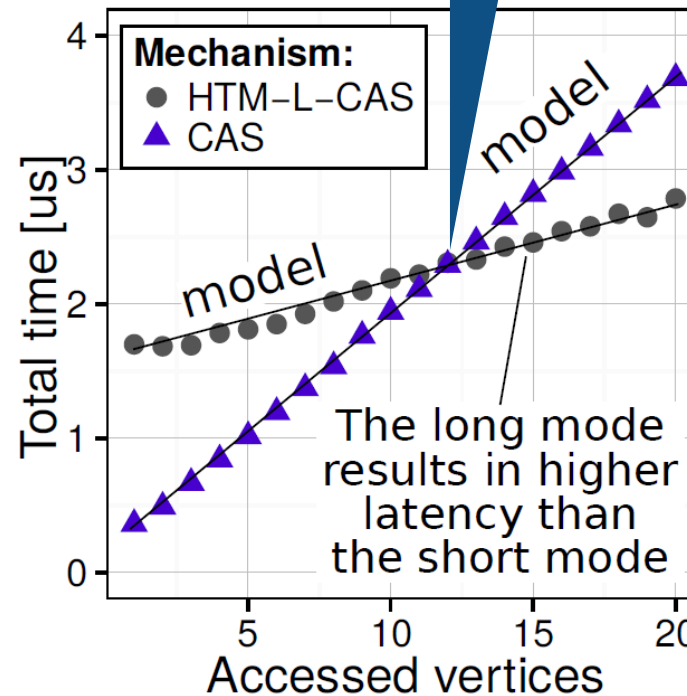
- Can we amortize HTM startup/commit overheads with larger transaction sizes?

### Haswell



Yes, we can!

### Blues



Yes, we can!

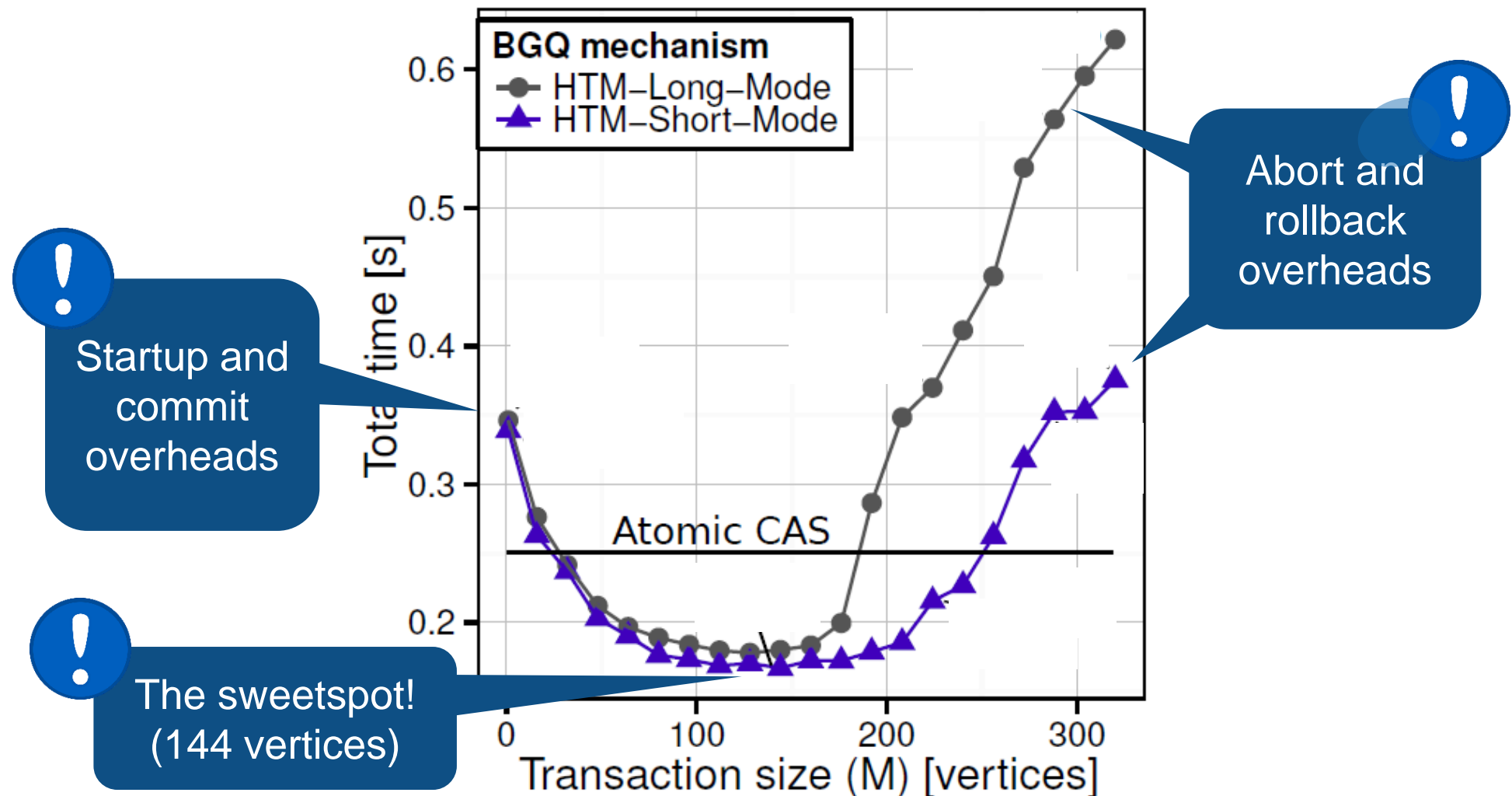
The long mode results in higher latency than the short mode



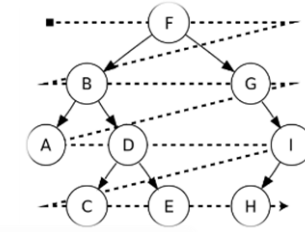


# MULTI-VERTEX TRANSACTIONS IN A BFS (GRAPH 500)

## MARKING VERTICES AS VISITED



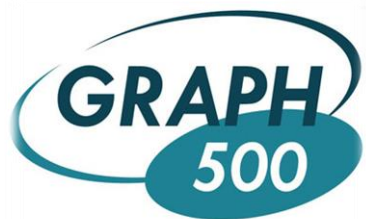
# REAL-GRAPH PERFORMANCE



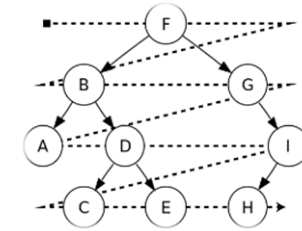
Input graph properties					BG/Q analysis				Haswell analysis				
Type	ID	Name	$V$	$E$	$S$ over $g500$ ( $M = 24$ )	$M$	$S$ over $g500$	$S$ over $g500$ ( $M = 2$ )	$S$ over Galois ( $M = 2$ )	$M$	$S$ over $g500$	$S$ over Galois	$S$ over HAMA
Comm. networks (RN)	eWT	wiki-Talk	2.4M	5M	2.82	45	1.35	0.91	1.2	8	0.96	1.28	344
	EDU	email-EuAll	265k	120k	3.07	32	1.36	0.76	0.88	5	0.97	1.12	1448
Social networks (SN)	slm	slm-Eu	1.8M	89M	1.18	19	1.88	1.05	1.1	8	1.05	1.13	$> 10^4$
	slm	slm-NA	1.8M	89M	1.18	19	1.88	1.05	1.1	8	1.05	1.13	$> 10^4$
	sDB	comp-DB	1.8M	89M	1.18	19	1.88	1.05	1.1	8	1.05	1.13	$> 10^4$
	sAM	comp-AM	1.8M	89M	1.18	19	1.88	1.05	1.1	8	1.05	1.13	$> 10^4$
Purchase network (PN)	pAM	amazon	1.8M	89M	1.18	19	1.88	1.05	1.1	8	1.05	1.13	$> 10^4$
	rCA	roadNet-CA	1.9M	3.5M	$> 1$	2	1.89	1.93	1.74	8	1.98	1.80	$> 10^4$
	rTX	roadNet-TX	1.3M	3.8M	$> 1$	2	1.83	1.29	1.89	6	1.42	2.08	$> 10^4$
Road networks (RN)	rPA	roadNet-PA	1M	3M	$> 1$	2	1.83	$> 1$	2.00	9	1.07	2.10	$> 10^4$
	cIP	cit-Patents	3.7M	16.5M	1.16	8	1.67	1.01	1.26	2	1.01	1.26	1875
Citation graphs (CG)	wGL	web-Google	875k	5.1M	1.75	12	2.08	0.98	1.26	6	1.06	1.35	305
	wDS	web-Berkeley	685k	7.6M	1.91	24	1.91	0.93	1.31	5	1.07	1.40	755
	wSP	web-Stanford	281k	2.3M	1.80	24	1.80	0.98	1.54	5	1.07	1.58	1077
Web graphs (WG)													

😊 No, you don't have to read it.

😊 Here: just a summary.



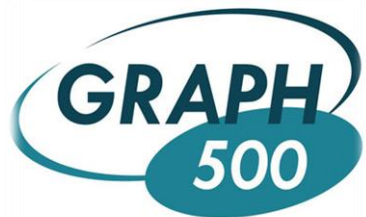
# REAL-GRAPH PERFORMANCE



Average overall speedup (geomean) over Graph 500: 1.07,  
Galois [1]: 1.40, HAMA: ~1000



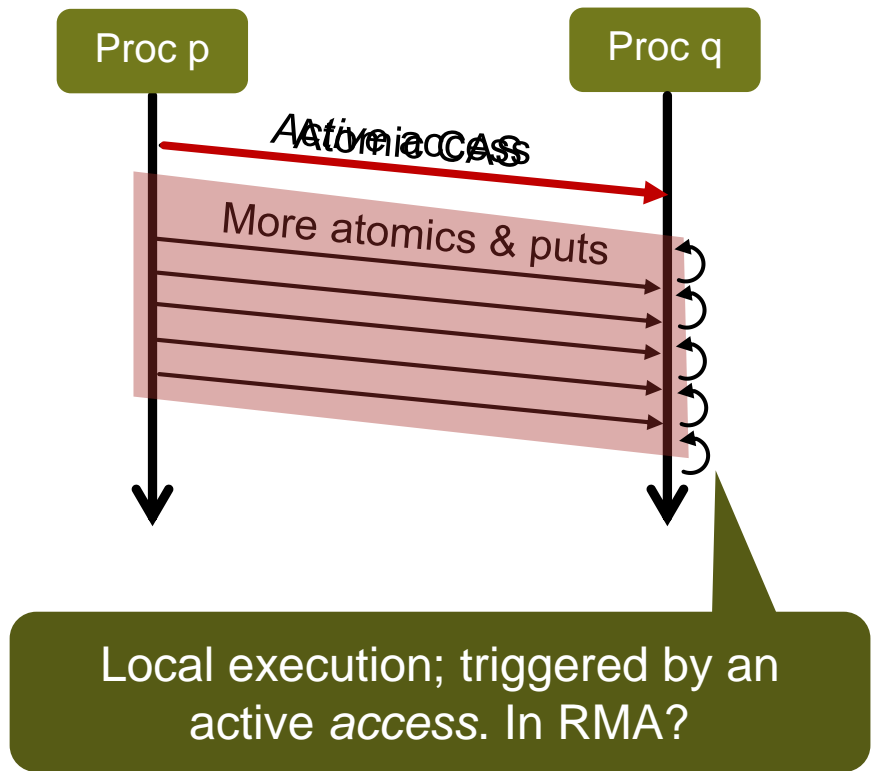
1.85x on average, up to 4.3x



# Remote Invocation



# IMAGINE A SIMPLE DISTRIBUTED HASH-TABLE



No collision:

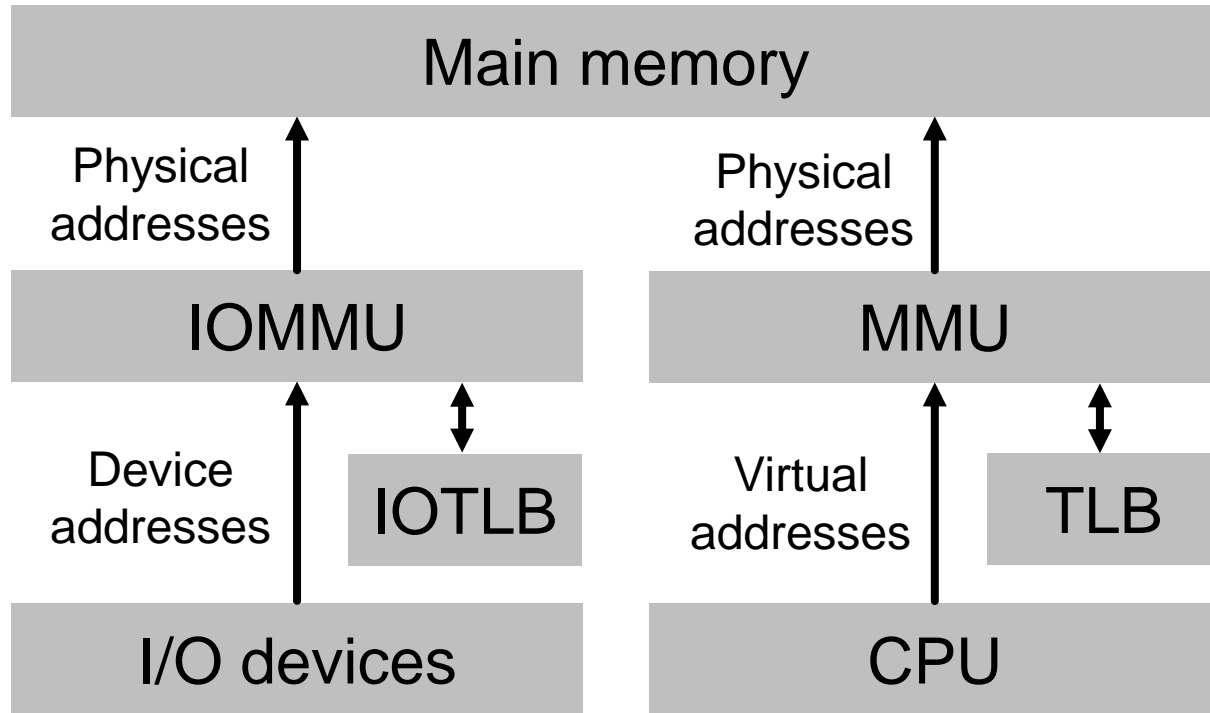
- 1 remote atomic
- Up to 5x speedup over MP [1]

A collision:

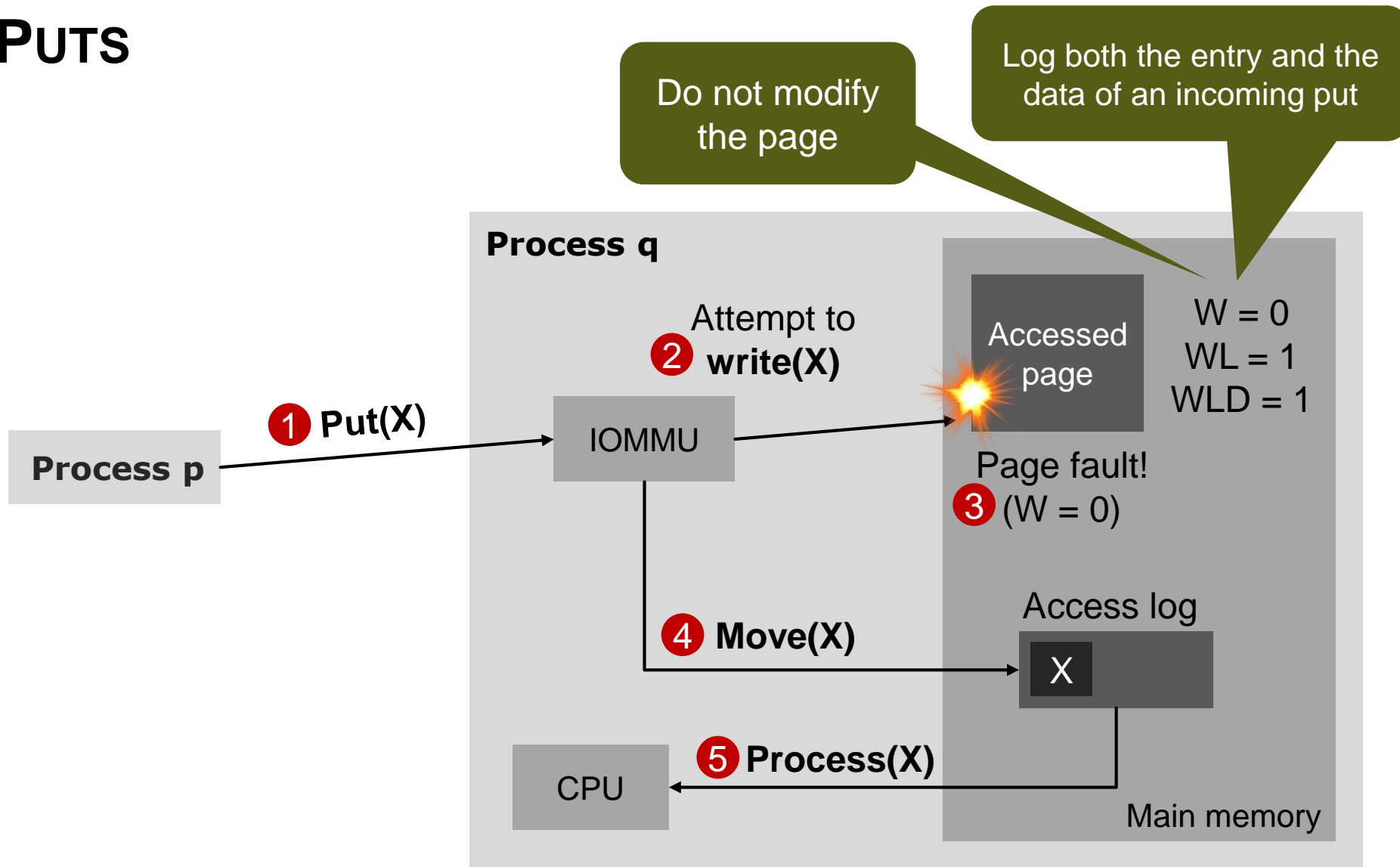
- 4 remote atomics + 2 remote puts
- Significant performance drops

Processes	3% Collisions	7% Collisions	14% Collisions	25% Collisions
0	0	0	0	0
250	~30	~20	~15	~10
500	~60	~40	~30	~20
1000	~100	~70	~50	~35

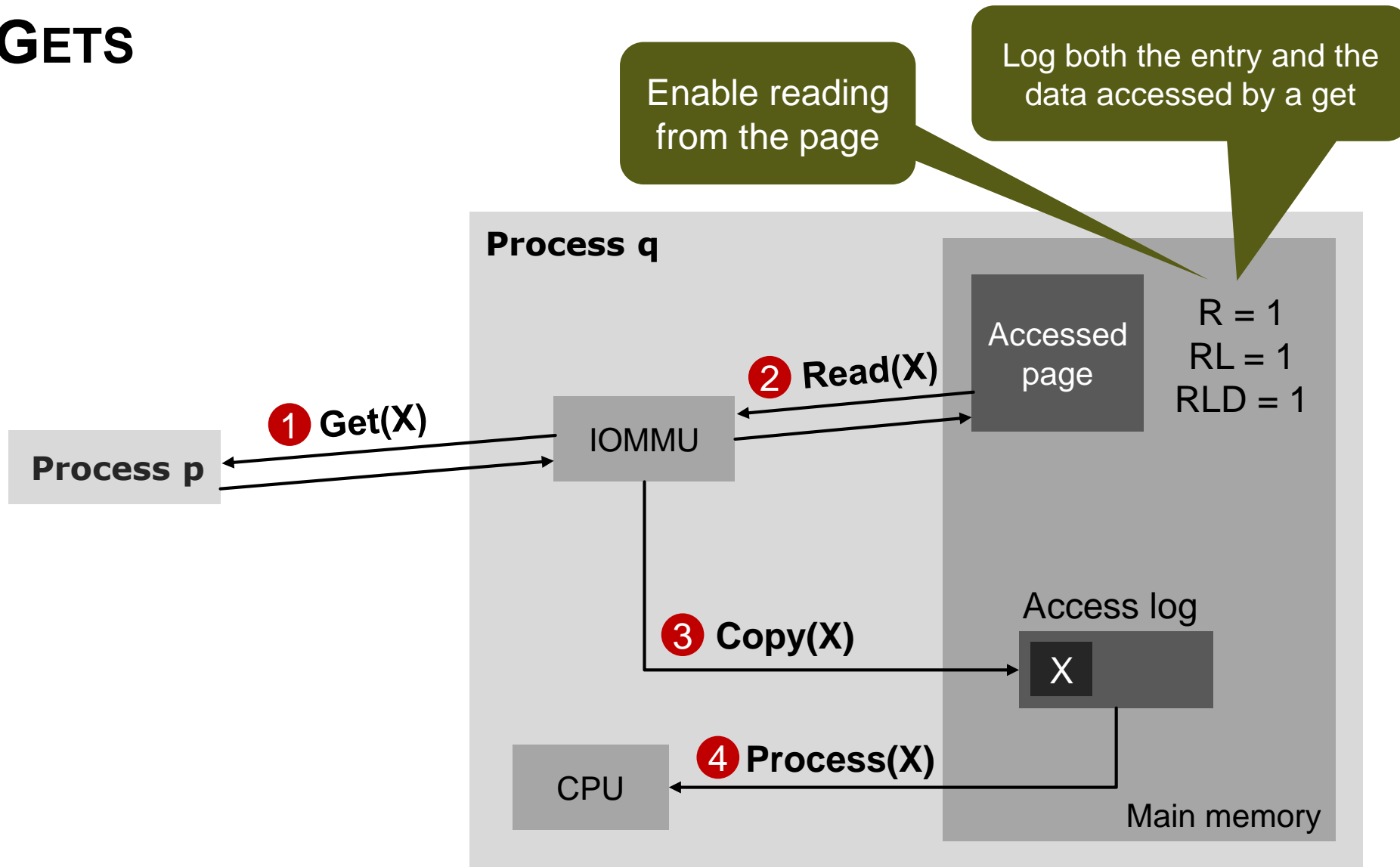
# USE INPUT/OUTPUT MEMORY MANAGEMENT UNITS



# ACTIVE PUTS

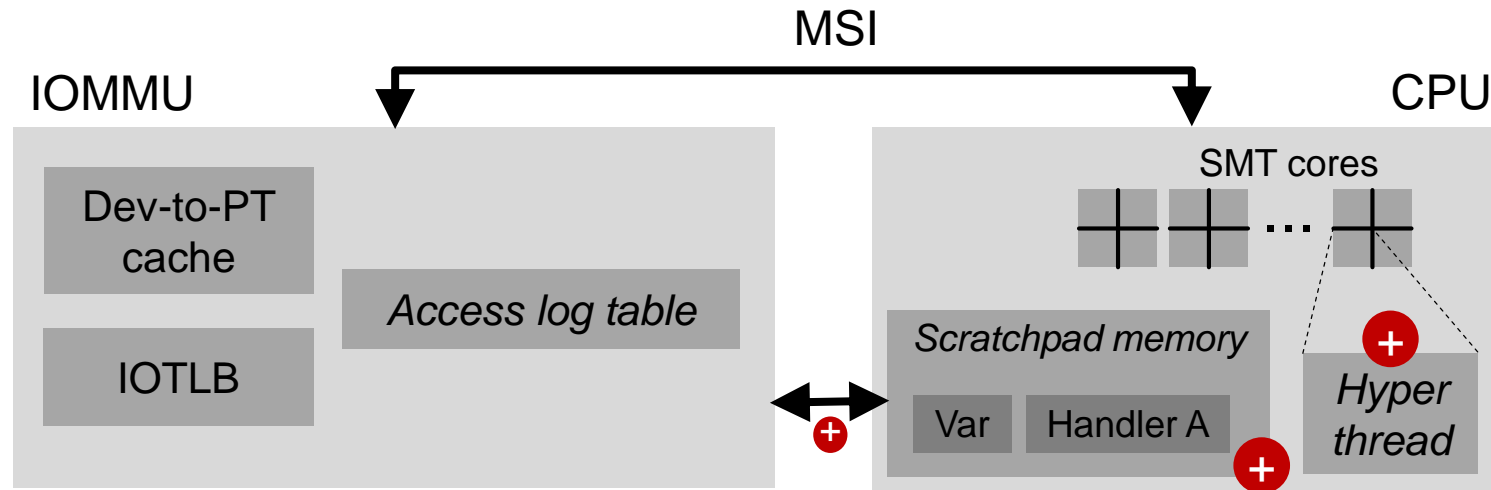


# ACTIVE GETS





# INTERACTIONS WITH THE CPU



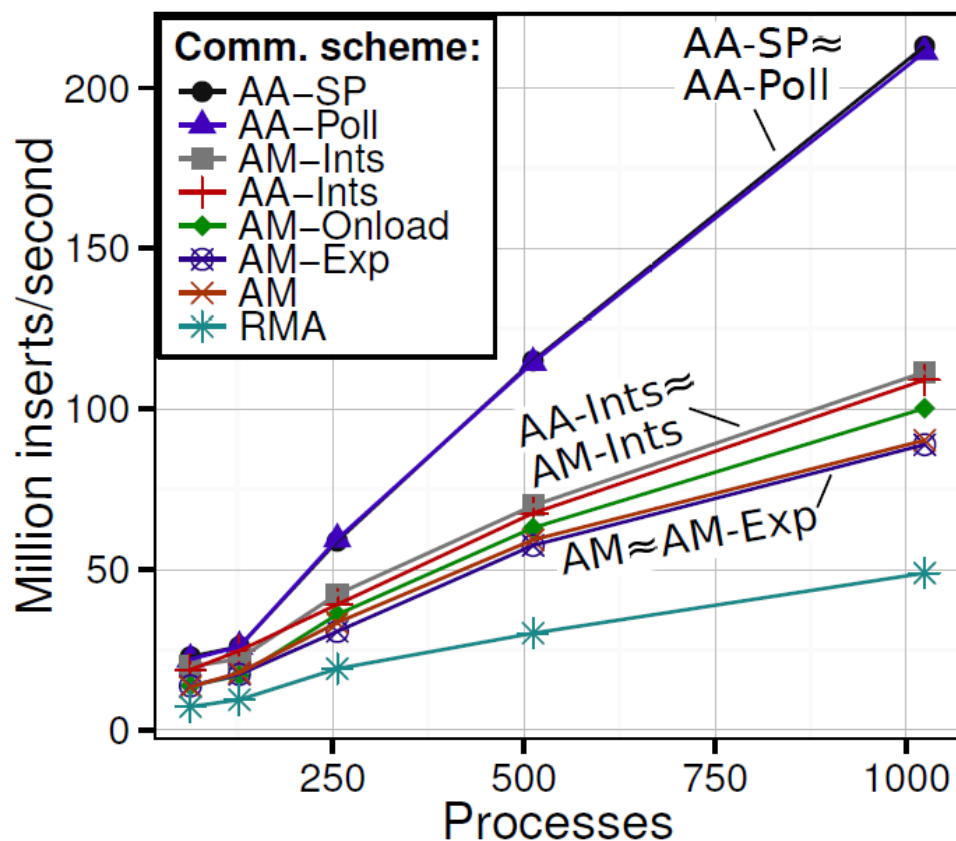
- Interrupts
- Polling
- Direct notifications via scratchpads

# PERFORMANCE: LARGE-SCALE CODES

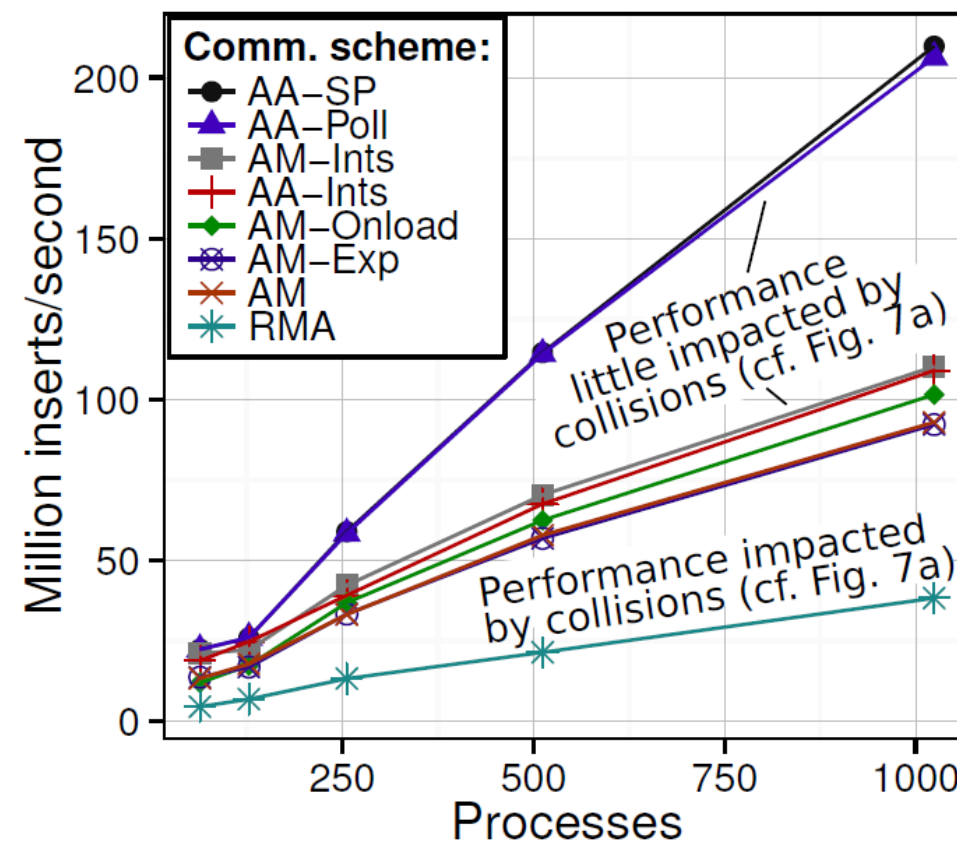
## DISTRIBUTED HASHTABLE



Collisions: 5%

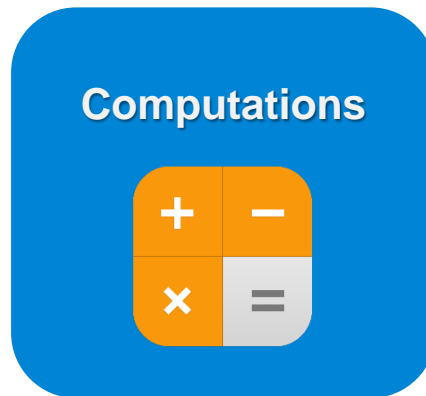


Collisions: 25%



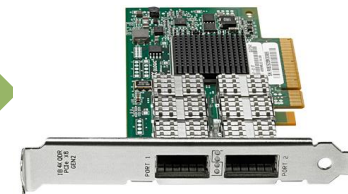
# Towards a Network Instruction Set Architecture (NISA)

## An example for offloading



```
L0: recv a from P1;
L1: b = compute f(buff, a);
L2: send b to P1;
L0 and CPU-> L1
L1 -> L2
```

CPU



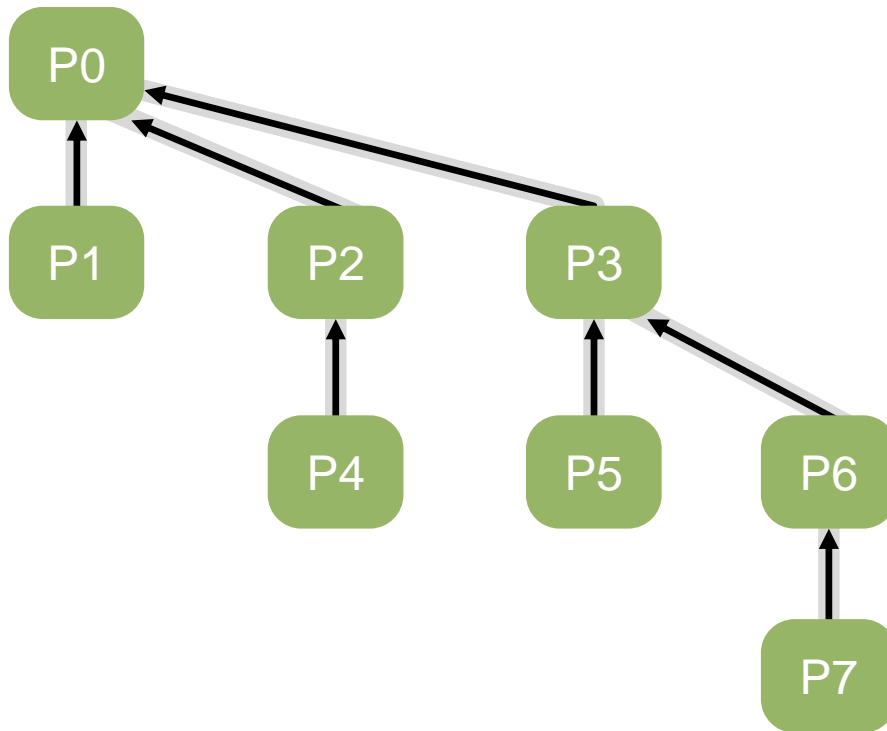
Offload Engine



# Fully Offloaded Collectives

**Collective communication:** A communication that involves a group of processes

**Non-blocking collective:** Once initiated the operation may progress independently of any computation or other communication at participating processes





# Fully Offloaded Collectives

**Collective communication:** A communication that involves a group of processes

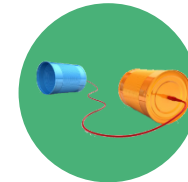
**Non-blocking collective:** Once initiated the operation may progress independently of any computation or other communication at participating processes

P0

## Fully Offloading:

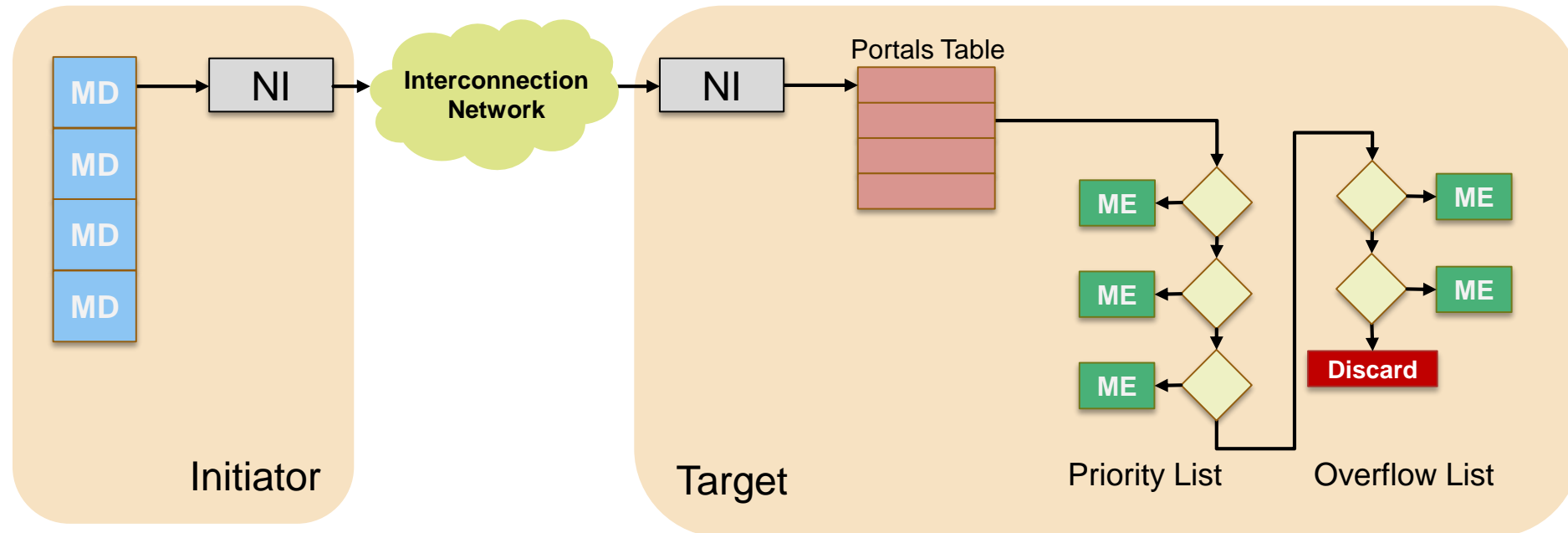
1. *No synchronization* is required in order to start the collective operation
2. Once a collective operation is started, *no further CPU intervention* is required in order to progress or complete it.

P7



# A Case Study: Portals 4

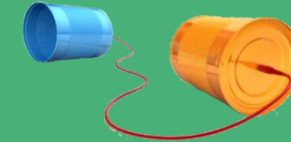
- Based on the one-sided communication model
- Matching/Non-Matching semantics can be adopted



# A Case Study: Portals 4

## Communication primitives

- Put/Get operations are natively supported by Portals 4
- One-sided + matching semantic



## Atomic operations

- Operands are the data specified by the MD at the initiator and by the ME at the target
- Available operators: *min*, *max*, *sum*, *prod*, *swap*, *and*, *or*, ...



## Counters

- Associated with MDs or MEs
- Count specific events (e.g., operation completion)

## Triggered operations

- Put/Get/Atomic associated with a counter
- Executed when the associated counter reaches the specified threshold



# FFlib: An Example

Proof of concept library implemented on top of Portals 4

```
ff_schedule_h sched = ff_schedule_create(...);
```

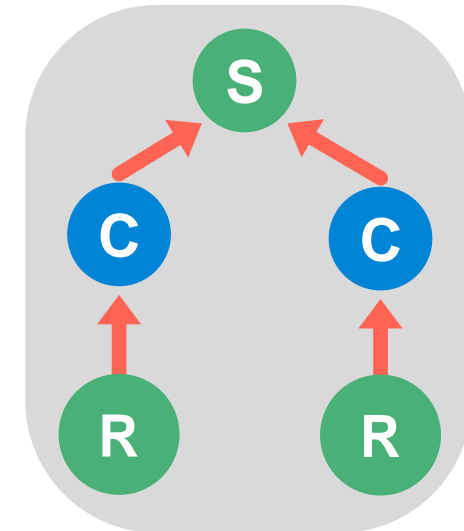
```
ff_op_h r1 = ff_op_create_recv(tmp + blocksize, blocksize, child1, tag);  
ff_op_h r2 = ff_op_create_recv(tmp + 2*blocksize, blocksize, child2, tag);
```

```
ff_op_h c1 = ff_op_create_computation(rbuff, blocksize, tmp + blocksize, blocksize, operator, datatype, tag)  
ff_op_h c2 = ff_op_create_computation(rbuff, blocksize, tmp + 2*blocksize, blocksize, operator, datatype, tag)
```

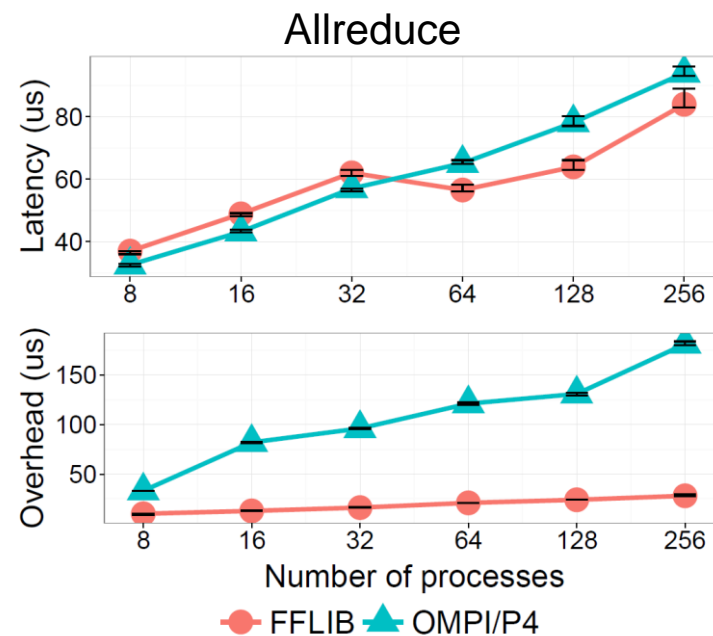
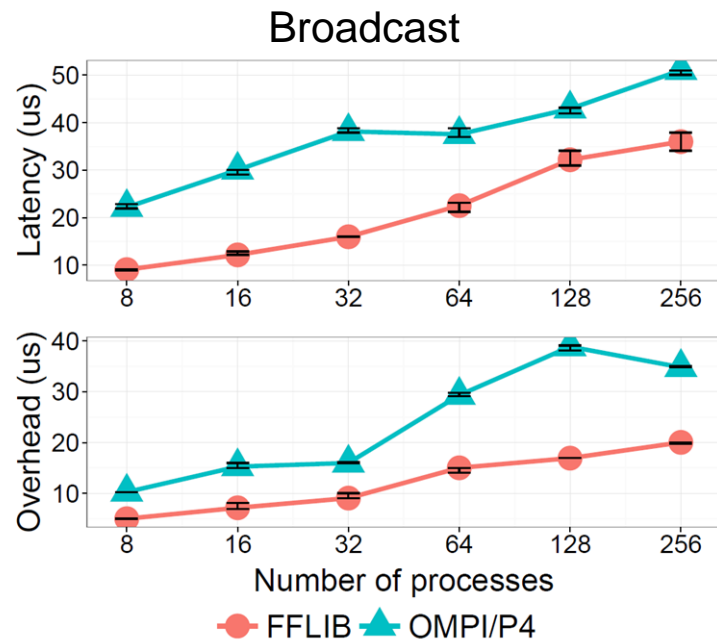
```
ff_op_h s = ff_op_create_send(rbuff, blocksize, parent, tag)
```

```
ff_op_hb(r1, c1)  
ff_op_hb(r2, c2)  
ff_op_hb(c1, s)  
ff_op_hb(c2, s)
```

```
ff_schedule_add(sched, r1)  
ff_schedule_add(sched, r2)  
ff_schedule_add(sched, c1)  
ff_schedule_add(sched, c2)  
ff_schedule_add(sched, s)
```



# Experimental Results: Latency/Overhead



Target machine: Curie

5,040 nodes

2 eight-core Intel Sandy Bridge processors

Full fat-tree Infiniband QDR

OMPI/P4: Open MPI 1.8.4 + Portals 4 RL

FFLIB: proof of concept library

More about FFLIB at : [http://spcl.inf.ethz.ch/Research/Parallel\\_Programming/FFlib/](http://spcl.inf.ethz.ch/Research/Parallel_Programming/FFlib/)



# Active RDMA – what could it be?

## NISA: Process the data while it moves!

Remote Synchronization  
[IPDPS'15]

- Extend RMA semantics
- Fully one-sided (in HW)
- Synchronization

Network  
Instruction Set  
Architecture  
(NISA)

- Remote Transactions  
[HPDC'15]
- Similar to HTM
  - Extend across nodes
  - Think active memory

Remote Invocation

[ICS'15]

- Utilizes IOMMUs
- Control transfer
- Active memory