

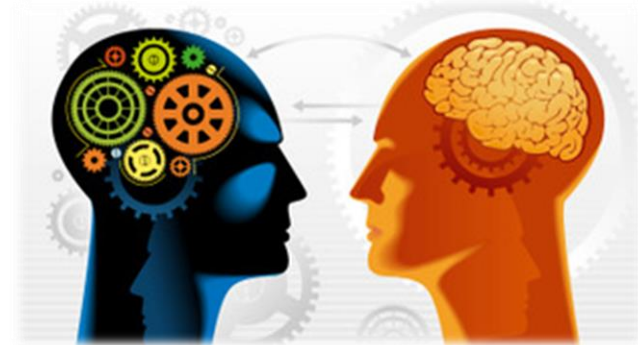
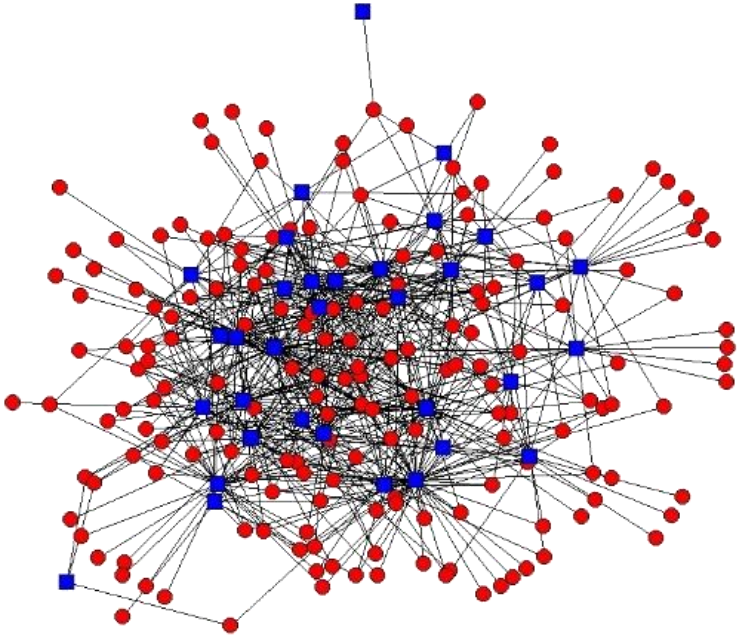
Towards scalable RDMA locking on a NIC

TORSTEN HOEFLER

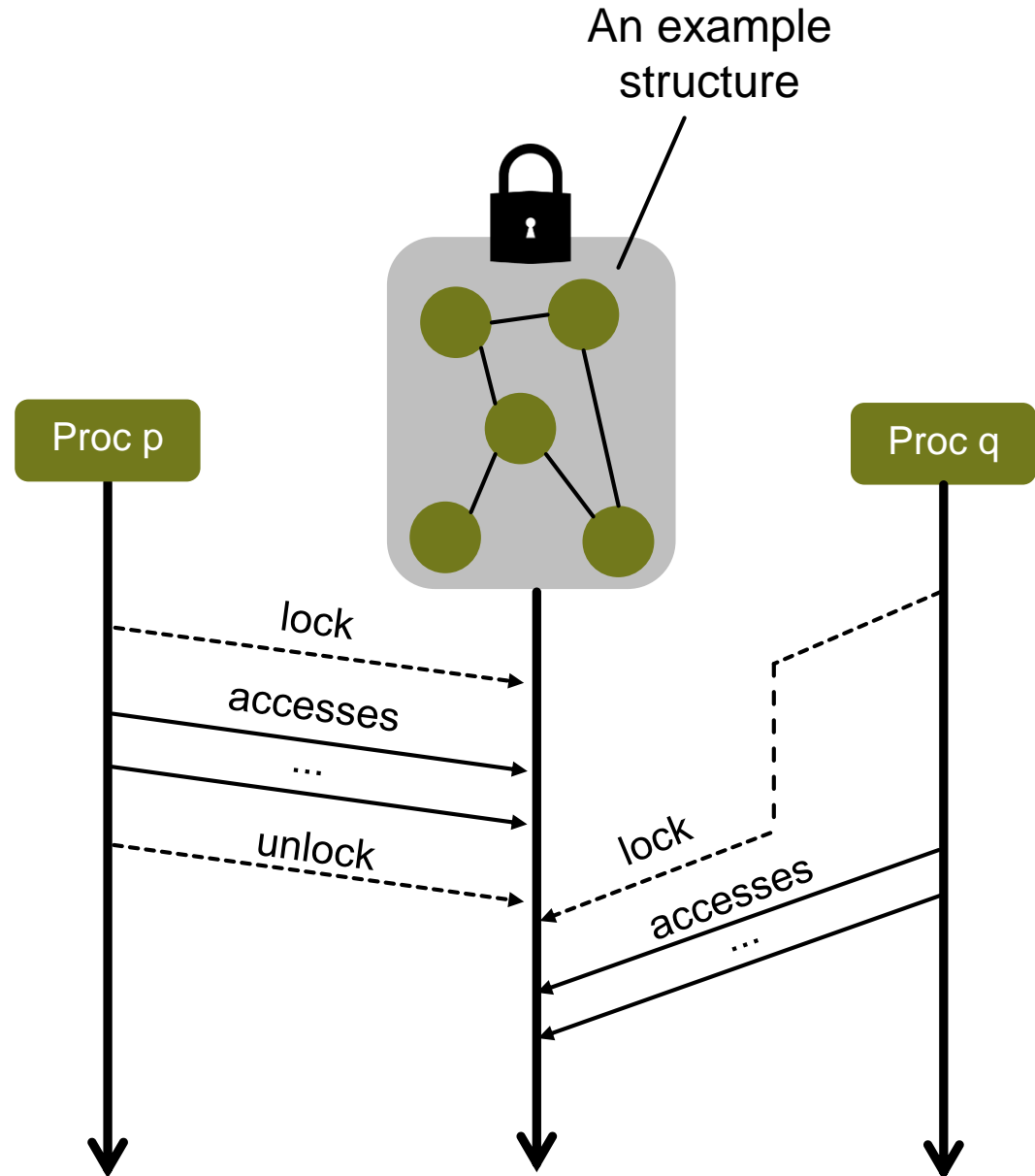
with support of Patrick Schmid, Maciej Besta, Salvatore di Girolamo @ SPCL
presented at HP Labs, Palo Alto, CA, USA



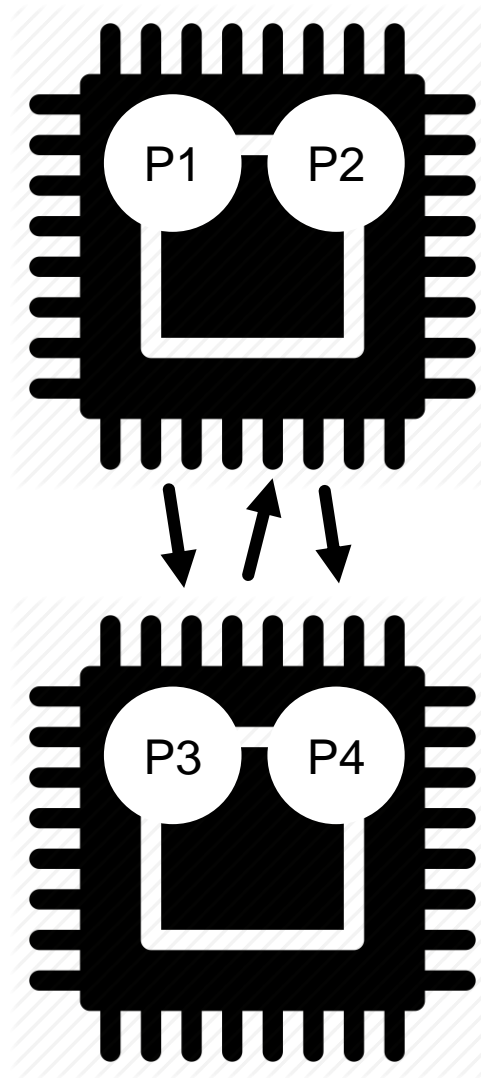
NEED FOR EFFICIENT LARGE-SCALE SYNCHRONIZATION



LOCKS



LOCKS: CHALLENGES



LOCKS: CHALLENGES



We need intra- and inter-node topology-awareness



We need to cover arbitrary topologies



LOCKS: CHALLENGES



We need to distinguish
between readers and writers

Reader

Reader

Reader

Writer



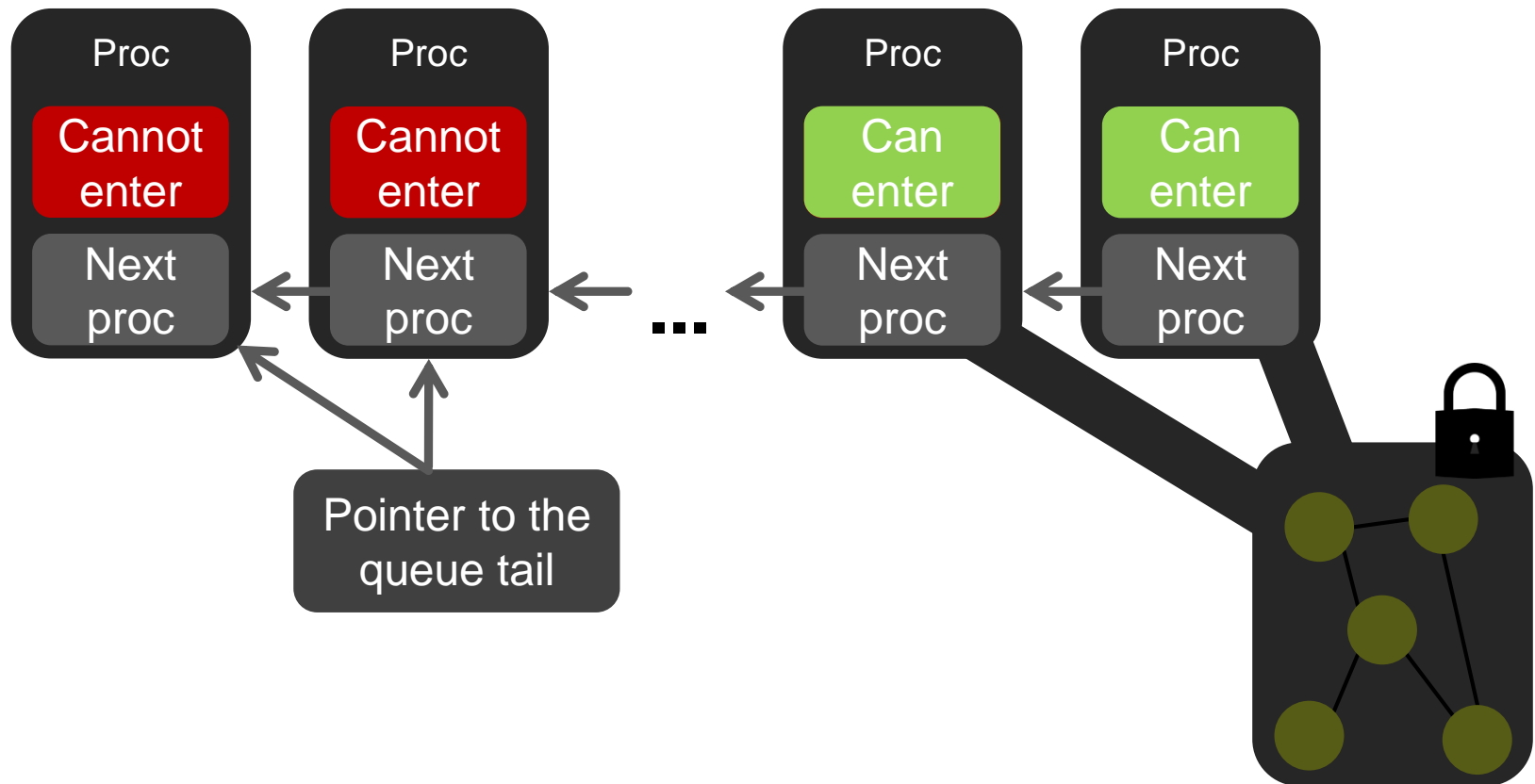
We need flexible
performance for both types
of processes



What will we use in the design?

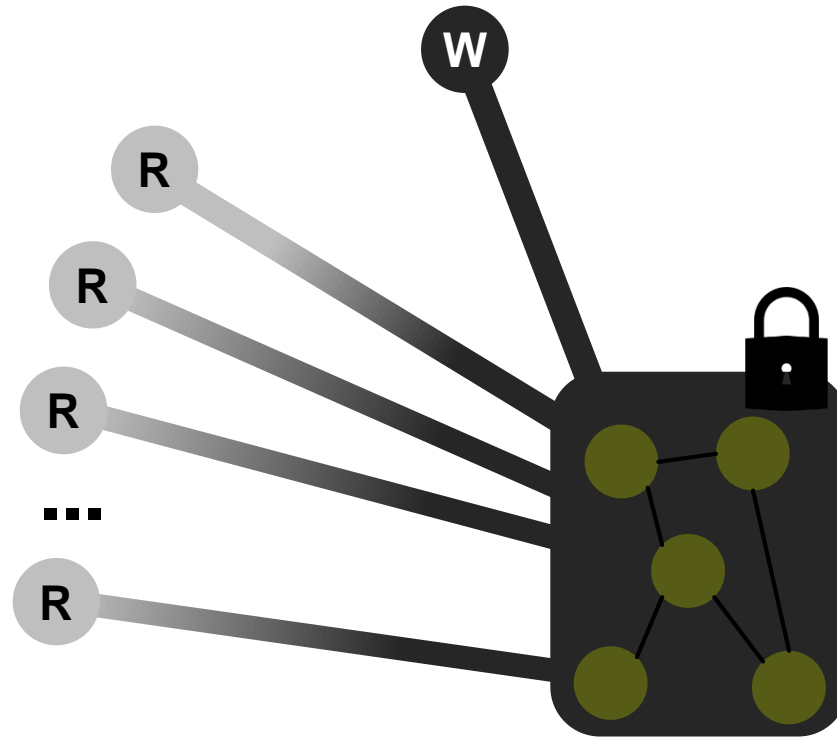
WHAT WE WILL USE

MCS Locks



WHAT WE WILL USE

Reader-Writer Locks





How to manage the design complexity?

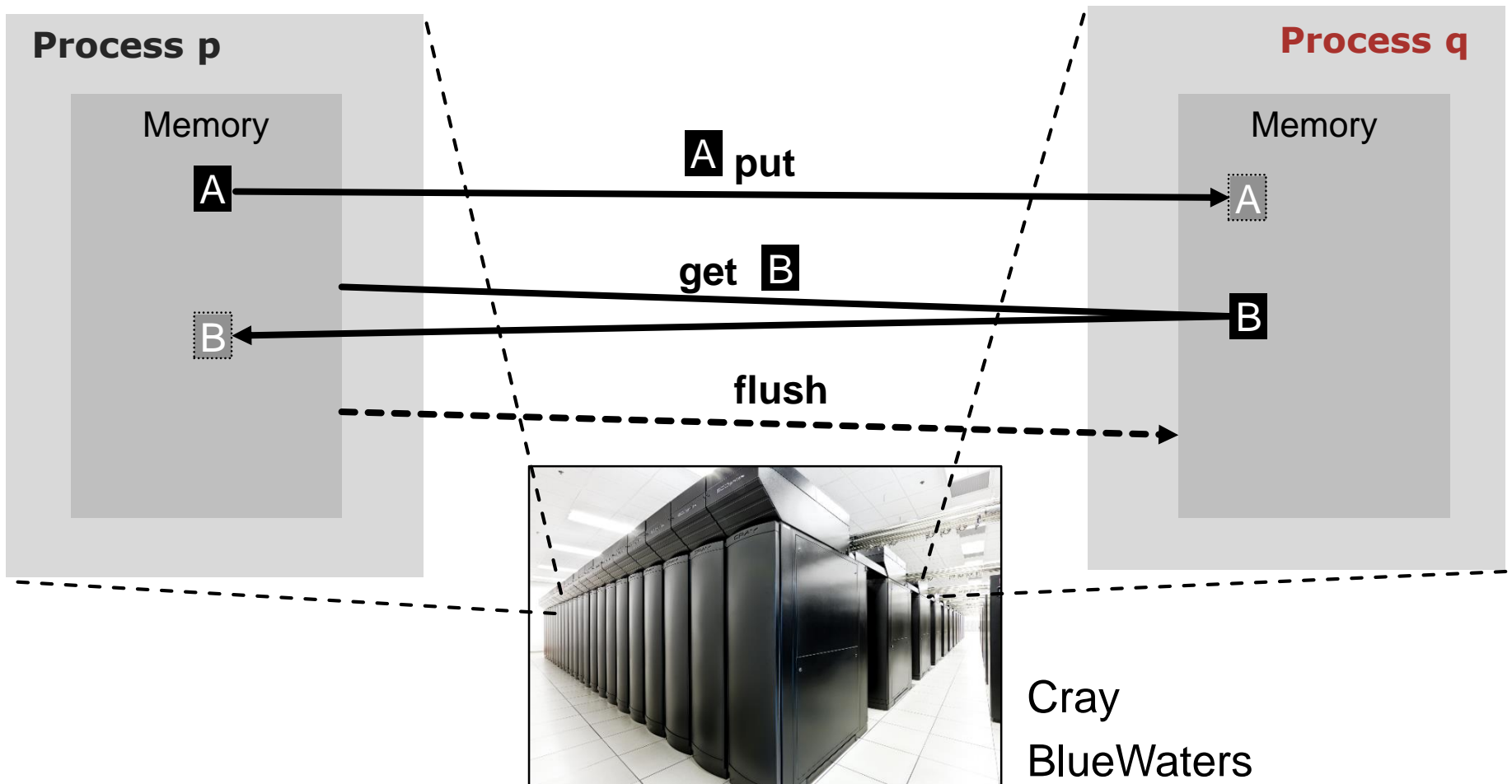


How to ensure tunable performance?



What mechanism to use for efficient implementation?

REMOTE MEMORY ACCESS (RMA) PROGRAMMING



REMOTE MEMORY ACCESS PROGRAMMING

- Implemented in hardware in NICs in the majority of HPC networks support RDMA





How to manage the design complexity?



How to ensure tunable performance?



What mechanism to use for efficient implementation?

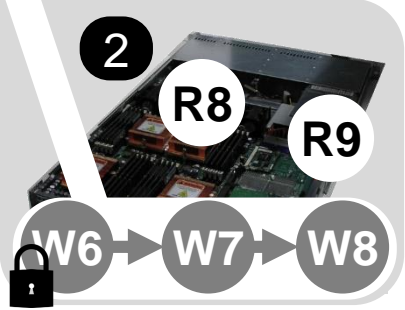
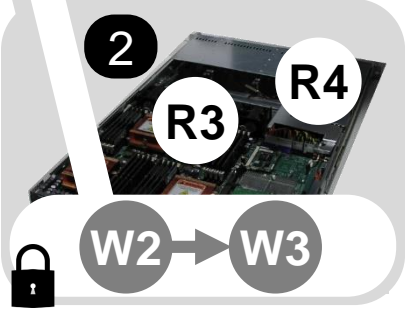
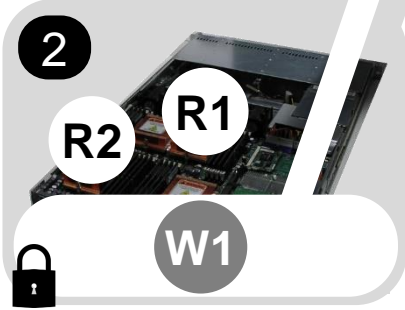
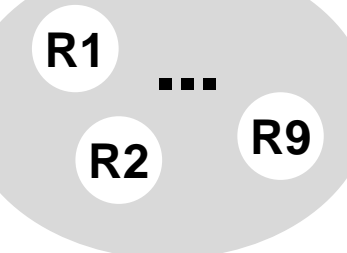
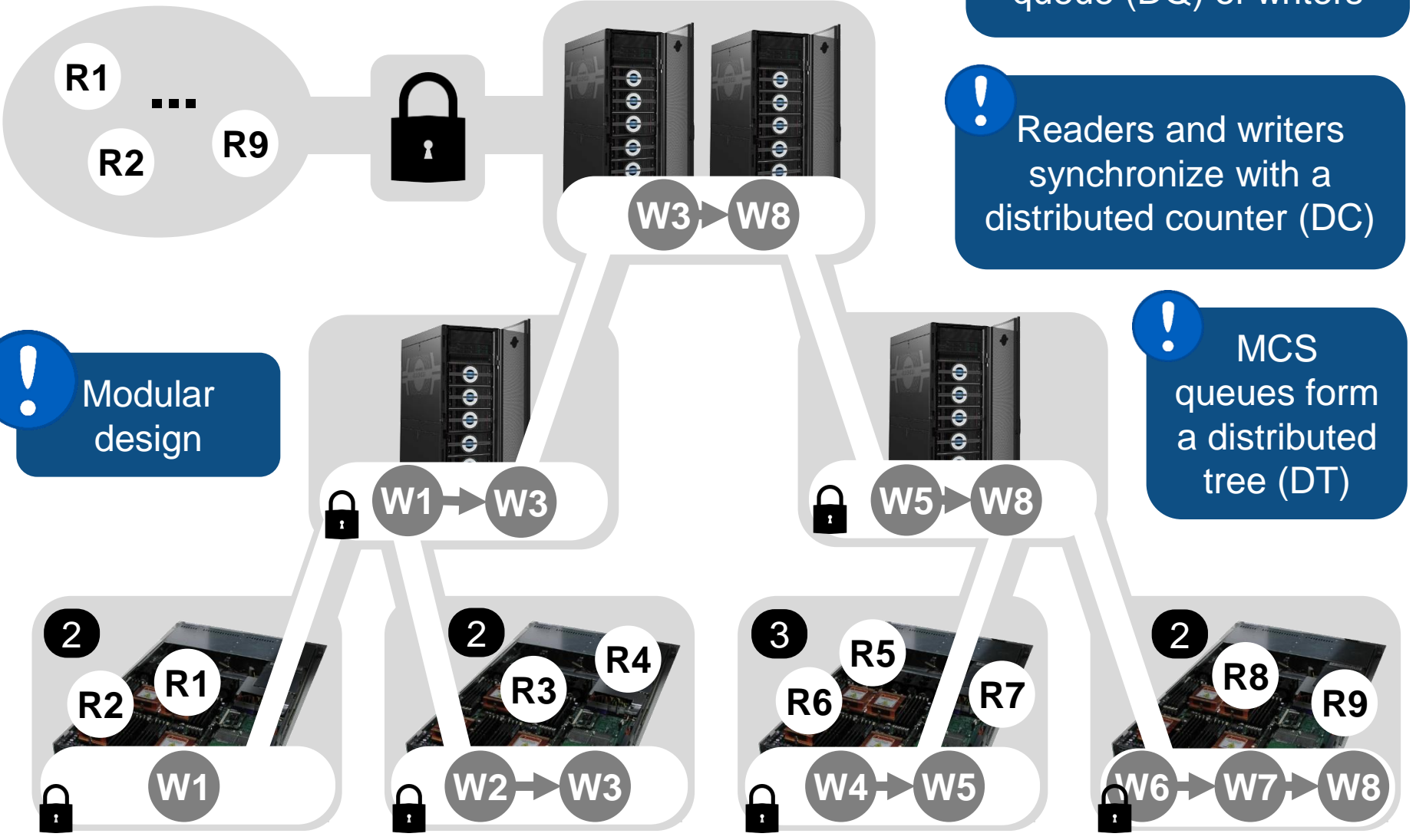
? How to manage the design complexity?

! Each element has its own distributed MCS queue (DQ) of writers

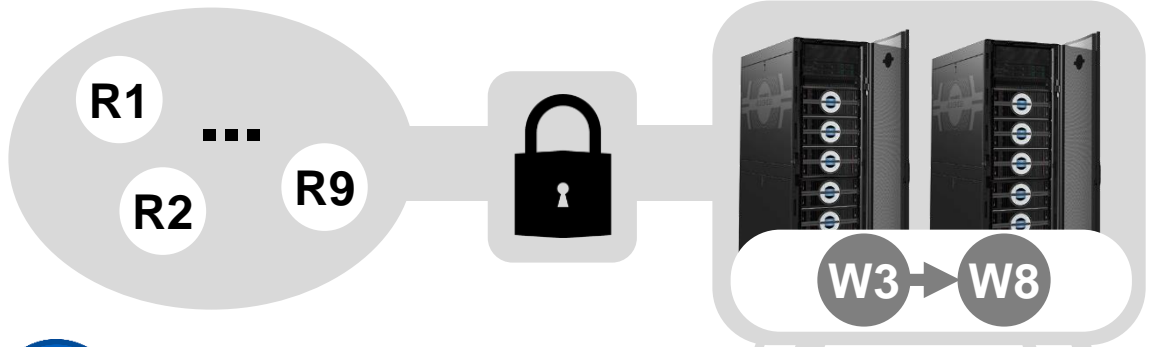
! Readers and writers synchronize with a distributed counter (DC)

! MCS queues form a distributed tree (DT)

! Modular design



? How to ensure tunable performance?

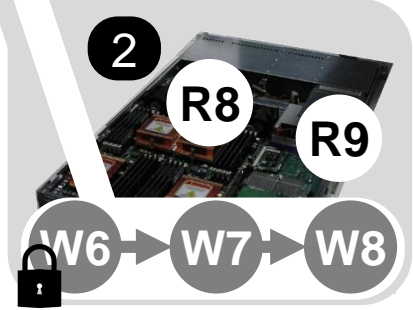
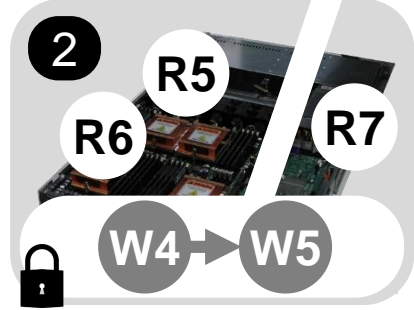
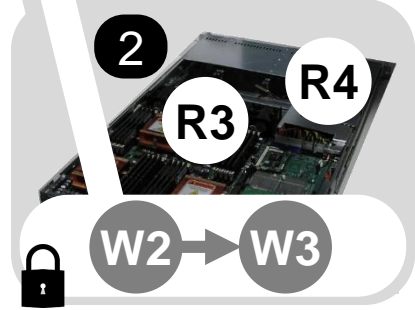
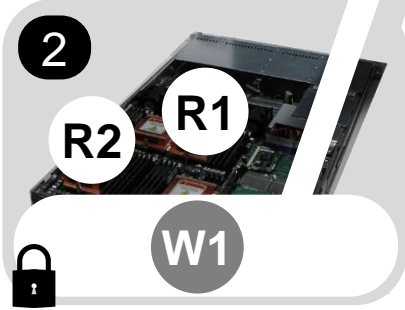


! Each DQ: fairness vs throughput of writers

! DC: a parameter for the latency of readers vs writers

! A tradeoff parameter for every structure

! DT: a parameter for the throughput of readers vs writers



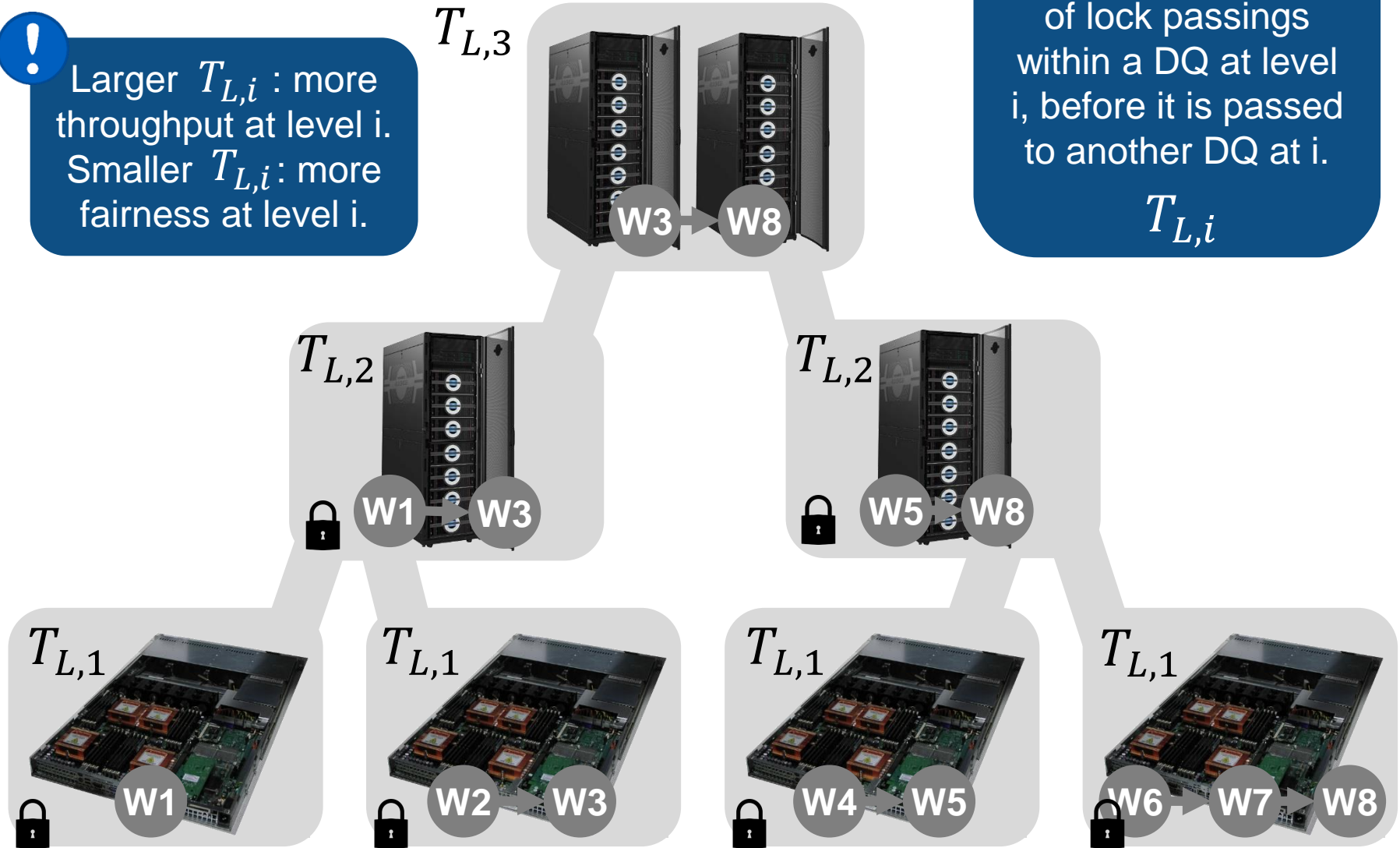
DISTRIBUTED MCS QUEUES (DQs)

Throughput vs Fairness

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

! Each DQ: The maximum number of lock passings within a DQ at level i , before it is passed to another DQ at i .

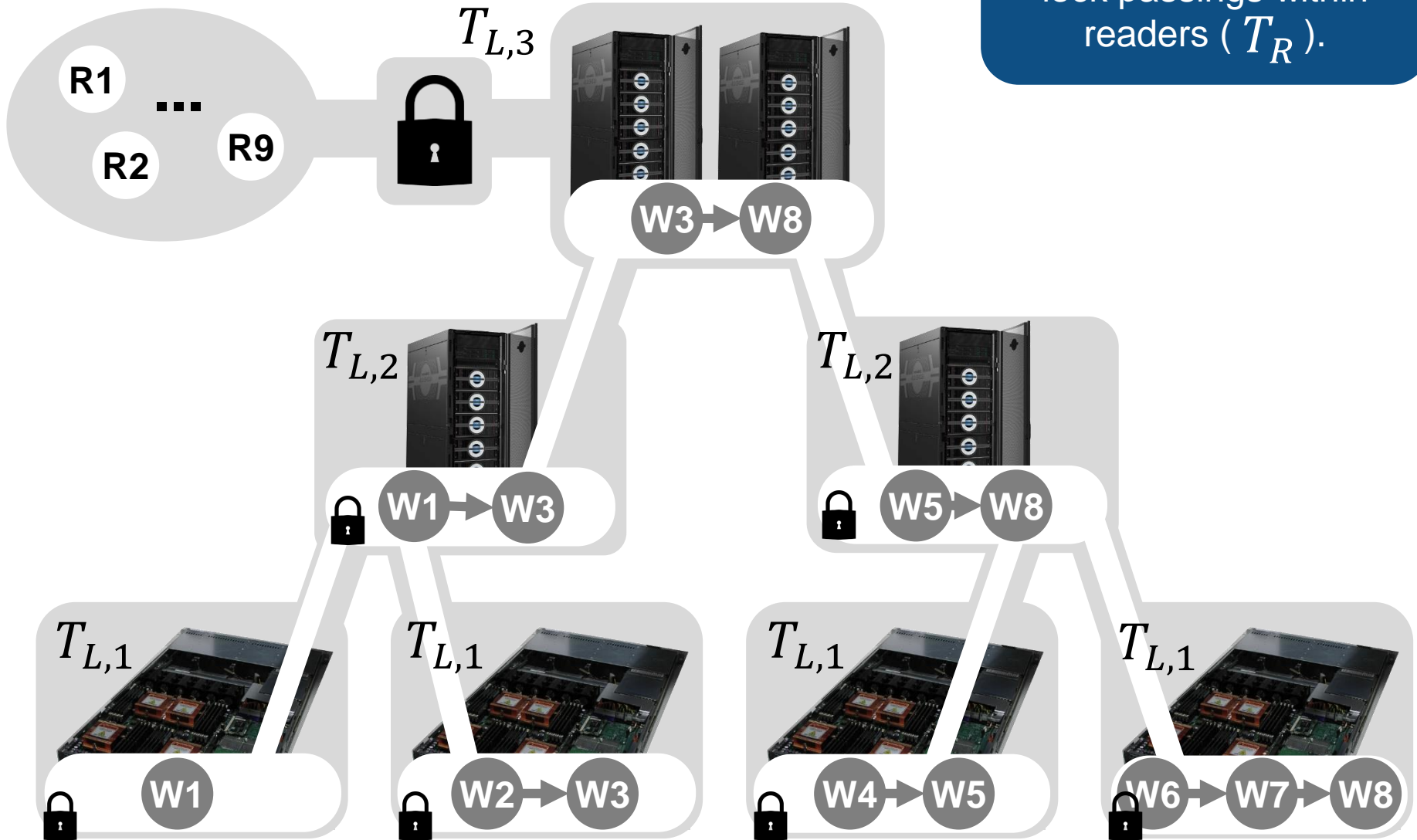
$T_{L,i}$



DISTRIBUTED TREE OF QUEUES (DT)

Throughput of readers vs writers

! DT: The maximum number of consecutive lock passings within readers (T_R).



DISTRIBUTED COUNTER (DC)

Latency of readers vs writers

DC: every k th compute node hosts a partial counter, all of which constitute the DC.

! $k = T_{DC}$



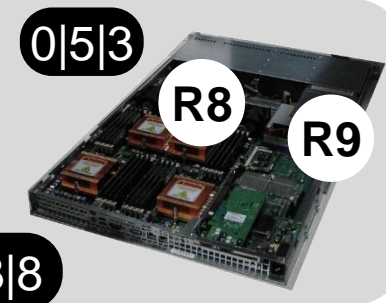
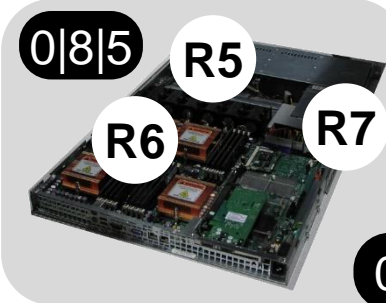
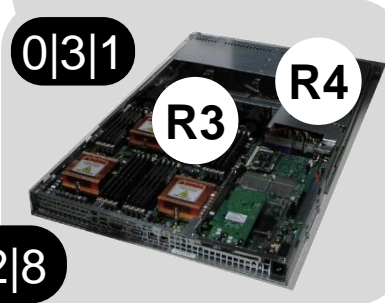
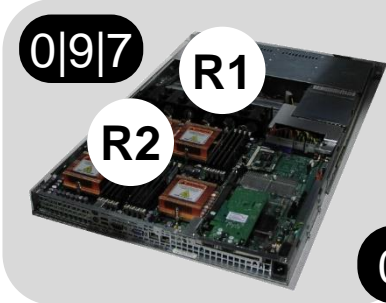
A writer holds the lock $b|x|y$

Readers that arrived at the CS

Readers that left the CS

$T_{DC} = 1$

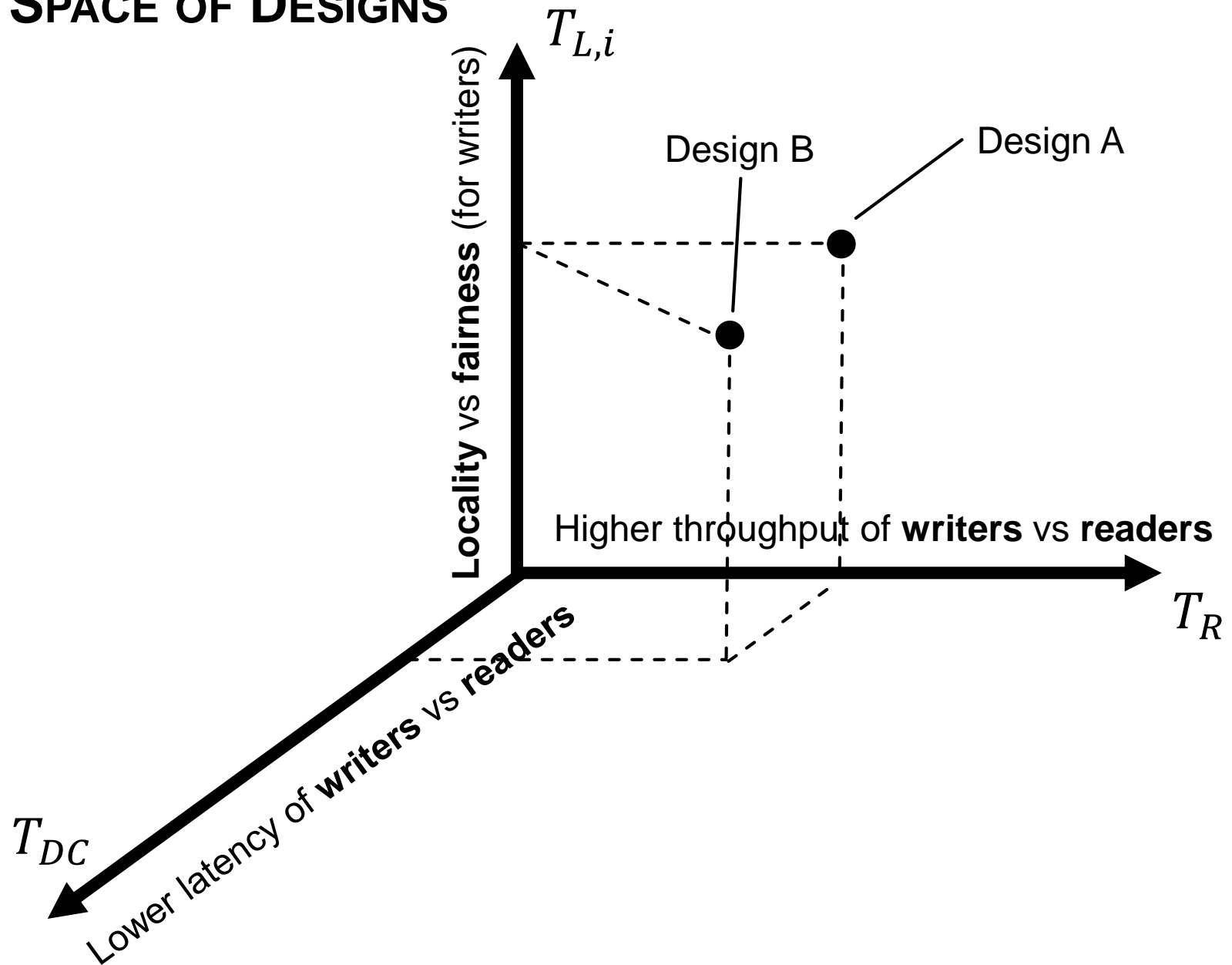
$T_{DC} = 2$



0|12|8

0|13|8

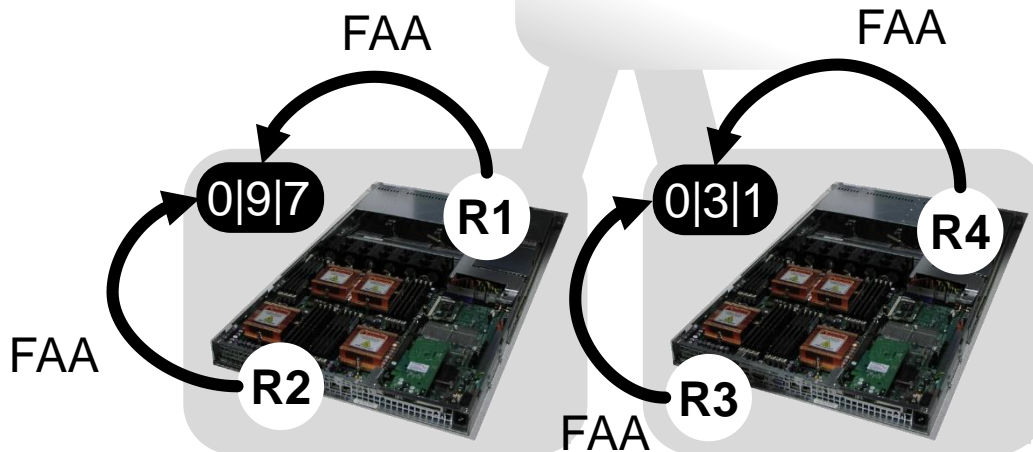
THE SPACE OF DESIGNS



LOCK ACQUIRE BY READERS



A lightweight acquire protocol for readers: only one atomic fetch-and-add (FAA) operation



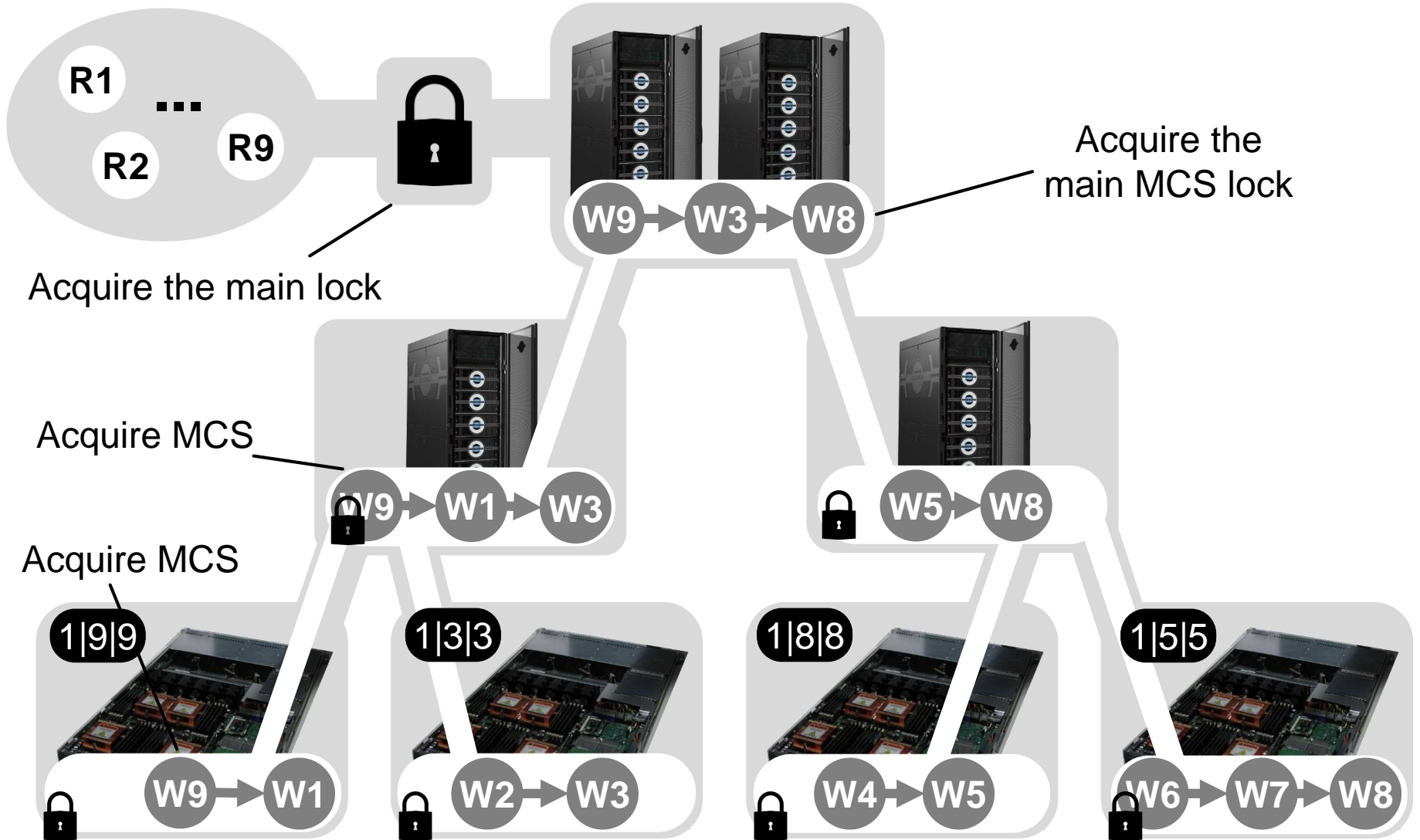
A writer holds
the lock

b|x|y

Readers that
arrived at the CS

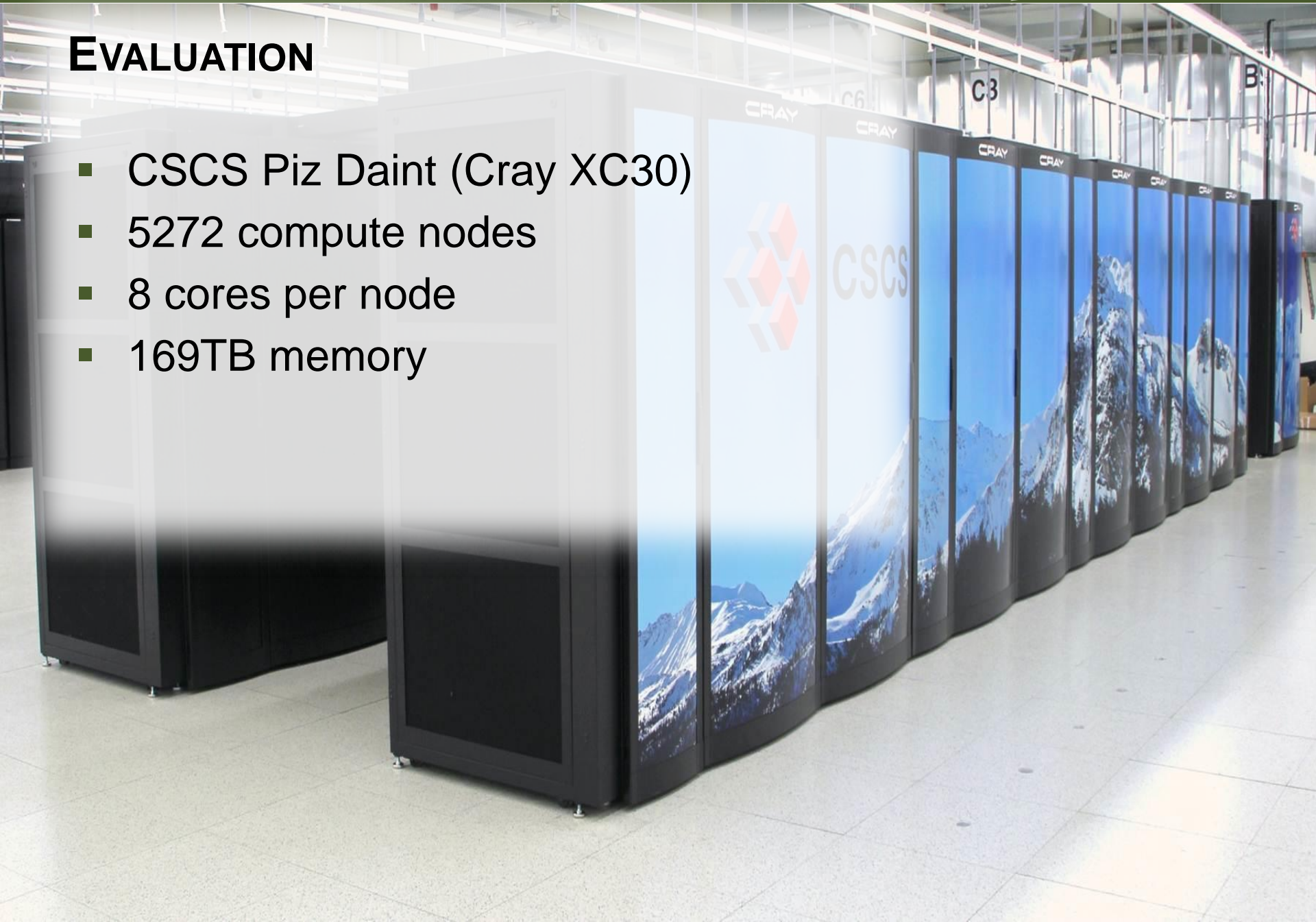
Readers that
left the CS

LOCK ACQUIRE BY WRITERS



EVALUATION

- CSCS Piz Daint (Cray XC30)
- 5272 compute nodes
- 8 cores per node
- 169TB memory



EVALUATION

CONSIDERED BENCHMARKS

The **latency**
benchmark

DHT

Distributed
hashtable
evaluation

Throughput
benchmarks:

Empty-critical-section

Single-operation

Wait-after-release

Workload-critical-section

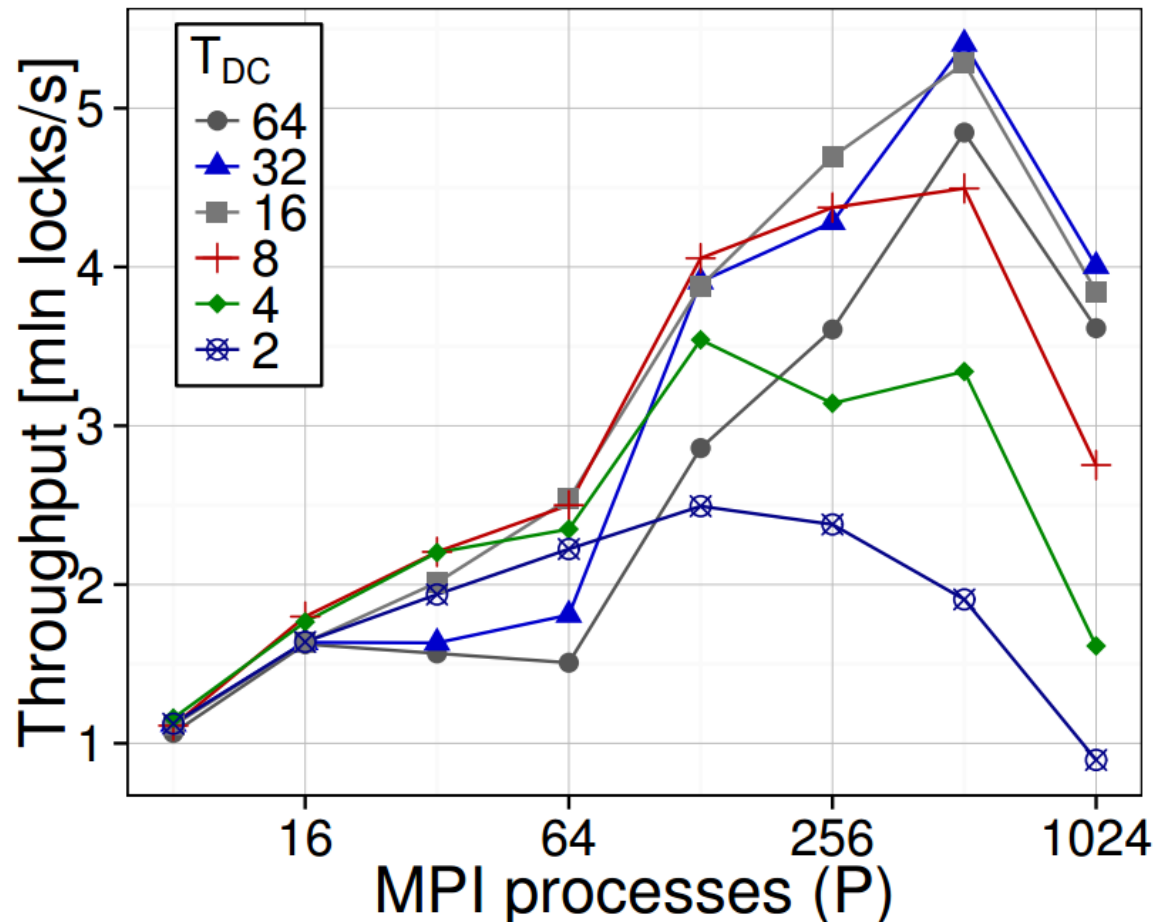
EVALUATION

DISTRIBUTED COUNTER ANALYSIS

0|9|7

0|3|1

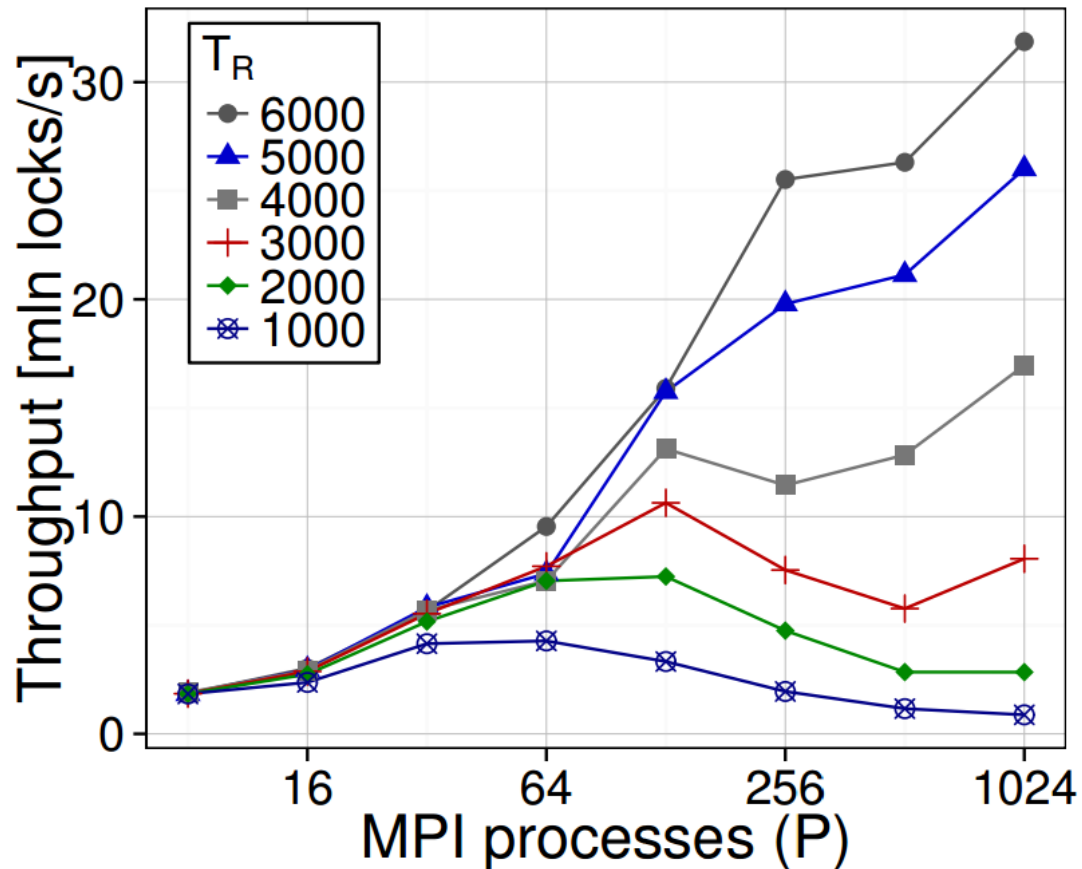
0|12|8

Throughput, 2% writers
Single-operation benchmark

EVALUATION

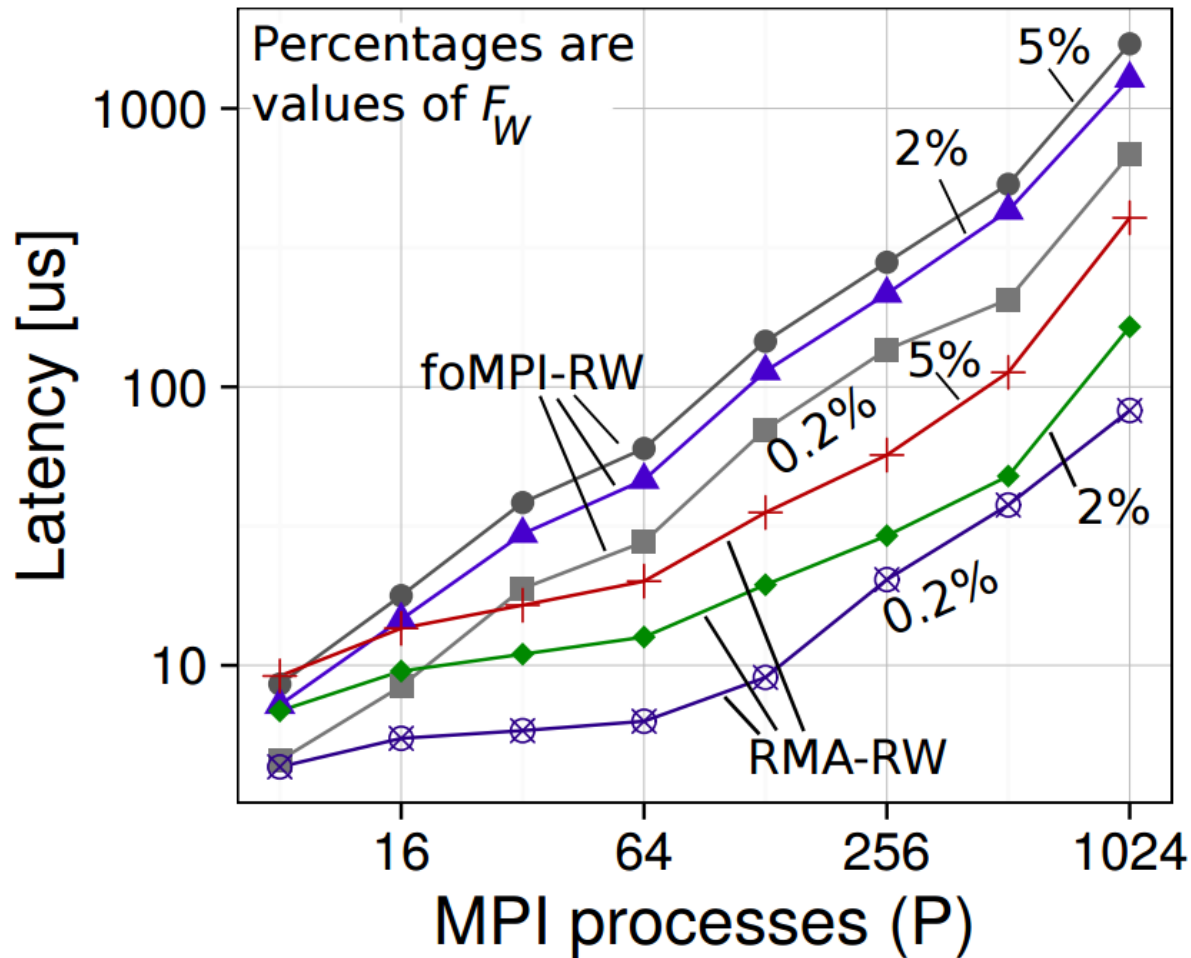
READER THRESHOLD ANALYSIS

Throughput, 0.2% writers,
Empty-critical-section benchmark



EVALUATION

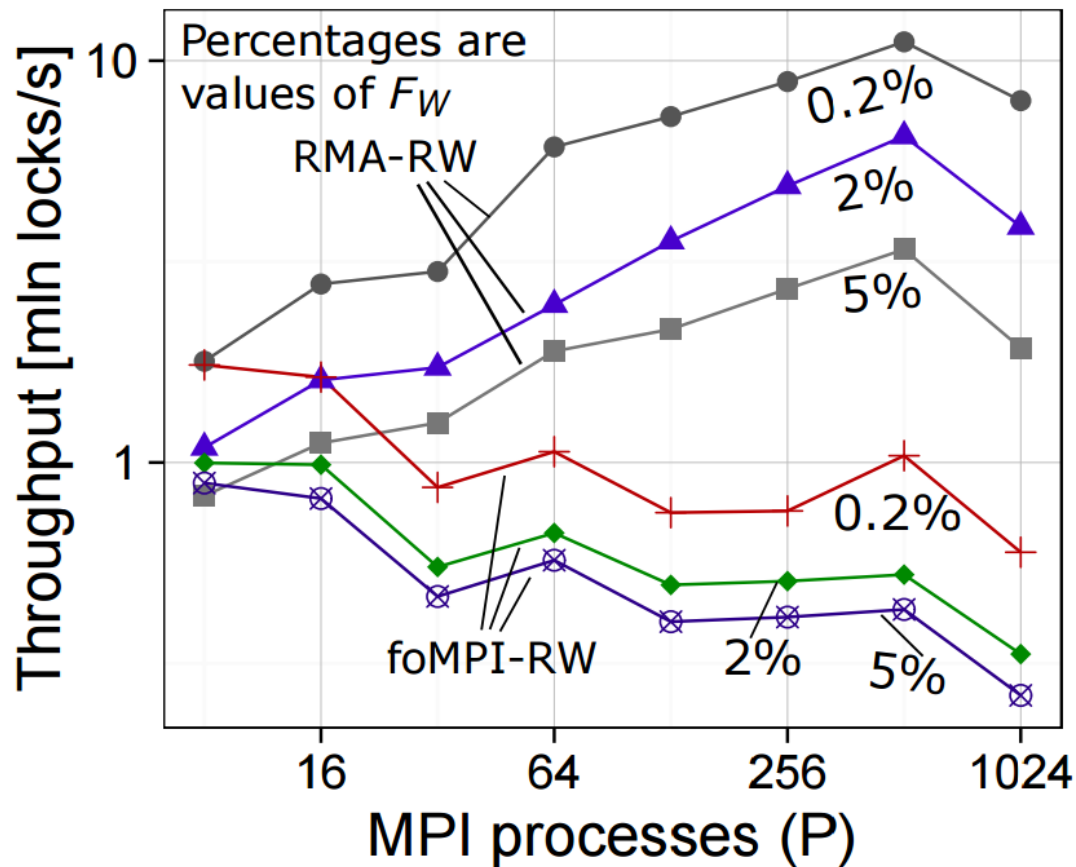
COMPARISON TO THE STATE-OF-THE-ART



EVALUATION

COMPARISON TO THE STATE-OF-THE-ART

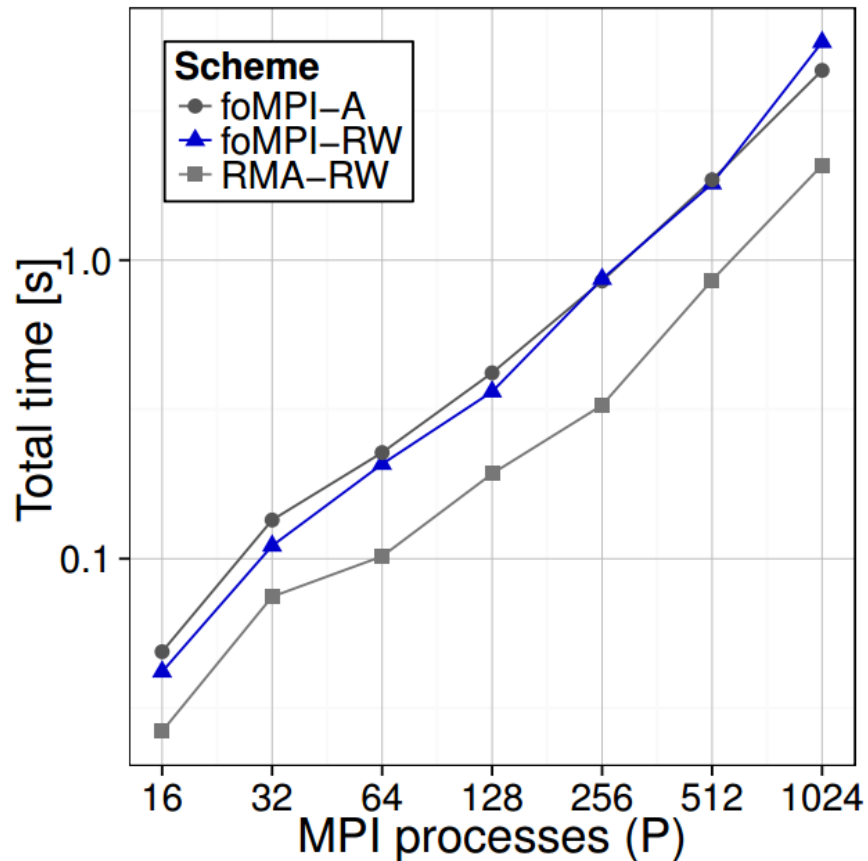
Throughput, single-operation benchmark



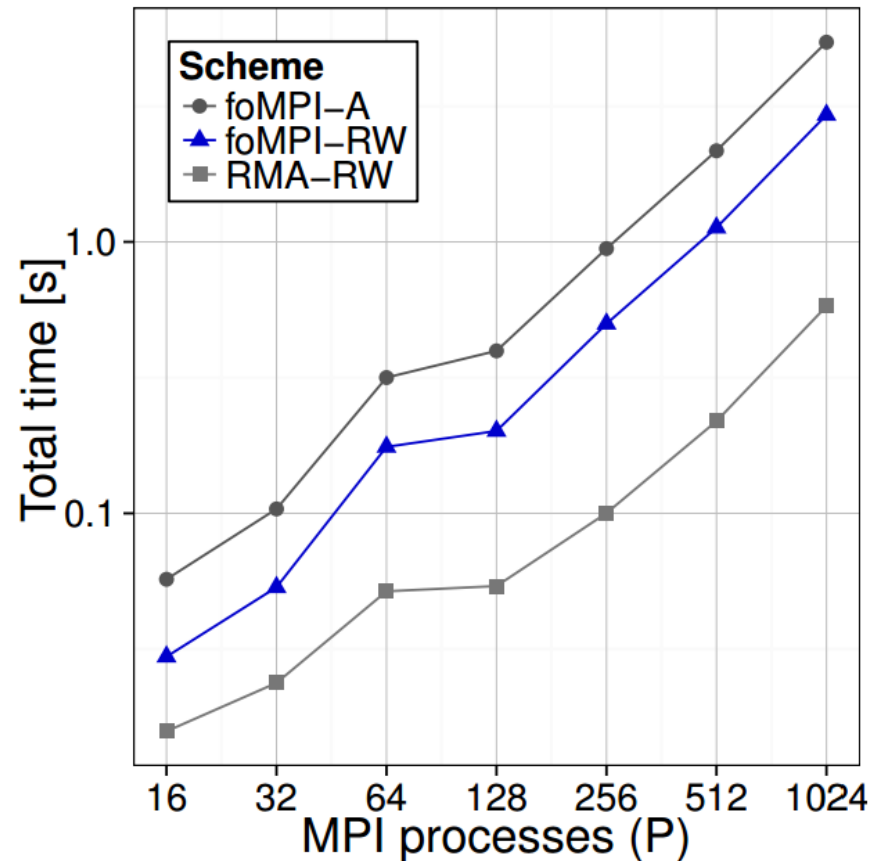
EVALUATION

DISTRIBUTED HASHTABLE

20% writers



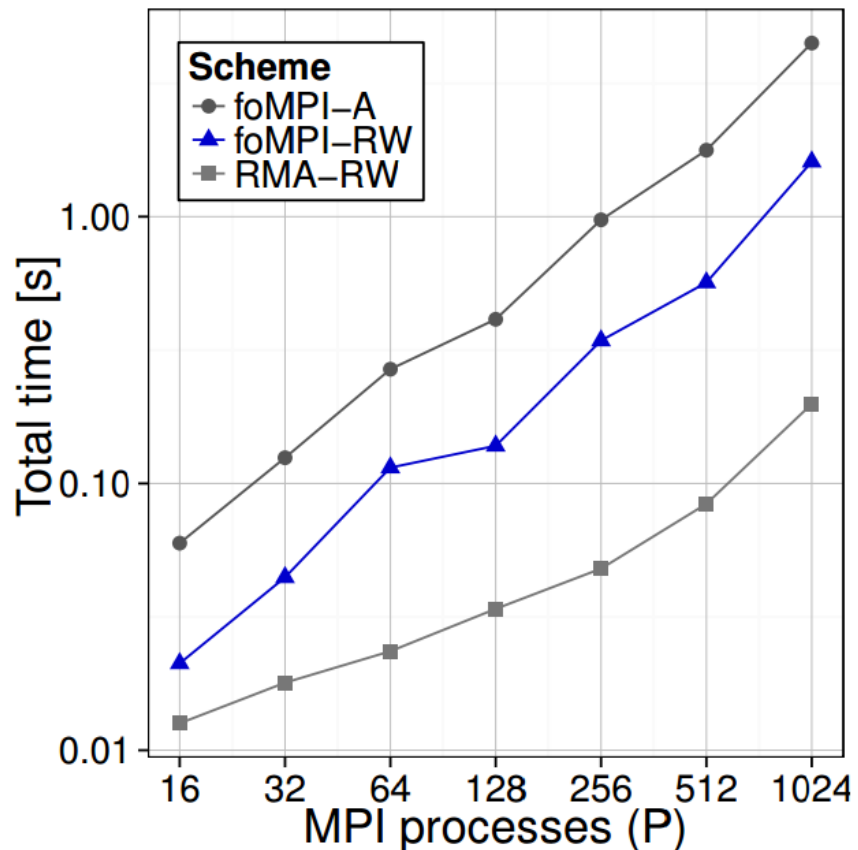
10% writers



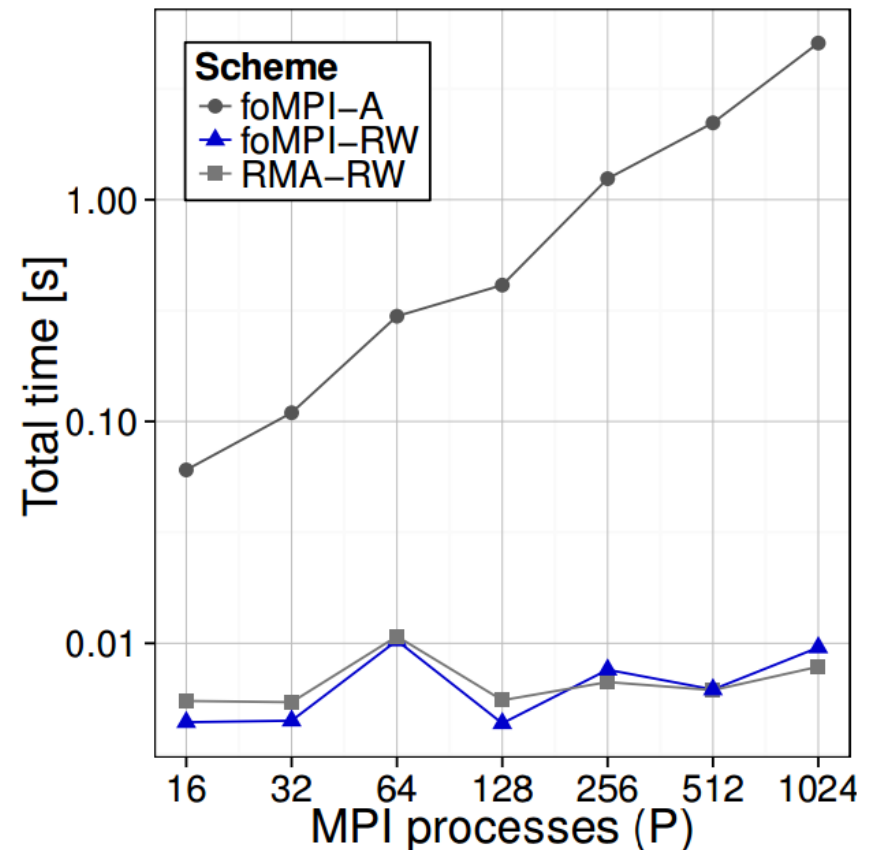
EVALUATION

DISTRIBUTED HASHTABLE

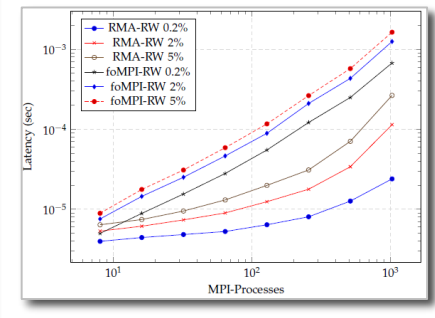
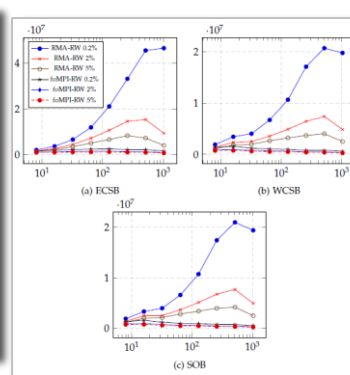
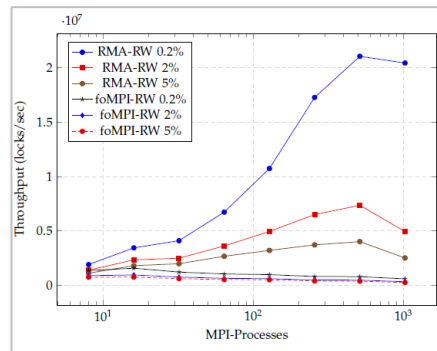
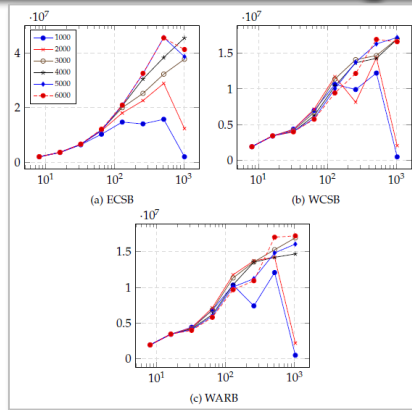
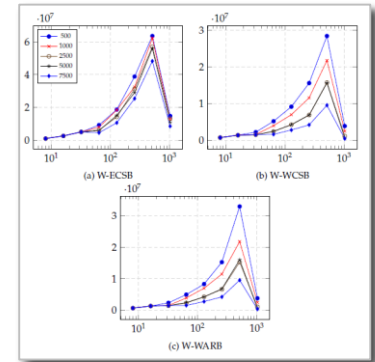
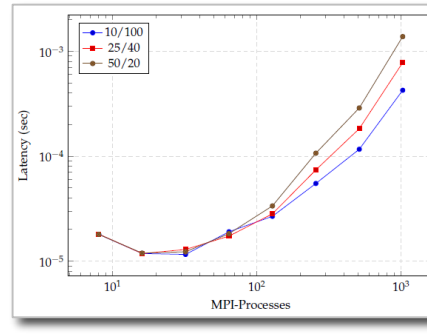
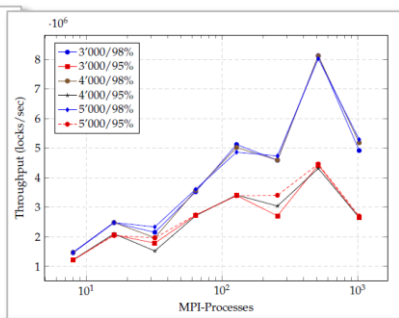
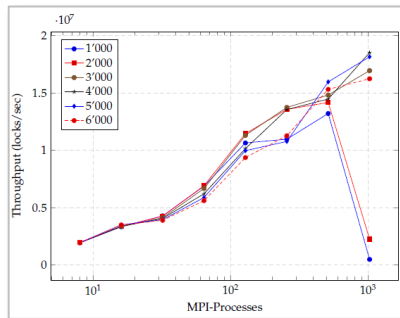
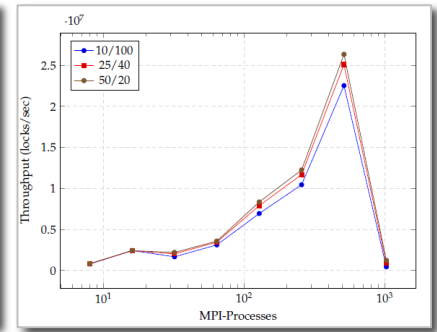
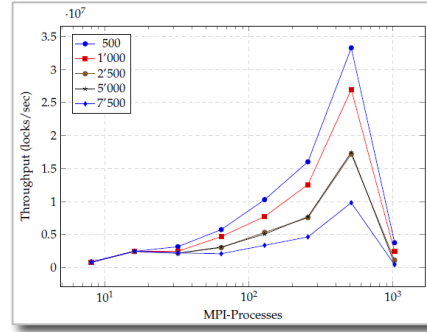
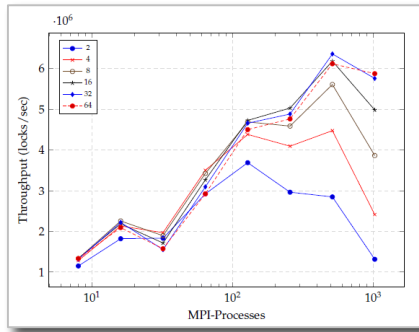
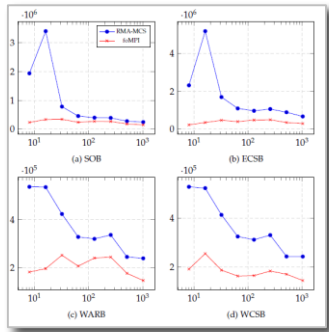
2% of writers



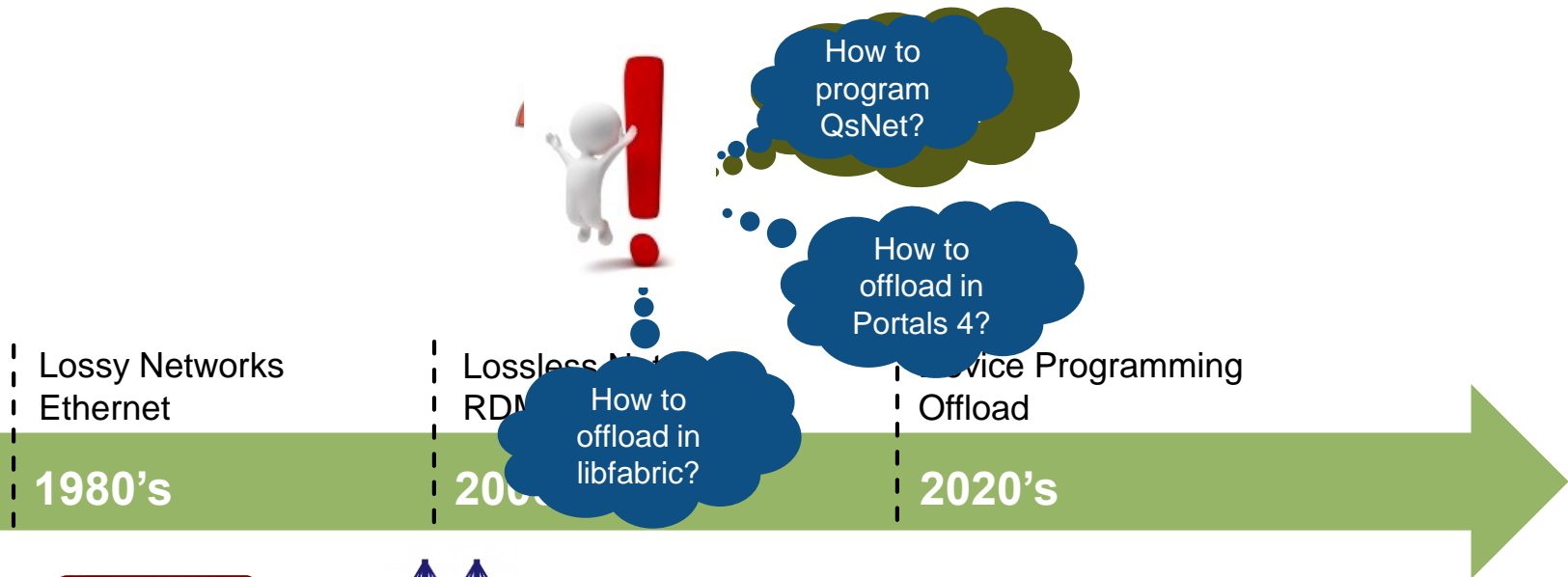
0% of writers

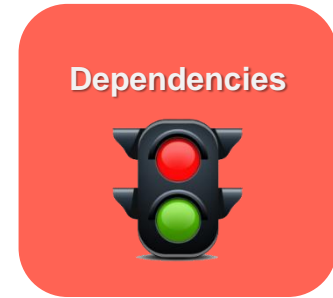


OTHER ANALYSES



But why stop at RDMA -- A brief history





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

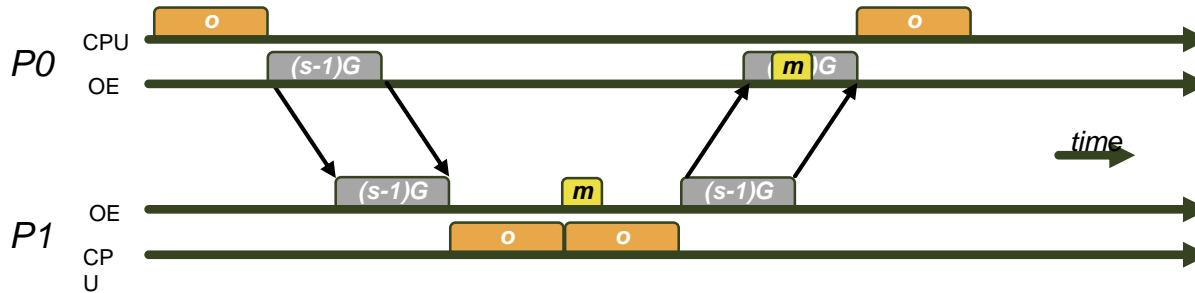
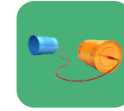
CPU



Offload Engine



Performance Model



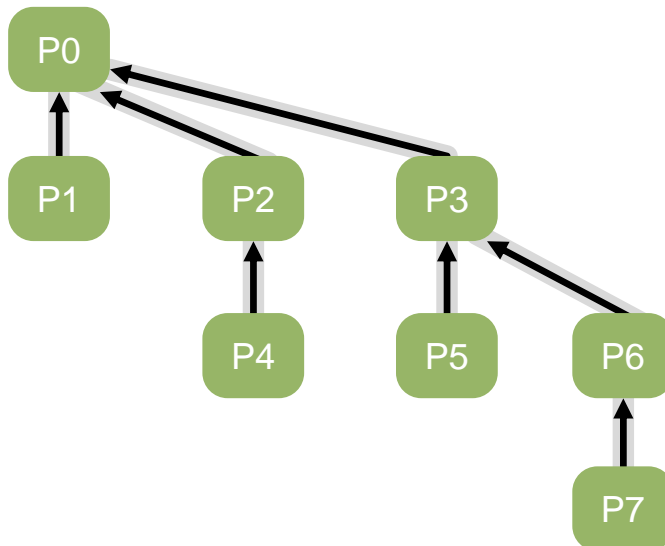
```
P0{
  L0: recv m1 from P1;
  L1: send m2 to P1;
}
```

```
P1{
  L0: recv m1 from P1;
  L1: send m2 to P1;
  L0 -> L1
}
```

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

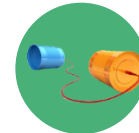
Non-blocking collective: Once initiated the operation may progress independently of any computation or other communication at participating processes

P0

Fully Offloading:

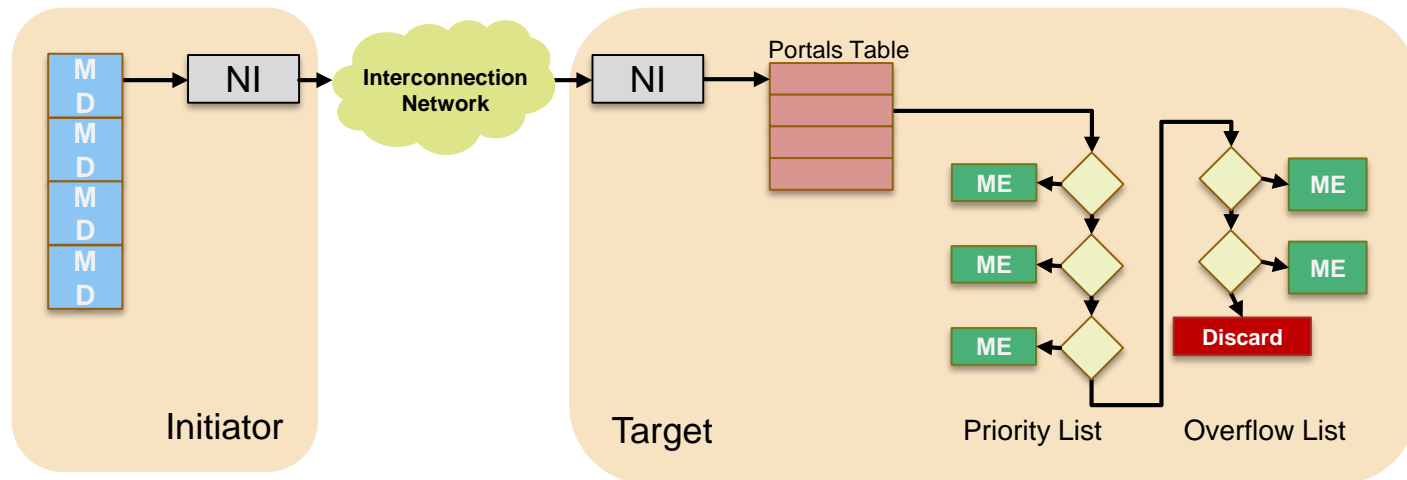
1. *No synchronization* is required in order to start the collective operation
2. Once a collective operation is started, *no further CPU intervention* is required in order to progress or complete it.

P7



A Case Study: Portals 4

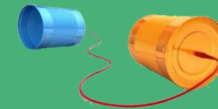
- Based on the one-sided communication model
- Matching/Non-Matching semantics can be adopted



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(...);
```

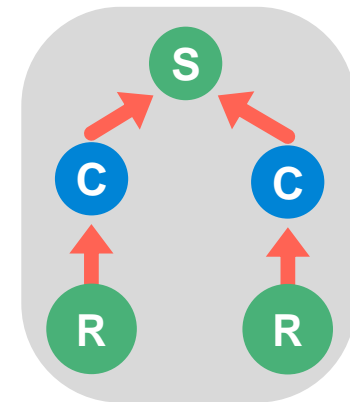
```
ff_op_h r1 = ff_op_create_recv(tmp + blocksize, blocksize, child1, tag);  
ff_op_h r2 = ff_op_create_recv(tmp + 2*blocksize, blocksize, child2, tag);
```

```
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

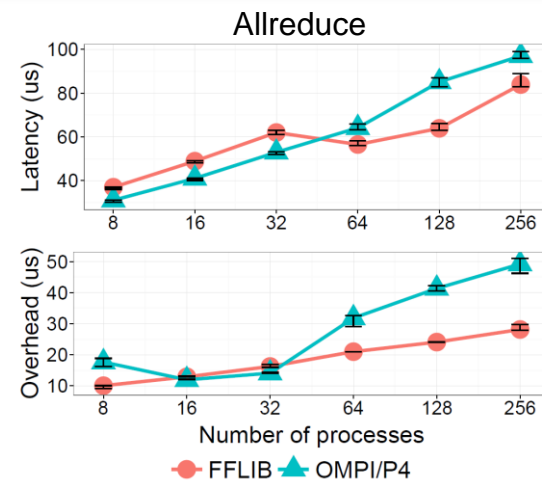
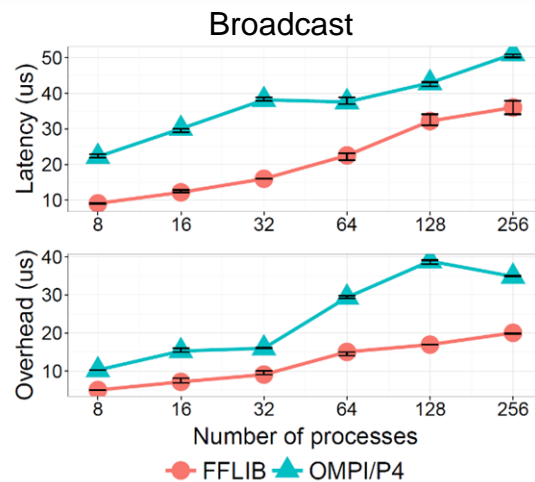
Target machine: Curie

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/

Experimental Results: Latency/Overhead



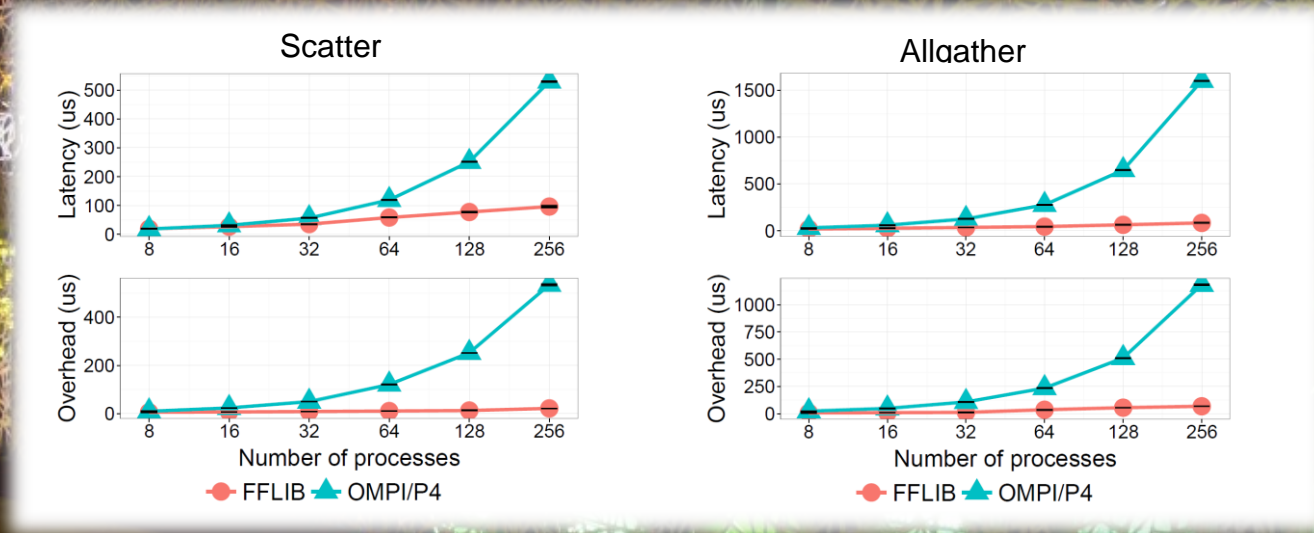
Target machine: Curie

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/

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/

Experimental Results: Micro-Benchmarks

3DFFT

PGMRES

Target machine: Curie

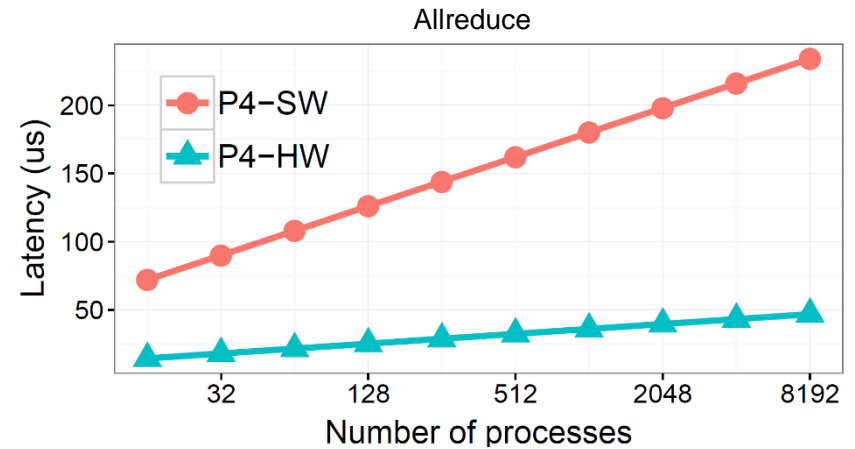
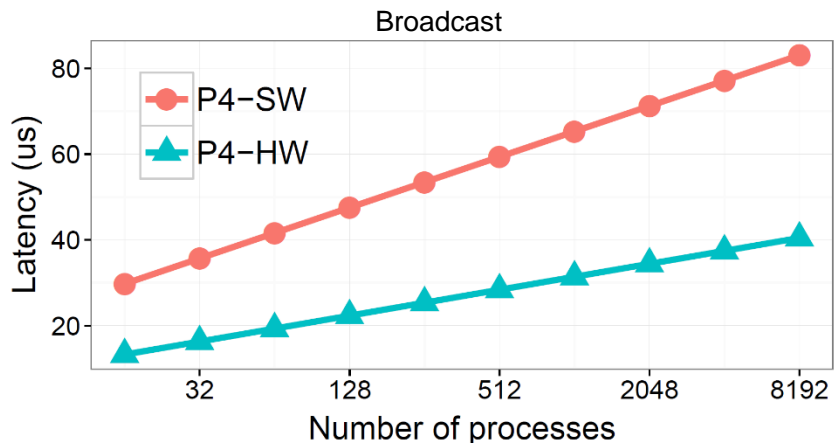
5,040 nodes
2 eight-core Intel Sandy Bridge processors
Full fat-tree Infiniband QDR

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/

Simulations

- **Why?** To study offloaded collectives at large scale
- **How?** Extending the LogGOPSim to simulate Portals 4 functionalities



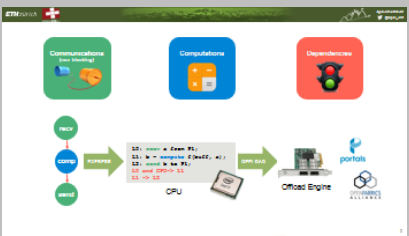
	L	o	g	G	m
P4-SW	$5\mu s$	$6\mu s$	$6\mu s$	$0.4ns$	$0.9ns$
P4-HW	$2.7\mu s$	$1.2\mu s$	$0.5\mu s$	$0.4ns$	$0.3ns$ [4]

[3] T. Hoefler, T. Schneider, A. Lumsdaine. "LogGOPSim - Simulating Large-Scale Applications in the LogGOPS Model", In *Proceedings of the 19th ACM International Symposium on High Performance Distributed Computing (HPDC '10)*. ACM, 2010.

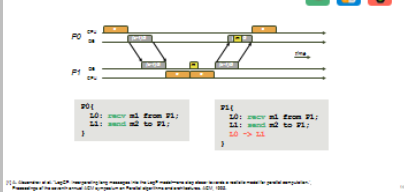
[4] Underwood et al., "Enabling Flexible Collective Communication Offload with Triggered Operations", *IEEE 19th Annual Symposium on High Performance Interconnects (HOTI '11)*. IEEE, 2011.



Abstract Machine Model



Performance Model



Offloading Collectives

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 of participating processes

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.

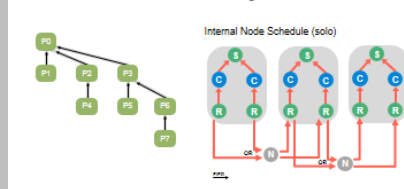


Solo Collectives

Solo Collectives

- Synchronized collectives lead to the synchronization of the participating nodes
- Solo collectives never wait: it is executed as soon as one node (the initiator) starts its own schedule

Solo Collectives: Multi-Version Scheduling



Mapping to Portals 4

A Case Study: Portals 4

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

FFlib

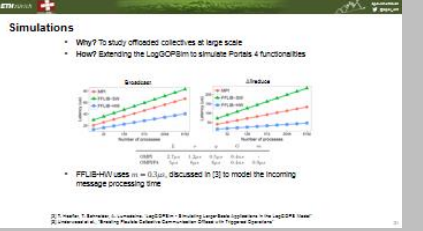
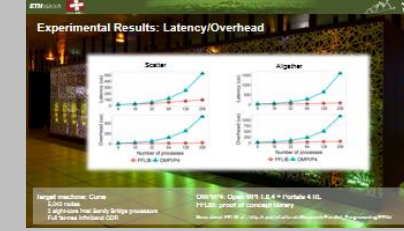
FFlib: An Example

Proof of concept library implemented on top of Portals 4

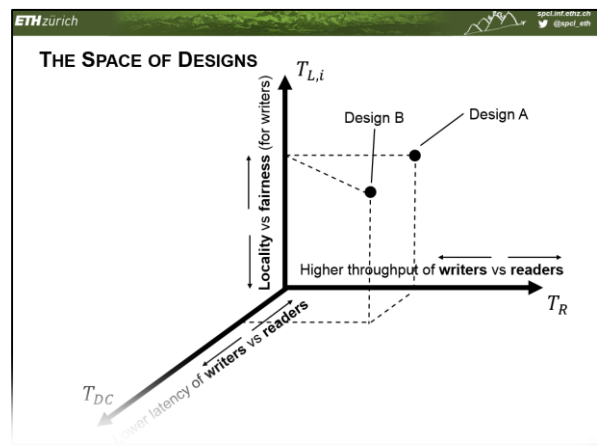
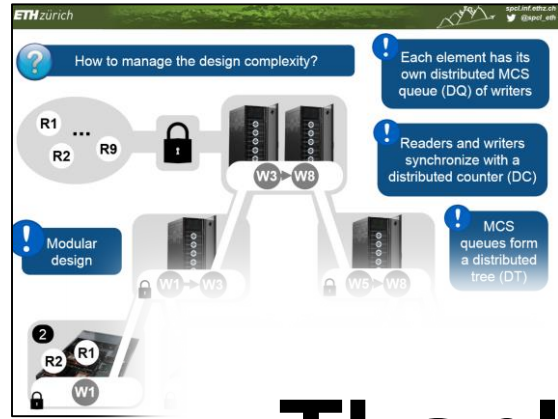
```

fflib_collective_send = ff_collective_send(-);
fflib_collective_send(1);
fflib_collective_send(2);
fflib_collective_send(3);
fflib_collective_send(4);
fflib_collective_send(5);
fflib_collective_send(6);
fflib_collective_send(7);
fflib_collective_send(8);
fflib_collective_send(9);
fflib_collective_send(10);
fflib_collective_send(11);
fflib_collective_send(12);
fflib_collective_send(13);
fflib_collective_send(14);
fflib_collective_send(15);
fflib_collective_send(16);
fflib_collective_send(17);
fflib_collective_send(18);
fflib_collective_send(19);
fflib_collective_send(20);
fflib_collective_send(21);
fflib_collective_send(22);
fflib_collective_send(23);
fflib_collective_send(24);
fflib_collective_send(25);
fflib_collective_send(26);
fflib_collective_send(27);
fflib_collective_send(28);
fflib_collective_send(29);
fflib_collective_send(30);
fflib_collective_send(31);
fflib_collective_send(32);
fflib_collective_send(33);
fflib_collective_send(34);
fflib_collective_send(35);
fflib_collective_send(36);
fflib_collective_send(37);
fflib_collective_send(38);
fflib_collective_send(39);
fflib_collective_send(40);
fflib_collective_send(41);
fflib_collective_send(42);
fflib_collective_send(43);
fflib_collective_send(44);
fflib_collective_send(45);
fflib_collective_send(46);
fflib_collective_send(47);
fflib_collective_send(48);
fflib_collective_send(49);
fflib_collective_send(50);
fflib_collective_send(51);
fflib_collective_send(52);
fflib_collective_send(53);
fflib_collective_send(54);
fflib_collective_send(55);
fflib_collective_send(56);
fflib_collective_send(57);
fflib_collective_send(58);
fflib_collective_send(59);
fflib_collective_send(60);
fflib_collective_send(61);
fflib_collective_send(62);
fflib_collective_send(63);
fflib_collective_send(64);
fflib_collective_send(65);
fflib_collective_send(66);
fflib_collective_send(67);
fflib_collective_send(68);
fflib_collective_send(69);
fflib_collective_send(70);
fflib_collective_send(71);
fflib_collective_send(72);
fflib_collective_send(73);
fflib_collective_send(74);
fflib_collective_send(75);
fflib_collective_send(76);
fflib_collective_send(77);
fflib_collective_send(78);
fflib_collective_send(79);
fflib_collective_send(80);
fflib_collective_send(81);
fflib_collective_send(82);
fflib_collective_send(83);
fflib_collective_send(84);
fflib_collective_send(85);
fflib_collective_send(86);
fflib_collective_send(87);
fflib_collective_send(88);
fflib_collective_send(89);
fflib_collective_send(90);
fflib_collective_send(91);
fflib_collective_send(92);
fflib_collective_send(93);
fflib_collective_send(94);
fflib_collective_send(95);
fflib_collective_send(96);
fflib_collective_send(97);
fflib_collective_send(98);
fflib_collective_send(99);
fflib_collective_send(100);
    
```

Result S



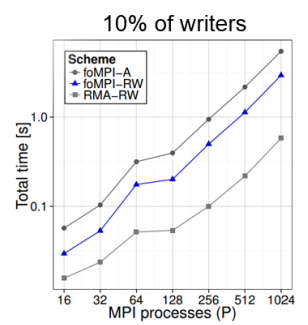
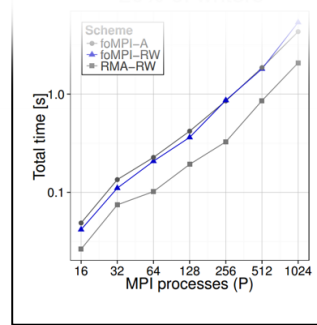
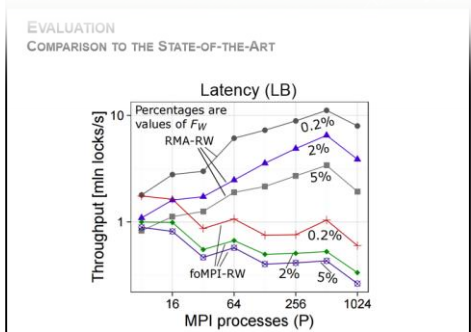
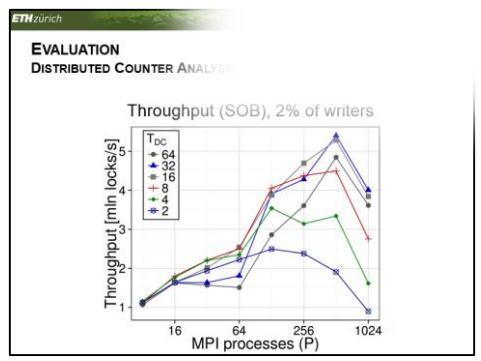
CONCLUSIONS



Thank you for your attention

Modular design
correct

able
ges



Improves latency and throughput over state-of-the-art

Enables high-performance distributed hashtable