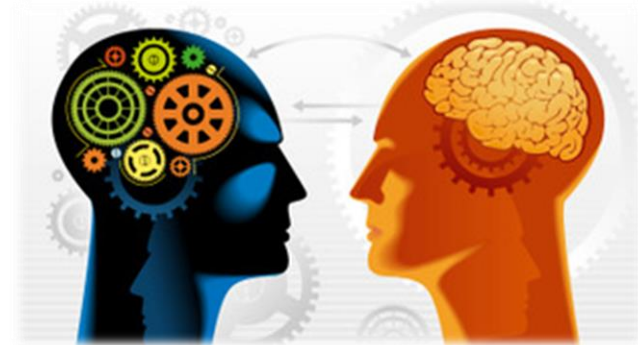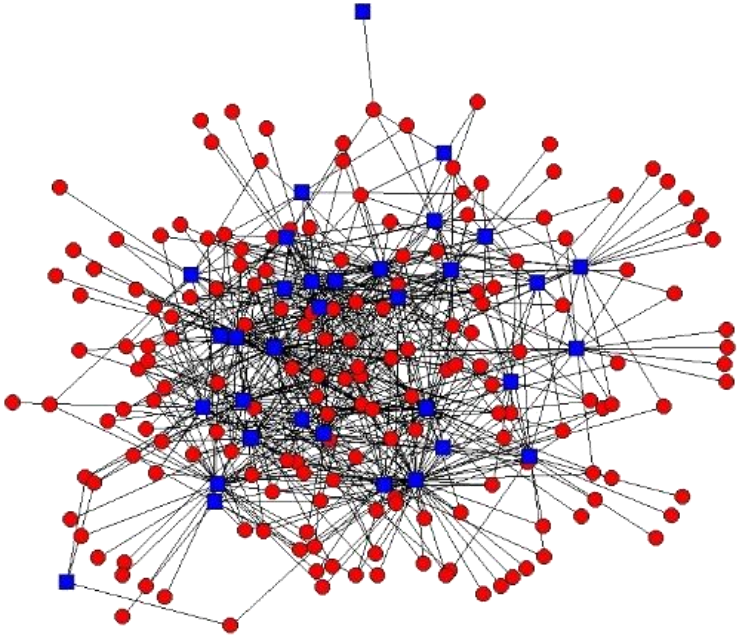**ETH** *zürich*

# High-Performance Distributed RMA Locks

**PATRICK SCHMID, MACIEJ BESTA, TORSTEN HOEFLER**
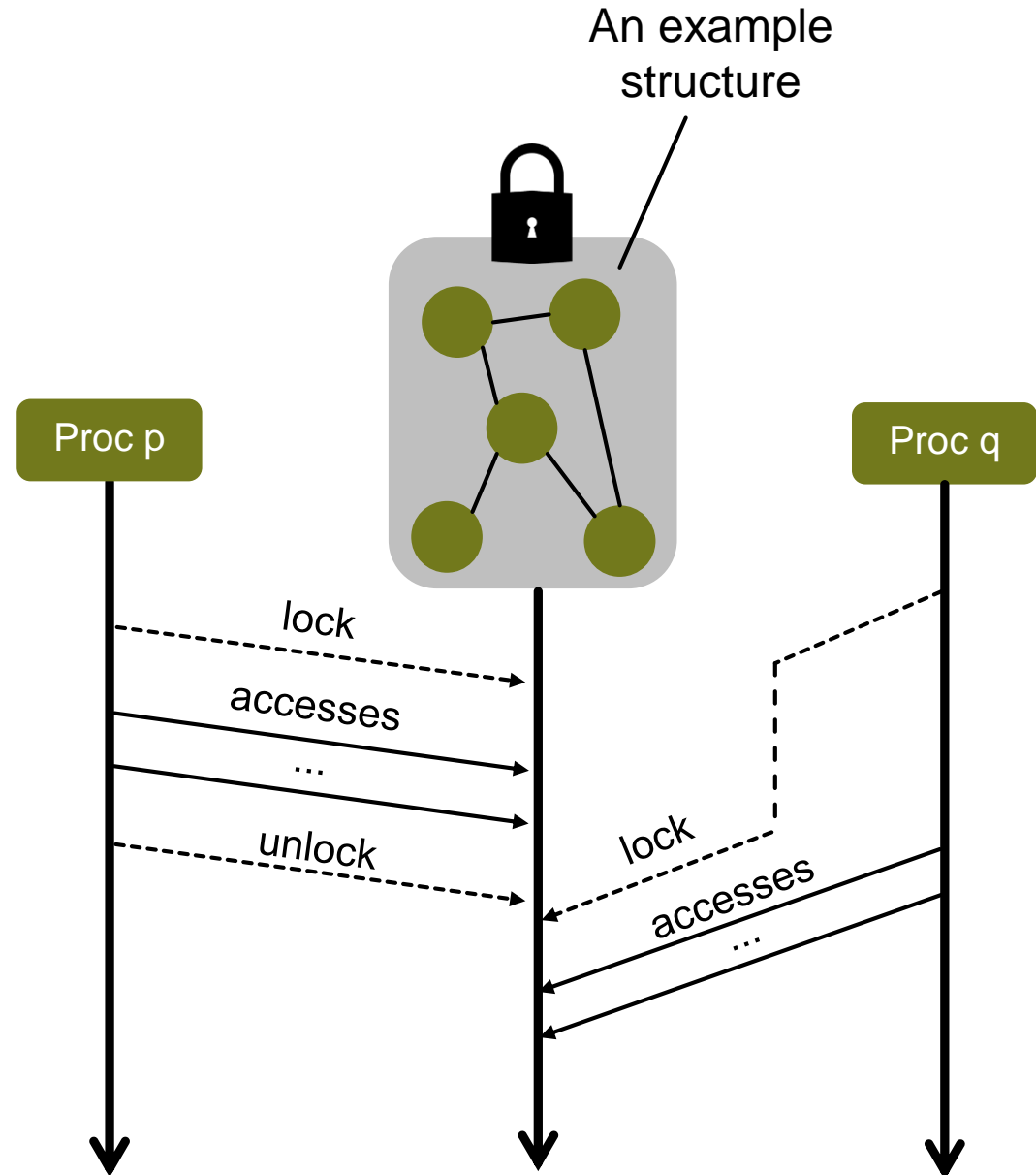
# NEED FOR EFFICIENT LARGE-SCALE SYNCHRONIZATION

# LOCKS
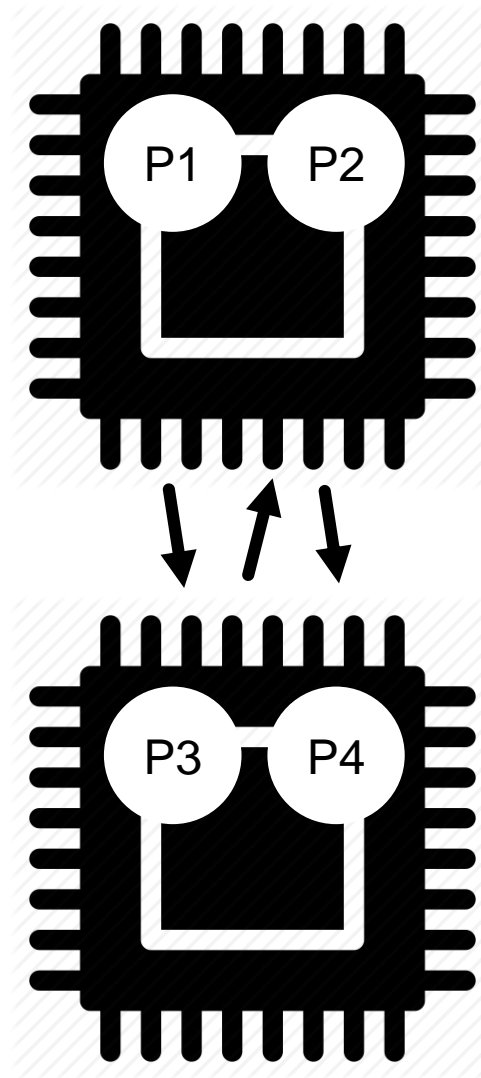
An example structure



Proc p

Proc q

Inuitive semantics

*lock*

accesses

...

*unlock*

*lock*

accesses

...

Various performance penalties

# LOCKS: CHALLENGES

# LOCKS: CHALLENGES

> ! We need intra- and inter-node topology-awareness

> ! We need to cover arbitrary topologies

We need to distinguish between readers and writers

We need flexible performance for both types of processes

Reader

Reader

Reader

Writer

[1] V. Venkataramani et al. Tao: How facebook serves the social graph. SIGMOD'12.

What will we use in the design?
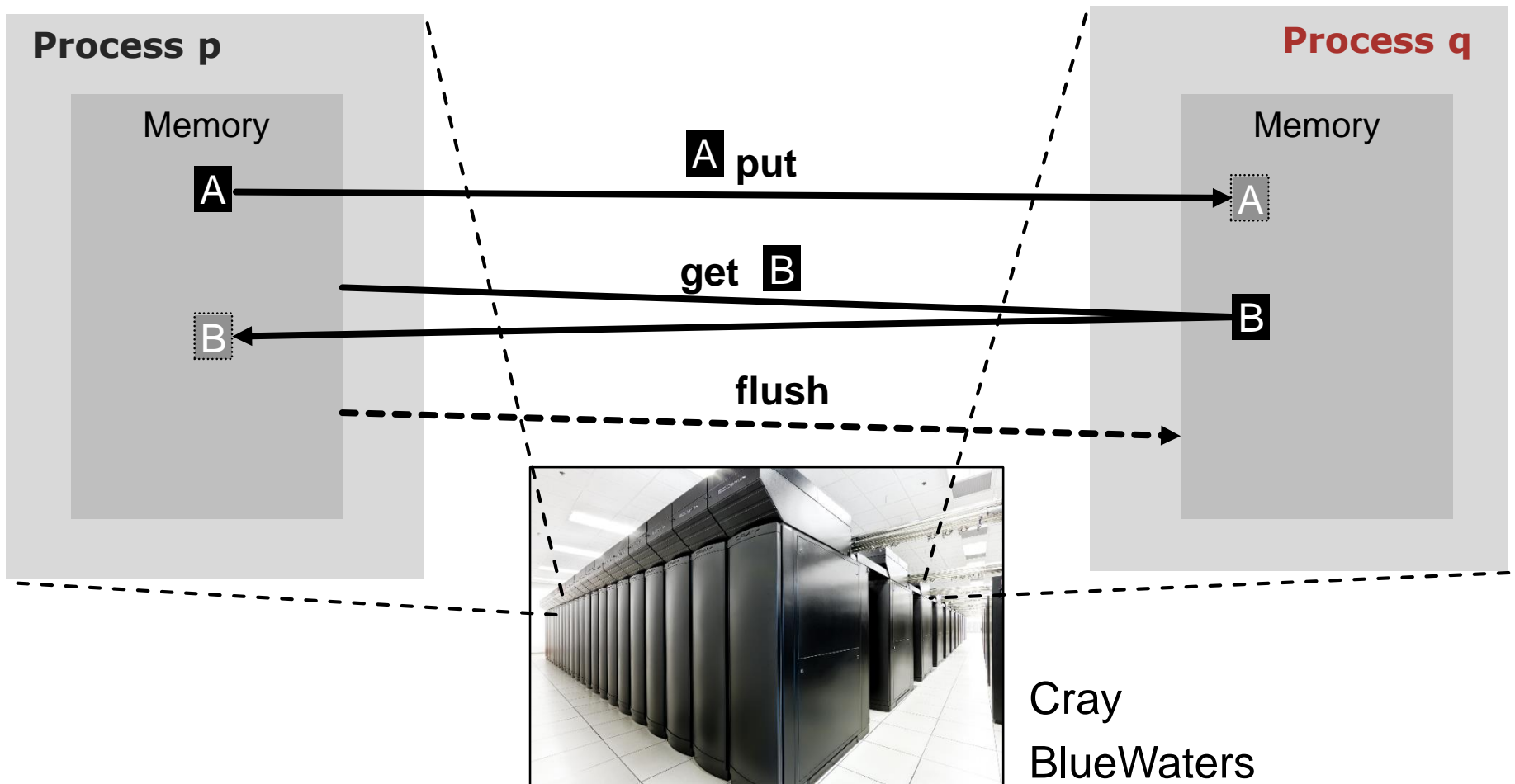
# WHAT WE WILL USE
## MCS Locks

How to manage the design complexity?

How to ensure tunable performance?

What mechanism to use for efficient implementation?

# REMOTE MEMORY ACCESS (RMA) PROGRAMMING



Cray BlueWaters

# 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?

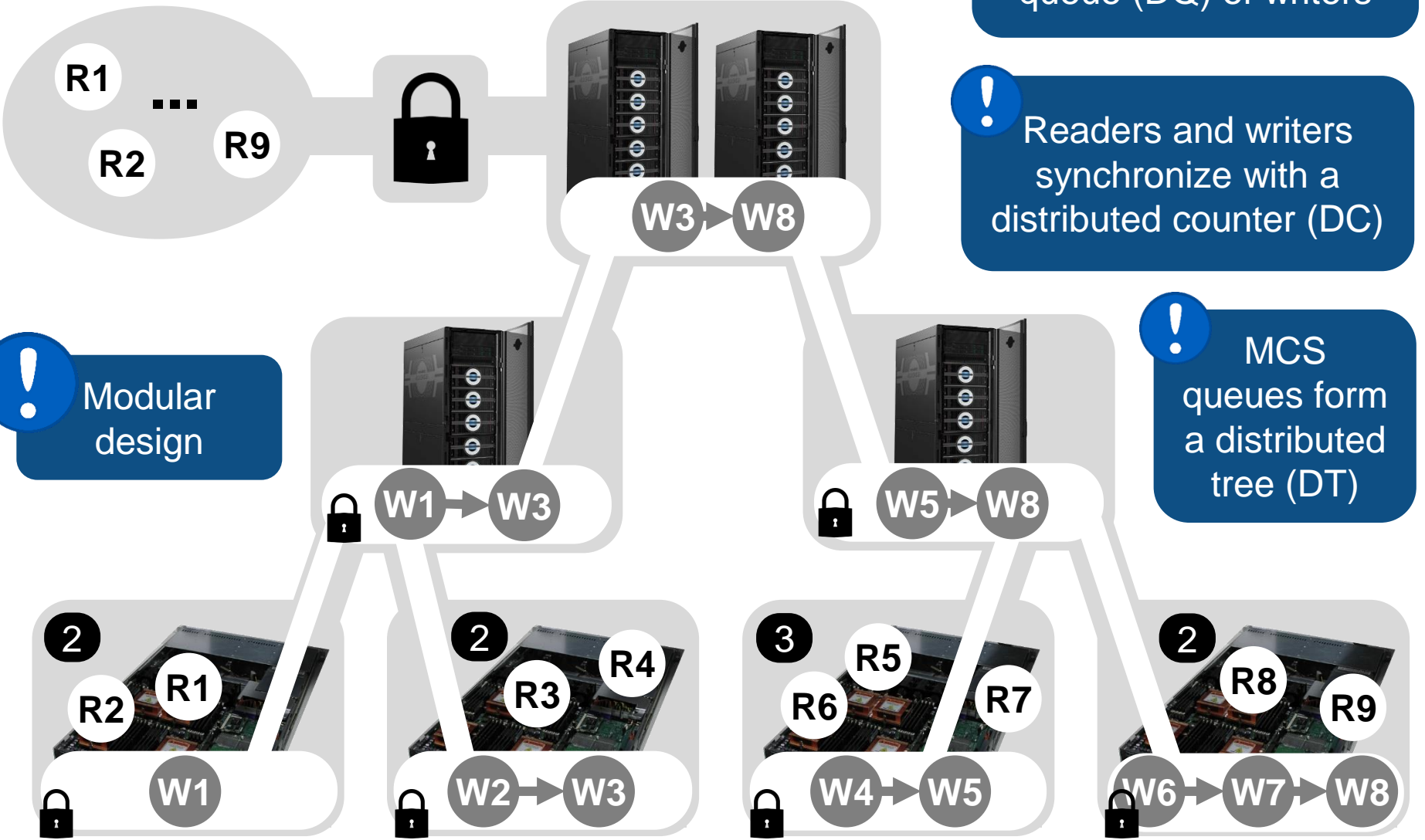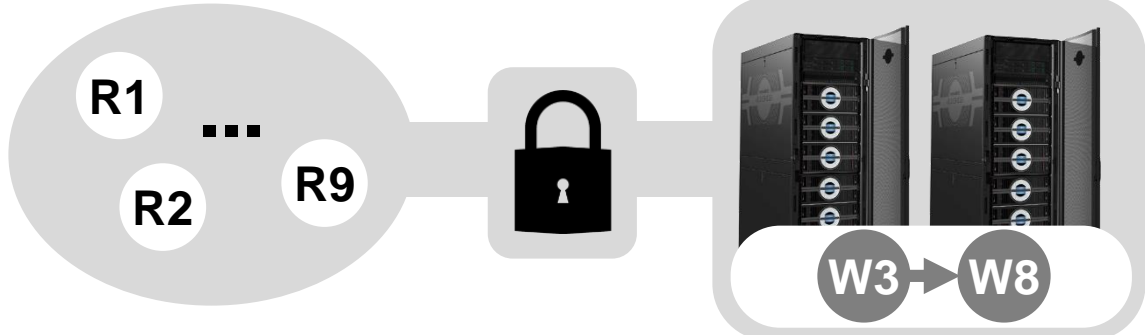What mechanism to use for efficient implementation?
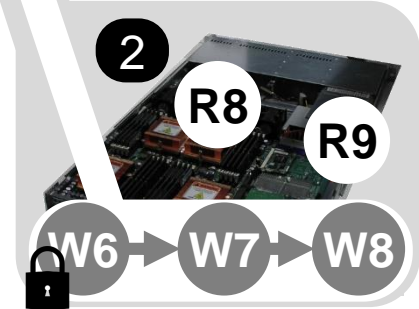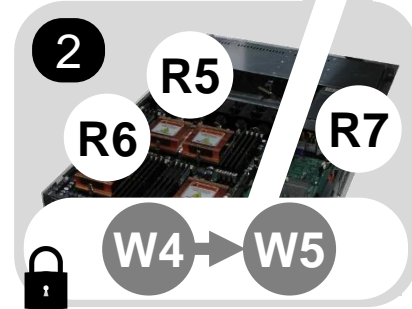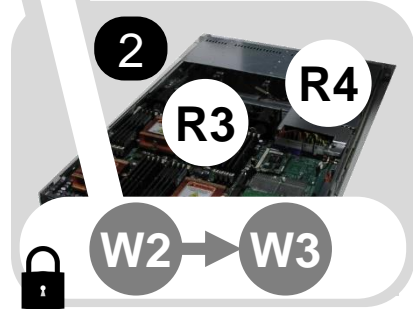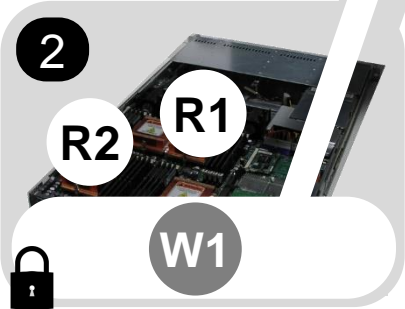
How to ensure tunable performance?

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

R1 ... R9 R2

W3 → W8

W1 → W3

W5 → W8

2  R2 R1

W1

2  R3 R4

W2 → W3

2  R5 R6 R7

W4 → W5

2  R8 R9

W6 → W7 → W8

# DISTRIBUTED MCS QUEUES (DQs)
## Throughput vs Fairness

$T_{L,3}$

$T_{L,2}$

$T_{L,2}$

$T_{L,1}$

$T_{L,1}$

$T_{L,1}$

$T_{L,1}$

**W3 → W8**

**W1 → W3**

**W5 → W8**

**W1**

**W2 → W3**

**W4 → W5**

**W6 → W7 → W8**

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

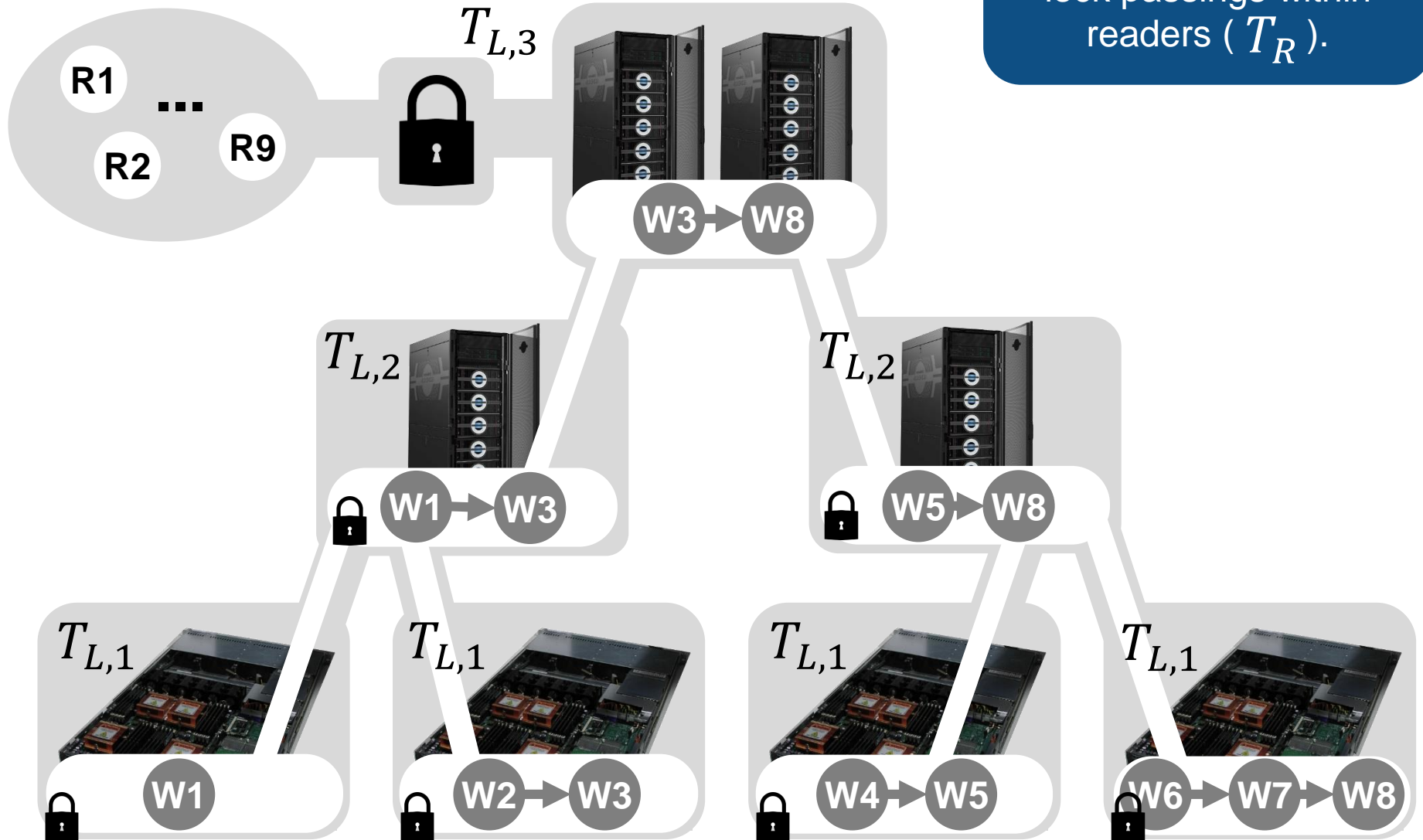Each DQ: The maximum number of lock passings within a DQ at level i, before it is passed to another DQ at i.

$T_{L,i}$

# DISTRIBUTED TREE OF QUEUES (DT)
## Throughput of readers vs writers

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

$T_{L,3}$

R1
...
R2    R9

W3 → W8

$T_{L,2}$

$T_{L,2}$

W1 → W3

W5 → W8

$T_{L,1}$

$T_{L,1}$

$T_{L,1}$

$T_{L,1}$

W1

W2 → W3

W4 → W5

W6 → W7 → W8

# DISTRIBUTED COUNTER (DC)
## Latency of readers vs writers

DC: every *k*th compute node hosts a partial counter, all of which constitute the DC.
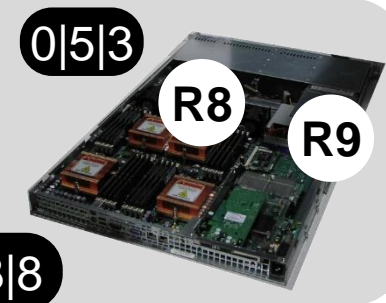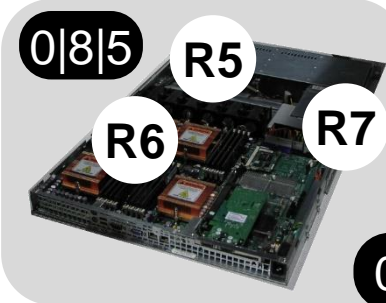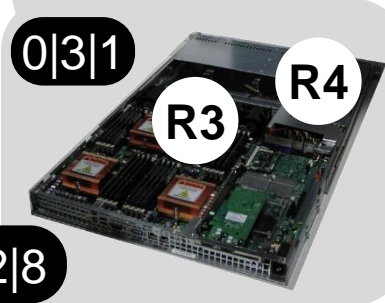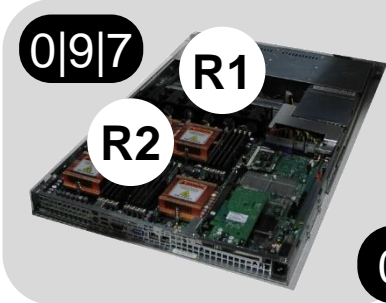
$$k = T_{DC}$$

A writer holds the lock

b|x|y

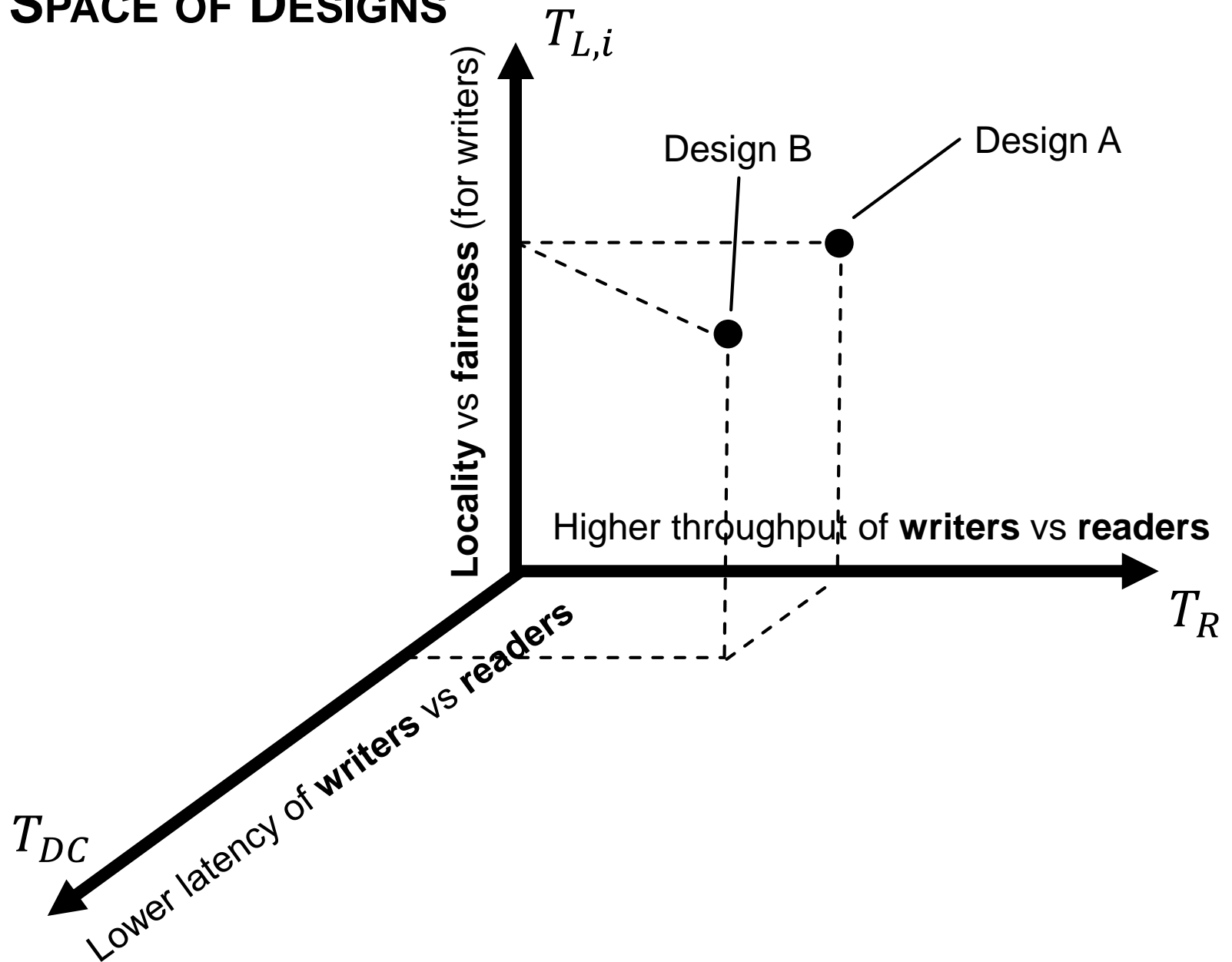Readers that arrived at the CS

Readers that left the CS

$$T_{DC} = 1$$
$$T_{DC} = 2$$

0|9|7  R1
R2
0|12|8

0|3|1  R4
R3

0|8|5  R5
R6  R7
0|13|8

0|5|3  R8
R9

# THE SPACE OF DESIGNS



$T_{L,i}$

Design B

Design A

**Locality** vs **fairness** (for writers)

Higher throughput of **writers** vs **readers**

$T_R$

$T_{DC}$

Lower latency of **writers** vs **readers**
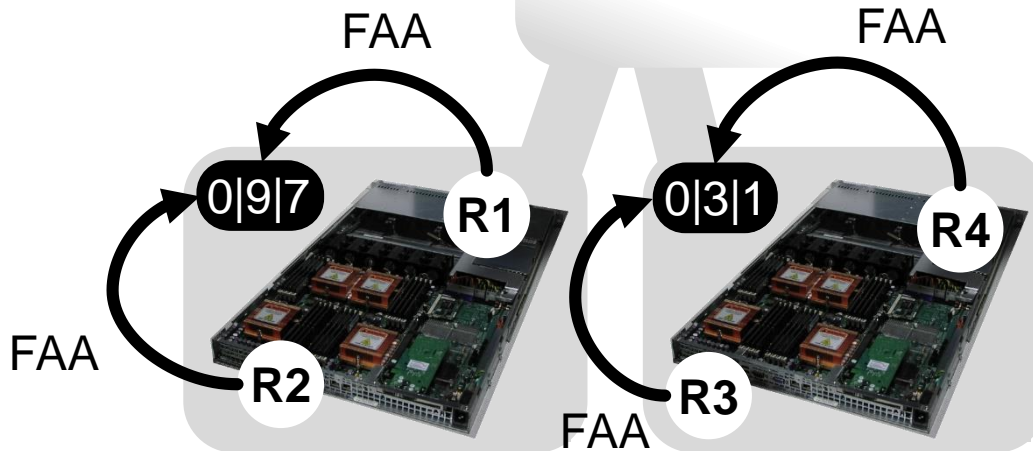
# LOCK ACQUIRE BY READERS

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

A writer holds the lock — b|x|y
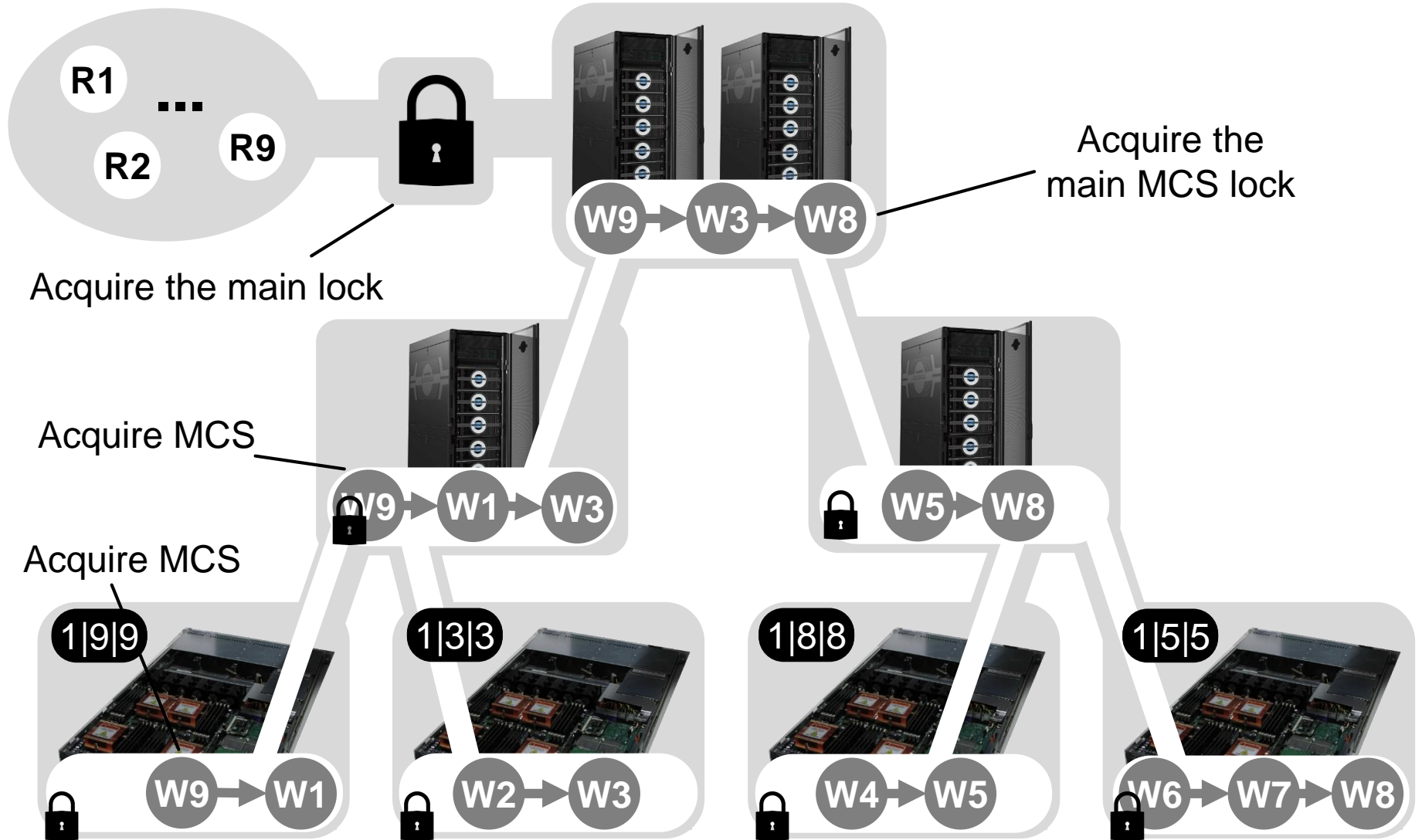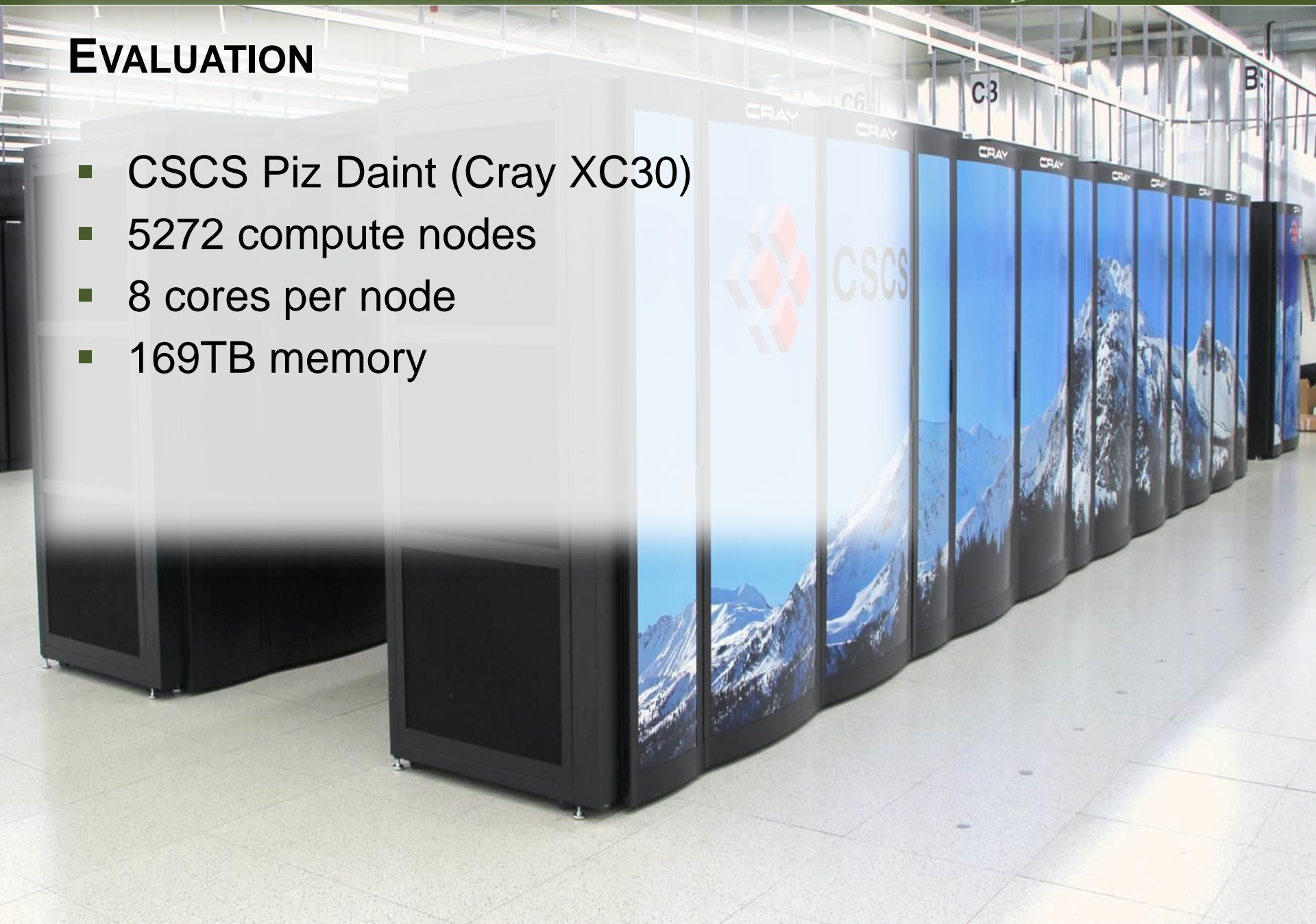
Readers that arrived at the CS

Readers that left the CS

FAA        FAA

0|9|7    R1      0|3|1    R4

FAA

R2        R3

FAA

# LOCK ACQUIRE BY WRITERS

Acquire the main MCS lock

Acquire the main lock

Acquire MCS

Acquire MCS

1|9|9

1|3|3

1|8|8

1|5|5

# EVALUATION
## CONSIDERED BENCHMARKS

The **latency** benchmark

**Throughput benchmarks:**

Empty-critical-section

Single-operation

Wait-after-release

Workload-critical-section

DHT

Distributed hashtable evaluation

# EVALUATION
## DISTRIBUTED COUNTER ANALYSIS

0|9|7   0|3|1

0|12|8

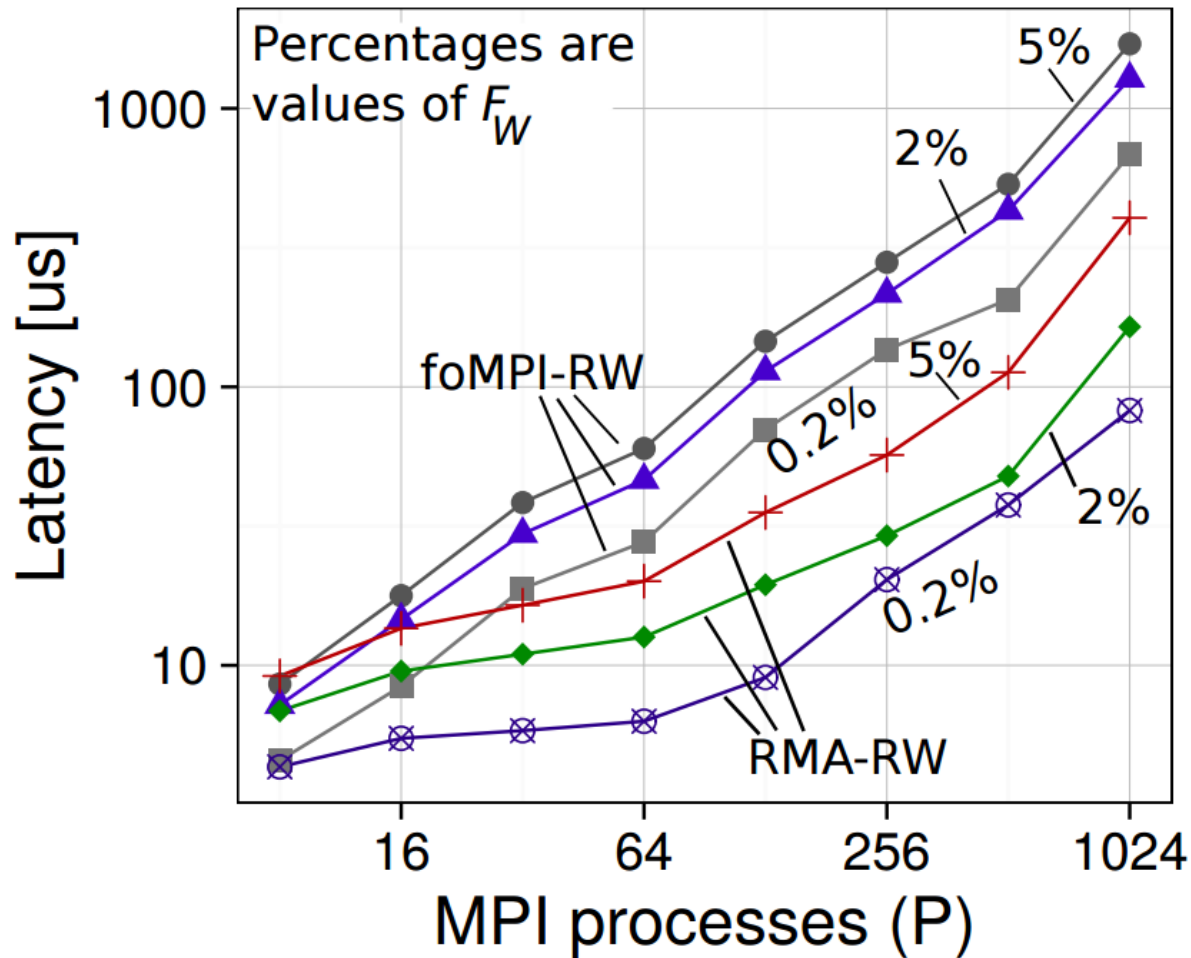Throughput, 2% writers

Single-operation benchmark

# EVALUATION
## READER THRESHOLD ANALYSIS

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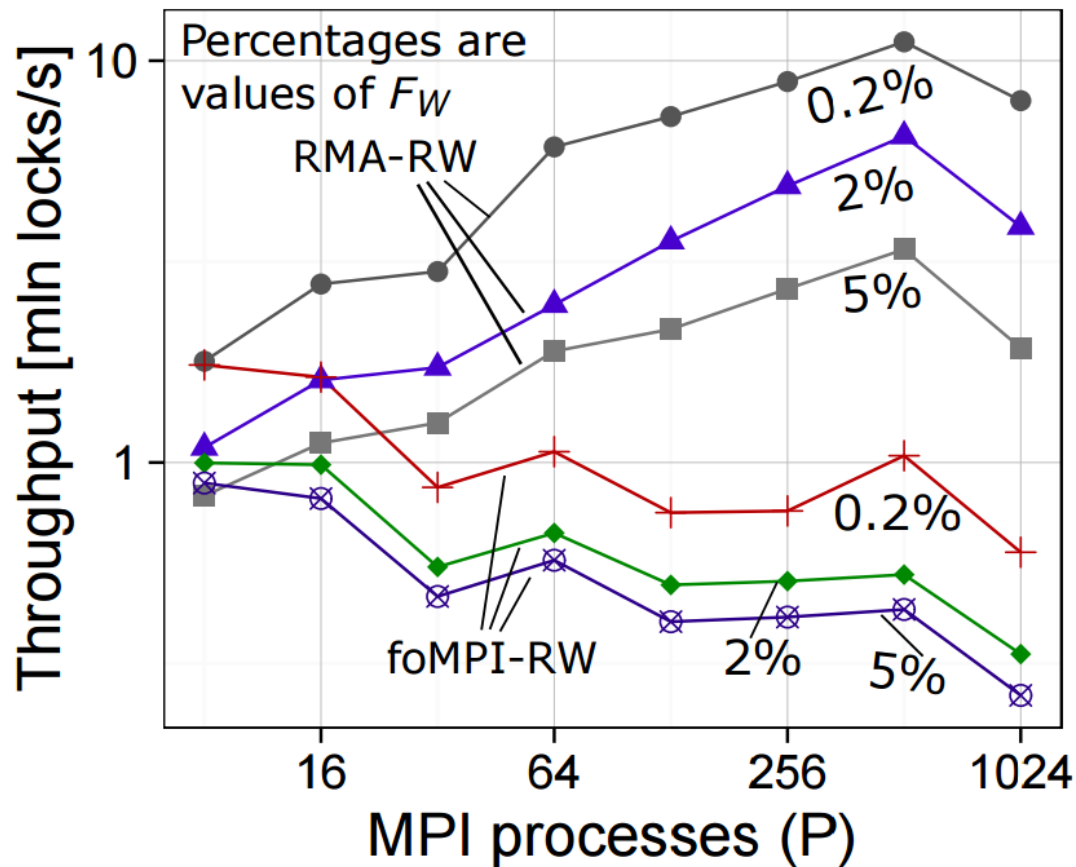
Empty-critical-section benchmark

# EVALUATION
## COMPARISON TO THE STATE-OF-THE-ART



[1] R. Gerstenberger et al. Enabling Highly-scalable Remote Memory Access Programming with MPI-3 One Sided. ACM/IEEE Supercomputing 2013.

# EVALUATION
## COMPARISON TO THE STATE-OF-THE-ART

Throughput, single-operation benchmark



[1] R. Gerstenberger et al. Enabling Highly-scalable Remote Memory Access Programming with MPI-3 One Sided. ACM/IEEE Supercomputing 2013.

# EVALUATION
## DISTRIBUTED HASHTABLE



20% writers

10% writers

[1] R. Gerstenberger et al. Enabling Highly-scalable Remote Memory Access Programming with MPI-3 One Sided.
ACM/IEEE Supercomputing 2013.

# EVALUATION
## DISTRIBUTED HASHTABLE

## 2% of writers

## 0% of writers



[1] R. Gerstenberger et al. Enabling Highly-scalable Remote Memory Access Programming with MPI-3 One Sided.
ACM/IEEE Supercomputing 2013.

# OTHER ANALYSES

# Thank you for your attention

# DISTRIBUTED TREE OF QUEUES (DT)
**Throughput of readers vs writers**

DT: The maximum number of consecutive lock passings within writers ($T_W$) and readers ($T_R$).

$$T_W = \prod_{i=1}^{N} T_{L,i}$$

$T_{L,2}$  $T_{L,2}$

$T_{L,1}$  $T_{L,1}$  $T_{L,1}$  $T_{L,1}$

W3 → W8

W1 → W3   W5 → W8

W1   W2 → W3   W4 → W5   W6 → W7 → W8

R1 ... R2 R9

36

# THE SPACE OF DESIGNS

# EVALUATION
## D-MCS VS OTHERS

## Latency (LB)

## Throughput (ECSB)

# EVALUATION
## WRITER THRESHOLD ANALYSIS

Throughput, 25% of writers

Single-operation benchmark

# EVALUATION
## FAIRNESS VS THROUGHPUT ANALYSIS

Throughput, 25% of writers,

Single-operation benchmark
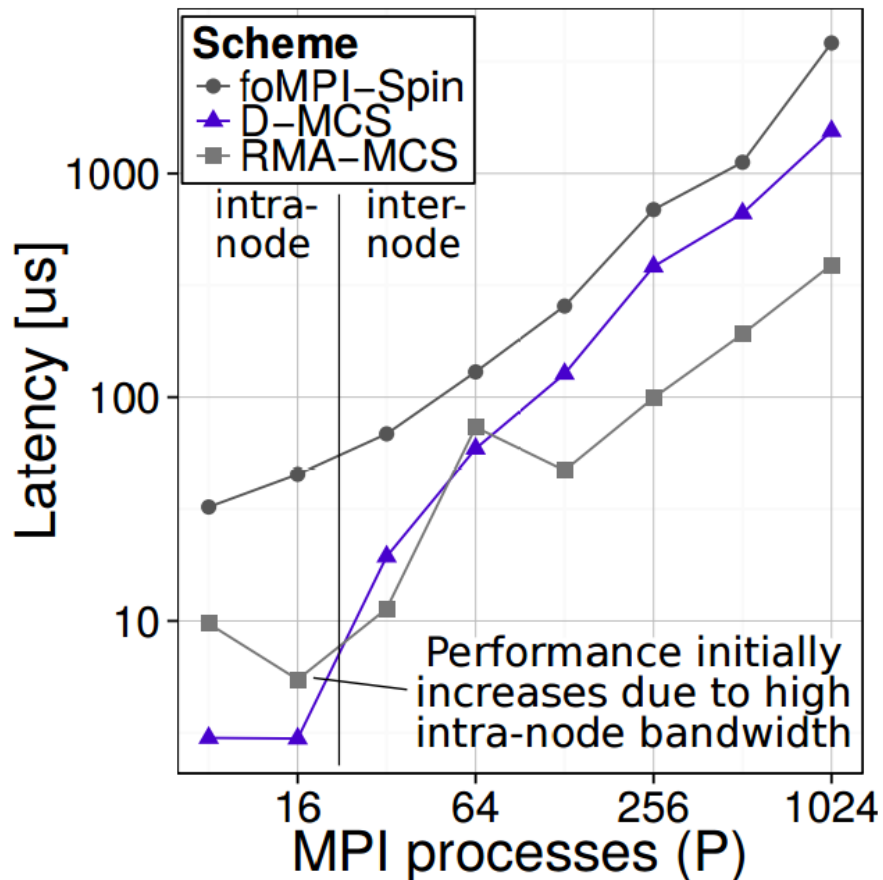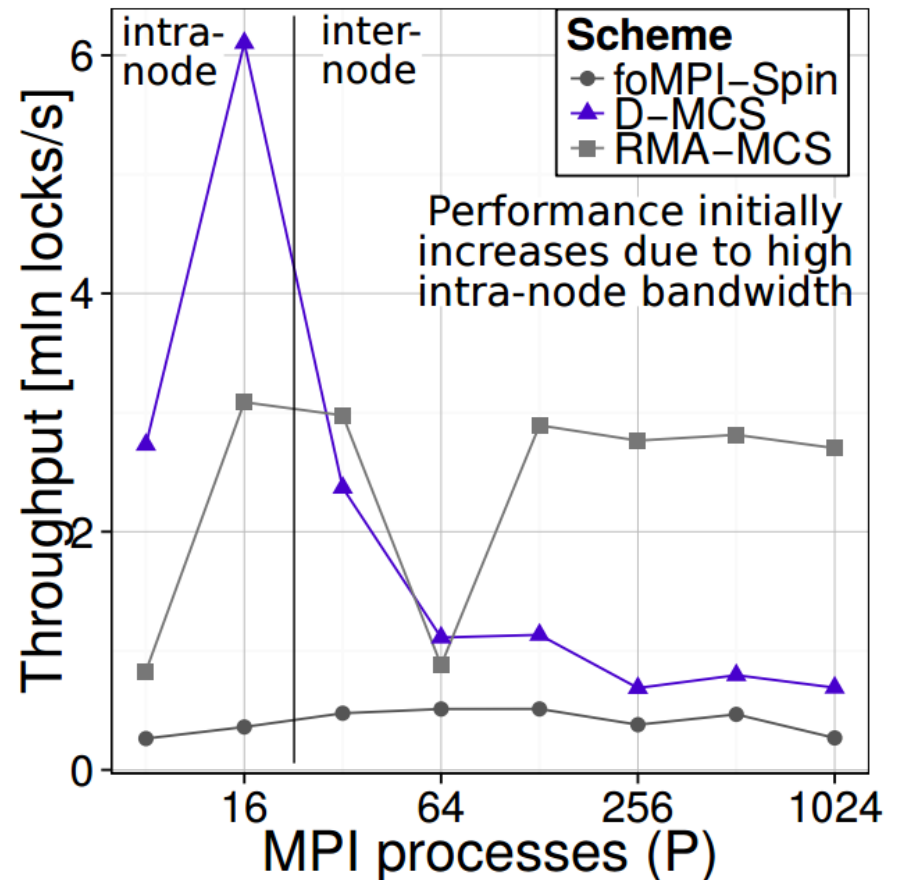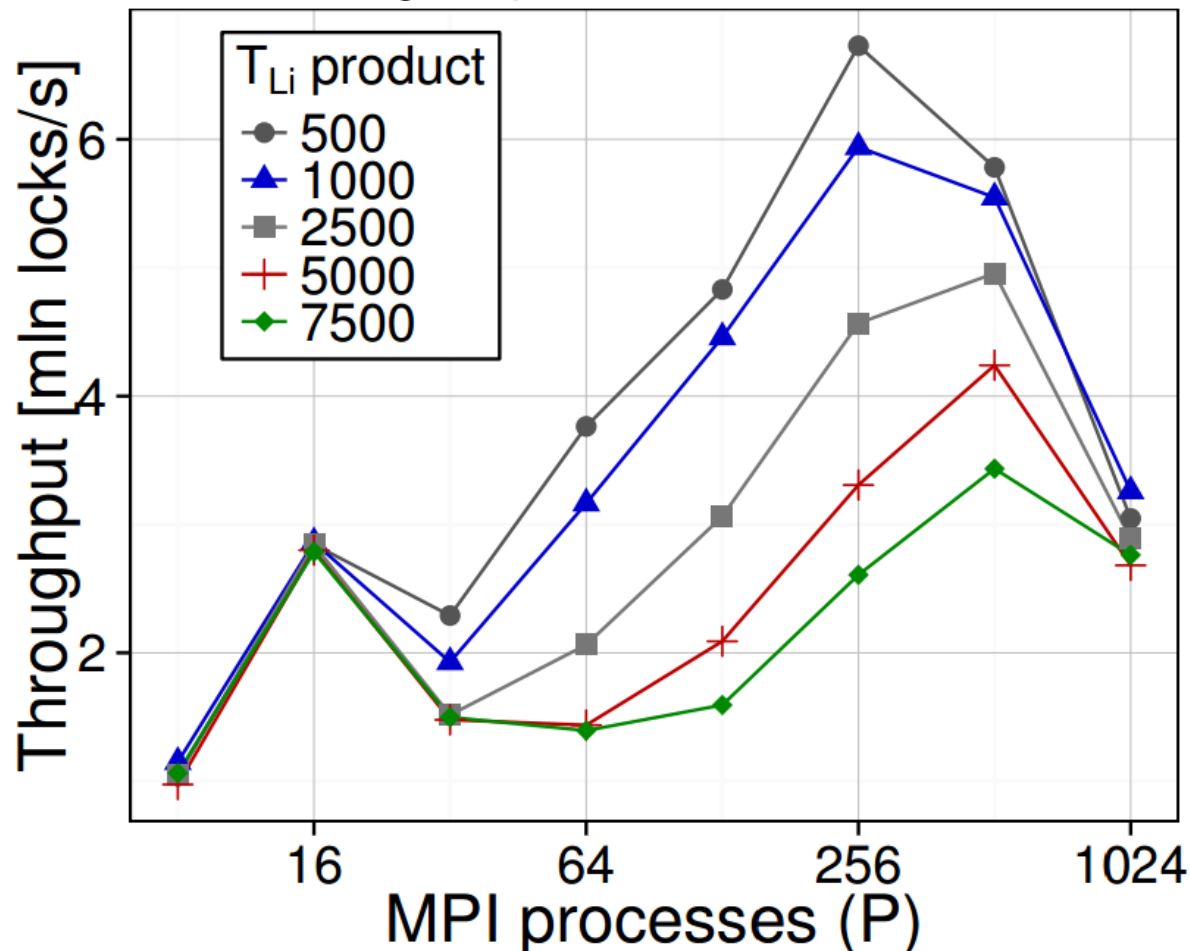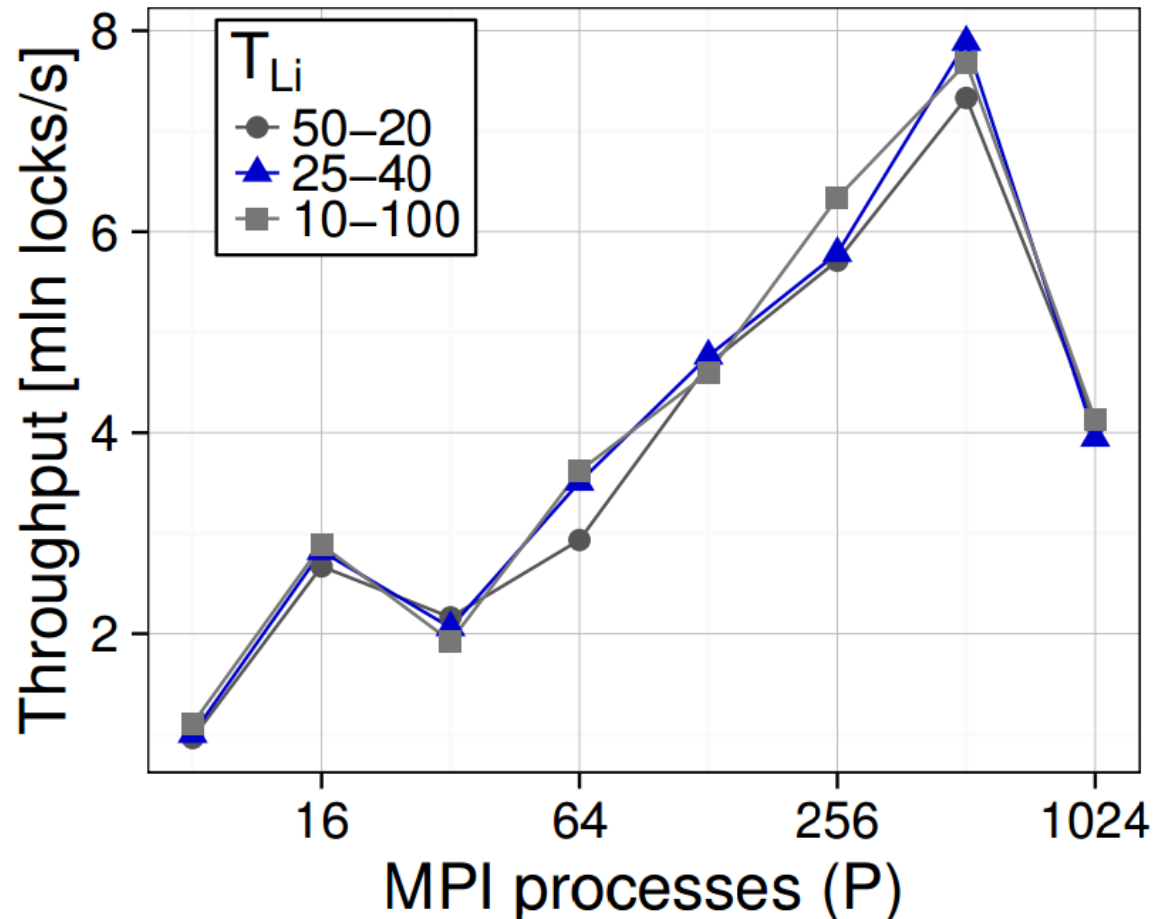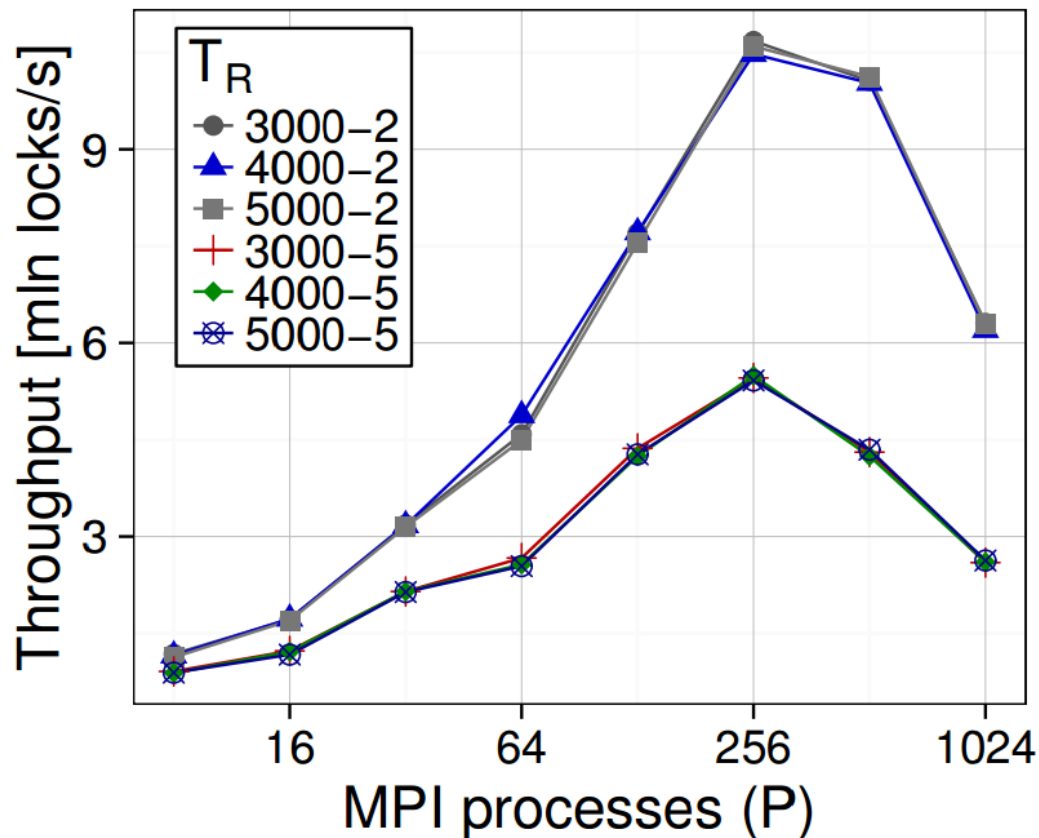
# EVALUATION
## READER THRESHOLD ANALYSIS

Throughput, 2% and 5% writers,
Empty-critical-section benchmark

# FEASIBILITY ANALYSIS

|  | UPC (standard) [44] | Berkeley UPC [1] | SHMEM [4] |
|---|---|---|---|
| Put | UPC_SET | bupc_atomicX_set_RS | shmem_swap |
| Get | UPC_GET | bupc_atomicX_read_RS | shmem_mswap |
| Accumulate | UPC_INC | bupc_atomicX_fetchadd_RS | shmem_fadd |
| FAO (SUM) | UPC_INC, UPC_DEC | bupc_atomicX_fetchadd_RS | shmem_fadd |
| FAO (REPLACE) | UPC_SET | bupc_atomicX_swap_RS | shmem_swap |
| CAS | UPC_CSWAP | bupc_atomicX_cswap_RS | shmem_cswap |

|  | Fortran 2008 [27] | Linux RDMA/IB [33, 43] | iWARP [18, 41] |
|---|---|---|---|
| Put | atomic_define | MskCmpSwap | masked CmpSwap |
| Get | atomic_ref | MskCmpSwap | masked CmpSwap |
| Accumulate | atomic_add | FetchAdd | FetchAdd |
| FAO (SUM) | atomic_add | FetchAdd | FetchAdd |
| FAO (REPLACE) | atomic_define* | MskCmpSwap | masked CmpSwap |
| CAS | atomic_cas | CmpSwap | CmpSwap |