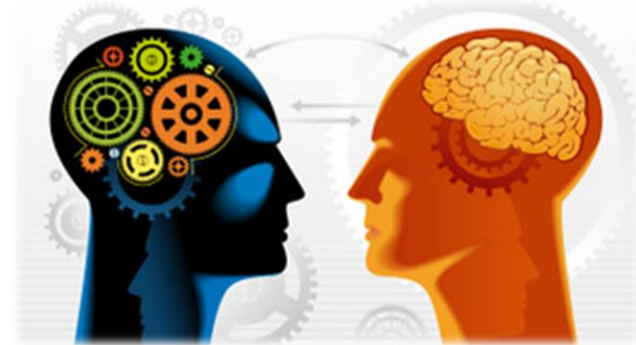**ETH** *zürich*

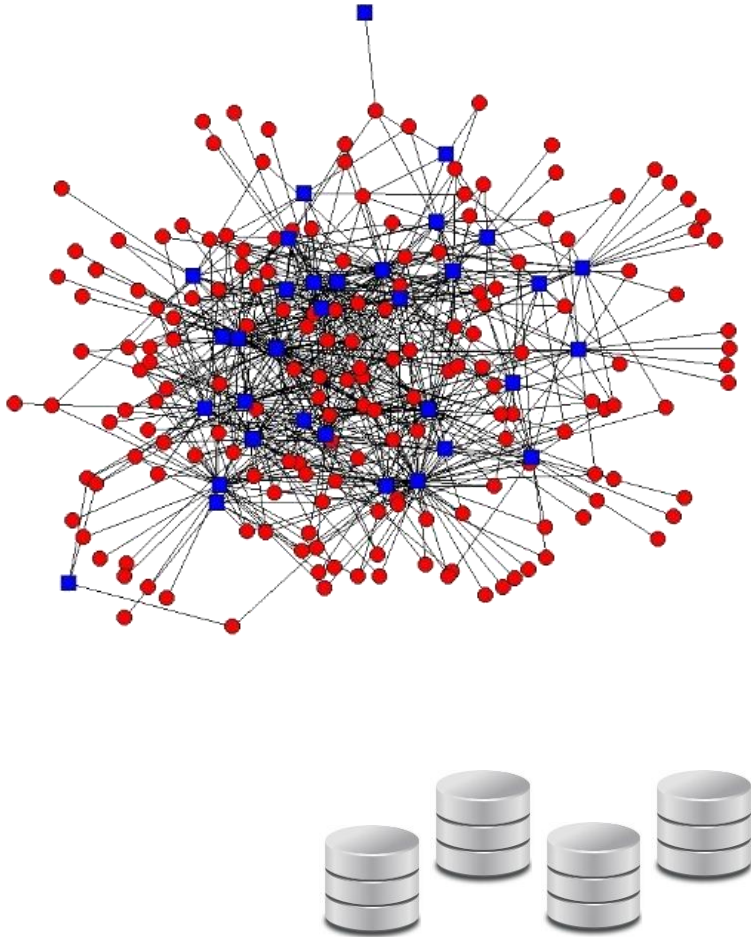# High-Performance Distributed RMA Locks

## PATRICK SCHMID, MACIEJ BESTA, TORSTEN HOEFLER

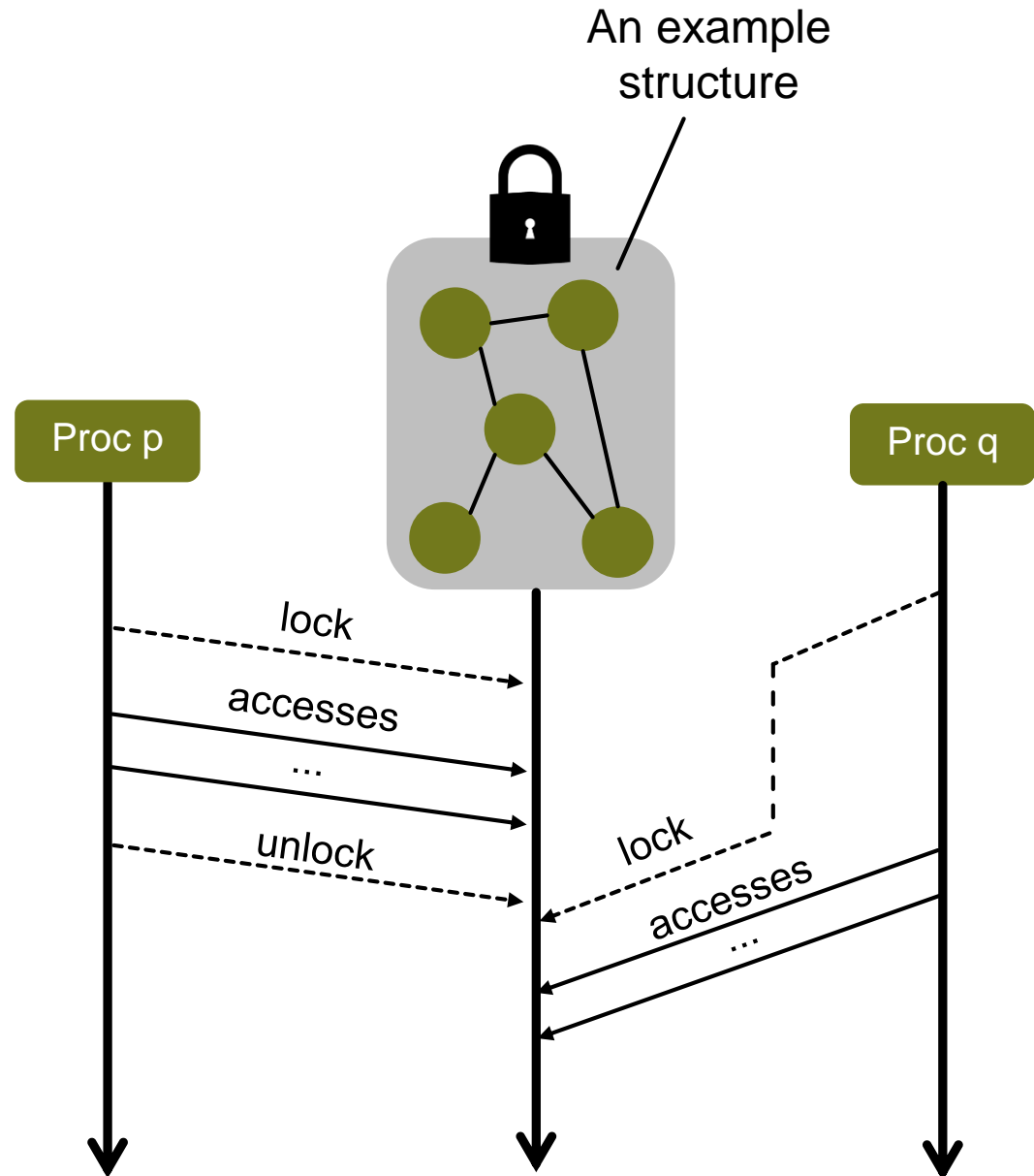# NEED FOR EFFICIENT LARGE-SCALE SYNCHRONIZATION

# LOCKS
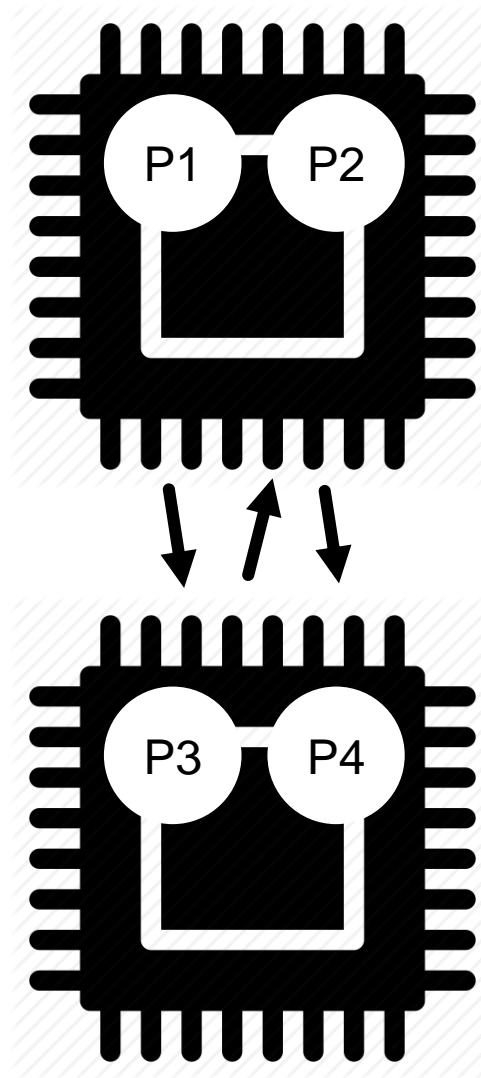


An example structure

✔ Inuitive semantics

✖ Various performance penalties

# LOCKS: CHALLENGES



Calciu et al.: NUMA-aware reader-writer locks, PPoPP'13

# 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
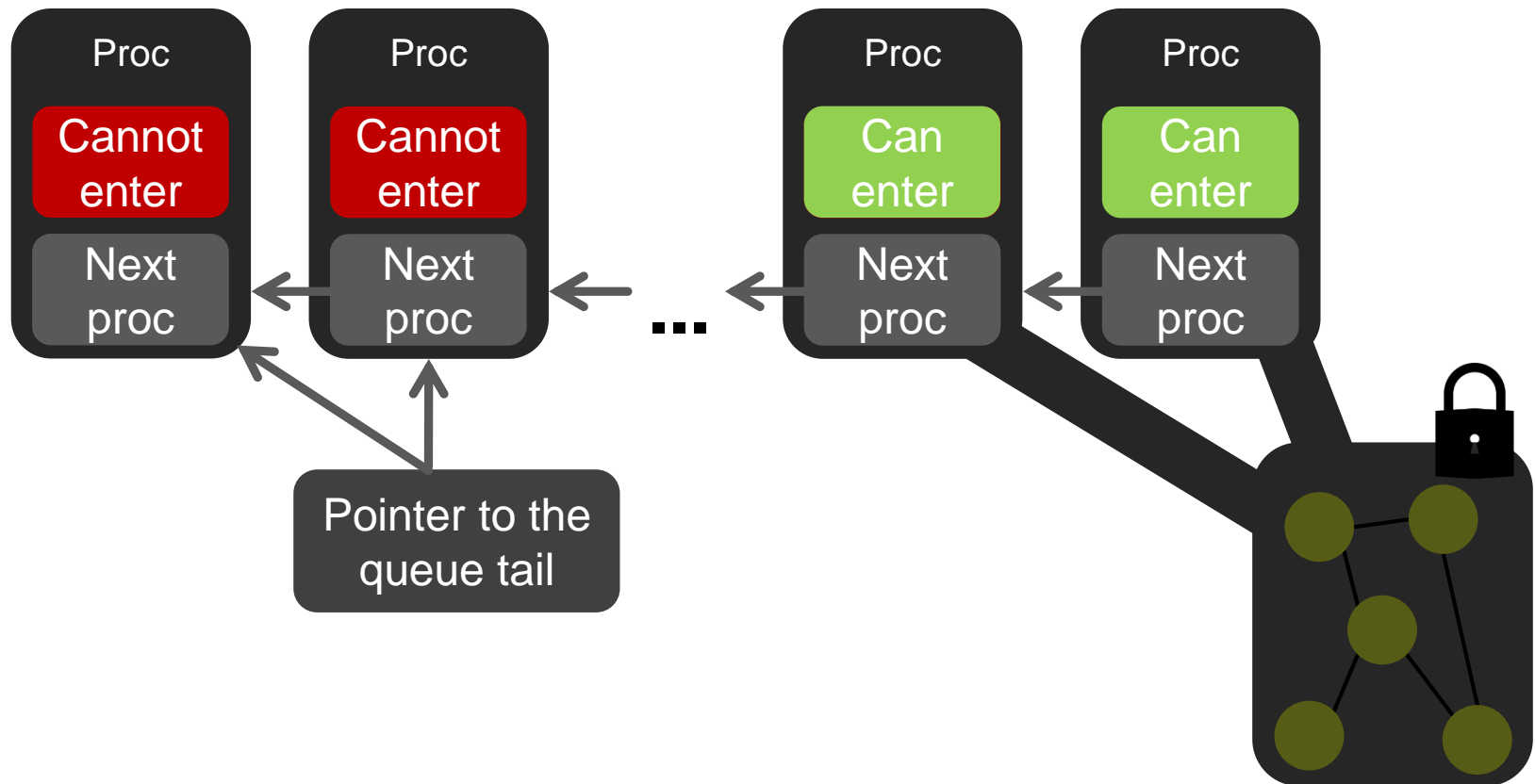
Reader

Reader

Reader

Writer

We need flexible performance for both types of processes

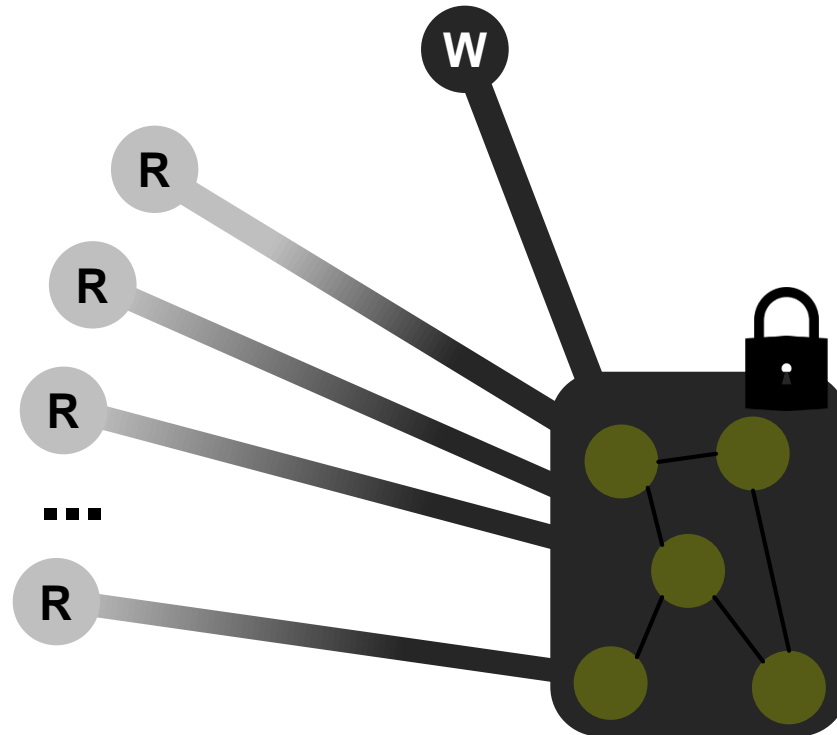[1] V. Venkataramani et al. Tao: How facebook serves the social graph. SIGMOD'12.

What will we use in the design?

# WHAT WE WILL USE
## MCS Locks



Mellor-Crummey and Scott: Algorithms for Scalable Synchronization on Shared-Memory Multiprocessors, ACM TOCS'91

# WHAT WE WILL USE
**Reader-Writer Locks**

# REMOTE MEMORY ACCESS (RMA) PROGRAMMING



Cray BlueWaters

TH, J. Dinan, R. Thakur, B. Barrett, P. Balaji, W. Gropp, K. Underwood: Remote Memory Access Programming in MPI-3, ACM TOPC'15

# REMOTE MEMORY ACCESS PROGRAMMING

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

How to manage the design complexity?

How to ensure tunable performance?

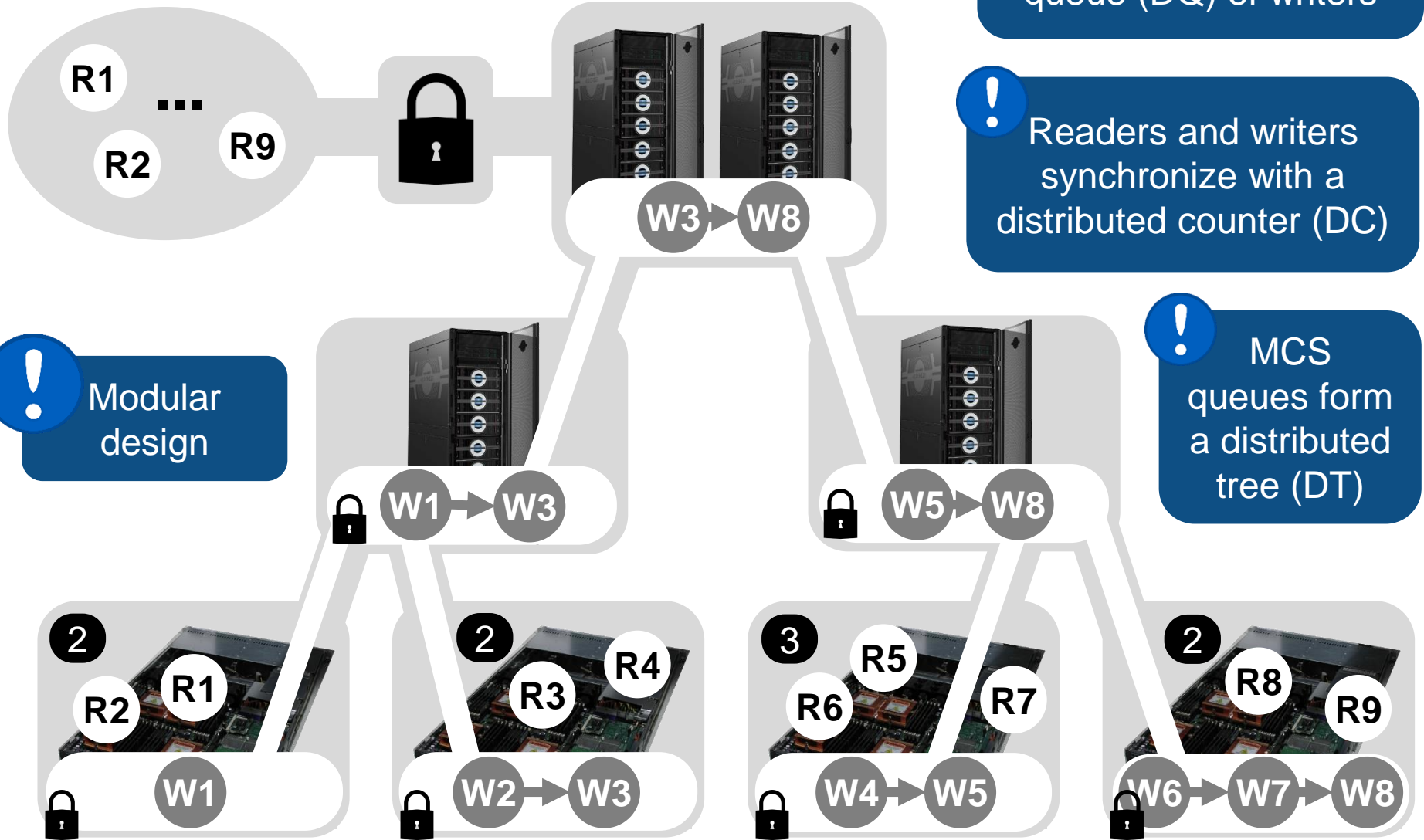What mechanism to use for efficient implementation?

How to manage the design complexity?

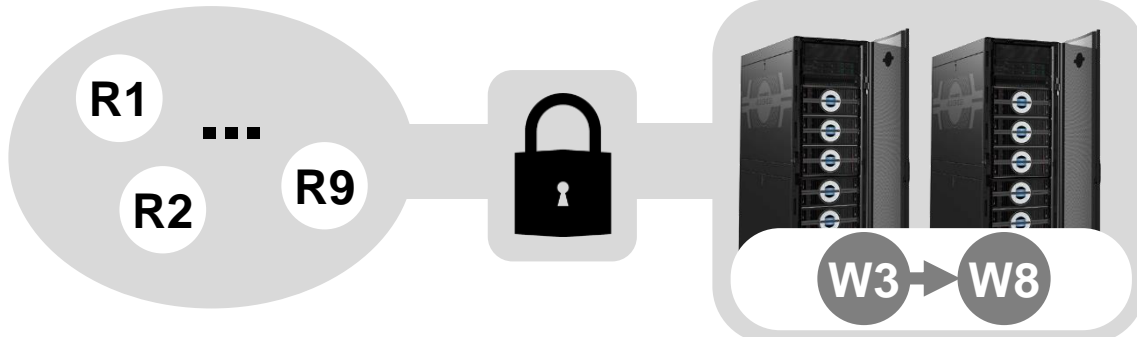Each element has its own distributed MCS queue (DQ) of writers

Readers and writers synchronize with a distributed counter (DC)

Modular design

MCS queues form a distributed tree (DT)

R1 ... R9 R2

W3 → W8

W1 → W3          W5 → W8

2  R2 R1          2  R3 R4          3  R5 R6 R7          2  R8 R9

W1          W2 → W3          W4 → W5          W6 → W7 → W8

P. Schmid, M. Besta, TH: High-Performance Distributed RMA Locks, ACM HPDC'16, best paper

**How to ensure tunable performance?**
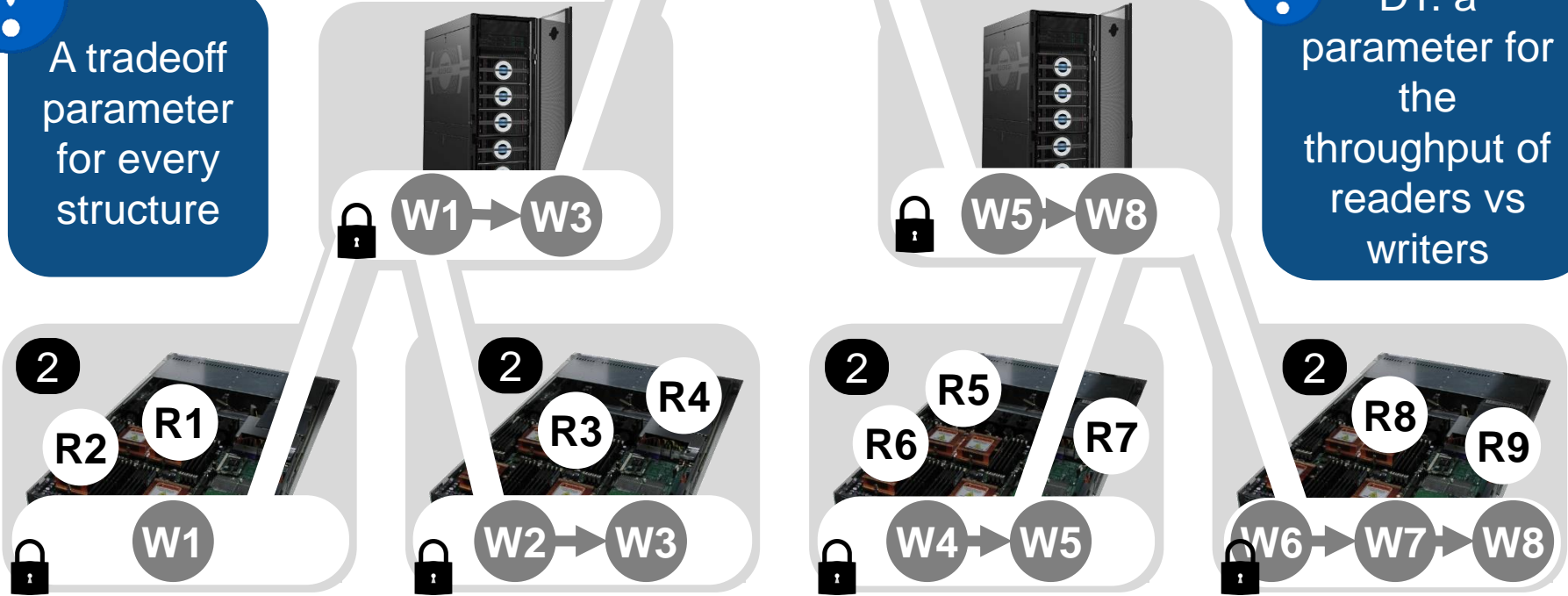
Each DQ: fairness vs throughput of writers

DC: a parameter for the latency of readers vs writers

DT: a parameter for the throughput of readers vs writers

A tradeoff parameter for every structure

R1 ... R2 R9

W3 → W8

W1 → W3

W5 → W8

2   R2 R1

2   R3 R4

2   R5 R6 R7

2   R8 R9

W1

W2 → W3

W4 → W5

W6 → W7 → W8

P. Schmid, M. Besta, TH: High-Performance Distributed RMA Locks, ACM HPDC'16, best paper

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

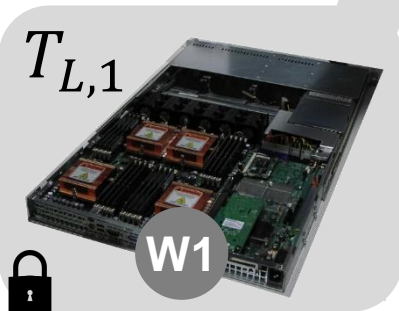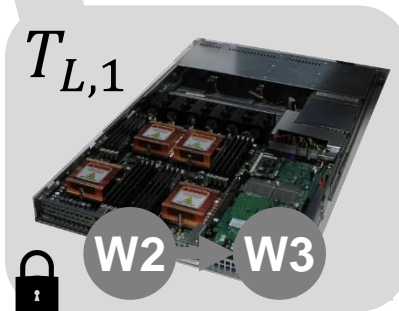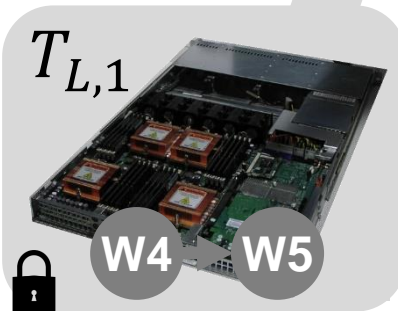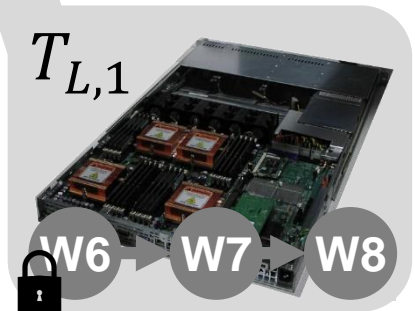Larger $T_{L,i}$ : more throughput at level i. Smaller $T_{L,i}$ : more fairness at level i.

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

$T_{L,i}$

$T_{L,3}$

W3 → W8

$T_{L,2}$
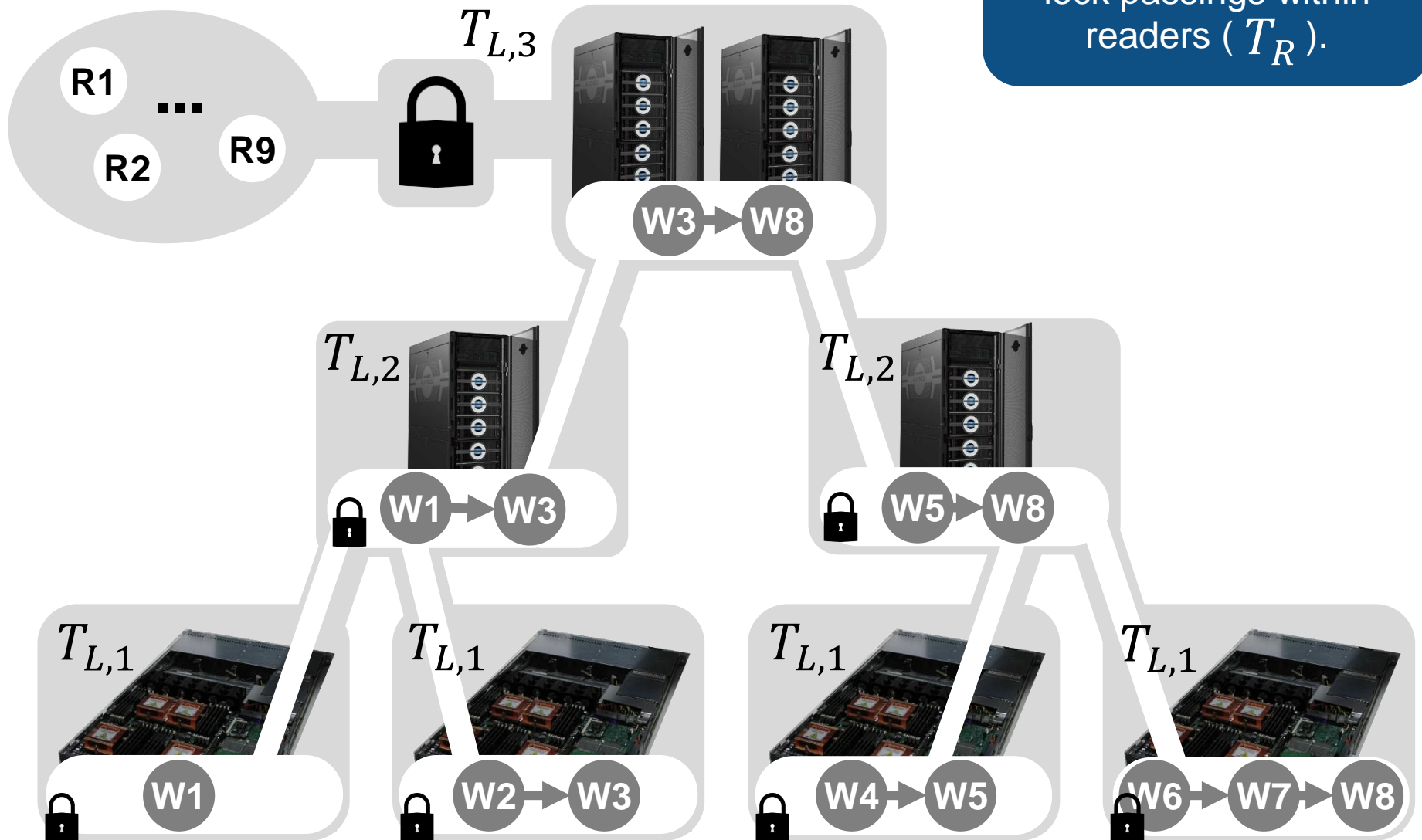
W1 → W3

$T_{L,2}$

W5 → W8

$T_{L,1}$

W1

$T_{L,1}$

W2 → W3

$T_{L,1}$

W4 → W5

$T_{L,1}$

W6 → W7 → W8

P. Schmid, M. Besta, TH: High-Performance Distributed RMA Locks, ACM HPDC'16, best paper

# 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 ... R9 R2

W3 → W8

$T_{L,2}$ W1 → W3

$T_{L,2}$ W5 → W8

$T_{L,1}$ W1

$T_{L,1}$ W2 → W3

$T_{L,1}$ W4 → W5

$T_{L,1}$ W6 → W7 → W8

P. Schmid, M. Besta, TH: High-Performance Distributed RMA Locks, ACM HPDC'16, best paper

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

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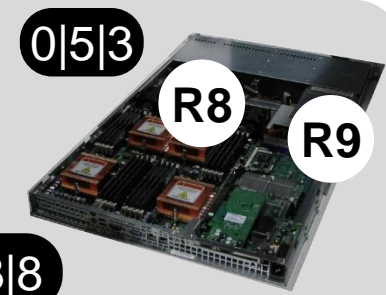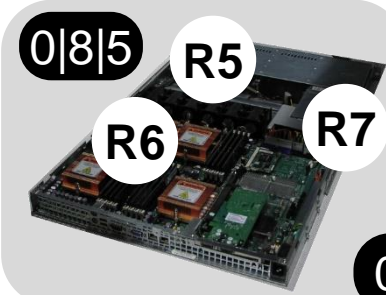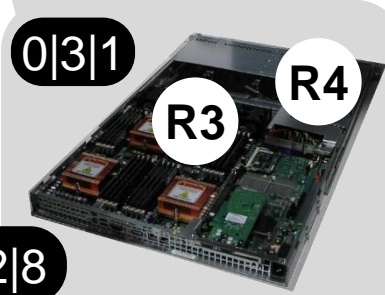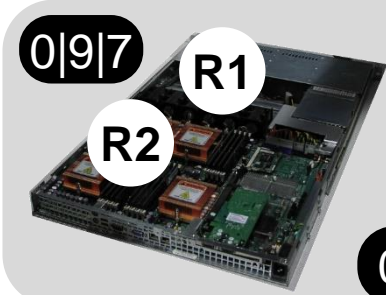
$$k = T_{DC}$$

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

P. Schmid, M. Besta, TH: High-Performance Distributed RMA Locks, ACM HPDC'16, best paper

# THE SPACE OF DESIGNS



$T_{L,i}$

Locality vs **fairness** (for writers)

Design B

Design A

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

$T_R$

$T_{DC}$

Lower latency of **writers** vs **readers**

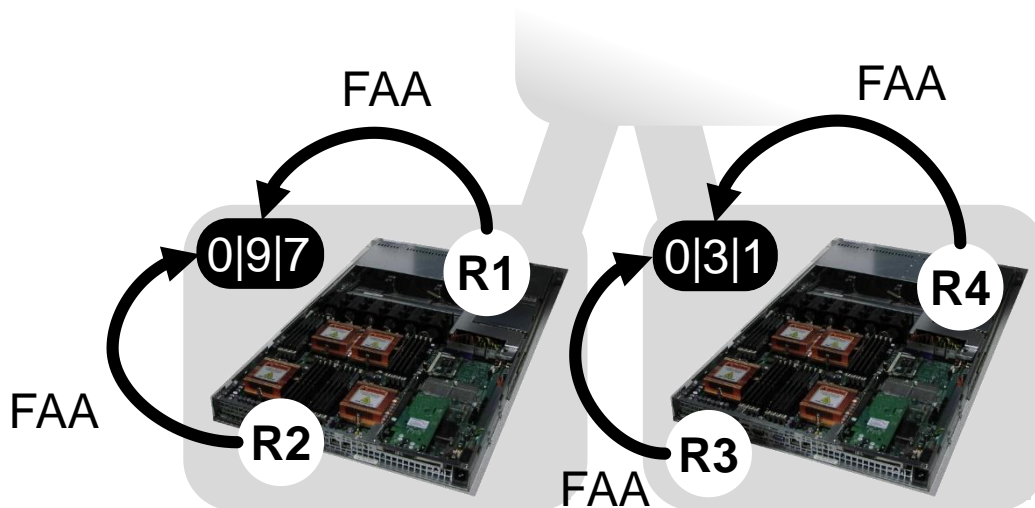P. Schmid, M. Besta, TH: High-Performance Distributed RMA Locks, ACM HPDC'16, best paper

# 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
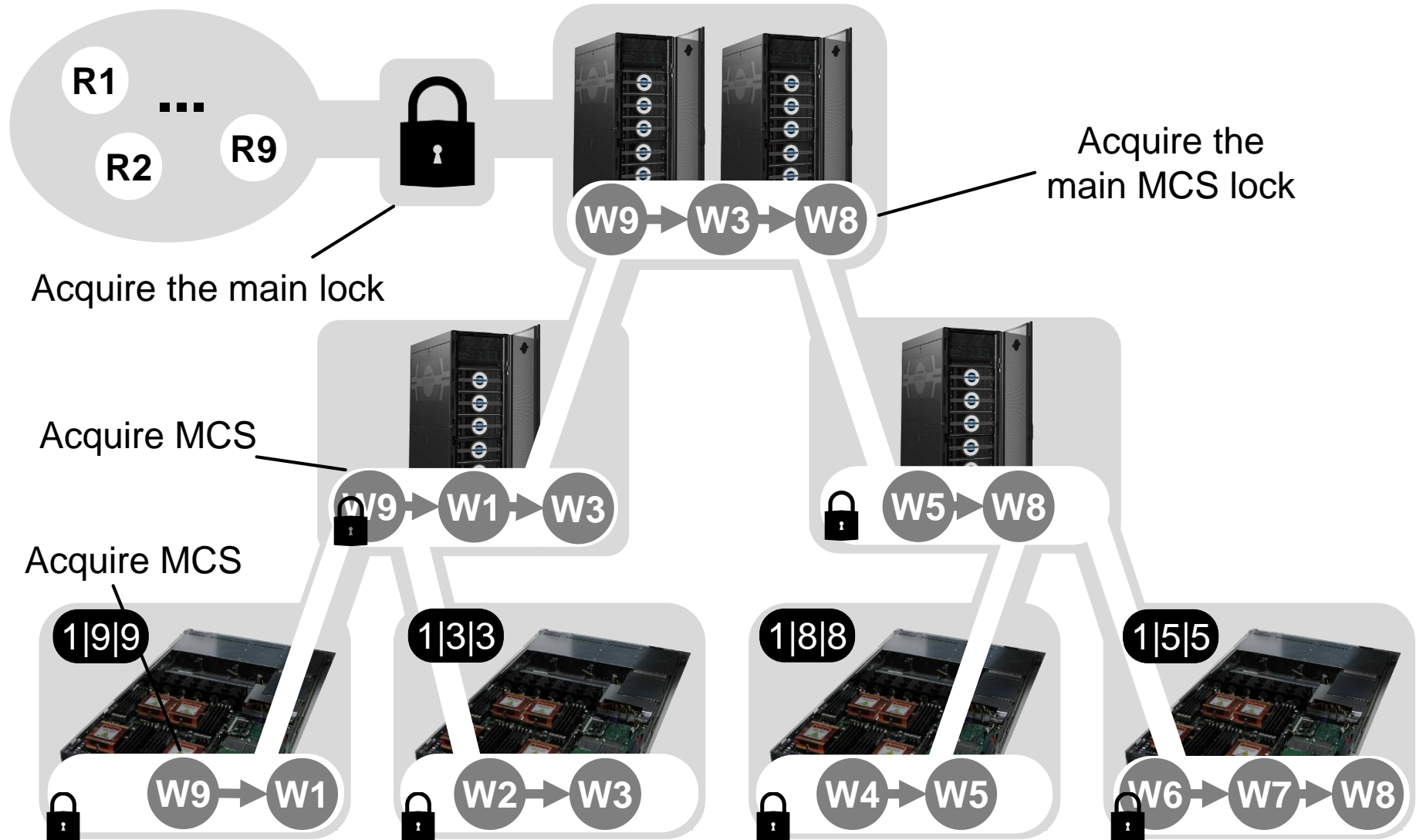
Readers that arrived at the CS

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 lock

Acquire the main MCS lock

Acquire MCS

Acquire MCS

R1 ... R2 R9

W9 → W3 → W8

W9 → W1 → W3

W5 → W8

1|9|9    1|3|3    1|8|8    1|5|5

W9 → W1    W2 → W3    W4 → W5    W6 → W7 → W8

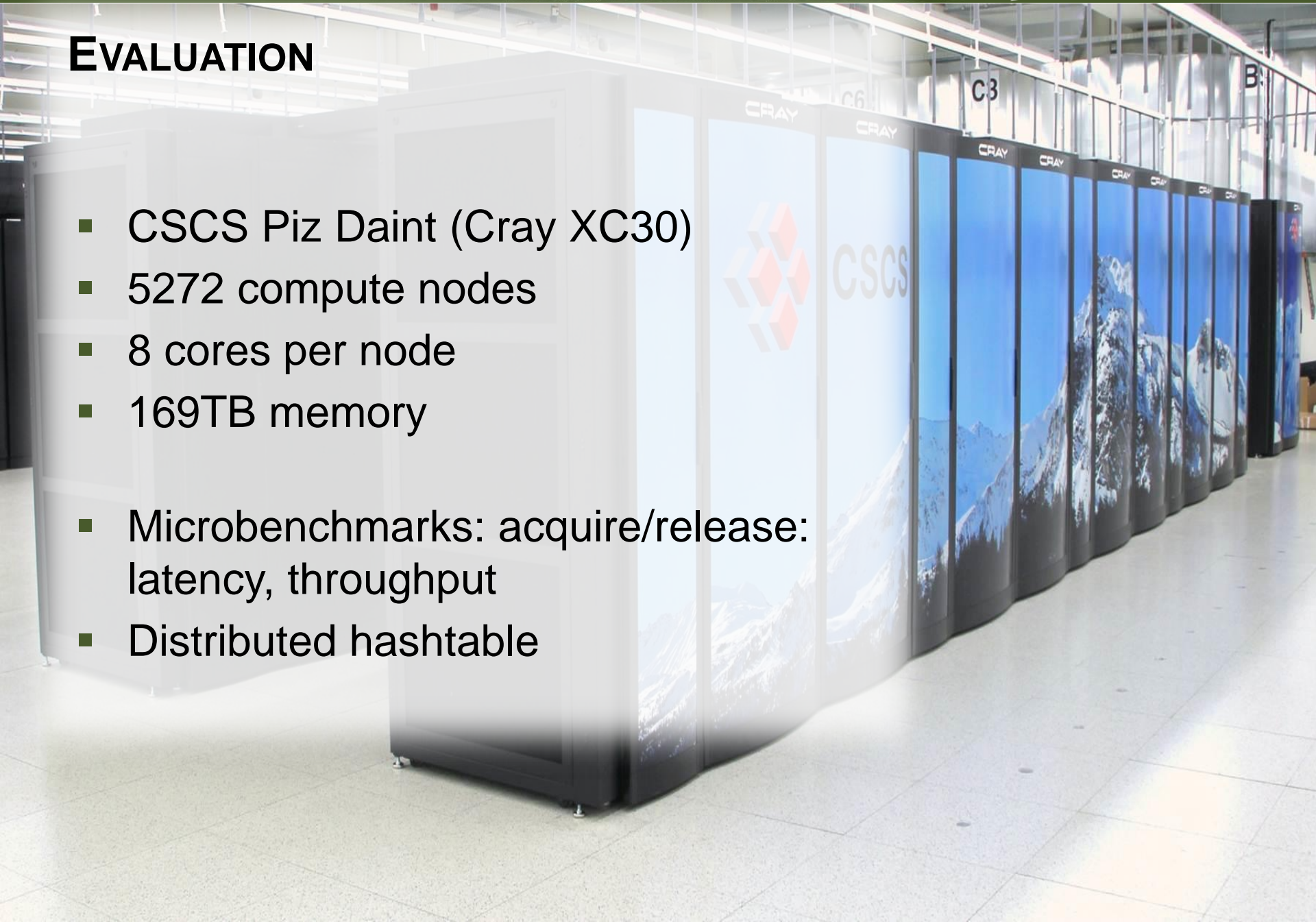P. Schmid, M. Besta, TH: High-Performance Distributed RMA Locks, ACM HPDC'16, best paper

# EVALUATION

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

- Microbenchmarks: acquire/release: latency, throughput
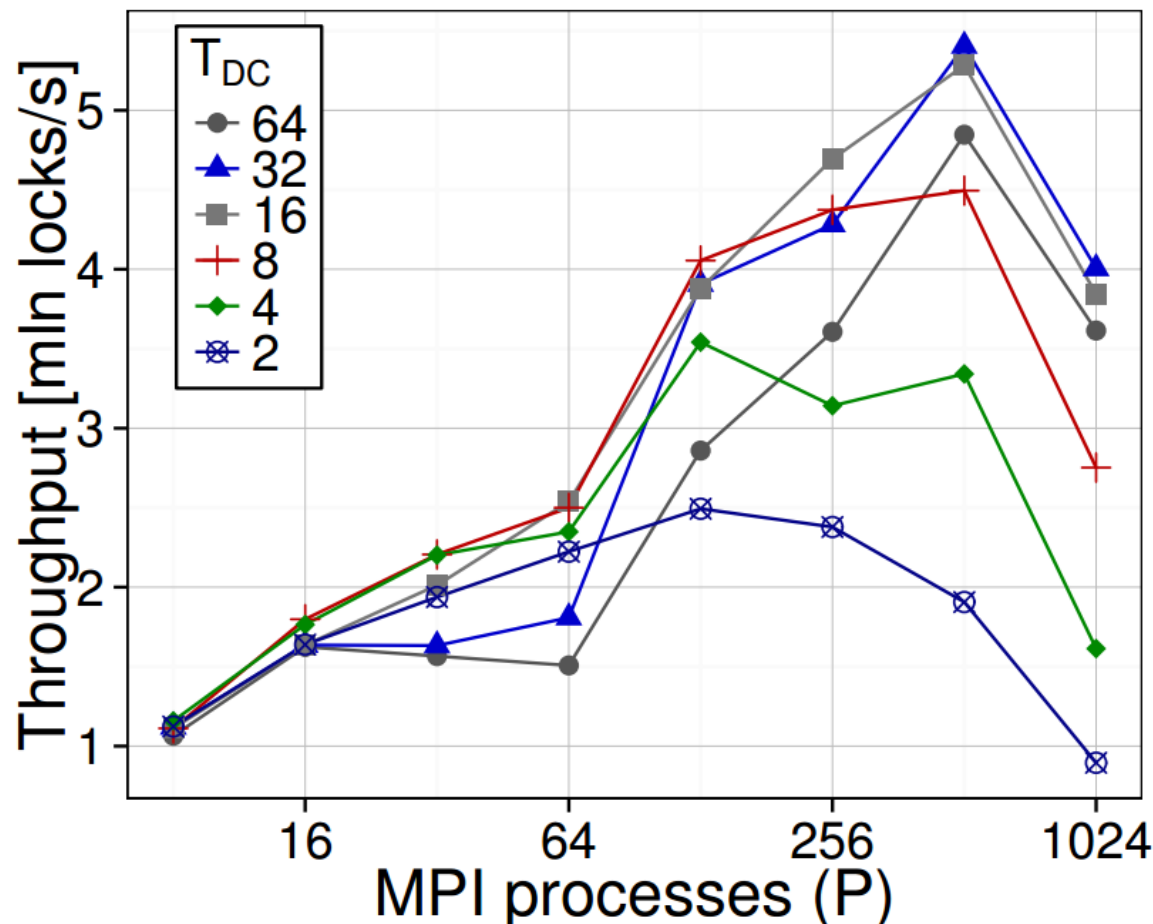- Distributed hashtable

# EVALUATION
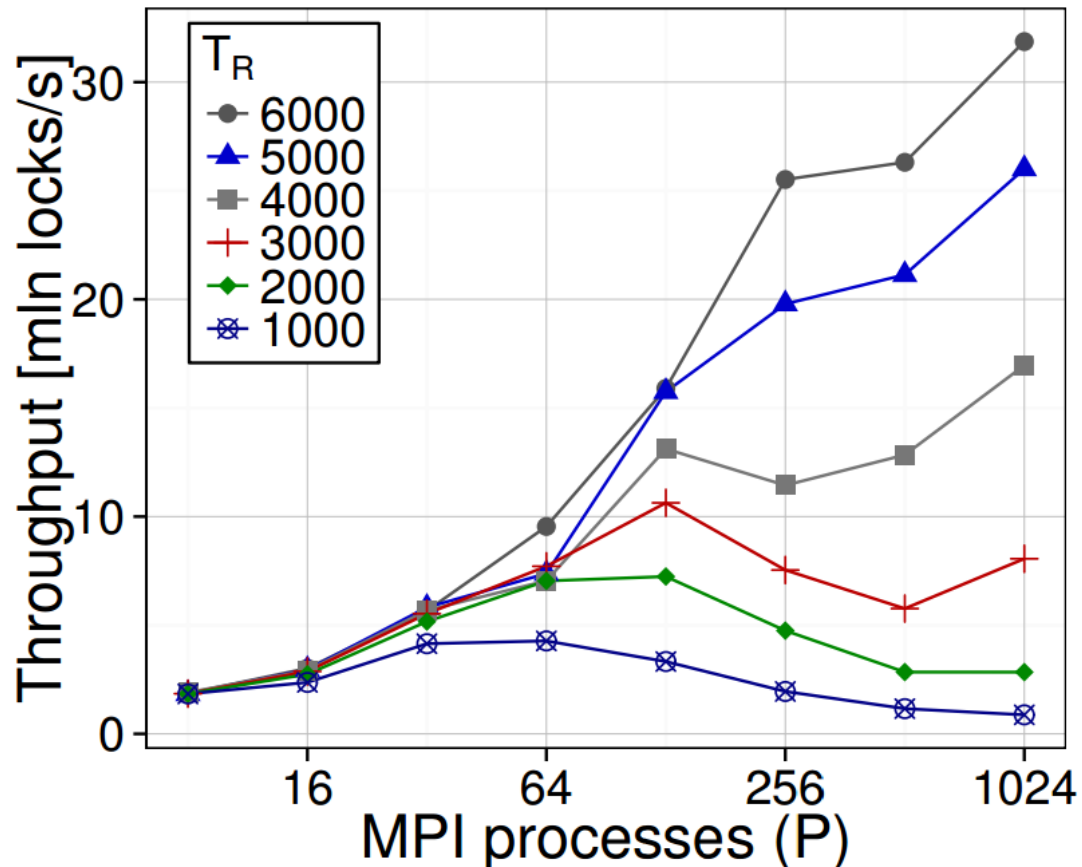## DISTRIBUTED COUNTER ANALYSIS

0|9|7   0|3|1

0|12|8
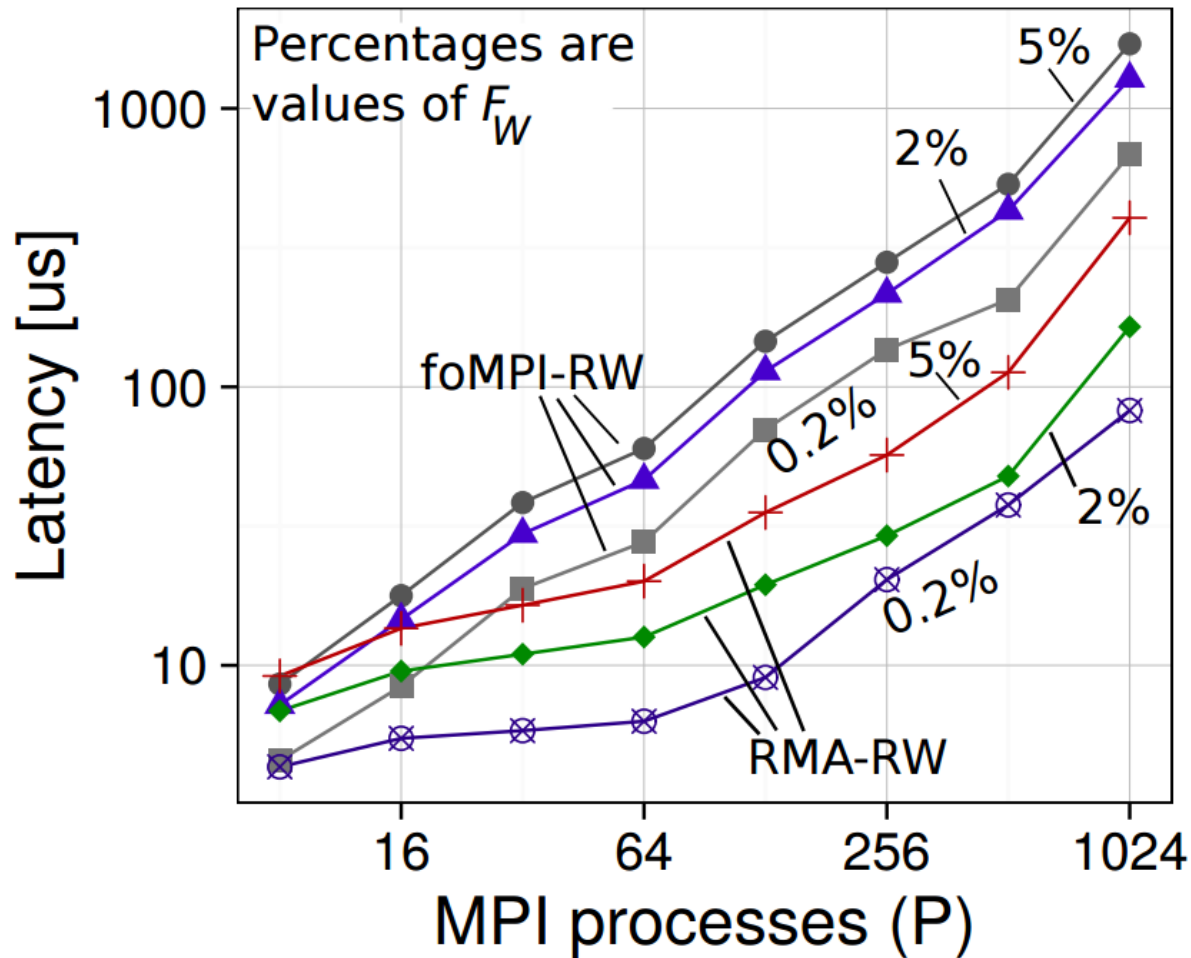


Throughput, 2% writers

Single-operation benchmark

P. Schmid, M. Besta, TH: High-Performance Distributed RMA Locks, ACM HPDC'16, best paper

# 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



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

# Another application area - Databases

- **MPI-RMA for distributed databases?**

**Hash-Join**

**Sort-Join**



C. Barthels, et al., TH: Distributed Join Algorithms on Thousands of Cores presented in Munich, Germany, VLDB Endowment, Aug. 2017

**ETH** *zürich*
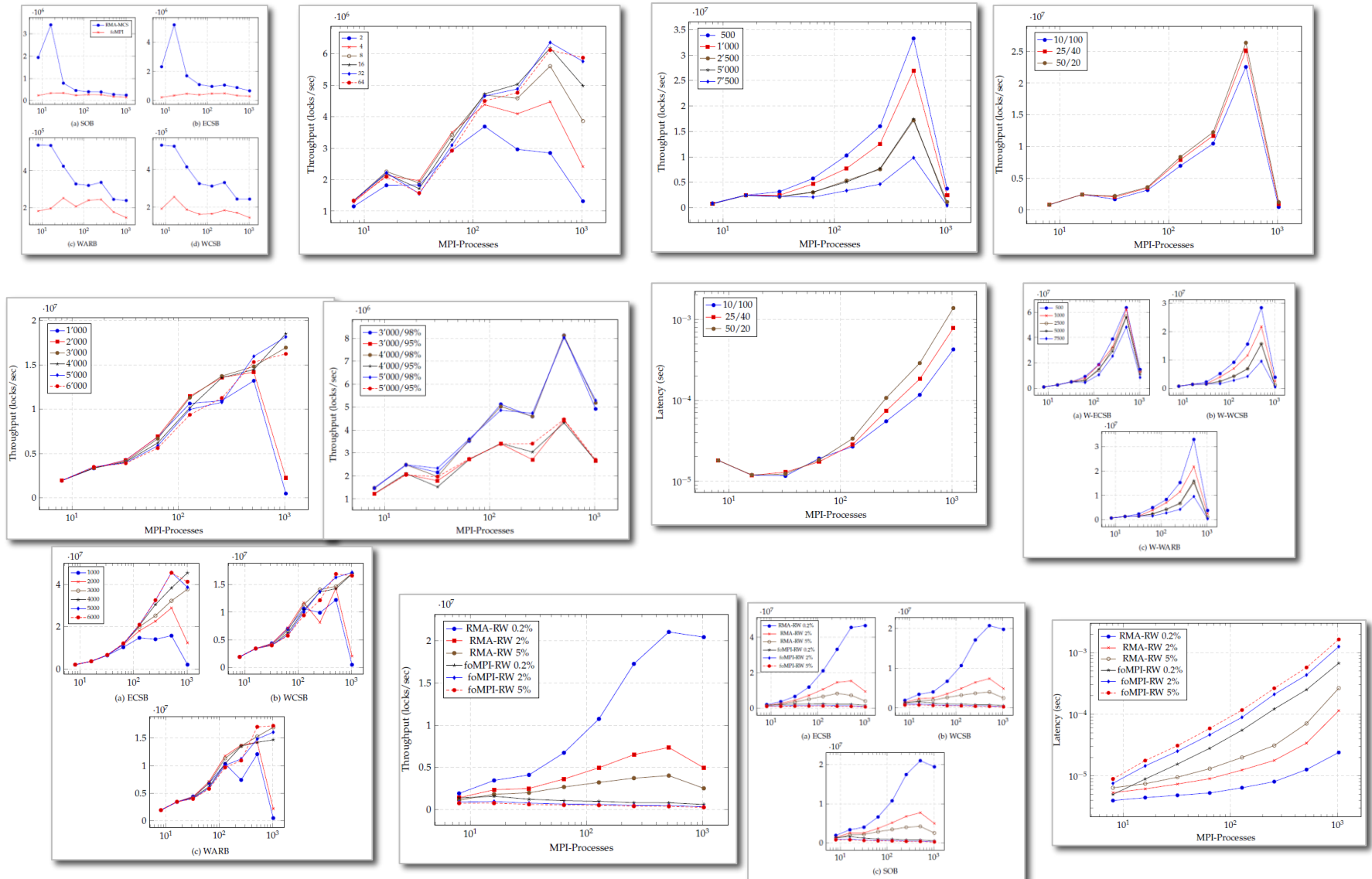
# Another application area - Databases

- **MPI-RMA for distributed databases on Piz Daint**



Network dominating

Radix hash join (MPI) with data compression

Sort-merge join (MPI) with data compression

Compute dominating

C. Barthels, et al., TH: Distributed Join Algorithms on Thousands of Cores presented in Munich, Germany, VLDB Endowment, Aug. 2017
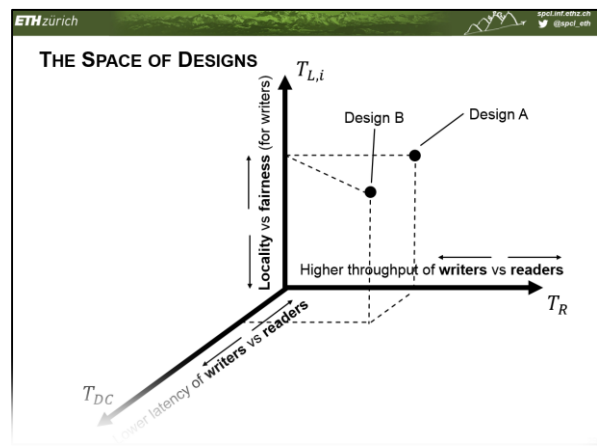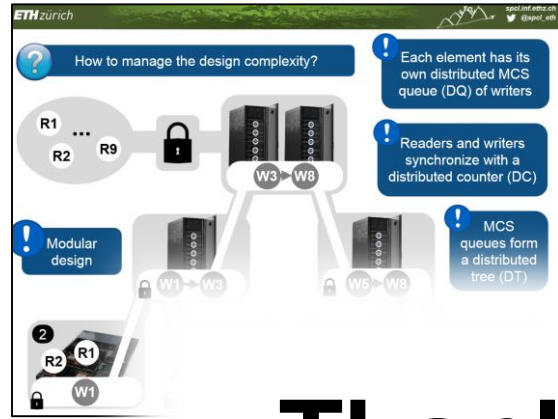
# Another application area - Databases

- **MPI-RMA for distributed databases on Piz Daint**

# OTHER ANALYSES



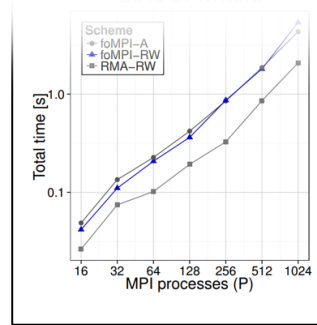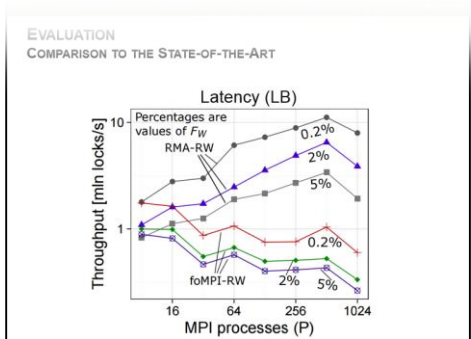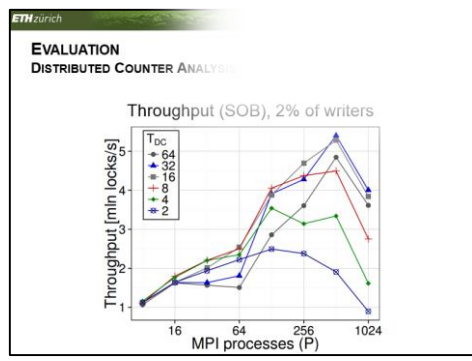P. Schmid, M. Besta, TH: High-Performance Distributed RMA Locks, ACM HPDC'16, best paper

# CONCLUSIONS

Modular design correc...

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

Accelerates distributed hashtabled

# Thank you for your attention

P. Schmid, M. Besta, TH: High-Performance Distributed RMA Locks, ACM HPDC'16, best paper