High-Performance Distributed RMA Locks

Patrick Schmid* Department of Computer Science ETH Zurich patrick.schmid@ieffects.com bestam@inf.ethz.ch

Maciej Besta* Department of Computer Science ETH Zurich

Torsten Hoefler Department of Computer Science ETH Zurich htor@inf.ethz.ch

ABSTRACT

We propose a topology-aware distributed Reader-Writer lock that accelerates irregular workloads for supercomputers and data centers. The core idea behind the lock is a modular design that is an interplay of three distributed data structures: a counter of readers/writers in the critical section, a set of queues for ordering writers waiting for the lock, and a tree that binds all the queues and synchronizes writers with readers. Each structure is associated with a parameter for favoring either readers or writers, enabling adjustable performance that can be viewed as a point in a three dimensional parameter space. We also develop a distributed topologyaware MCS lock that is a building block of the above design and improves state-of-the-art MPI implementations. Both schemes use non-blocking Remote Memory Access (RMA) techniques for highest performance and scalability. We evaluate our schemes on a Cray XC30 and illustrate that they outperform state-of-the-art MPI-3 RMA locking protocols by 81% and 73%, respectively. Finally, we use them to accelerate a distributed hashtable that represents irregular workloads such as key-value stores or graph processing.

Code: spcl.inf.ethz.ch/Research/Parallel_Programming/RMALocks

INTRODUCTION

The scale of today's data processing is growing steadily. For example, the size of Facebook's social graph is many petabytes [7, 45] and graphs processed by the well-known HPC benchmark Graph500 [37] can have trillions of vertices. Efficient analyses of such datasets require distributedmemory (DM) machines with deep Non-Uniform Memory Access (NUMA) hierarchies.

Locks are among the most effective synchronization mechanisms used in codes for such machines [6]. On one hand, if used improperly, they may cause deadlocks. Yet, they have intuitive semantics and they often outperform other schemes such as atomic operations or transactions.

Designing efficient locks for machines with deep hierarchical memory systems is challenging. Consider four pro-

Publication rights licensed to ACM. ACM acknowledges that this contribution was authored or co-authored by an employee, contractor or affiliate of a national government. As such, the Government retains a nonexclusive, royalty-free right to publish or reproduce this article, or to allow others to do so, for Government purposes only.

HPDC'16, May 31-June 04, 2016, Kyoto, Japan

ACM ISBN 978-1-4503-4314-5/16/05...\$15.00 DOI: http://dx.doi.org/10.1145/2907294.2907323 cesses competing for the same lock. Assume that two of them (A and B) run on one socket and the remaining two (C and D) execute on the other one. Now, in a naive lock design oblivious to the memory hierarchy, the lock may be passed between different sockets up to three times, degrading performance (e.g., if the order of the processes entering the critical section (CS) is A, C, B, and D). Recent advances [9,15] tackle this problem by reordering processes acquiring the lock to reduce inter-socket communication. Here, the order of A, B, C, and D entails only one inter-socket lock transfer, trading fairness for higher throughput. Extending such schemes to DM machines with weak memory models increases complexity. Moreover, expensive inter-node data transfers require more aggressive communication-avoidance strategies than those in intra-node communication [17]. To our best knowledge, no previous lock scheme addresses these challenges.

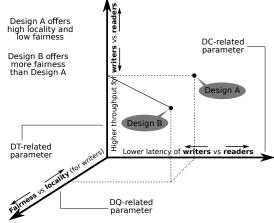


Figure 1: The space of parameters of the proposed Reader-Writer lock.

Another property of many large-scale workloads is that they are dominated by reads (e.g., they constitute 99.8% of requests to the Facebook graph [45]). Here, simple locks would entail unnecessary overheads. Instead, the Reader-Writer (RW) lock [35] can be used to reduce the overhead among processes that only perform reads in the critical section (CS). Initial RW NUMA-aware designs have recently been introduced [8], but they do not address DM machines.

In this work, we develop a lock that addresses the above challenges. Its core concept is a modular design for adjusting performance to various types of workloads. The lock consists of three key data structures. First, the distributed counter (DC) indicates the number of readers or the presence of a writer in the CS. Second, the distributed queue

^{*} Both authors contributed equally to this work.

(DQ) synchronizes writers belonging to a given element of the memory hierarchy (e.g., a rack). Finally, the distributed tree (DT) binds together all queues at different levels of the memory hierarchy and synchronizes writers with readers. Each of these three structures offers an adjustable performance tradeoff, enabling high performance in various settings. DC can lower the latency of lock acquire/release performed by either readers or writers, DQ can be biased towards improving either locality or fairness, and DT can increase the throughput of either readers or writers. The values of these parameters constitute a three dimensional space that is illustrated in Figure 1. Each point is a specific lock design with selected performance properties.

Most DM machines offer Remote Direct Memory Access (RDMA) [39], a hardware scheme that removes the OS and the CPU from the inter-node communication path. RDMA is the basis of many Remote Memory Access (RMA) [17] programming models. Among others, they offer a Partitioned Global Address Space (PGAS) abstraction to the programmer and enable low-overhead direct access to remote memories with put/get communication primitives. RMA principles are used in various HPC languages and libraries: Unified Parallel C (UPC) [44], Fortran 2008 [27], MPI-3 [36], or SHMEM [4]. We will illustrate how to utilize RMA in the proposed locks for DM machines, addressing the abovementioned challenges. In the following, we use MPI-3 RMA but we keep our protocols generic and we discuss (§ 6) how other RMA languages and libraries can also be used.

In summary, our key contributions are as follows:

- We develop a topology-aware distributed Reader-Writer lock that enables various tradeoffs between fairness, throughput, latency, and locality.
- We offer a topology-aware distributed MCS lock that accelerates the state-of-the-art MPI-3 RMA codes [17].
- We illustrate that our designs outperform the state-ofthe-art in throughput/latency (7.2x/6.8x on average) and that they accelerate distributed hashtables used in key-value (KV) stores or graph processing.

2 RMA AND LOCKS

We start by discussing RMA (§ 2.1), our tool to develop the proposed locks. Next, we present traditional (§ 2.2) and state-of-the-art (§ 2.3, § 2.4) locks that we use and extend.

Notation/Naming: We denote the number of processes as P; we use the notion of a process as it occurs frequently in DM codes such as MPI [36]. Still, our schemes are independent of whether heavyweight processes or lightweight threads are incorporated. Each process has a unique ID called the $rank \in \{1,...,P\}$. A process in the CS is called active. A null pointer is denoted as \emptyset . Then, N is the number of levels of the memory hierarchy of the used machine. Here, the selection of the considered levels depends on the user. For example, one can only focus on the nodes connected with a network and racks that contain nodes and thus N = 3 (three levels: the nodes, the racks, and the whole machine). We refer to a single considered machine part (e.g., a node) as an element. We refer to a node that is a sharedmemory cache-coherent domain connected to other such domains with a non-coherent network as a compute node (or just node). One compute node may contain smaller elements that

are cache-coherent and together offer *non-uniform memory access* (*NUMA*). We refer to such elements as *NUMA nodes*; an example NUMA node is a socket with a local DRAM. We present symbols used in the paper in Table 1.

Table 1: Symbols used in the paper.

2.1 RMA Programming

In RMA programming, processes communicate by directly accessing one another's memories. Usually, RMA is built over OS-bypass RDMA hardware for highest performance. RMA non-blocking *puts* (writes to remote memories) and *gets* (reads from remote memories) offer low latencies, potentially outperforming message passing [17]. Remote *atomics* such as compare-and-swap [21,36] are also available. Finally, RMA *flushes* ensure the consistency of data by synchronizing respective memories. RDMA is provided in virtually all modern networks (e.g., IBM PERCS [3], IBM's on-chip Cell, InfiniBand [43], iWARP [18], and RoCE [26]). Moreover, numerous libraries and languages offer RMA features. Examples include MPI-3 RMA [36], UPC [44], Titanium [22], Fortran 2008 [27], X10 [11], or Chapel [10]. The number of RMA codes is growing steadily.

RMA Windows: In RMA, each process explicitly exposes an area of its local memory as shared. In MPI, this region is called a *window*. Once shared, a window can be accessed with puts/gets/atomics and synchronized with flushes. We will refer to such an exposed memory in any RMA library/language as a window.

RMA Functions: We describe the syntax/semantics of the used RMA calls in Listing 1. All ints are 64-bit. For clarity, we also use the bool type and assume it to be an int that can take the 0 (false) or 1 (true) values, respectively. Values returned by Get/FAO/CAS are only valid after the subsequent Flush. The syntax is simplified for clarity: we omit a pointer to the accessed window (we use a single window). We use an *origin/a target* to refer to a process that issues or is targeted by an RMA call.

```
1 /* Common parameters: target: target's rank; offset: an offset
2 * into target's window that determines the location of the
3 * targeted data; op: an operation applied to a remote piece of
4 * data (either an atomic replace (REPLACE) or a sum (SUM));
5 * oprd: the operand of an atomic operation op.*/
6
7 /* Place atomically src_data in target's window.*/
8 void Put(int src_data, int target, int offset);
9
10 /* Fetch and return atomically data from target's window.*/
11 int Get(int target, int offset);
12
13 /* Apply atomically op using oprd to data at target.*/
14 void Accumulate(int oprd, int target, int offset, MPI_Op op);
15
16 /* Atomically apply op using oprd to data at target
17 * and return the previous value of the modified data.*/
18 int FAO(int oprd, int target, int offset, MPI_Op op);
19
20 /* Atomically compare cmp_data with data at target and, if
21 * equal, replace it with src_data; return the previous data.*/
22 int CAS(int src_data, int cmp_data, int target, int offset);
23
24 /* Complete all pending RMA calls started by the calling process
25 * and targeted at target.*/
26 void Flush(int target);
```

Listing 1: The syntax/semantics of the utilized RMA calls.

2.2 Traditional Hardware-Oblivious Locks

We now present hardware-oblivious locks used in this work.

2.2.1 Reader-Writer (RW) Locks

Reader-Writer (RW) locks [12] distinguish between processes that only perform reads when in the CS (*readers*) and those that issue writes (*writers*). Here, multiple readers may simultaneously enter a given CS, but only one writer can be granted access at a time, with no other concurrent readers or writers. RW locks are used in OS kernels, databases, and present in various HPC libraries such as MPI-3 [36].

2.2.2 MCS Locks

Unlike RW locks, the MCS lock (due to Mellor-Crummey and Scott) [34,40,42] does not distinguish between readers or writers. Instead, it only allows one process p at a time to enter the CS, regardless of the type of memory accesses issued by p. Here, processes waiting for the lock form a queue, with a process at the head holding the lock. The queue contains a single global pointer to its tail. Moreover, each process in the queue maintains: (1) a local flag that signals if it can enter the CS and (2) a pointer to its successor. To enter the queue, a process p updates both the global pointer to the tail and the pointer at its predecessor so that they both point to p. A releasing process notifies its successor by changing the successor's local flag. The MCS lock reduces the amount of coherence traffic that limits the performance of spinlocks [2]. Here, each process in the queue spin waits on its local flag that is modified once by its predecessor.

2.3 State-of-the-Art NUMA-Aware Locks

We now discuss lock schemes that use the knowledge of the NUMA structure of the underlying machine for more performance. We will combine and extend them to DM domains, and enrich them with a family of adjustable parameters for high performance with various workloads.

2.3.1 NUMA-Aware RW Locks

Many traditional RW locks (§ 2.2.1) entail performance penalties in NUMA systems as they usually rely on a centralized structure that becomes a bottleneck and entails high latency when accessed by processes from remote NUMA elements. Calciu et al. [8] tackle this issue with a flag on each NUMA node that indicates if there is an active reader on that node. This reduces contention due to readers (each reader only marks a local flag) but may entail additional overheads for writers that check for active readers.

2.3.2 Hierarchical MCS Locks

Hierarchical locks tackle expensive lock passing described in § 1. They trade fairness for higher throughput by ordering processes that enter the CS to reduce the number of such passings. Most of the proposed schemes address two-level NUMA machines [9,14,31,38]. Chabbi et al. consider a multi-level NUMA system [9]. Here, each NUMA hierarchy element (e.g., a socket) entails a separate MCS lock. To acquire the global lock, a process acquires an MCS lock at each machine level. This increases locality but reduces fairness: processes on the same NUMA node acquire the lock consecutively even if processes on other nodes are waiting.

2.4 Distributed RMA MCS Locks

Finally, we present a distributed MCS (D-MCS) lock based on an MPI-3 MCS lock [19]. We will use it to accelerate state-of-the-art MPI RMA library foMPI [17] and as a building

block of the proposed distributed topology-aware RW and MCS locks (§ 3).

2.4.1 Summary and Key Data Structures

Here, processes that wait for the D-MCS lock form a queue that may span multiple nodes. Each process maintains several globally visible variables. A naive approach would use one window per variable. However, this would entail additional memory overheads (one window requires $\Omega(P)$ storage in the worst case [17]). Thus, we use one window with different offsets determining different variables: a pointer to the next process in the MCS queue (offset NEXT, initially \varnothing) and a flag indicating if a given process has to spin wait (offset WAIT, initially false). A selected process (rank tail_rank) also maintains a pointer to a process with the queue tail (offset TAIL, initially \varnothing).

2.4.2 Lock Protocols

We now describe the protocols for acquire/release. We refer to respective variables using their offsets in the window.

Lock Acquire (Listing 2) First, *p* atomically modifies TAIL with its own rank and fetches the predecessor rank (Line 6). If there is no predecessor, it proceeds to the CS. Otherwise, it enqueues itself (Line 10) and waits until its local WAIT is set to false. Flushes ensure the data consistency.

```
1 void acquire() {
2     /* Prepare local fields. */
3     Put(Ø, p, NEXT);
4     Put(frue, p, STATUS);
5     /* Enter the tail of the MCS queue and get the predecessor. */
6     int pred = FAO(p, tail_rank, TAIL, REPLACE);
7     Flush(tail_rank); /* Ensure completion of FAO. */
8     if(pred != Ø) { /* Check if there is a predecessor. */
9          /* Make the predecessor see us. */
10     Put(p, pred, NEXT); Flush(pred);
11     bool waiting = true;
12     do { /* Spin locally until we get the lock. */
13          waiting = Get(p, WAIT); Flush(p);
14     ) while(waiting == true); } }
```

Listing 2: Acquiring D-MCS.

Lock Release (Listing 3) First, p checks if it has a successor in the queue (Line 3). If there is none, it atomically verifies if it is still the queue tail (Line 5); if yes, it sets TAIL to \emptyset . Otherwise, p waits for a process that has modified TAIL to update its NEXT field (Lines 9-11). If there is a successor, the lock is passed with a single Put (Line 14).

Listing 3: Releasing D-MCS.

3 DISTRIBUTED RMA RW LOCKS

We now present a distributed *topology-aware* RW lock (RMA-RW) for scalable synchronization and full utilization of parallelism in workloads dominated by reads. For brevity, we focus on the RW semantics and only briefly discuss the topology-aware MCS lock (§ 3.5); full algorithm listings for

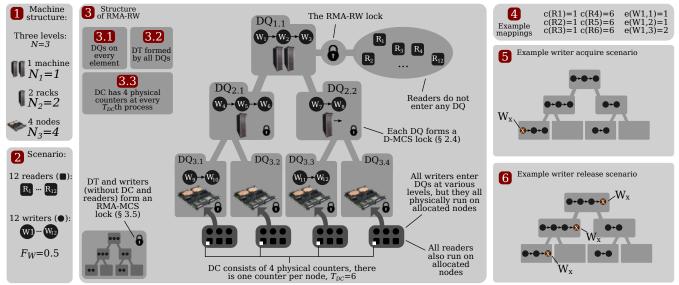


Figure 2: An example RMA-RW on a three-level system.

the latter can be found in an extended technical report¹. Symbols specific to RMA-RW are presented in Table 2.

Lock Abbreviations We always refer to the proposed topology-aware distributed RW and MCS lock as RMA-RW and RMA-MCS, respectively. Both RMA-RW and RMA-MCS use as their building block a simple distributed topology-oblivious MCS lock (§ 2.4) denoted as D-MCS.

Example In the whole section, we will use the example shown in Figure 2. Here, N=3 and the considered levels are: compute nodes, racks, and the whole machine.

```
\begin{array}{ll} T_{DC} & \text{The $Distributed Counter} \text{ threshold } (\S~3.2.1). \\ T_{L,i} & \text{The $Locality$ threshold at level } i~(\S~3.2.2). \\ T_{R} & \text{The $Reader$ threshold } (\S~3.2.3). \\ T_{W} & \text{The $Writer$ threshold; } T_{W} = \prod_{i=1}^{N} T_{L,i}~(\S~3.2.3). \\ c(p) & \text{Mapping from a process $p$ to its physical counter } (\S~3.2.1). \\ e(p,i) & \text{Mapping from a process $p$ to its home machine element at level } i~(\S~3.2.2). \\ F_{W} & \text{The fraction of writers in a given workload (the fraction of readers: } 1-F_{W}). \\ \end{array}
```

Table 2: Symbols used in RMA-RW.

3.1 Design Summary and Intuition

As explained in § 1, RMA-RW consists of three types of core data structures: distributed queues (DQs), a distributed tree (DT), and a distributed counter (DC). They are illustrated in Figure 2. First, every machine element (at each considered level) has an associated DQ and thus a D-MCS lock local to this element (as opposed to the global RMA-RW lock). In our example, every node, rack, and the whole machine have their own DQ (and thus a local MCS lock). Note that some DQs that are associated with elements such as nodes are not necessarily distributed, but we use the same name for clarity. Second, all the DQs form a DT that corresponds to the underlying memory hierarchy, with one DQ related to one tree vertex. For example, DQs associated with nodes that belong to a given rack r constitute vertices that are children of a vertex associated with a DQ running on rack r. Third, DC counts active readers and writers and consists of several physical counters located on selected processes. DT on its own (without DC and any readers) constitutes RMA-MCS.

Writers A writer that wants to acquire a lock starts at a leaf of DT located at the lowest level N (a node in our

example). At any level i ($2 \le i \le N$), it acquires a local D-MCS lock that corresponds to a subtree of D-MCS locks (and thus DQs) rooted at the given element. Here, it may compete with other writers. When it reaches level 1, it executes a different protocol for acquiring the whole RMA-RW lock. Here, it may also compete with readers. RMA-RW's locality-aware design enables a *shortcut*: some writers stop before reaching level 1 and directly proceed to the CS. This happens if a lock is passed within a given machine element.

Readers Readers do not enter DQs and DT and thus have a single acquire protocol. This design reduces synchronization overhead among readers.

3.2 Key Data Structures

We now present the key structures in more detail.

3.2.1 Distributed Counter (DC)

DC maintains the number of active readers or writers. It enables an adjustable performance tradeoff that accelerates readers or writers. For this, one DC consists of multiple physical counters, each maintained by every T_{DC} th process; T_{DC} is a parameter selected by the user. To enter the CS, a reader p increments only one associated physical counter while a writer must check each one of them. Thus, selecting more physical counters (smaller T_{DC}) entails lower reader latency (as each reader can access a counter located on a closer machine element) and contention (as each counter is accessed by fewer readers). Yet, higher T_{DC} entails lower latency for a writer that accesses fewer physical counters.

A physical counter associated with a reader p is located at a rank c(p); $c(\cdot) \in \{1,...,P\}$ can be determined at compile-or run-time. In a simple hardware-oblivious scheme, one can fix $c(p) = \lceil p/T_{DC} \rceil$. For more performance, the user can locate physical counters in a topology-aware way. For example, if the user allocates x processes/node and a node s hosts processes with x successive ranks starting from (s-1)x+1, then setting $T_{DC}=kx$ in the above formula results in storing one physical counter every kth node. This can be generalized to any other machine element.

To increase performance, we implement each physical counter as two 64-bit fields that count the readers (assigned to this counter) that arrived and departed from the CS, re-

 $^{^{1}} http://spcl.inf.ethz.ch/Research/Parallel_Programming/RMALocks$

spectively. This facilitates obtaining the number of readers that acquired the lock since the last writer and reduces contention between processes that acquire and release the lock. We dedicate one bit of the field that counts arriving readers to indicate whether the CS of RMA-RW is in the READ mode (it contains readers) or the WRITE mode (it contains a writer).

RMA Design of DC: Each physical counter occupies two words with offsets ARRIVE (for counting arriving readers) and DEPART (for counting departing readers); physical counters together constitute an RMA window.

3.2.2 Distributed Queue (DQ)

DQ orders writers from a single element of the machine that attempt to enter the CS. DQs from level i have an associated threshold $T_{L,i}$ that determines the maximum number of lock passings between writers running on a machine element from this level before the lock is passed to a process from a different element. We use a separate threshold $T_{L,i}$ for each i because some levels (e.g., racks) may need more locality (a higher threshold) than others (e.g., nodes) due to expensive data transfers. This design enables an adjustable tradeoff between fairness and throughput at each level.

DQ extends D-MCS in that the local flag that originally signals whether a process can enter the CS now becomes an integer that carries (in the same RMA operation) the number of past lock acquires within a given machine element. We use this value to decide whether to pass the lock to a different element at a given level i (if the value reaches $T_{L,i}$) or not (if the value is below $T_{L,i}$).

RMA Design of DQ: All DQs at a given level constitute an RMA window. Respective offsets in the window are as follows: NEXT (a rank of the next process in the queue), STATUS (an integer that both signals whether to spin wait and carries the number of past lock acquires in the associated machine element), and TAIL (a rank of the process that constitutes the current tail of the queue). TAIL in DQ at level *i* associated with *j*th element is stored on a process tail_rank[*i*, *j*].

3.2.3 Distributed Tree of Queues (DT)

DT combines DQs at different memory hierarchy levels into a single structure. This enables p to make progress in acquiring/releasing RMA-RW by moving from level N to level 1. Then, at the tree root, writers synchronize with readers. Specifically, the lock is passed from writers to readers (if there are some waiting) when the total number of lock passings between writers reaches a threshold T_W . In our design, $T_W = \prod_{i=1}^N T_{L,i}$. To avoid starvation of writers, we also introduce a threshold T_R that is the maximum number of readers that can enter the CS consecutively before the lock is passed to a writer (if there is one waiting). Increasing T_R or T_W improves the throughput of readers or writers because more processes of a given type can enter the CS consecutively.

While climbing up DT, a writer must determine the next DQ (and thus D-MCS) to enter. This information is encoded in a mapping $e(\cdot,\cdot)$ and structure $tail_rank[i,j]$. $e(p,i) \in \{1,...,N_i\}$ returns the ID of a machine element associated with a process p at level i. An expression $tail_rank[i,e(p,i)]$ returns the rank of a process that points to the tail of a DQ at level i within a machine element assigned to p. This enables p to enter D-MCS at the next level on the way to the CS. Similarly to c(p), e(p,i) can be determined statically or dynamically.

Depending on $T_{L,i}$, some writers do not have to climb all DT levels and can proceed directly to the CS. Thus, we fur-

ther extend the STATUS field used in DQ with one more special value ACQUIRE_PARENT. This indicates that *p* cannot directly enter the CS and should continue up DT.

3.2.4 Discussion on the Status Field

A central part of DQ and DT is the STATUS field that enables processes to exchange various additional types of information in a single RMA communication action, including: (1) if a lock mode changed (e.g., from READ to WRITE), (2) if a given process should acquire a lock at a higher DT level, (3) if a given process can enter the CS, and (4) the number of past consecutive lock acquires. Two selected integer values are dedicated to indicate (1) and (2). All the remaining possible values indicate that the given process can enter the CS (3); at the same time the value communicates (4).

3.3 Distributed Reader-Writer Protocol

We now illustrate how the above data structures play together in the acquire and release protocols. A writer starts at the leaf of DT (level N) both for acquiring and releasing. At any level i ($2 \le i \le N$), it proceeds up the tree executing a protocol for a partial acquire/release of the respective part of the tree (§ 3.3.1, § 3.3.2). At level 1, it executes a different protocol for locking or releasing the whole lock (§ 3.3.3, § 3.3.4). Readers do not follow such a hierarchy and thus have single acquire (§ 3.3.5) and release (§ 3.3.6) protocols.

3.3.1 Writer Lock Acquire: Level N to 2 (Listing 4)

Intuition: *p* enters the DQ associated with a given level *i* and its home element e(p,i); it then waits for the update from its predecessor. If the predecessor does not have to hand over the lock to a process from another element (i.e., has not reached the threshold $T_{L,i}$), the lock is passed to p that immediately enters the CS. Otherwise, p moves to level i-1. **Details:** p first modifies its NEXT and STATUS to reflect it spin waits at the DQ tail (Lines 2-3). Then, it enqueues itself (Line 5). If there is a predecessor at this level, p makes itself visible to it with a Put (Line 8) and then waits until it obtains the lock. While waiting, p uses Gets and Flushes to check for any updates from the predecessor. If the predecessor reached $T_{L,i}$ and released the lock to the parent level, p must itself acquire the lock from level i-1 (Line 23). Otherwise, it can directly enter the CS as the lock is simply passed to it (Line 18). If there is no predecessor at level *i*, *p* also proceeds to acquire the lock for level i - 1 (Line 23).

```
void writer-acquire<i>() {
        Put(\emptyset, p, NEXT);

Put(WAIT, p, STATUS); Flush(p);

/* Enter the DQ at level i and in this machine element. */

int pred = FAO(p, tail_rank[i,e(p,i)], TAIL, REPLACE);
        Flush(tail_rank[i,e(p,i)]);

if(pred != \emptyset) {

Put(p, pred, NEXT); Flush(pred); /* pred sees us. */
            rut(p, pred, NEXI); Flush(pred); /* pred set
int status = WAIT;
do { /* Wait until pred passes the lock. */
    status = Get(p, STATUS); Flush(p);
} while(status == WAIT);
10
11
12
             /* Check if pred released the lock to the parent level. This
13
            would happen if T_{L,i} is reached. */
if(status != ACQUIRE_PARENT) {
14
15
                /* T_{L,i} is not reached. Thus, the lock is passed to p that directly proceeds to the CS. */
16
18
                 return; /* The global lock is acquired. */
19
        }
21
        /* Start to acquire the next level of the tree.*/
Put(ACQUIRE_START, p, STATUS); Flush(p);
        writer-acquire\langle i-1\rangle();}
```

Listing 4: Acquiring the RMA-RW lock by a writer; levels N to 2.

3.3.2 Writer Lock Release: Level N to 2 (Listing 5)

Intuition: p passes the lock within e(p,i) if there is a successor and $T_{L,i}$ is not yet reached. Otherwise, it releases the lock to the parent level i-1, leaves the DQ, and informs any new successor that it must acquire the lock at level i-1.

Details: p first finds out whether it has a successor. If there is one and $T_{L,i}$ is not yet reached, the lock is passed to it with a Put (Line 8). If $T_{L,i}$ is reached, p releases the lock for this level and informs its successor (if any) that it has to acquire the lock at level i-1. If there is no known successor, it checks atomically if some process has already entered the DQ at level i (Line 15). If so, the releaser waits for the successor to make himself visible before it is notified to acquire the lock at level i-1.

Listing 5: Releasing an RMA-RW lock by a writer; levels N to 2.

3.3.3 Writer Lock Acquire: Level 1 (Listing 7)

Intuition: This scheme is similar to acquiring the lock at lower levels (§ 3.3.1). However, the predecessor may notify p of the *lock mode change* that enabled readers to enter the CS, forcing p to acquire the lock from the readers.

Details: p first tries to obtain the lock from a predecessor (Lines 2-18). If there is one, p waits until the lock is passed. Still, it can happen that the predecessor hands the lock over to the readers (Line 14). Here, p changes the mode back to WRITE before entering the CS (Line 16); this function checks each counter to verify if there are active readers. If not, it switches the value of each counter to WRITE (see Listing 6). If there is no predecessor (Line 19), p tries to acquire the lock from the readers by changing the mode to WRITE (Line 21).

3.3.4 Writer Lock Release: Level 1 (Listing 8)

Intuition: p first checks if it has reached T_W and if there is a successor waiting at level 1. If any case is true, it passes the lock to the readers and notifies any successor that it must acquire the lock from them. Otherwise, the lock is handed over to the successor.

Details: First, if T_W is reached, p passes the lock to the readers by resetting the counters (Line 6). Then, if it has no successor, it similarly enables the readers to enter the CS (Line 12). Later, p appropriately modifies the tail of the DQ and verifies if there is a new successor (Line 17). If necessary, it passes the lock to the successor with a Put (line 23) and simultaneously (using next_stat) notifies it about a possible lock mode change.

```
1 /***** Change all physical counters to the WRITE mode *****/
 2 void set_counters_to_WRITE() {
     /* To simplify, we use one counter every T_{DC}th process.*/ for(int p = 0; p < P; p += T_{DC}) {    /* Increase the arrival counter to block the readers.*/
        {\tt Accumulate(INT64\_MAX/2, \ p, \ ARRIVE, \ SUM); \ Flush(p);}
  /******* Reset one physical counter ***********/
10 void reset_counter(int rank) {
11  /* Get the current values of the counters.*/
12  int arr_cnt = Get(rank, ARRIVE), dep_cnt = Get(rank, DEPART);
      /* Prepare the values to be subtracted from the counters.*/
     int sub_arr_cnt = -dep_cnt, sub_dep_cnt = -dep_cnt;
16
17
      /* Make sure that the WRITE is reset if it was set.*/
     if(arr_cnt >= INT64_MAX/2)
19
        sub_arr_cnt -= INT64_MAX/2;
20
21
      /* Subtract the values from the current counters.*/
     Accumulate(sub_arr_cnt, rank, ARRIVE, SUM);
Accumulate(sub_dep_cnt, rank, DEPART, SUM); Flush(rank);
22
24 }
25
   /******* Reset all physical counters ***********/
27 void reset_counters() {
     for(int p = 0; p < P; p \leftarrow T_{DC}) { reset_counter(p); } }
```

Listing 6: Functions that manipulate counters.

```
1 void writer-acquire<1>() {
      Put(\emptyset, p, NEXT); Put(WAIT, p, STATUS); Flush(p); /* Prepare to enter the DQ.*/ /* Enqueue oneself to the end of the DQ at level 1.*/ int pred = FAO(p, tail_rank[1,e(p,1)], TAIL, REPLACE); Flush(tail_rank[1,e(p,1)]);
       if(pred != \emptyset) { /* If there is a predecessor...*/
         Put(p, pred, NEXT); Flush(pred);
int curr_stat = WAIT;
10
11
          do { /* Wait until pred notifies us.*/
  curr_stat = Get(p, STATUS); Flush(p);
} while (curr_stat == WAIT);
12
13
14
          if(curr_stat == MODE_CHANGE) { /* The lock mode changed...*/
15
             \slash * The readers have the lock now; try to get it back.*/
16
             set_counters_to_WRITE();
17
             Put(ACQUIRE_START, p, STATUS); Flush(p);
18
19
       else { /* If there is no predecessor...*/
20
          /* Change the counters to WRITE as we have the lock now.*/
set_counters_to_WRITE();
21
          Put(ACQUIRE_START, p, STATUS); Flush(p); } }
```

Listing 7: Acquiring an RMA-RW lock by a writer; level $1. \,$

3.3.5 Reader Lock Acquire (Listing 9)

Intuition: Here, p first spin waits if there is an active writer or if p's arrival made its associated counter c(p) exceed T_R . Then, it can enter the CS. If $c(p) = T_R$, then p resets DC.

Details: In the first part, p may spin wait on a boolean barrier variable (Line 5), waiting to get the lock from a writer. Then, p atomically increments its associated counter and checks whether the count is below T_R . If yes, the lock mode is READ and p enters the CS. Otherwise, either the lock mode is WRITE or T_R is reached. In case of the latter, p checks if there are any waiting writers (Line 17). If there are none, p resets the DC (Line 20) and re-attempts to acquire the lock. If there is a writer, p sets the local barrier and waits for DC to be reset by the writer.

3.3.6 Reader Lock Release (Listing 10)

Releasing a reader lock only involves incrementing the departing reader counter.

```
1 void reader-release() { 2    Accumulate(1, c(p), DEPART, SUM); Flush(c(p)); }
```

Listing 10: Releasing an RMA-RW reader lock.

3.4 Example

Consider the scenario from Figure 2. Here, there are three machine levels, 12 readers, and 12 writers ($F_W = 0.5$).

```
1 void writer-release<1>(){
       bool counters_reset = false;
/* Get the count of consecutive lock acquires (level 1).*/
        int next_stat = Get(p, STATUS); Flush(p);
       if(++next\_stat == T_W) { /* Pass the lock to the readers.*/
reset\_counters();/* See Listing 6.*/
           next_stat = MODE_CHANGE; counters_reset = true;
 8
       int succ = Get(p, NEXT); Flush(p);
if(succ == Ø) { /* No known successor.*/
   if(!counters_reset) { /* Pass the lock to the readers.*/
     reset_counters(); next_stat = MODE_CHANGE;/* Listing 6.*/
11
12
13
            /* Check if some process has already entered the DO.*/
          \begin{array}{ll} \textbf{int} \  \, \text{curr\_rank} = \  \, \text{CAS}(\varnothing, \ p, \ \text{tail\_rank[1}, e(p, 1)], \ \text{TAIL}); \\ \text{Flush(tail\_rank[1}, e(p, 1)]); \end{array}
15
16
          if(p = curr_rank) { return; } /* No successor...*/ do { /* Wait until the successor makes itself visible.*/
          succ = Get(p, NEXT); Flush(p);
} while (succ == \emptyset);
19
20
21
        /* Pass the lock to the successor.*/
22
       Put(next_stat, succ, STATUS); Flush(succ); }
```

Listing 8: Releasing an RMA-RW lock by a writer; level 1.

```
1 \; {f void} \; {f reader-acquire()} \; \{
      bool done = false; bool barrier = false;
       while(!done) {
         int curr stat = 0:
         if(barrier) {
             \begin{array}{lll} {\rm curr\_stat} \; = \; {\rm Get}(c(p), \; {\rm ARRIVE}) \, ; \; \; {\rm Flush}(c(p)) \, ; \\ {\rm Supple}(curr\_stat \; >= \; T_R) \, ; \\ \end{array} 
10
          /* Increment the arrival counter.*/
          curr_stat = FAO(1, c(p), ARRIVE, SUM); Flush(c(p)); if(curr_stat >= T_R) { /* T_R has been reached...*/
13
            barrier = true; if (curr\_stat == T_R) {/* We are the first to reach T_R.*/
15
                /* Pass the lock to the writers if there are any.*/
16
                int curr_tail = Get(tail_rank[1,e(p,1)], TAIL);
               Flush(tail_rank[1,e(p,1)]);
if(curr_tail == \emptyset) { /* There are no waiting writers.*/
18
19
20
21
                   reset_counter(c(p)); barrier = false;/* Listing 6.*/
               }
23
             /* Back off and try again.*/
              \label{eq:accumulate}  \mbox{Accumulate(-1, $c(p)$, $ARRIVE, $SUM)$; $Flush(c(p))$;}
```

Listing 9: Acquiring an RMA-RW lock by a reader.

Writer Acquire Assume a new writer W_x running on a node related to DQ_{3.1} attempts to acquire RMA-RW (Figure 2, Part 5). First, it enters DQ_{3.1} (Listing 4). As this queue is not empty, W_x spins locally (Lines 10-12) until its predecessor W_9 modifies W_x 's STATUS. Now, if W_9 has not yet reached $T_{L,3}$, W_x gets the lock and immediately proceeds to the CS (Lines 15-19). Otherwise, it attempts to move to level 2 by updating its STATUS (Line 22) and calling writer-acquire $\langle i-1 \rangle$ (). Thus, it enters DQ_{2.1} and takes the same steps as in $DQ_{3,1}$: it spins locally until W_4 changes its STATUS and it either directly enters the CS or it proceeds up to level 1. Assuming the latter, W_x enters $DQ_{1,1}$ and waits for W_1 to change its STATUS (Listing 7, Lines 10-12). If STATUS is different from MODE_CHANGE (Line 17), W_x can enter the CS. Otherwise, the lock was handed over to the readers and W_x calls set_counters_to_WRITE() to change all physical counters to the WRITE mode (Line 15), which blocks new incoming readers. At some point, the readers reach the T_R threshold and hand the lock over to W_x .

Writer Release Assume writer W_x occupies the CS and starts to release RMA-RW (Figure 2, Part 6). It begins with level 3 (Listing 5). Here, it first checks if it has a successor in DQ_{3.1} and if $T_{L,3}$ is not yet reached (Line 5). Its successor is W_{10} and assume that the latter condition is true. Then, W_x passes the lock to W_{10} by updating its STATUS so that it contains the number of lock acquires within the given element.

If $T_{L,3}$ is reached, W_x releases the lock at level 2 (Line 12). Here, it repeats all the above steps (its successor is W_6) and then starts to release the lock at level 1 (Listing 8). Here it hands the lock over to the readers if T_W is reached (Lines 5-8). Finally, it notifies its successors at each level (N to 2) to acquire the lock at the parent level (Listing 5, Line 23).

Reader Acquire A reader R_x that attempts to acquire RMA-RW first increments $c(R_x)$ (Listing 9, Line 12) and checks if T_R is reached (in the first attempt Lines 6-8 are skipped). If yes, it sets barrier (Line 14), backs off (Line 24), and reattempts to acquire the lock. In addition, if R_x is the first process to reach T_R , it also checks if there are any waiting writers (Lines 15-21). If not, it resets $c(R_x)$ and sets barrier to false so that it can enter the CS even if T_R was reached. Then, it reexecutes the main loop (Line 3); this time it may enter the loop in Lines 6-8 as the lock was handed over to a writer (if T_R was reached). In that case, R_x waits until its $c(R_x)$ is reset (Listing 9, Lines 6-8).

Reader Release This is a straightforward scenario in which R_x only increments DEPART at $c(R_x)$.

3.5 RMA-RW vs. RMA-MCS

We now outline the design of RMA-MCS; the details are in the technical report. RMA-MCS consists of DQs and DT but not DC. T_R and T_W are excluded as the are no readers. Similarly, $T_{L,1}$ is not applicable because there is no need to hand the lock to readers. The acquire/release protocols are similar to the ones in Listings 4 and 5 for any $i \in \{1, ..., N\}$.

4 CORRECTNESS ANALYSIS

We now discuss how RMA-RW ensures three fundamental correctness properties: mutual exclusion (ME), deadlock freedom (DF), and starvation freedom (SF) [21]. At the end of this section, we show how we use model checking to verify the design.

4.1 Mutual Exclusion

ME is violated if two writers or a reader and a writer enter the CS concurrently. We now discuss both cases.

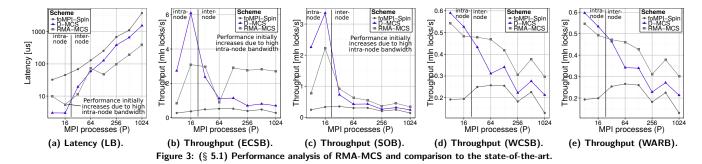
Writer & Writer: We distinguish between writers that are in the same DQ (case A) or in different ones (case B). In case A, they operate on the same TAIL. Thus, they could only violate ME if both writers do not see any predecessor. This is prevented by using FAO for atomically modifying TAIL. In case B, two writers competing in different DQs have a common DQ in DT where they or their predecessor compete for the lock. Similarly as above, the MCS lock must be acquired at each DT level. If a predecessor has to compete for the lock, a writer waits until he gets notified by its predecessor and thus does not interfere in the lock acquiring process.

Reader & Writer: A reader and a writer can be active at the same time if the lock mode is READ and about to change to WRITE. This is because the reader on its own cannot change the mode and as a consequence cannot acquire a lock while a writer is active. However, a writer can alter the mode to WRITE while a reader is active. This is prevented by a writer that checks each counter again for active readers after changing all of them.

4.2 Deadlock Freedom

Here, we also differentiate two base cases: two writers dead-lock or a reader and a writer deadlock.

Writer & Writer The only way how writers deadlock is if there is a cycle in a queue. For two writers it means that



one becomes the predecessor of the other. Therefore, both wait on the other to get notified. This cannot happen as the processes use an atomic FAO to obtain their predecessor. As explained, this function is atomic and thus we can order the uses of FAO in a timeline. This contradicts that the writers have a cycle in their waiting queue.

Reader & Writer A reader may deadlock after T_R is reached (case A) or the mode goes into WRITE (case B). In case A, either there is no writer active and the reader resets the DC or a writer is waiting and a reader backs off. Thus, the writer changes the mode to WRITE after all readers back off which is done in a finite time. As writers do not deadlock and the last writer changes the mode back to READ, no reader will deadlock in case B either.

4.3 Starvation Freedom

Finally, we show that no writer or reader can starve.

Writers A writer may starve while other writers or readers are active. We prevent it with different thresholds. First, there is $T_{L,i}$ at each DT level i. After reaching $T_{L,i}$, writers in one of the associated DQs at i release the lock to the next DQ at the same level. Thus, we only need to show that one DQ is starvation-free which is already provided by the underlying MCS queue lock design. Yet, there is the T_R threshold that regulates the number of lock acquires by readers for one counter before the readers associated to the counter back off. We already showed that the readers make progress. Thus, at some point, all counters have reached T_R and a writer changes the mode to WRITE.

Readers There are two ways how readers could starve. First, other readers are active while processes associated with a certain counter back off to let writers acquire the lock. However, there is the T_R threshold for each counter after which the readers associated with this counter back off. Thus, eventually, all readers wait on the writers to take over. This leads us to the second case where the writers have the lock and do not pass it to the waiting readers. This is not possible since there is the $T_{L,i}$ threshold at each level of the writer hierarchy and at most after $T_W = \prod_{i=1}^N T_{L,i}$ lock passings between writers the lock goes to readers; we have also already illustrated that the writers will make progress until this threshold is reached.

4.4 Model Checking

To confirm that RMA-RW provides the desired correctness properties, we also conduct model checking with SPIN [24] (v6.4.5), a software tool for the formal verification of multithreaded codes. The input to SPIN is constructed in PROMELA, a verification modeling language that allows for the dynamic creation of concurrent processes to model, for example, distributed systems. We evaluate RMA-RW for up

to $N \in \{1, ..., 4\}$ and a maximum of 256 processes. The machine elements on each level of the simulated system have the same number of children. Thus, for N=3 and four subelements per machine element, the system would consist of 4^3 processes. Each process is defined randomly either as a reader or a writer at the beginning and after that, it tries to acquire the lock 20 times. We choose this value as it generates a feasible number of cases that SPIN has to check even for a high count of processes. During the execution of a test, we use a designated process that verifies that either only one writer or multiple readers hold a lock. All the tests confirm mutual exclusion and deadlock freedom.

5 EVALUATION

We now illustrate performance advantages of RMA-MCS and RMA-RW over state-of-the-art distributed locks from the foMPI implementation of MPI-3 RMA [17].

Comparison Targets We consider D-MCS and both foMPI locking schemes: a simple spin-lock (foMPI-Spin) that enables mutual exclusion, and an RW lock (foMPI-RW) that provides both shared and exclusive accesses to the CS.

Selection of Benchmarks We conduct six series of experiments. The latency benchmark (LB) measures the latency of both acquiring and releasing a lock; an important performance metric in workloads such as real-time queries. Four other analyses obtain throughput under varying conditions and parameters. The empty-critical-section benchmark (ECSB) derives the throughput of acquiring an empty lock with no workload in the CS. The singleoperation benchmark (SOB) measures the throughput of acquiring a lock with only one single operation (one memory access) in the CS; it represents irregular parallel workloads such as graph processing with vertices protected by fine locks. Next, the workload-critical-section benchmark (WCSB) covers variable workloads in the CS: each process increments a shared counter and then spins for a random time (1-4 μ s) to simulate local computation. The wait-afterrelease benchmark (WARB) varies lock contention: after release, processes wait for a random time (1-4 μ s) before the next acquire. The throughput experiments represent dataand communication-intensive workloads. Finally, we integrate and evaluate the proposed locks with a distributed hashtable (DHT) to cover real codes such as key-value stores.

Varied Parameters To evaluate various scenarios, we vary: T_{DC} , $T_{L,i}$, and T_R . Unless stated otherwise, we set the fraction of writers $F_W = 0.2\%$ as it reflects Facebook workloads [45]; however, we also evaluate other values.

Experimentation Methodology To calculate the latency, we derive the arithmetic mean of 100,000 operations per process (for each latency benchmark). Throughput is the aggre-

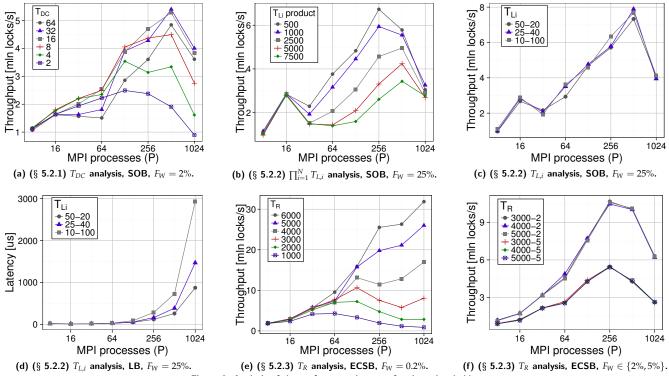


Figure 4: Analysis of the performance impact of various thresholds.

gate count of lock acquires or releases divided by the total time to run a given benchmark. 10% of the first measurements are discarded (warmup). All time measurements are taken using a high precision rdtsc timer [23].

Experimental Setup We conduct experiments on CSCS Piz Daint (Cray XC30). Each node has an 8-core HT-enabled Intel Xeon E5-2670 CPU with 32 GiB DDR3-1600 RAM. The interconnection is based on Cray's Aries and it implements the Dragonfly topology [16, 28]. The batch system is slurm 14.03.7. We use C++ and the GNU 5.2.40 g++ compiler with -O3 optimizations. The utilized Cray DMAPP is 7.0.1-1.0501.8315.8.4.ari. Unless stated otherwise, we use all the compute resources and run one MPI process per one HT resource (16 processes per one compute node).

Machine Model We consider two levels of the hierarchy: the whole machine and compute nodes, thus N=2.

Implementation Details We use the *libtopodisc* [20] library for discovering the structure of the underlying compute nodes and for obtaining MPI communicators that enable communication within each node. We group all the locking structures in MPI allocated windows to reduce the memory footprint [17].

5.1 Performance Analysis of RMA-MCS

We present the results in Figure 3. The latency of RMA-MCS is lower than any other target. For example, for P=1,024, it is $\approx 10 \mathrm{x}$ and $\approx 4 \mathrm{x}$ lower than foMPI-Spin and D-MCS, respectively. This is because foMPI-Spin entails lock contention that limits performance. In addition, both foMPI-Spin and D-MCS are topology-oblivious. Then, the throughput analysis confirms the advantages of RMA-MCS across all the considered benchmarks. The interesting spike in ECSB and SOB is because moving from P=8 to P=16 does not entail internode communication, initially increasing RMA-MCS's and

D-MCS's throughput. We conclude that RMA-MCS consistently outperforms the original foMPI design and D-MCS.

5.2 Performance Analysis of RMA-RW

We now proceed to evaluate RMA-RW. First, we analyze the impact of various design parameters (Figure 4) and then compare it to the state-of-the-art (Figure 5). Due to space constraints, we only present a subset of the results, all remaining plots follow similar performance patterns and are included in the technical report².

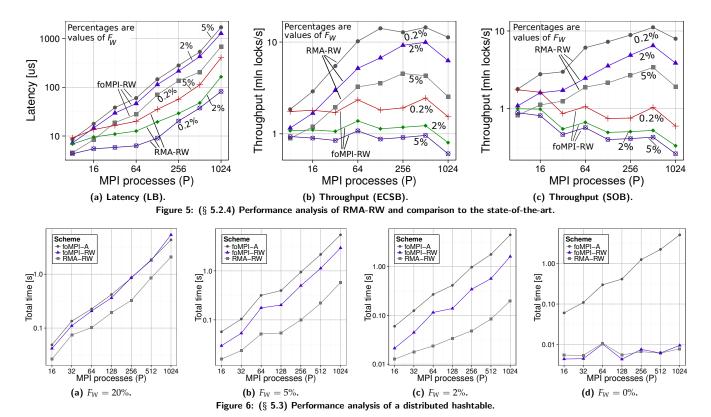
5.2.1 Influence of T_{DC}

We first discuss how different T_{DC} values impact performance. We consider $T_{DC} \in \{1, 2, 4\}$ (one physical counter on each compute node and every 2nd and 4th compute node, respectively). We also vary the number of counters on one node (1, 2, 4, 8). The results are presented in Figure 4a. First, lower T_{DC} entails more work for writers that must access more counters while changing the lock mode. This limits performance, especially for high P, because of the higher total number of counters. Larger T_{DC} increases throughput (less work for writers), but at some point (e.g., P = 512a counter on every 2nd node) the overhead due to readers (contention and higher latency) begins to dominate. We conclude that selecting the proper T_{DC} is important for high performance of RMA-RW, but the best value depends on many factors and should be tuned for a specific machine. For example, higher T_{DC} might entail unpredictable performance penalties on Cray XE because the job scheduler does not enforce contiguous job allocations [5].

5.2.2 Influence of $T_{L,i}$

Next, we analyze the performance impact of $T_{L,i}$ in the considered system $i \in \{1,2\}$. We fix $F_W = 25\%$ to en-

²http://spcl.inf.ethz.ch/Research/Parallel_Programming/RMALocks



sure that there are multiple writers per machine element on each level. We start with various $\prod_{i=1}^{N} T_{L,i}$: the maximal number of writer acquires before the lock is passed to the readers; see Figure 4b. As expected, smaller product increases throughput because more readers can enter the CS, but reduces fairness as writers wait longer. In the second step, we analyze how varying each $T_{L,i}$ impacts performance. We first fix $\prod_{i=1}^{N} T_{L,i} = 1000$. As N = 2, we use $T_{L,2} \in (10,25,50)$ and $T_{L,1} \in (100,40,20)$. The outcome is shown in Figure 4c. When more writers consecutively acquire the lock within one node (higher $T_{L,2}$), the throughput increases. Still, the differences between the considered options are small (up to 25% of the relative difference), especially for lower P. This is because of smaller amounts of inter-node communication. Interestingly, options that increase throughput (e.g., 50-20) also increase latency, see Figure 4d. We conjecture this is due to improved fairness caused by smaller $T_{L,2}$ (more processes from different nodes can acquire the lock). However, the average latency increases because other writers have to wait for a longer time.

5.2.3 Influence of T_R

Next, we analyze the impact of T_R ; see Figure 4e. We first use $F_W = 0.2\%$. The throughput for $T_R \in \{1,000 ; 2,000\}$ drops significantly for P > 512 due to the higher overhead of writers. Contrarily, increasing T_R improves the throughput significantly. This is because the latency of readers is lower than that of writers and a higher T_R entails a preference of readers. However, the larger T_R the longer the waiting time for writers is. Finally, we analyze the relationship between T_R and F_W in more detail; see Figure 4f. Here, we vary $F_W \in \{2\%, 5\%\}$. The results indicate no consistent significant advantage (<1% of relative difference for most P) of one threshold over others within a fixed F_W .

5.2.4 Comparison to the State-of-the-Art

We now present the advantages of RMA-RW over the state-of-the-art foMPI RMA library [17]; see Figure 5. Here, we consider different F_W rates. As expected, any RW distributed lock provides the highest throughput for $F_W=0.2\%$. This is because readers have a lower latency for acquiring a lock than writers and they can enter the CS in parallel. The maximum difference between the rates $F_W=0.2\%$ and $F_W=2\%$ is 1.8x and between $F_W=0.2\%$ and $F_W=5\%$ is 4.4x. We then tested other values of F_W up to 100% to find out that for $F_W>30\%$ the throughput remains approximately the same. At such rates, the throughput is dominated by the overhead of writers that enter the CS consecutively.

In each case, RMA-RW always outperforms foMPI by >6x for $P \ge 64$. One reason for this advantage is the topology-aware design. Another one is the presence of $T_{L,i}$ and T_R that prevent one type of processes to dominate the other one resulting in performance penalties.

5.3 Case Study: A Distributed Hashtable

We now illustrate how RMA-RW accelerates a distributed hashtable (DHT) that represents irregular codes. Our DHT stores 64-bit integers and it consists of parts called local volumes. Each local volume consists of a table of elements and an overflow heap for elements after hash collisions. The table and the heap are constructed with fixed-size arrays. Every local volume is managed by a different process. Inserts are based on atomic CASes. If a collision happens, the losing thread places the element in the overflow list by atomically incrementing the next free pointer. In addition, a pointer to the last element is also updated with a second CAS. Flushes are used to ensure memory consistency.

We illustrate a performance analysis in Figure 6. In the benchmark, P-1 processes access a local volume of a se-

	UPC (standard) [44]	Berkeley UPC [1]	SHMEM [4]	Fortran 2008 [27]	Linux RDMA/IB [33, 43]	iWARP [18, 41]
Put	UPC_SET	bupc_atomicX_set_RS	shmem_swap	atomic_define	MskCmpSwap	masked CmpSwap
Get	UPC_GET	<pre>bupc_atomicX_read_RS</pre>	shmem_mswap	atomic_ref	MskCmpSwap	masked CmpSwap
Accumulate	UPC_INC	<pre>bupc_atomicX_fetchadd_RS</pre>		atomic_add	FetchAdd	FetchAdd
FAO (SUM)	UPC_INC, UPC_DEC	<pre>bupc_atomicX_fetchadd_RS</pre>	shmem_fadd	atomic_add	FetchAdd	FetchAdd
FAO (REPLACE)	UPC_SET	<pre>bupc_atomicX_swap_RS</pre>	shmem_swap	atomic_define*	MskCmpSwap	masked CmpSwap
CAS	UPC_CSWAP	<pre>bupc_atomicX_cswap_RS</pre>	shmem_cswap	atomic_cas	CmpSwap	CmpSwap

Table 3: Illustration of the feasibility of using libraries/languages other than MPI RMA for RMA-MCS/RMA-RW. * indicates the lack of an atomic swap in Fortran 2008, suggesting that some of RMA-RW protocols that depend on it would have to be adjusted to a different set of available atomics.

lected process with a specified number of inserts and reads targeted at random hashtable elements. We compare the total execution time of foMPI-A (a variant that only synchronizes accesses with CAS/FAO), foMPI-RW, and RMA-RW. For $F_W \in \{2\%, 5\%, 20\%\}$ RMA-RW outperforms both the remaining variants. For $F_W = 0\%$, foMPI-RW and RMA-RW offer comparable performance.

6 DISCUSSION

Using Different RMA Libraries/Languages In our implementation, we use MPI RMA. Still, the proposed schemes are generic and can be implemented using several other existing RMA/PGAS libraries/languages that support the required operations described in Listing 1. We illustrate this in Table 3 (we omit the distinction between blocking and non-blocking operations as any type can be used in the proposed locks). The analysis indicates that RMA-MCS and RMA-RW can be used in not only traditional HPC domains (by utilizing UPC, SHMEM, or RDMA/IB), but also in TCP/IP-based settings (by using iWARP).

Selecting RMA-RW Parameters To set the parameters, we first find an appropriate value for T_{DC} . This is because our performance analysis indicates that T_{DC} has on average the highest impact on performance of both readers and writers. Here, our evaluation indicates that placing one counter per compute node results in a reasonable balance between reader throughput and writer latency. In the second step, we further influence the reader/writer performance tradeoff by manipulating with T_R and $T_{L,i}$. To reduce the parameter space, we fix T_W as indicated in Table 2. Selecting $T_{L,i}$ depends on the hardware hierarchy and would ideally incorporate several performance tests before fixing final numbers. One rule of the thumb is to reserve larger values for $T_{L,i}$ associated with components with higher inter-component communication costs, such as racks; this may reduce fairness, but increases throughput.

7 RELATED WORK

Queue-Based Locks The well-known traditional examples of this family are CLH [13, 32] and MCS [34]. Yet, they are oblivious to the memory hierarchy and cannot use this knowledge to gain performance. More recently, Radovic and Hagersten [38] proposed a hierarchical backoff lock that exploits memory locality: a thread reduces its backoff delay if another thread from the same cluster owns the lock. This increases the chance to keep the lock within the cluster, but introduces the risk of starvation. Luchangco et al. [31] improved this scheme by introducing a NUMA-aware CLH queue that ensures no starvation. Yet, it considers only two levels of the memory hierarchy. Chabbi et al. [9] generalized it to any number of memory hierarchy levels. Similarly to our scheme, they introduce an MCS lock for each level. Yet, they do not target DM machines. None of these protocols can utilize the parallelism of miscellaneous workloads where the majority of processes only read the data.

RW Locks There exist various traditional RW proposals [25,29]. Recently, Courtois et al. [12] introduced different preference schemes that favor either readers (a reader can enter the CS even if there is a writer waiting) or writers (a writer can enter the CS before waiting readers). Yet, this protocol neither prevents starvation nor scales well. Mellor-Crummey and Scott [35] extended their MCS lock to distinguish between readers and writers. This algorithm however does not scale well under heavy read contention. Next, Krieger et al. [29] use a double-linked list for more flexibility in how processes traverse the queue. Yet, there is still a single point of contention. Hsieh and Weihl [25] overcome this by trading writer throughput for reader throughput. In their design, each thread has a private mutex; the readers acquire the lock by acquiring their private mutex but the writers need to obtain all mutex objects. This introduces a massive overhead for the writers for large thread counts. Other approaches incorporate elaborate data structures like the Scalable Non-Zero Indicator (SNZI) tree [30] that traces readers in the underlying NUMA hierarchy for more locality. Yet, writers remain NUMA-oblivious. Calciu et al. [8] extend this approach with an RW lock in which both readers and writers are NUMA-aware. This design improves memory locality but it only considers two levels in a NUMA hierarchy. None of these schemes address DM environments.

Distributed Locks To the best of our knowledge, little research has been performed into locks for DM systems. Simple spin-lock protocols for implementing MPI-3 RMA synchronization were proposed by Gerstenberger et al. [17]. Some other RMA languages and libraries (e.g., UPC) also offer locks, but they are not RW, their performance is similar to that of foMPI, and they are hardware-oblivious.

We conclude that our work offers the first lock for DM systems that exploits the underlying inter-node structure and utilizes the RW parallelism present in various data- and communication-intensive workloads.

8 CONCLUSION

Large amounts of data in domains such as graph computations require distributed-memory machines for efficient processing. Such machines are characterized by weak memory models and expensive inter-node communication. These features impact the performance of topology-oblivious locks or completely prevent a straightforward adoption of existing locking schemes for shared-memory systems.

In this work, we propose a distributed topology-aware Reader-Writer (RMA-RW) and MCS lock that outperform the state-of-the-art. RMA-RW offers a modular design with three parameters that offer performance tradeoffs in selected parts of the lock. These are: higher lock fairness or better locality, larger throughput of readers or writers, and lower latency of readers or writers. This facilitates performance tuning for a specific workload or environment. RMA-RW could also be extended with adaptive schemes for a runtime

selection and tuning of the values of the parameters. This might be used in accelerating dynamic workloads.

Microbenchmark results indicate that the proposed locks outperform the state-of-the-art in both latency and throughput. Finally, RMA-RW accelerates a distributed hashtable that represents irregular workloads such as key-value stores.

Acknowledgements

This work was supported by Microsoft Research through its Swiss Joint Research Centre. We thank our shepherd Patrick G. Bridges and the anonymous reviewers for their insightful comments. We thank the CSCS team granting access to the Piz Dora and Daint machines, and for their excellent technical support. MB is supported by the 2013 Google European Doctoral Fellowship in Parallel Computing.

9 References

- [1] Berkeley UPC User's Guide version 2.22.0. http://upc.lbl.gov/docs/user/.
- [2] T. E. Anderson. The performance of spin lock alternatives for shared-memory multiprocessors. *IEEE Trans. Parallel Distrib. Syst.*, 1(1):6–16, Jan. 1990.
- [3] B. Arimilli et al. The PERCS High-Performance Interconnect. In Proc. of the IEEE Symp. on High Perf. Inter., HOTI '10, pages 75–82, 2010.
- [4] R. Barriuso and A. Knies. SHMEM user's guide for C, 1994.
- [5] A. Bhatele et al. There goes the neighborhood: performance degradation due to nearby jobs. In *Proc. of the ACM/IEEE Supercomputing*, page 41. ACM, 2013.
- [6] C. Bienia. Benchmarking Modern Multiprocessors. PhD thesis, Princeton University, January 2011.
- [7] N. Bronson et al. TAO: Facebook's Distributed Data Store for the Social Graph. In USENIX Annual Technical Conference, pages 49–60, 2013.
- [8] I. Calciu et al. NUMA-aware Reader-writer Locks. In Proc. of the ACM Symp. on Prin. and Prac. of Par. Prog., PPoPP '13, pages 157–166, 2013.
- [9] M. Chabbi, M. Fagan, and J. Mellor-Crummey. High Performance Locks for Multi-level NUMA Systems. In Proc. of the ACM Symp. on Prin. and Prac. of Par. Prog., PPoPP 2015, pages 215–226, 2015.
- [10] B. Chamberlain, S. Deitz, M. B. Hribar, and W. Wong. Chapel. Technical report, Cray Inc., 2005.
- [11] P. Charles et al. X10: an Object-Oriented Approach to Non-Uniform Cluster Computing. SIGPLAN Not., 40(10):519–538, Oct. 2005.
- [12] P. J. Courtois, F. Heymans, and D. L. Parnas. Concurrent control with "readers" and "writers". Commun. ACM, 14(10):667–668, Oct. 1971.
- [13] T. S. Craig. Building FIFO and Priority-Queuing Spin Locks from Atomic Swap. Technical report, 1993.
- [14] D. Dice, V. J. Marathe, and N. Shavit. Flat-combining NUMA Locks. In *Proc. of the ACM Symp. on Par. in Alg. and Arch.*, SPAA '11, pages 65–74, 2011.
- [15] D. Dice, V. J. Marathe, and N. Shavit. Lock Cohorting: A General Technique for Designing NUMA Locks. In Proc. of the ACM Symp. on Prin. and Prac. of Par. Prog., PPoPP '12, pages 247–256, 2012.
- [16] G. Faanes et al. Cray cascade: a scalable HPC system based on a Dragonfly network. In *Proc. of the ACM/IEEE* Supercomputing, page 103, 2012.
- [17] R. Gerstenberger, M. Besta, and T. Hoefler. Enabling Highly-scalable Remote Memory Access Programming with MPI-3 One Sided. In *Proc. of ACM/IEEE Supercomputing*, SC '13, pages 53:1–53:12, 2013.
- [18] R. Grant, M. Rashti, A. Afsahi, and P. Balaji. RDMA Capable iWARP over Datagrams. In Par. Dist. Proc. Symp. (IPDPS), 2011 IEEE Intl., pages 628–639, 2011.

- [19] W. Gropp, T. Hoefler, R. Thakur, and E. Lusk. Using Advanced MPI: Modern Features of the Message-Passing Interface. MIT Press, 2014.
- [20] W. D. Gropp. Personal exchange, 2013.
- [21] M. Herlihy and N. Shavit. The Art of Multiprocessor Programming. Morgan Kaufmann Publishers Inc., 2008.
- [22] P. N. Hilfinger et al. Titanium Language Reference Manual, version 2.19. Technical report, UC Berkeley Tech Rep. UCB/EECS-2005-15, 2005.
- [23] T. Hoefler et al. Netgauge: A Network Performance Measurement Framework. In *Proc. of High Perf. Comp. and Comm.*, *HPCC'07*, volume 4782, pages 659–671, 2007.
- [24] G. J. Holzmann. The Model Checker SPIN. *IEEE Trans. Softw. Eng.*, 23(5):279–295, May 1997.
- [25] W. C. Hsieh and W. W. Weihl. Scalable reader-writer locks for parallel systems. In *Proc. of Par. Proc. Symp.*, pages 656–659, Mar 1992.
- [26] InfiniBand Trade Association. Supplement to InfiniBand Architecture Spec., Vol. 1, Rel. 1.2.1. Annex A16: RDMA over Converged Ethernet (RoCE). 2010.
- [27] ISO Fortran Committee. Fortran 2008 Standard (ISO/IEC 1539-1:2010). 2010.
- [28] J. Kim, W. J. Dally, S. Scott, and D. Abts. Technology-driven, highly-scalable dragonfly topology. In ACM SIGARCH Comp. Arch. News, volume 36, pages 77–88, 2008.
- [29] O. Krieger, M. Stumm, R. Unrau, and J. Hanna. A fair fast scalable reader-writer lock. In *In Proc. of the Intl. Conf. on Par. Proc.*, pages 201–204, 1993.
- [30] Y. Lev, V. Luchangco, and M. Olszewski. Scalable reader-writer locks. In *Proc. of the Symp. on Par. in Alg. and Arch.*, SPAA '09, pages 101–110, 2009.
- [31] V. Luchangco, D. Nussbaum, and N. Shavit. A Hierarchical CLH Queue Lock. In W. Nagel, W. Walter, and W. Lehner, editors, Euro-Par 2006 Par. Proc., volume 4128 of Lecture Notes in Computer Science, pages 801–810. 2006.
- [32] P. S. Magnusson, A. Landin, and E. Hagersten. Queue Locks on Cache Coherent Multiprocessors. In *Proc. of the Intl. Symp.* on Par. Proc., pages 165–171, 1994.
- [33] Mellanox Technologies. Mellanox OFED for Linux User Manual, 2015.
- [34] J. M. Mellor-Crummey and M. L. Scott. Algorithms for scalable synchronization on shared-memory multiprocessors. *ACM Trans. Comput. Syst.*, 9(1):21–65, Feb. 1991.
- [35] J. M. Mellor-Crummey and M. L. Scott. Scalable reader-writer synchronization for shared-memory multiprocessors. In *Proc.* of the ACM SIGPLAN Symp. on Prin. and Prac. of Par. Prog., PPOPP '91, pages 106–113, 1991.
- [36] MPI Forum. MPI: A Message-Passing Interface Standard. Ver. 3, 2012.
- [37] R. C. Murphy, K. B. Wheeler, B. W. Barrett, and J. A. Ang. Introducing the graph 500. Cray User's Group (CUG), 2010.
- [38] Z. Radovic and E. Hagersten. Hierarchical backoff locks for nonuniform communication architectures. In Proc. of the Intl. Symp. on High-Perf. Comp. Arch., HPCA '03, pages 241–, 2003.
- [39] R. Recio et al. A remote direct memory access protocol specification, Oct 2007. RFC 5040.
- [40] M. L. Scott and W. N. Scherer. Scalable Queue-based Spin Locks with Timeout. In Proc. of the ACM SIGPLAN Symp. on Prin. and Prac. of Par. Prog., PPoPP '01, pages 44–52, 2001.
- [41] R. Sharp et al. Remote Direct Memory Access (RDMA) Protocol Extensions. 2014.
- [42] H. Takada and K. Sakamura. Predictable spin lock algorithms with preemption. In Real-Time Operating Systems and Software. RTOSS '94, Proc., IEEE Workshop on, pages 2–6, 1994.
- [43] The InfiniBand Trade Association. *Infiniband Architecture Spec.* Vol. 1-2, Rel. 1.3. InfiniBand Trade Association, 2004.
- [44] UPC Consortium. UPC language spec., v1.3. Technical report, Lawrence Berkeley National Laboratory, 20013. LBNL-6623E.
- [45] V. Venkataramani et al. TAO: How Facebook Serves the Social Graph. In *Proc. of the ACM Intl. Conf. on Manag. of Data*, SIGMOD '12, pages 791–792, 2012.