

A Verified High-Performance Composable Object Library for Remote Direct Memory Access

GUILLAUME AMBAL*, Imperial College London, UK

GEORGE HODGKINS*, University of Colorado, Boulder, USA

MARK MADLER, University of Colorado, Boulder, USA

GREGORY CHOCKLER, University of Surrey, UK

BRIJESH DONGOL, University of Surrey, UK

JOSEPH IZRAELEVITZ, University of Colorado, Boulder, USA

AZALEA RAAD, Imperial College London, UK

VIKTOR VAFEIADIS, MPI-SWS, Germany

Remote Direct Memory Access (RDMA) is a memory technology that allows remote devices to directly write to and read from each other's memory, bypassing components such as the CPU and operating system. This enables low-latency high-throughput networking, as required for many modern data centres, HPC applications and AI/ML workloads. However, baseline RDMA comprises a highly permissive weak memory model that is difficult to use in practice and has only recently been formalised.

In this paper, we introduce the Library of Composable Objects (LOCO), a formally verified library for building multi-node objects on RDMA, filling the gap between shared memory and distributed system programming. LOCO objects are well-encapsulated and take advantage of the strong locality and the weak consistency characteristics of RDMA. They have performance comparable to custom RDMA systems (e.g. distributed maps), but with a far simpler programming model amenable to formal proofs of correctness.

To support verification, we develop a novel modular declarative verification framework, called Mowgli, that is flexible enough to model multinode objects and is independent of a memory consistency model. We instantiate Mowgli with the RDMA memory model, and use it to verify correctness of LOCO libraries.

CCS Concepts: • **Theory of computation** → **Axiomatic semantics**; **Distributed computing models**; • **Software and its engineering** → **Distributed programming languages**.

Additional Key Words and Phrases: RDMA, Distributed Computing, Declarative Semantics, Verification

ACM Reference Format:

Guillaume Ambal, George Hodgkins, Mark Madler, Gregory Chockler, Brijesh Dongol, Joseph Izraelevitz, Azalea Raad, and Viktor Vafeiadis. 2026. A Verified High-Performance Composable Object Library for Remote Direct Memory Access. *Proc. ACM Program. Lang.* 10, POPL, Article 71 (January 2026), 32 pages. <https://doi.org/10.1145/3776713>

*co-first authors.

Authors' Contact Information: [Guillaume Ambal](mailto:g.ambal@imperial.ac.uk), Imperial College London, UK, g.ambal@imperial.ac.uk; [George Hodgkins](mailto:George.Hodgkins@colorado.edu), University of Colorado, Boulder, USA, George.Hodgkins@colorado.edu; [Mark Madler](mailto:Mark.Madler@colorado.edu), University of Colorado, Boulder, USA, Mark.Madler@colorado.edu; [Gregory Chockler](mailto:Gregory.Chockler@surrey.ac.uk), University of Surrey, UK, g.chockler@surrey.ac.uk; [Brijesh Dongol](mailto:Brijesh.Dongol@surrey.ac.uk), University of Surrey, UK, b.dongol@surrey.ac.uk; [Joseph Izraelevitz](mailto:Joseph.Izraelevitz@colorado.edu), University of Colorado, Boulder, USA, Joseph.Izraelevitz@colorado.edu; [Azalea Raad](mailto:Azalea.Raad@imperial.ac.uk), Imperial College London, UK, azalea.raad@imperial.ac.uk; [Viktor Vafeiadis](mailto:Viktor.Vafeiadis@mpi-sws.org), MPI-SWS, Germany, viktor@mpi-sws.org.



This work is licensed under a Creative Commons Attribution 4.0 International License.

© 2026 Copyright held by the owner/author(s).

ACM 2475-1421/2026/1-ART71

<https://doi.org/10.1145/3776713>

1 Introduction

The *remote direct memory access* (RDMA) protocol provides a load/store interface, allowing a machine to access the memory of a remote machine across a network without communicating with the remote processor. The memory accesses are performed directly by the *network interface card* (NIC), bypassing the software networking stack on both ends of the connection. As such, RDMA achieves low-latency, high-throughput communication, making it a key technology in many production-grade data centres such as those at Microsoft [Zhu et al. 2015], Google [Lu et al. 2018], Alibaba [Wang et al. 2023b], and Meta [Gangidi et al. 2024].

Despite its memory-like interface, RDMA is a hardware-accelerated networking protocol, and has traditionally been programmed as such—not as shared memory. This has resulted in a very weak memory model with out-of-order behaviours visible even in a sequential setting [Ambal et al. 2024]. Consider, for example, the following program, where all memories are zero-initialised.

```
 $\bar{z} := x;$  // RDMA put: write the value of local variable  $x$  to remote location  $z$   
 $x := 1$  // update local variable  $x$  to 1
```

Somewhat counterintuitively, this program can result in z getting the value 1, with the following execution steps: (1) the put instruction ($\bar{z} := x$) is offloaded to the NIC; (2) the CPU executes $x := 1$ updating the value of x in the local memory; and (3) the NIC executes the put instruction, fetching the *new* value of x from local memory before performing the remote write.

Since programming RDMA directly is challenging, prior work has developed custom RDMA libraries. Most existing libraries are monolithic: they encapsulate a useful distributed protocol (such as consensus [Aguilera et al. 2020] or distributed storage [Dragojević et al. 2014; Wang et al. 2022]) as a single, global entity—not one that can be reused by other RDMA libraries. Some other libraries (e.g. [Cai et al. 2018; Wang et al. 2020]) provide a simple high-level memory abstraction that hides all the complexities of a highly non-uniform, weakly consistent network memory, but also loses a lot of the performance that can be achieved by knowing the system layout [Liu and Mellor-Crummey 2014; Majo and Gross 2017; Tang et al. 2013]. Other intermediate layers, such as MPI [Message Passing Interface Forum 2023] or NCCL [NVIDIA Corporation 2020] are designed explicitly for networks and present a message passing interface that is ideal for embarrassingly parallel or task-oriented workflows, but ill-suited for irregular and data-dependent workloads, such as data stores or stateful transactional systems, for which shared-memory solutions excel [Liu et al. 2021]. Although these library implementations are impressive engineering artefacts and have often been carefully tuned to achieve very good performance, they are almost impossible to verify formally due to their lack of modularity.

In this paper, we argue for a new way for programming RDMA applications—and more generally systems with non-uniform weakly consistent memories—with *flexible* libraries that can expose the non-uniform memory aspects and that support formal verification. Key to our approach is *composability*—namely, the ability to put together smaller/simpler objects to build larger ones—and this composability is reflected both in the design and implementation of our library as well as in the formal proofs about its correctness.

LOCO. As a first contribution, we introduce the *Library of Composable Objects* (LOCO). A LOCO object is a concurrent object as in Herlihy and Wing [1990], exposing a collection of methods, but storing its state in a distributed fashion across all participating nodes. Familiar examples include cross-node locks, barriers, queues, and maps. LOCO objects provide encapsulation and can be composed together to build other LOCO objects. We define objects encapsulating the underlying RDMA operations and the local CPU instructions, and use them to build intermediate objects, such as ring buffers, which in turn are used to build larger objects, such as a key-value store (see Fig. 1).

The source code for LOCO is available on Github [9]. Additionally, a previous preprint version of this paper focused on the LOCO library was published on arXiv [Hodgkins et al. 2025].

For concreteness, we implement and verify LOCO objects over RDMA^{TSO} (which combines an RDMA networking fabric with Intel x86-TSO nodes), making use of an existing formalisation by Ambal et al. [2024]. RDMA^{TSO} is, however, too low-level for our purposes: it does not provide a compositional way to wait for RDMA operations to complete, making it impossible to encapsulate it as a LOCO object. For this reason, we introduce $\text{RDMA}^{\text{WAIT}}$, a thin layer of abstraction over RDMA^{TSO} , that attaches identifiers to RDMA operations and allows threads to perform a `Wait` operation to wait for all RDMA operations with a given identifier to finish executing. To attain good performance, the practical implementation of the `Wait` operation in LOCO is quite involved. Nevertheless, we prove the correctness of a simplified version over the underlying RDMA^{TSO} model.

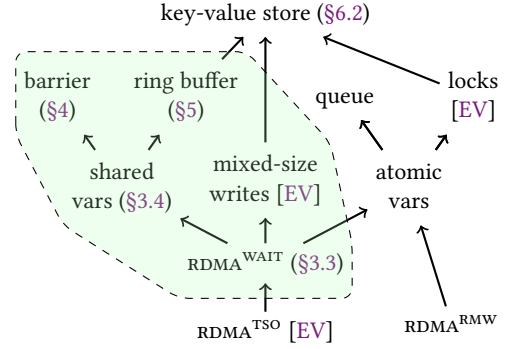


Fig. 1. LOCO libraries with their dependencies

MOWGLI. As a second contribution, we introduce MOWGLI (MODular Weak Graph-based Libraries), a generic, modular framework for modelling and verifying weak libraries. MOWGLI is *generic* in that it makes no assumptions about the underlying memory model (e.g. RDMA or TSO) in its core theory; and it is *modular* in that it allows proof decomposition at library interfaces and reasoning about individual components without referring to the internals of other components.

We instantiate MOWGLI with $\text{RDMA}^{\text{WAIT}}$ and establish the correctness of all the LOCO libraries that do not depend on atomic read-modify-write (RMW) RDMA operations because the latter are not covered by the existing RDMA^{TSO} model. The verified libraries are highlighted in Fig. 1.

MOWGLI represents program executions as graphs, whose nodes are called *events* and represent either a simple operation like a read or a write, or a more complex operation such as a method call. Following the declarative approach of Raad et al. [2019] and Stefanescu et al. [2024], we specify each concurrent object with a set of axioms (i.e., consistency predicates) over events. As we shall see in §2.3, however, events are too coarse-grained to model the intricate synchronisation guarantees of RDMA operations.

We therefore introduce the novel notion of a *subevent*, allowing one to split complex library operations into multiple subevents, each with a different *stamp* (representing, e.g., the node affected by the subevent). Stamps are meta-categories of behaviours, shared by all libraries, and are independent from programs. Stamps are then used to induce ordering among (sub)events. Within a thread, they are used to define the *preserved program order* (ppo) [Alglave et al. 2014], which relates (sub)events executed by a thread that may not be reordered. Across threads and nodes, stamps are used to define the *synchronisation order* (so) [Dongol et al. 2018] between methods calls of the same library. Together ppo and so are used to define the happens-before relation.

Our main result supporting modular proofs in MOWGLI is a new locality result that decomposes proving correctness of a system into proofs about the correctness of its individual components. This is akin to the notion of compositionality for linearisability [Herlihy and Wing 1990], but generalised to a partially ordered setting. In our verification of LOCO, this allows us to verify a library, then use the *specification* of the library in any program that uses the library. Moreover, we show that our locality result supports both *horizontal composition*, where a library is used within a

$x=0$	$z=0$
$\bar{z} := x$	
$\text{Poll}(2)$	
$x := 1$	

(a) $z=0$ ✓ $z=1$ ✗

$x=0$	$z=0$
$\bar{z} := x$	
$\bar{z} := x$	
$\text{Poll}(2)$	
$x := 1$	

(b) $z=0$ ✓ $z=1$ ✓

$x=0$	$z=0$
$\bar{z} := x$	
$\bar{z} := x$	
$\text{Poll}(2)$	
$\text{Poll}(2)$	
$x := 1$	

(c) $z=0$ ✓ $z=1$ ✗

$x=0$	$z=0$
$\bar{z} :=^d x$	
$\text{Wait}(d)$	
$x := 1$	

(a) $z=0$ ✓ $z=1$ ✗

$x=0$	$z=0$
$\bar{z} :=^e x$	
$\bar{z} :=^d x$	
$\text{Wait}(d)$	
$x := 1$	

(b) $z=0$ ✓ $z=1$ ✗Fig. 2. Polling under RDMA^{TSO} Fig. 3. Waiting under $\text{RDMA}^{\text{WAIT}}$

client program, and *vertical composition*, where a library is developed from other libraries via a series of abstractions.

Contributions. In summary, we make the following contributions:

- We define a new consistency model, $\text{RDMA}^{\text{WAIT}}$, that supports a `Wait` operation that allows CPUs to wait for the confirmation (by the NIC) for a *specific* group of remote operations. We verify the correctness of the $\text{RDMA}^{\text{WAIT}}$ implementation over the existing RDMA^{TSO} model.
- We develop LOCO [10], a flexible, modular object library for RDMA, and demonstrate its compositionality by using simpler objects to build more advanced objects: e.g., a barrier, a ring buffer, a linearisable key-value store, a transactional locking scheme, and a distributed DC/DC converter.
- We introduce a new modular formal framework, MOWGLI, for specifying and verifying concurrent libraries over weakly consistent memory and distributed architectures.
- We instantiate MOWGLI to verify correctness of the aforementioned LOCO libraries.
- We benchmark LOCO's barrier and ring buffer objects and show that they outperform the highly-tuned OpenMPI implementations of the same objects.

2 Overview of LOCO and MOWGLI

In this section, we provide an informal, more detailed overview of LOCO and MOWGLI. We present LOCO's base memory model, $\text{RDMA}^{\text{WAIT}}$, in §2.1, then discuss the key libraries that we consider. In §2.3, we provide an overview of our MOWGLI verification framework.

2.1 The $\text{RDMA}^{\text{WAIT}}$ Memory Model

We start by informally describing LOCO's base memory model, $\text{RDMA}^{\text{WAIT}}$, and contrast it to RDMA^{TSO} [Ambal et al. 2024] via a set of simple examples. Both models provide put operations ($\bar{x} := y$) for writing to remote memory and get operations ($y := \bar{x}$) for reading from remote memory, which are executed asynchronously. The models differ in how a thread can wait for these asynchronous operations to terminate.

In RDMA^{TSO} , waiting is achieved with the `Poll` primitive. Consider the programs in Fig. 2, which comprise two nodes, with a variable x in node 1 and a variable z in node 2. In the first program, Fig. 2a, node 1 comprises a single thread that first puts the value of x to the remote location z (located in node 2), and then polls node 2, which causes the thread to wait until the put has been executed, and finally updates x to 1. This means that the final value of z is 0, and not 1. Note that in the absence of the `Poll` operation, the final outcome $z = 1$ would be permitted since the instruction $\bar{z} := x$ could simply be offloaded to the NIC, followed by the update of x to 1. When $\bar{z} := x$ is later executed by the NIC, it will load the value 1 for x .

Synchronisation via `Poll` is however brittle, and sensitive to the number of instructions occurring before the `Poll`. For example, as shown in Fig. 2b, the final outcome $z = 1$ is once again permitted because the `Poll` only waits for the earliest unpollled operation to be executed at node 2. In particular,

although Poll does wait for the first put instruction, the second put may be offloaded to the NIC and the local write $x := 1$ executed before the second put ($\bar{z} := x$) is executed. This weak behaviour is also allowed if we replace the first operation with *any* RDMA operation, even unrelated to locations x and z . This demonstrates that RDMA^{TSO} programs are *not compositional*: we cannot reason about a property (e.g. the final value of z) by focusing only on the part of the program that seems relevant; only monolithic analysis of the full program is possible. To prevent the weak behaviour of Fig. 2b, one must add a second Poll operation as shown in Fig. 2c.

In $\text{RDMA}^{\text{WAIT}}$, synchronisation is performed with the Wait operation. RDMA operations are associated with a work identifier, e.g. d in Fig. 3a, which can be waited upon with a Wait operation. Thus, unlike Poll, which waits for the first unpolled operation, $\text{RDMA}^{\text{WAIT}}$ can wait for a specific put or get operation. This improves robustness since the Wait is independent of the number of instructions that have been executed by each thread. For example, in Fig. 3b, the Wait can target the *second* put instruction using the work identifier d and exclude the unintended outcome $z = 1$.

While Wait makes targeting a remote operation easier, it does not provide more synchronisation guarantees than the Poll operation. Waiting for a put operation ($\bar{z} := x$) only guarantees that the local value of x has been read, not that the remote location z has been modified. Thus, as shown in Fig. 4, the store buffering behaviour across nodes is possible even if we wait for every remote operation. In contrast, waiting for a get operation ($x := \bar{z}$) does guarantee it has fully completed, i.e. that z has been read and x modified. This can be exploited to prevent the store buffering behaviour. RDMA ordering rules ensure that later gets execute after previous puts towards the same remote node. Thus, waiting for a (seemingly unrelated) get operation can be used to ascertain the completion of previous remote writes.

$y=0$	$x=0$
$\bar{x} :=^d 1$	$\bar{y} :=^e 1$
Wait(d)	Wait(e)
$a := y$	$b := x$
$(a, b) = (0, 0) \checkmark$	

$y, w=0$	$x, z=0$
$\bar{x} := 1$	$\bar{y} := 1$
$c :=^d \bar{z}$	$d :=^e \bar{w}$
Wait(d)	Wait(e)
$a := y$	$b := x$
$(a, b) = (0, 0) \times$	

Fig. 4. Preventing RDMA store buffering

We present the formal definitions of $\text{RDMA}^{\text{WAIT}}$ in §3.3 using a declarative style. Although, like RDMA^{TSO} , it is also possible to derive an equivalent operational model, we elide these details since the proof technique that we use (see §2.3) directly uses the declarative semantics.

Note that the actual implementation of Wait in LOCO is non-trivial, relying on a highly optimised code path to track outstanding operations and match them to an associated Wait. This extension to the RDMA interface is described within Section 2.2, with a more complete treatment in the extended version [EV].

2.2 LOCO Libraries

LOCO provides a set of commonly used distributed objects, which we call *channels*, built on top of $\text{RDMA}^{\text{WAIT}}$. Channels are *named* and *composable*. To communicate over a channel, each participating node constructs a local channel object, or *channel endpoint*, with the same name. Each channel endpoint allocates zero or more named local regions of network memory when constructed, and delivers the metadata necessary to access these local memory regions to the other endpoints during the setup process.

Channels make it easy to develop RDMA applications and prove their correctness, for minimal performance loss. A LOCO application will usually consist of many channels (objects) of many different channel types (classes). In addition, each channel can itself instantiate member sub-channels. For instance, a key-value store might include several mutexes as sub-channels to synchronise access to its contents.

Shared Variable Library (sv, §3.4). One of the most basic components of LOCO is the *shared variable* library. Each shared variable is replicated across all (participating) nodes in the network and supports Write_{sv} and Read_{sv} operations, which only access the *local* copy of the variable. Any updates to the variable may be pushed to the other replicas by the modifying node with a Bcast_{sv} operation.¹ We provide examples in §2.3, Fig. 9.

$y = 0$	$x = 0$
$\bar{x} := 1$	$\bar{y} := 1$
$\text{GF}_{\text{sv}}(2)$	$\text{GF}_{\text{sv}}(1)$
$a := y$	$b := x$
$(a, b) = (0, 0)$ ✗	

Fig. 5. Using GF_{sv}

The shared variable library also provides a mechanism for synchronising different nodes using a *global fence* (GF_{sv}) operation. GF_{sv} takes the node(s) on which the fence should be performed as a parameter and causes the executing thread to wait until all prior operations executed by the thread towards the given nodes have fully completed. This is stronger than using the Wait primitive, as the global fence also ensures the remote write parts have completed. An example program using a GF_{sv} is the store buffering setting given in Fig. 5, which disallows the final outcome $(a, b) = (0, 0)$, but allows all other combinations for a and b with values from $\{0, 1\}$. As can be guessed from the similarity with Fig. 4, this global fence can be implemented by submitting get operations and waiting for them.

Barrier Library (BAL, §4). A commonly used object in distributed systems is a barrier, which provides a stronger synchronisation guarantee than global fences. All threads synchronising on a barrier must finish their operations before execution continues. For example, consider the program in Fig. 6, which only allows the final outcome $(a, b) = (1, 1)$ and forbids all other outcomes. Here, nodes 1 and 2 synchronise on the barrier z , and hence nodes 1 and 2 both wait until both writes to x and y have completed.

$y = 0$	$x = 0$
$\bar{x} := 1$	$\bar{y} := 1$
$\text{BAR}_{\text{BAL}}(z)$	$\text{BAR}_{\text{BAL}}(z)$
$a := y$	$b := x$
$(a, b) = (1, 1)$ ✓	

Fig. 6. Using BAR_{BAL}

Ring Buffer Library (RBL, §5). Similarly useful is a ring buffer, which allows one to develop producer-consumer systems. LOCO's ring buffer supports a one-to-many broadcast, and is the most sophisticated of the libraries that we consider.

Mixed-Size Writes (msw, in extended version).

The final library we consider is the mixed-size write library, which allows safe transmission of data spanning multiple words. Here, due to the asynchrony between the CPU and the NIC, it is possible for corrupted data to be transmitted that does not correspond to any write performed by the CPU. There are multiple solutions to this problem; we consider a simple solution that transmits a hash alongside the data.

LOCO API Example. As an example of the LOCO C++ API, Fig. 7 shows our implementation of a barrier object, based on Gupta et al. [2002]. The class uses an array (arr) of shared variables as a sub-object [Jha et al. 2019, 2017], demonstrating composition. As with a traditional shared memory barrier, it is used to synchronise all participants at a certain point in execution. For each use of the

```

1 class barrier : public loco::channel {
2   unsigned count;
3   loco::var_array<unsigned> arr;
4   public:
5   void waiting() {
6     // complete outstanding RDMA ops
7     loco::global_fence();
8     count++; // increment our counter
9     arr[loco::my_node()].store(count);
10    arr[loco::my_node()].push_broadcast
11    (); //and push
12    bool waiting = true;
13    while(waiting){ // wait for others
14      waiting = false; // to match
15      for (auto& i : arr) {
16        if (i.load() < count){
17          waiting = true;
18          break;}
19    } } } };
```

Fig. 7. Complete C++ code for the LOCO barrier

¹It is also possible for replicas to pull the new value from a source node when a shared variable is modified, but we do not model this aspect because it is not used in the libraries we consider. Moreover, LOCO also defines a stronger form of a shared variable called an *owned variable*, which provides a mechanism for defining a variable's owner that provides a single authoritative version of the variable (describing its true value), defining a single-writer multi-reader register.

barrier, participants increment their local count variable, then broadcast the new value to others using their index in the array. They then wait locally, leaving the barrier only when all participants have a count in the array not less than their own. This code is a near-complete implementation of a single-threaded barrier in LOCO, missing only a boilerplate constructor.

Implementing $\text{RDMA}^{\text{WAIT}}$. In general, RDMA operations are assigned a unique ID at initialisation. Subsequent queries to a corresponding *completion queue* (i.e. via the `Poll` operation) indicate the oldest ID that has been received at the remote node and acknowledged. As mentioned in Section 2.1, this default system results in non-local effects.

In contrast, LOCO's backend allows for a practical implementation of $\text{RDMA}^{\text{WAIT}}$ with a high-performance and composable system for tracking RDMA operations. LOCO uses a dedicated *polling thread* to query the completion queue and notify the application of tracked RDMA operations. If the application wishes to monitor the progress of a single RDMA operation (or a set of them, e.g. for a broadcast to all remote nodes), it creates a special `ack_key` object with the associated operation IDs. In $\text{RDMA}^{\text{WAIT}}$, `ack_key` objects are abstracted by work identifiers (see §2.1). When all associated IDs have been dequeued by the polling thread, the `ack_key` object is marked completed. The `ack_key` object exports methods to check its status, i.e. the `Wait` operation simply looks for a completed status. Communication between the application and the polling thread for outstanding operation IDs is managed via a single-writer, multiple-reader lock-free queue [Morrison and Afek 2013]. A full description of this system can be found in the extended version [EV].

Additional LOCO Libraries. In addition to the proven libraries that are the focus of this paper, LOCO [Hodgkins et al. 2025] contains a number of additional objects that rely on an RDMA read-modify-write primitive currently missing from the formalisation provided in RDMA^{TSO} . These include an atomic variable library for accessing these operations, a set of locks (both ticket and test-and-set with optional local flat-combining [Hendler et al. 2010]), and a shared FiFo queue porting the cyclic ring queue [Morrison and Afek 2013]. We intend to fully prove the correctness of these libraries in future work.

LOCO-Based Applications. As mentioned earlier, LOCO enables one to quickly build distributed applications. We demonstrate this by using LOCO to construct a linearisable key-value store (§6.2), a transactional locking scheme [EV], and a distributed DC/DC converter [EV]. Additional obvious targets for LOCO include distributed shared memories [Kaxiras et al. 2015; Keleher et al. 1994], distributed communication collectives [Graham et al. 2006], and other HPC communication library backends (e.g. global arrays [Nieplocha et al. 1994; Zheng et al. 2014]).

2.3 Towards a Modular Verification Framework for LOCO

To support reasoning about LOCO libraries, we develop a modular verification framework for $\text{RDMA}^{\text{WAIT}}$ programs. Our point of departure is the Yacovet framework [Raad et al. 2019; Stefanescu et al. 2024] that was used to reason about weak shared memory *within* a single node. Yacovet, however, is not expressive enough to model $\text{RDMA}^{\text{WAIT}}$ programs, and so we need to develop a framework that can take into account both sources of weak consistency: shared-memory concurrency (TSO) and distribution (RDMA). This poses three main challenges.

Lack of Causality. $\text{RDMA}^{\text{WAIT}}$ assumes the TSO memory model [Alglave et al. 2014; Owens et al. 2009] for each CPU within each node. This means that well-known effects such as store buffering (see Fig. 8a) are possible, where both reads in the two threads read from the initial state. Despite this weakness, TSO guarantees causal consistency, i.e. message passing (see Fig. 8b), where the right thread reading the new value 1 for y guarantees that it also reads 1 for x . Formally, this is due to a relation known as *preserved program order* (ppo) between the read of w , the write of this value to x ,

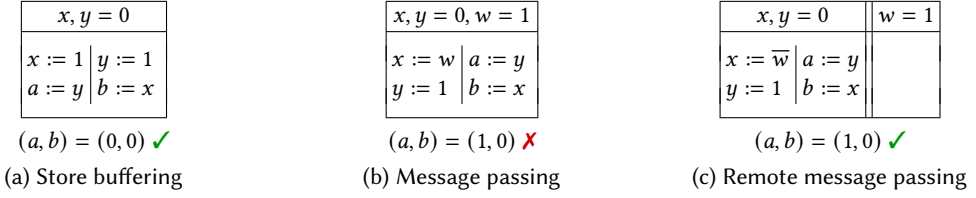
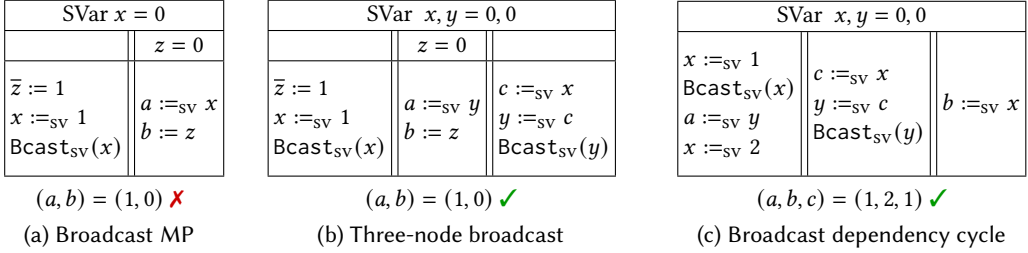
Fig. 8. TSO effects of $\text{RDMA}^{\text{WAIT}}$ 

Fig. 9. Broadcast synchronisation

and the write to y . However, under $\text{RDMA}^{\text{WAIT}}$, when interacting with the NIC, causal consistency is no longer guaranteed (see Fig. 8c). This leads to our first modelling challenge: $\text{RDMA}^{\text{WAIT}}$ has a much weaker **ppo** relation than TSO [Alglave et al. 2014]. Here, compositionality is critical to ensure proofs for scalability; we offer this through our locality result (Theorem 3.14).

Fine-Grained Synchronisation. A second challenge in specifying RDMA libraries is that the *same* method call may interact with different library methods in different ways. To make this problem concrete, consider a version of message passing in Fig. 9a, where node 1 updates the remote variable z (located in node 2), and then broadcasts a new value of a shared variable x to signal that the remote value has changed. In Fig. 9a, when node 2 sees the new value of x , it means that the (earlier) write to z must have also taken effect. To represent this, we require that $\bar{z} := 1$ happens before $\text{Bcast}_{\text{sv}}(x)$ and that $\text{Bcast}_{\text{sv}}(x)$ happens before $a :=_{\text{sv}} x$. These orders *must* be part of the declarative semantics, in some shape or form, to disallow the behaviour $(a, b) = (1, 0)$.

However, naively specifying broadcast in this way is problematic. Consider the example in Fig. 9b, where node 1 behaves as before, but the “signal variable” x is picked up by node 3 and a new signal using y is broadcast by node 3. This time, when node 2 receives the signal on y (i.e. $a = 1$), there is actually no guarantee that the write on z has completed. The outcome $(a, b) = (1, 0)$ is allowed, as communication between each pair of nodes is independent. Thus we *must not* have a happens-before dependency between the write to z (from node 1) and the read on z .

For an even more precarious example, consider Fig. 9c, which is a possible behaviour of LOCO’s broadcast library. The final outcome $(a, b, c) = (1, 2, 1)$ is only possible if node 1 broadcasts $x = 1$ to node 2, and $x = 2$ to node 3 with a single broadcast. The broadcast is allowed to pick up the later value 2 since the CPU might run the command $x :=_{\text{sv}} 2$ before the NIC reads the value of x . As mentioned above, reading the result of a broadcast *must* create happens-before order so that we can preclude behaviours like in Fig. 9a. In this example, we thus need a sequence of dependencies: $x :=_{\text{sv}} 1 \rightarrow \text{Bcast}_{\text{sv}}(x) \rightarrow c :=_{\text{sv}} x \rightarrow y :=_{\text{sv}} c \rightarrow \text{Bcast}_{\text{sv}}(y) \rightarrow a :=_{\text{sv}} y \rightarrow x :=_{\text{sv}} 2 \rightarrow \text{Bcast}_{\text{sv}}(x) \rightarrow b :=_{\text{sv}} x$. This sequence seemingly contains a dependency cycle from $\text{Bcast}_{\text{sv}}(x)$ to itself, and thus any reasonable system of dependencies on events would not allow this valid behaviour.

We fix this apparent cycle by splitting the broadcast event into its four basic components called subevents: (1) reading x to send to node 2 (stamp aNLR_2); (2) writing x on node 2 (stamp aNRW_2); (3) reading x to send to node 3 (stamp aNLR_3); (4) writing x on node 3 (stamp aNRW_3). With this we can create a more fine-grain sequence of dependencies: $x :=_{\text{sv}} 1 \rightarrow \langle \text{Bcast}_{\text{sv}}(x), \text{aNLR}_2 \rangle \rightarrow \langle \text{Bcast}_{\text{sv}}(x), \text{aNRW}_2 \rangle \rightarrow c :=_{\text{sv}} x \rightarrow \dots \rightarrow x :=_{\text{sv}} 2 \rightarrow \langle \text{Bcast}_{\text{sv}}(x), \text{aNLR}_3 \rangle \rightarrow \langle \text{Bcast}_{\text{sv}}(x), \text{aNRW}_3 \rangle \rightarrow b :=_{\text{sv}} x$. For each remote node the broadcast reads before writing, and we have a dependency between writing on node 2 and reading for node 3, but this does not create a dependency cycle at the level of the subevents and we can authorise the behaviour of Fig. 9c.

Stamps are shared by all libraries and also allow us to precisely define **ppo**, i.e. which pairs of effects are required to execute in order, even across libraries. For instance in example Fig. 9a we have a dependency $\langle \bar{z} := 1, \text{aNRW}_2 \rangle \xrightarrow{\text{ppo}} \langle \text{Bcast}_{\text{sv}}(x), \text{aNRW}_2 \rangle$ guaranteeing that the contents of z and x on node 2 are modified in order. However, note that this is more subtle than a dependency between events as the location x might still be read by the broadcast *before* the content of z is modified, i.e. $\langle \text{Bcast}_{\text{sv}}(x), \text{aNLR}_2 \rangle \rightarrow \langle \bar{z} := 1, \text{aNRW}_2 \rangle \rightarrow \langle \text{Bcast}_{\text{sv}}(x), \text{aNRW}_2 \rangle$, as is allowed by the semantics of RDMA.

Modularity. A final challenge in developing Mowgli is to support modularity through both horizontal composition (the use of libraries in a client program) and vertical composition (the development of libraries using other libraries as a subcomponent). Mowgli presents a generic framework that is independent of a memory model to support such proofs through a locality theorem. It allows the simultaneous use of multiple libraries within a single program, and defines a semantics when the specification of a library is used in place of an implementation. Finally, it provides local methods for proving that a library implementation satisfies its specification.

3 The Mowgli Framework and the Shared Variable Library

In this section we define Mowgli’s meta-language and general theory for modelling weak memory libraries, as well as its notion of compositionality that enables modular proofs. We note that our language and theory is generic and could be applied to other memory models. We present the syntax and semantics of Mowgli in §3.1 and model for formalising libraries in §3.2. Throughout the section, we use the shared variable library (sv) as a running example and define its consistency in §3.4. Then we present library abstraction in §3.5 and our main locality result in §3.6.

3.1 Syntax and Semantics

In this section, we present the syntax and semantics of our basic programming language. Our language is inspired by Cminor [Appel and Blazy 2007] and Yacovet [Stefanescu et al. 2024].

Programs. We assume a type Val of values, a type $\text{Loc} \subseteq \text{Val}$ of locations², and a type Method of methods. The syntax of sequential programs is given by the following grammar:

$$\begin{aligned} v, v_i &\in \text{Val} & m &\in \text{Method} & f &\in \text{Val} \rightarrow \text{SeqProg} & k &\in \mathbb{N}^+ \\ \text{SeqProg} &\ni p ::= v \mid m(v_1, \dots, v_k) \mid \text{let } p \text{ f} \mid \text{loop } p \mid \text{break}_k v \end{aligned}$$

A method call is parameterised by a sequence of input values. In later sections, we will instantiate m to basic operations such as read and write, as well as operations corresponding to method calls of a high-level library.

For a function f mapping values to sequential programs, the syntax $\text{let } p \text{ f}$ denotes the execution of p with an output that is then used as an input for f . This constructor is a generalisation of

²In Mowgli, every argument of a method call is a value. Thus identifiers (x, y, \dots) are called “locations” by the libraries but are seen as values by the meta-language.

the more standard let-in syntax, and for a program p_2 with a free meta-variable x we can define $\text{let } x = p_1 \text{ in } p_2$ as $\text{let } p_1 (\lambda v. p_2[x := v])$. We can also model sequential composition, i.e. $p_1; p_2$, as syntactic sugar for $\text{let } p_1 (\lambda_. p_2)$ using a constant function that discards its input. The syntax $\text{let } p \text{ f}$ also allows programs to perform branching and pattern-matching, via a function mapping different kinds of values to different continuations. In particular, $\text{if } v \text{ then } p_1 \text{ else } p_2$ can be taken as syntactic sugar for $\text{let } v \{ \text{true} \mapsto p_1, \text{false} \mapsto p_2 \}$.

Finally, our syntax includes $\text{loop } p$ that infinitely executes the program p , as well as the $\text{break}_k v$ construct which exits k levels of nested loops and returns v . While uncommon, these constructs can be used to define usual while and for loops.

We assume top-level concurrency. We assume a fixed number T of threads and let $\text{Tid} \triangleq \{1, 2, \dots, T\}$ be the set of all threads. A concurrent program is thus given by a tuple $\tilde{p} = \langle p_1, \dots, p_T \rangle$, where each thread t corresponds to a program $p_t \in \text{SeqProg}$. Note that we allow libraries to discriminate threads, and so the position of a program in \tilde{p} matters, e.g. the program $\langle p_1, \dots, p_T \rangle$ is *not* equivalent to $\langle p_T, \dots, p_1 \rangle$. For instance, a pair of RDMA threads have different interactions depending on whether they run on the same node or not.

Example 3.1 (Shared Variables). For our RDMA libraries, we assume a set of nodes, Node , of fixed size. Each thread t is associated to a node $n(t)$. The sv library uses the following methods:

$$m(\bar{v}) ::= \text{Write}_{\text{sv}}(x, v) \mid \text{Read}_{\text{sv}}(x) \mid \text{Bcast}_{\text{sv}}(x, d, \{n_1, \dots, n_k\}) \mid \text{Wait}_{\text{sv}}(d) \mid \text{GF}_{\text{sv}}(\{n_1, \dots, n_k\})$$

$\text{Write}_{\text{sv}}(x, v)$ writes a new value v to the location x of the current node. $\text{Read}_{\text{sv}}(x)$ reads the location x of the current node and returns its value. $\text{Bcast}_{\text{sv}}(x, d, \{n_1, \dots, n_k\})$ broadcasts the local value of x and overwrites the values of the copies of x on the nodes $\{n_1, \dots, n_k\}$, which might include the local node. $\text{Wait}_{\text{sv}}(d)$ waits for previous broadcasts of the thread marked with the same work identifier $d \in \text{Wid}$. As mentioned in the overview, this operation only guarantees that the local values of the broadcasts have been read, but not that remote copies have been modified. Finally, the global fence operation $\text{GF}_{\text{sv}}(\{n_1, \dots, n_k\})$ ensures every previous operation of the thread towards one of the nodes in the argument is fully finished, including the writing part of broadcasts.

Plain Executions. The semantics of a program is given by an execution, which is a graph over events. Each event has a label taken from the set $\text{Lab} \triangleq \text{Method} \times \text{Val}^* \times \text{Val}$, i.e. a triple comprising the method, the input values, and the output value. Labels are used to define events, which are elements of the set $\text{Event} \triangleq \text{Tid} \times \text{EventId} \times \text{Lab}$, where $\text{EventId} \triangleq \mathbb{N}$. For each event $\langle t, \iota, l \rangle \in \text{Event}$, we have that $t \in \text{Tid}$ is the thread that executes the label $l \in \text{Lab}$, and ι is a unique identifier for the event. For an event $e = \langle t, \iota, l \rangle$, we note $\text{t}(e) \triangleq t$.

Definition 3.2. We say that $\langle E, \text{po} \rangle$ is a *plain execution* iff $E \subseteq \text{Event}$, $\text{po} \subseteq E \times E$, and $\text{po} = \bigcup_{t \in \text{Tid}} \text{po}|_t$ where every $\text{po}|_t$ (i.e. po restricted to the events of thread t) is a total order.

Here, po represents *program order* i.e. $\langle e_1, e_2 \rangle \in \text{po}$ iff e_1 is executed before e_2 by the same thread.

We write $\emptyset_G \triangleq \langle \emptyset, \emptyset \rangle$ for the empty execution and $\{e\}_G \triangleq \langle \{e\}, \emptyset \rangle$ for the execution with a single event e . Given two executions, $G_1 = \langle E_1, \text{po}_1 \rangle$ and $G_2 = \langle E_2, \text{po}_2 \rangle$, with disjoint sets of events (i.e. $E_1 \cap E_2 = \emptyset$), we define their sequential composition, $G_1; G_2$, by ordering all events of G_1 before those of G_2 . Similarly, we define their parallel composition, $G_1 \parallel G_2$, by taking the union of G_1 and G_2 . That is,

$$G_1; G_2 \triangleq \langle E_1 \cup E_2, \text{po}_1 \cup \text{po}_2 \cup (E_1 \times E_2) \rangle \quad G_1 \parallel G_2 \triangleq \langle E_1 \cup E_2, \text{po}_1 \cup \text{po}_2 \rangle$$

The plain semantics of a program p executed by a thread t is given by $\llbracket p \rrbracket_t$, which is a set of pairs of the form $\langle r, G \rangle$, where r is the output and G is a plain execution. This set represents all conceivable unfoldings of the program into method calls, even those that will be rejected by the semantics of

the corresponding libraries. Each output is a pair $\langle v, k \rangle$, where v is a value and k a break number, indicating the program terminates by requesting to exit k nested loops and returning the value v .

$$\begin{aligned}
\llbracket v \rrbracket_t &\triangleq \{ \langle \langle v, 0 \rangle, \emptyset_G \rangle \} & \llbracket \text{break}_k v \rrbracket_t &\triangleq \{ \langle \langle v, k \rangle, \emptyset_G \rangle \} \\
\llbracket m(\widetilde{v}) \rrbracket_t &\triangleq \{ \langle \langle v', 0 \rangle, \{ \langle t, \iota, \langle m, \widetilde{v}, v' \rangle \} \rangle_G \mid v' \in \text{Val} \wedge \iota \in \text{EventId} \} \\
\llbracket \text{let } p \text{ f} \rrbracket_t &\triangleq \{ \langle \langle r, G_1; G_2 \rangle \mid \langle \langle v, 0 \rangle, G_1 \rangle \in \llbracket p \rrbracket_t \wedge \langle \langle r, G_2 \rangle \in \llbracket f v \rrbracket_t \} \\
&\quad \cup \{ \langle \langle v, k \rangle, G_1 \rangle \mid \langle \langle v, k \rangle, G_1 \rangle \in \llbracket p \rrbracket_t \wedge k \neq 0 \} \\
\llbracket \text{loop } p \rrbracket_t &\triangleq \bigcup_{j \in \mathbb{N}} \{ \langle \langle v, k \rangle, G_0; \dots; G_j \rangle \mid (\forall 0 \leq i < j. \langle \langle _, 0 \rangle, G_i \rangle \in \llbracket p \rrbracket_t) \wedge \langle \langle v, k+1 \rangle, G_j \rangle \in \llbracket p \rrbracket_t \}
\end{aligned}$$

The execution of a value v simply returns $\langle v, 0 \rangle$ with an empty graph. Similarly, the execution of $\text{break}_k v$ returns $\langle v, k \rangle$ with a non-zero break number and an empty graph.

The plain semantics of $\llbracket m(\widetilde{v}) \rrbracket_t$ considers every value v' as a possible output of the method call. For each, we can create a graph G with a single event $\langle t, _, \langle m, \widetilde{v}, v' \rangle \rangle$, and the corresponding output for the program is then $\langle v', 0 \rangle$ with a break number of 0.

The execution of $\text{let } p \text{ f}$ has two kinds of plain semantics. Either the execution of p requests a break, i.e. $\langle \langle v, k \rangle, G_1 \rangle \in \llbracket p \rrbracket_t$ with $k \neq 0$, in which case $\text{let } p \text{ f}$ breaks as well with the same output. Or p terminates with a break number of zero, and the output value v of p is given to f . In this second case, the plain execution of $\text{let } p \text{ f}$ is the sequential composition of the plain executions for p and $(f v)$, and its output value is the one of $(f v)$.

Finally, the execution of $\text{loop } p$ can be unfolded and corresponds to the execution of p any number $j+1$ of times. The first j times, p returns without requesting a break and its output value is ignored. The $(j+1)^{\text{th}}$ execution of p returns a value v and break number $k+1$, and $\text{loop } p$ propagates $\langle v, k \rangle$ with a decremented break number. The plain execution of the loop is then the sequential composition of the plain executions of the $j+1$ iterations of p .

We lift the plain semantics to the level of concurrent programs and define

$$\llbracket \widetilde{p} \rrbracket \triangleq \{ \langle \langle v_1, \dots, v_T \rangle, \parallel_{t \in \text{Tid}} G_t \rangle \mid \forall t \in \text{Tid}. \langle \langle v_t, 0 \rangle, G_t \rangle \in \llbracket p_t \rrbracket_t \}$$

Concurrent programs only properly terminate if each thread terminates with a break number of 0. In which case, the output of the concurrent program is the parallel composition of the values and plain executions of the different threads.

Executions. We generate executions from plain executions by (1) extending the model with subevents, then (2) introducing additional relations describing synchronisation and happens-before order. We will later define consistency conditions for executions in the context of libraries.

We assume a fixed set of stamps, $\text{Stamp} = \{a_1, \dots\}$, and a relation $\text{to} \subseteq \text{Stamp} \times \text{Stamp}$. We will use stamps to define subevents and to to define preserved program order over subevents within an execution.

Definition 3.3. We say that $\langle E, \text{po}, \text{stmp}, \text{so}, \text{hb} \rangle$ is an *execution* iff each of the following holds:

- $\langle E, \text{po} \rangle$ is a plain execution.
- $\text{stmp} : E \rightarrow \mathcal{P}(\text{Stamp})$ is a function that associates each event with a non-empty set of stamps and induces a set of *subevents*, $\text{SEvent} \triangleq \{ \langle e, a \rangle \mid e \in E \wedge a \in \text{stmp}(e) \}$.
- $\text{so} \subseteq \text{SEvent} \times \text{SEvent}$ and $\text{hb} \subseteq \text{SEvent} \times \text{SEvent}$ are relations on SEvent defining *synchronisation order* and *happens-before order*, respectively.

To define consistency, we must ultimately relate po , so , and hb . However, in many weak memory models such as RDMA, including all of po into hb is too restrictive. We therefore make use of a

			Second Stamp										
			single					families					
			1	2	3	4	5	6	7	8	9	10	11
First Stamp	single	to	aCR	aCW	aCAS	aMF	aWT	aNLR _n	aNRW _n	aNRR _n	aNLW _n	aRF _n	aGF _n
		A	aCR	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
		B	aCW	✗	✓	✓	✓	✓	✓	✓	✓	✓	✓
		C	aCAS	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
		D	aMF	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
		E	aWT	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
	families	F	aNLR _n	✗	✗	✗	✗	SN	SN	SN	SN	SN	SN
		G	aNRW _n	✗	✗	✗	✗	✗	SN	SN	SN	✗	SN
		H	aNRR _n	✗	✗	✗	✗	✗	✗	✗	SN	SN	SN
		I	aNLW _n	✗	✗	✗	✗	✗	✗	✗	SN	✗	SN
		J	aRF _n	✗	✗	✗	✗	✗	SN	SN	SN	SN	SN
		K	aGF _n	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓

Fig. 10. Stamp order **to** for the RDMA libraries. Lines indicate the earlier stamp, columns the later. A cell marked ✓ indicates that the stamps are ordered, and that the po ordering of subevents with these stamps is preserved. A cell marked ✗ indicates that the stamps are not ordered, and that such subevents can execute out of order. Finally, SN indicates the stamps are ordered iff they have the same node index.

weaker relation called *preserved program order*, $\text{ppo} \subseteq \text{SEvent} \times \text{SEvent}$, which we derive from po and **to** as follows:

$$\text{ppo} \triangleq \{ \langle \langle e_1, a_1 \rangle, \langle e_2, a_2 \rangle \rangle \mid \langle e_1, e_2 \rangle \in \text{po} \wedge a_1 \in \text{stmp}(e_1) \wedge a_2 \in \text{stmp}(e_2) \wedge \langle a_1, a_2 \rangle \in \text{to} \}$$

For our RDMA libraries, we define 11 kinds of stamps. We have aCR representing a CPU read; aCW representing a CPU write; aCAS for an atomic read-modify-write operation; aMF for a TSO memory fence; aWT for a wait t operation; aNLR_n for a NIC local read; aNRW_n for a NIC remote write; aNRR_n for a NIC remote read; aNLW_n for a NIC local write; aRF_n for a NIC remote fence; and aGF_n for a global fence operation. The last 6 are families of stamps, as we create a different copy for each node $n \in \text{Node}$.

The stamp order **to** we use is defined in Fig. 10. We note ✓ when two stamps are ordered, ✗ when they are not ordered, and SN when they are ordered iff they have the same index node. For instance, the ✗ in cell B1 indicates that when a CPU write is in program order before a CPU read, there is no ordering guarantee between the two operations, as we assume the CPUs follow the TSO memory model, and the read might execute first.

Example 3.4 (ppo for Shared Variables). For the sv library, we use the stamping function stmp_{sv} :

$$\text{stmp}_{\text{sv}}(\langle _, _, \langle \text{Write}_{\text{sv}}, _, _ \rangle \rangle) = \{\text{aCW}\}$$

$$\text{stmp}_{\text{sv}}(\langle _, _, \langle \text{Read}_{\text{sv}}, _, _ \rangle \rangle) = \{\text{aCR}\}$$

$$\text{stmp}_{\text{sv}}(\langle _, _, \langle \text{Wait}_{\text{sv}}, _, _ \rangle \rangle) = \{\text{aWT}\}$$

$$\text{stmp}_{\text{sv}}(\langle _, _, \langle \text{GF}_{\text{sv}}, (\{n_1, \dots, n_k\}), _ \rangle \rangle) = \{\text{aGF}_{n_1}, \dots, \text{aGF}_{n_k}\}$$

$$\text{stmp}_{\text{sv}}(\langle _, _, \langle \text{Bcast}_{\text{sv}}, (_, _, \{n_1, \dots, n_k\}), _ \rangle \rangle) = \{\text{aNLR}_{n_1}, \text{aNRW}_{n_1}, \dots, \text{aNLR}_{n_k}, \text{aNRW}_{n_k}\}$$

Broadcasts are associated with a NIC local read and NIC remote write stamp for each remote node they are broadcasting towards. Similarly, global fence operations are associated with a global fence stamp for each node.

With this, the stamp order is enough to enforce the behaviour of the global fence. If we have a program $\text{Bcast}_{\text{sv}}(x, d, \{\dots, n, \dots\})$; $\text{GF}_{\text{sv}}(\{\dots, n, \dots\})$, the plain execution has two events e_{BR} and e_{GF} , and the definitions of stmp_{sv} and **to** (cell G11 in Fig. 10) imply $\langle e_{\text{BR}}, \text{aNRW}_n \rangle \xrightarrow{\text{ppo}} \langle e_{\text{GF}}, \text{aGF}_n \rangle$.

3.2 Libraries

In this section, we describe how libraries and library consistency are modelled in our framework.

Definition 3.5. We say that a triple $\langle M, \text{loc}, C \rangle$ is a *library* iff each of the following holds.

- (1) $M \subseteq \text{Method}$ is a set of *methods*.
- (2) $\text{loc} : \text{Event}|_M \rightarrow \mathcal{P}(\text{Loc})$ is a function associating each method call to a set of locations accessed by the method call.
- (3) C is a *consistency predicate* over executions, respecting the following two properties.
 - **Monotonicity:** If $\langle E, \text{po}, \text{stmp}, \text{so}, \text{hb} \rangle \in C$ (i.e. is consistent), and $(\text{ppo} \cup \text{so})^+ \subseteq \text{hb}' \subseteq \text{hb}$, then $\langle E, \text{po}, \text{stmp}, \text{so}, \text{hb}' \rangle \in C$.
 - **Decomposability:** If $\langle (E_1 \uplus E_2), \text{po}, \text{stmp}, \text{so}, \text{hb} \rangle \in C$ and $\text{loc}(E_1) \cap \text{loc}(E_2) = \emptyset$, then $\langle E_1, \text{po}|_{E_1}, \text{stmp}|_{E_1}, \text{so}|_{E_1}, \text{hb}|_{E_1} \rangle \in C$.

Usually, including for all of the examples considered in this paper, the locations accessed by a method call are a subset of its arguments. E.g., we say that $\text{Write}(x, v)$ only accesses x . Monotonicity states that removing constraints cannot disallow a behaviour; this is trivially respected by all reasonable libraries. Decomposability states that method calls manipulating different locations can be considered independently. Crucially, this means combining independent programs *cannot* create additional behaviours; a prerequisite for modular verification. This holds for almost all libraries, and usually only breaks when programs have access to meta-information (e.g. the number of instructions of the whole program).

However, decomposability does *not* hold for RDMA^{TSO} . As show in Fig. 2, the program $\bar{z} := x; \text{Poll}(2); x := 1$ does not allow the outcome $z = 1$, while a combined program $p; \bar{z} := x; \text{Poll}(2); x := 1$ might, even when p seems independent (i.e. does not use locations z and x). This composition problem fundamentally prevents modular verification of RDMA^{TSO} programs. It is the reason we develop the alternative semantics of $\text{RDMA}^{\text{WAIT}}$, while ensuring the two semantics are as close as possible.

Notation. For a library L , we have $\text{Event}|_{L.M} = \{ \langle _, _, \langle m, _ \rangle \rangle \in \text{Event} \mid m \in L.M \}$. We use $\text{Event}|_L$ to refer to $\text{Event}|_{L.M}$. Moreover, $\text{loc}(e)$ is used to denote $L.\text{loc}(e)$, where L is the library containing e (i.e. $e \in \text{Event}|_L$) and for $E \subseteq \text{Event}$, we define $\text{loc}(E) \triangleq \bigcup_{e \in E} \text{loc}(e)$. From this, we can also define the locations $\text{loc}(\tilde{p})$ of a program \tilde{p} as $\text{loc}(\tilde{p}) \triangleq \bigcup_{\langle _, _, \langle E, _ \rangle \rangle \in [\tilde{p}]} \text{loc}(E)$.

Given a relation r and a set A , we write r^+ for the transitive closure of r ; r^* for its reflexive transitive closure; r^{-1} for the inverse of r ; $r|_A$ for $r \cap (A \times A)$; and $[A]$ for the identity relation on A , i.e. $\{ \langle a, a \rangle \mid a \in A \}$. We write $r_1; r_2$ for the relational composition of r_1 and r_2 : $\{ \langle a, b \rangle \mid \exists c. \langle a, c \rangle \in r_1 \wedge \langle c, b \rangle \in r_2 \}$.

Consistent Execution. Two libraries are *compatible* if their sets of methods are disjoint. We use Λ to denote a set of pairwise compatible libraries.

Definition 3.6. Let Λ be a set of pairwise compatible libraries. An execution $\langle E, \text{po}, \text{stmp}, \text{so}, \text{hb} \rangle$ is Λ -*consistent* iff each of the following holds.

- $(\text{ppo} \cup \text{so})^+ \subseteq \text{hb}$ and hb is a strict partial order (i.e. both irreflexive and transitive).
- $E = \bigcup_{L \in \Lambda} E|_L$ and $\text{so} = \bigcup_{L \in \Lambda} \text{so}|_L$.
- For all $L \in \Lambda$, we have $\langle E|_L, \text{po}|_L, \text{stmp}|_L, \text{so}|_L, \text{hb}|_L \rangle \in L.C$.

Although the definition of Λ -consistency allows hb relations that are bigger than $(\text{ppo} \cup \text{so})^+$, we usually have $\text{hb} = (\text{ppo} \cup \text{so})^+$ for the program executions we are interested in.

Given a concurrent program \tilde{p} using libraries Λ , we note $\text{outcome}_\Lambda(\tilde{p})$ the set of all output values of its Λ -consistent executions.

$$\text{outcome}_\Lambda(\tilde{p}) \triangleq \{\tilde{v} \mid \exists \langle E, \text{po}, \text{stmp}, \text{so}, \text{hb} \rangle \Lambda\text{-consistent. } \langle \tilde{v}, \langle E, \text{po} \rangle \rangle \in \llbracket \tilde{p} \rrbracket\}$$

3.3 The RDMA^{wait} Library

RDMA^{wait} is used as the lowest library of our tower of abstraction (Fig. 1). As mentioned in §3.4, it is the implementation target for the shared variable library (sv). It is an adaptation of RDMA^{ts} where the poll instruction is replaced by a more intuitive wait operation.

The RDMA^{wait} library uses the following 8 methods.

$$\begin{aligned} m(\tilde{v}) ::= & \text{Write}(x, v) \mid \text{Read}(x) \mid \text{CAS}(x, v_1, v_2) \mid \text{Mfence}() \\ & \mid \text{Get}(x, y, d) \mid \text{Put}(x, y, d) \mid \text{Wait}(d) \mid \text{Rfence}(n) \end{aligned}$$

The first line covers usual TSO operations: $\text{Write}(x, v)$ is a CPU write; $\text{Read}(x)$ is a CPU read; $\text{CAS}(x, v_1, v_2)$ is an atomic compare-and-swap operation that overwrites x to v_2 iff x contained v_1 , and returns the old value of x ; and $\text{Mfence}()$ is a TSO memory fence flushing the store buffer.

The second line covers RDMA-specific operations: $\text{Get}(x, y, d)$ (noted $x :=^d \bar{y}$ in our examples) is a get³ operation with work identifier d performing a NIC remote read on y and a NIC local write on x ; similarly $\text{Put}(x, y, d)$ (noted $\bar{x} :=^d y$) is a put operation with work identifier d performing a NIC local read on y and a NIC remote write on x ; $\text{Wait}(d)$ waits for previous operations with work identifier d ; and finally $\text{Rfence}(n)$ is an RDMA remote fence for the communication channel towards n that does not block the CPU.

We assume that each location x is associated with a specific node $n(x)$. From this, given $\langle E, \text{po} \rangle$, there is a single valid stamping function stmp_{RL} . Notably we have $\text{stmp}_{\text{RL}}(\text{Get}(x, y, d)) = \{\text{aNRr}_{n(y)}, \text{aNLW}_{n(y)}\}$ and $\text{stmp}_{\text{RL}}(\text{Put}(x, y, d)) = \{\text{aNLr}_{n(x)}, \text{aNRW}_{n(x)}\}$. Put and get operations perform both a NIC read and a NIC write, and as such are associated to two stamps, where the remote node can be deduced from the location. Also, a succeeding CAS has a single stamp aCAS, while a failing CAS has stamps {aMF, aCR}, as it behaves as both a memory fence (aMF) and a CPU read (aCR).

The formal semantics requires several functions and relations: v_R , v_W , rf , and mo , with roles similar to the semantics of sv (cf. §3.4), as well as the NIC-flush-order relation nfo representing the PCIe guarantees that NIC reads flush previous NIC writes. The consistency predicate for RDMA^{wait} is then stated from these relations and some derived relations, similarly to §3.4.

3.4 Example: Consistency for Shared Variables

As mentioned in Example 3.1, sv uses the methods $M = \{\text{Write}_{\text{sv}}, \text{Read}_{\text{sv}}, \text{Bcast}_{\text{sv}}, \text{Wait}_{\text{sv}}, \text{GF}_{\text{sv}}\}$. Since only the method and arguments matter for the location function, we use $\text{loc}(m(\tilde{v}))$ to denote $\text{loc}(\langle _, _, \langle m, \tilde{v}, _ \rangle \rangle)$, where $\text{loc}(\text{Write}_{\text{sv}}(x, _)) = \text{loc}(\text{Read}_{\text{sv}}(x)) = \text{loc}(\text{Bcast}_{\text{sv}}(x, _, _)) = \{x\}$ for events accessing a location x , and $\text{loc}(e) = \emptyset$ otherwise for methods Wait_{sv} and GF_{sv} .

Notation. For a subevent s , we note $s.e$ and $s.a$ its two components. Given an execution $\mathcal{G} = \langle E, \text{po}, \text{stmp}, \text{so}, \text{hb} \rangle$ and a stamp a , we write $\mathcal{G}.a$ for $\{s \in \mathcal{G}.\text{SEvent} \mid s.a = a\}$. For families, by abuse of notation, we also write e.g. $\mathcal{G}.\text{aNRr}$ for $\bigcup_{n \in \text{Node}} \mathcal{G}.\text{aNRr}_n$. We extend the notation loc to subevents by writing $\text{loc}(s)$ for $\text{loc}(s.e)$. We define the set of reads as $\mathcal{G}.\mathcal{R} \triangleq \mathcal{G}.\text{aCR} \cup \mathcal{G}.\text{aCAS} \cup \mathcal{G}.\text{aNLr} \cup \mathcal{G}.\text{aNRr}$ and writes as $\mathcal{G}.\mathcal{W} \triangleq \mathcal{G}.\text{aCW} \cup \mathcal{G}.\text{aCAS} \cup \mathcal{G}.\text{aNLW} \cup \mathcal{G}.\text{aNRW}$. We write $\mathcal{G}.\mathcal{W}_x \triangleq \{s \in \mathcal{G}.\mathcal{W} \mid \text{loc}(s) = \{x\}\}$ to constrain the set to writes on a specific location x . We also use

³In the RDMA specification, Get and Put are referred to as respectively “RDMA Read” and “RDMA Write” operations. We use the terms get and put to prevent confusion, as each of these perform both a read and a write subevents.

$|_t$ to restrict a set or relation to a specific thread. E.g. $E|_t = \{e \mid e \in E \wedge t(e) = t\}$ and $po|_t = [E|_t]; po; [E|_t]$.

For the sv library, we additionally define $\mathcal{G}.\mathcal{W}^n \triangleq \{\langle e, aCW \rangle \mid n(t(e)) = n\} \cup \mathcal{G}.aNRW_n$ as the set of write subevents occurring on node n . This includes CPU writes on the node, as well as broadcast writes towards n from all threads. We also note $\mathcal{G}.\mathcal{W}_x^n \triangleq \mathcal{G}.\mathcal{W}_x \cap \mathcal{G}.\mathcal{W}^n$ as expected. Similarly, $\mathcal{G}.\mathcal{R}^n \triangleq \{s \mid s \in \mathcal{G}.\mathcal{R} \wedge n(t(s)) = n\}$ covers reads occurring on n , either by a CPU read or as part of a broadcast.

Consistency. We now work towards a definition of consistency for shared variables.

Definition 3.7. For an execution $\mathcal{G} = \langle E, po, stmp_{sv}, _, _ \rangle$, we define the following:

- The *value-read* function $v_R : \mathcal{G}.\mathcal{R} \rightarrow \text{Val}$ that associates each read subevent with the value returned, if available, i.e. if $e = \langle _, _, \langle \text{Read}_{sv}, _, v \rangle \rangle$, then $v_R(e) = v$.
- The *value-written* function $v_W : \mathcal{G}.\mathcal{W} \rightarrow \text{Val}$ that associates each write subevent with a value \mathcal{G} , i.e. if $e = \langle _, _, \langle \text{Write}_{sv}, (_, v), _ \rangle \rangle$, then $v_W(e) = v$.
- A *reads-from* relation, $rf \triangleq \bigcup_n rf^n$, where each $rf^n \subseteq \mathcal{G}.\mathcal{W}^n \times \mathcal{G}.\mathcal{R}^n$ is a relation on subevents of the same location and node with matching values, i.e. if $\langle s_1, s_2 \rangle \in rf^n$ then $\text{loc}(s_1) = \text{loc}(s_2)$ and $v_W(s_1) = v_R(s_2)$.
- A *modification-order* relation $mo \triangleq \bigcup_{x,n} mo_x^n$ describing the order in which writes on x on node n reach memory.

We define *well-formedness* for rf and mo as follows. For each remote, a broadcast writes the corresponding read value: if $s_1 = \langle e, aNLR_n \rangle \in \mathcal{G}.SEvent$ and $s_2 = \langle e, aNRW_n \rangle \in \mathcal{G}.SEvent$, then $v_R(s_1) = v_W(s_2)$. Each rf^n is functional on its range, i.e. every read in $\mathcal{G}.\mathcal{R}^n$ is related to at most one write in $\mathcal{G}.\mathcal{W}^n$. If a read is not related to a write, it reads the initial value of zero, i.e. if $s_2 \in \mathcal{G}.\mathcal{R}^n \wedge \langle _, s_2 \rangle \notin rf^n$ then $v_R(s_2) = 0$. Finally, each mo_x^n is a strict total order on $\mathcal{G}.\mathcal{W}_x^n$.

We further define the *reads-from-internal* relation as $rf_i \triangleq [aCW]; (po \cap rf); [aCR]$ (which corresponds to CPU reads and writes using the same TSO store buffer), and the *reads-from-external* relation as $rf_e \triangleq rf \setminus rf_i$. As we shall see in Def. 3.8, rf_i does *not* contribute to synchronisation order, whereas rf_e does. Moreover, given an execution \mathcal{G} and well-formed rf and mo , we derive additional relations.

$$\begin{aligned}
 pf &\triangleq \left\{ \langle \langle e_1, aNLR_n \rangle, \langle e_2, aWT \rangle \rangle \mid \langle e_1, e_2 \rangle \in po \wedge \left(\begin{array}{l} \exists d. e_1 = \langle _, _, \langle \text{Bcast}_{sv}, (_, _, d), _ \rangle \rangle \\ \wedge e_2 = \langle _, _, \langle \text{Wait}_{sv}, (d), _ \rangle \rangle \end{array} \right) \right\} \\
 rb^n &\triangleq \left\{ \langle r, w \rangle \in \mathcal{G}.\mathcal{R}^n \times \mathcal{G}.\mathcal{W}^n \mid \begin{array}{l} \text{loc}(r) = \text{loc}(w) \\ \wedge (\langle r, w \rangle \in ((rf^n)^{-1}; mo^n) \vee r \notin \text{img}(rf^n)) \end{array} \right\} \quad rb \triangleq \bigcup_n rb^n \\
 iso &\triangleq \{ \langle \langle e, aNLR_n \rangle, \langle e, aNRW_n \rangle \rangle \mid e = \langle _, _, \langle \text{Bcast}_{sv}, (_, _, \{ \dots, n, \dots \}), _ \rangle \rangle \in E \}
 \end{aligned}$$

The *polls-from* relation pf states that a Wait_{sv} operation synchronises with the NIC local read subevents of previous broadcasts that use the same work identifier. The *reads-before* relation rb states that a read r executes before a specific write w on the same node and location. This is either because r reads the initial value of 0, or because r reads from a write that is mo -before w . Finally, the *internal-synchronisation-order* relation iso states that, within a broadcast, for each remote node the reading part occurs before the writing part.

We can then define the consistency predicate $sv.C$ as follows.

Definition 3.8 (sv-consistency). $\langle E, po, stmp, so, hb \rangle$ is sv-consistent if:

- $stmp = stmp_{sv}$ (defined in §3.1);
- there exists well-formed v_R, v_W, rf , and mo , such that $[aCR]; (po^{-1} \cap rb); [aCW] = \emptyset$ and $so = iso \cup rf_e \cup pf \cup rb \cup mo$.

It is straightforward to check that this consistency predicate satisfies monotonicity and decomposability. For CPU reads and writes, we ask that **rb** does not contradict the program order. E.g., a program $\text{Write}_{\text{sv}}(x, 1); \text{Read}_{\text{sv}}(x)$ must return 1 and cannot return 0, even if the semantics of TSO allows for the read to finish before the write.

There is no need to explicitly include conditions on **hb** in the consistency of the library, as the global consistency condition (cf. Def. 3.6) already enforces that $(\text{ppo} \cup \text{so} \cup \text{hb})^+$ is irreflexive.

3.5 Library Implementations

We now describe a mechanism for implementing the method calls of a library by an implementation. Our ideas build on Yacovet [Stefanescu et al. 2024], but have been adapted to our setting, which comprises a much weaker happens-before relation (based on **ppo** instead of **po**). In particular, MOWGLI's notions of implementation, soundness, and abstraction are similar to Yacovet (but simpler), but the notion of "local soundness" is more complicated due to the use of **ppo** and subevents.

An implementation for a library L is a function $I : (\text{Tid} \times L.M \times \text{Val}^*) \rightarrow \text{SeqProg}$ associating every method call of the library L to a sequential program.

Definition 3.9. We say that I is *well defined* for a library L using Λ iff for all $t \in \text{Tid}$, $m \in L.M$ and $\tilde{v} \in \text{Val}^*$, we have:

- (1) $L \notin \Lambda$, and $I(t, m, \tilde{v})$ only calls methods of the libraries of Λ .
- (2) $\langle \langle -, k+1 \rangle, - \rangle \notin \llbracket I(t, m, \tilde{v}) \rrbracket_t$, i.e. the implementation of a method call $m(\tilde{v})$ cannot return with a non-zero break number, and thus cannot cause a loop containing a call to $m(\tilde{v})$ to break inappropriately.
- (3) if $\langle \langle v, 0 \rangle, \langle E, \text{po} \rangle \rangle \in \llbracket I(t, m, \tilde{v}) \rrbracket_t$ then $E \neq \emptyset$, i.e. if an implementation successfully executes, it must contain at least one method call.

We note $\text{loc}(I)$ the set of all locations that can be accessed by the implementation of I : $\text{loc}(I) \triangleq \bigcup_{t, m, \tilde{v}} \bigcup_{\langle \langle -, k+1 \rangle, - \rangle \in \llbracket I(t, m, \tilde{v}) \rrbracket_t} \text{loc}(E)$. We then define a function $\llbracket _ \rrbracket_I$ to map an implementation I to a concurrent program as follows.

$$\begin{aligned} \llbracket v \rrbracket_{t,I} &\triangleq v & \llbracket m(v_1, \dots, v_k) \rrbracket_{t,I} &\triangleq \begin{cases} I(t, m, \langle v_1, \dots, v_k \rangle) & \text{if } m \in L.M \\ m(v_1, \dots, v_k) & \text{otherwise} \end{cases} \\ \llbracket \text{loop } p \rrbracket_{t,I} &\triangleq \text{loop } \llbracket p \rrbracket_{t,I} & \llbracket \text{let } p \text{ f } \rrbracket_{t,I} &\triangleq \text{let } \llbracket p \rrbracket_{t,I} (\lambda v. \llbracket f \ v \rrbracket_{t,I}) \\ \llbracket \text{break}_k \ v \rrbracket_{t,I} &\triangleq \text{break}_k \ v & \llbracket \langle p_1, \dots, p_T \rangle \rrbracket_I &\triangleq \langle \llbracket p_1 \rrbracket_{1,L}, \dots, \llbracket p_T \rrbracket_{T,L} \rangle \end{aligned}$$

As an example, we can define the implementation I_{SV} of the broadcast library into $\text{RDMA}^{\text{WAIT}}$. For each location x of the broadcast library, we create a location x_n for each node $n \in \text{Node}$. We also create a dummy location per node, \perp_n for $n \in \text{Node}$, and we use an additional dummy work identifier d_0 .

$$\begin{aligned} I_{\text{SV}}(t, \text{Write}_{\text{sv}}, (x, v)) &\triangleq \text{Write}(x_{n(t)}, v) \\ I_{\text{SV}}(t, \text{Read}_{\text{sv}}, (x)) &\triangleq \text{Read}(x_{n(t)}) \\ I_{\text{SV}}(t, \text{Bcast}_{\text{sv}}, (x, d, \{n_1, \dots, n_k\})) &\triangleq \text{Put}(x_{n_1}, x_{n(t)}, d); \dots; \text{Put}(x_{n_k}, x_{n(t)}, d) \\ I_{\text{SV}}(t, \text{Wait}_{\text{sv}}, (d)) &\triangleq \text{Wait}(d) \\ I_{\text{SV}}(t, \text{GF}_{\text{sv}}, (\{n_1, \dots, n_k\})) &\triangleq \text{Get}(\perp_{n(t)}, \perp_{n_1}, d_0); \dots; \text{Get}(\perp_{n(t)}, \perp_{n_k}, d_0); \text{Wait}(d_0) \end{aligned}$$

where $\{\text{Write}, \text{Read}, \text{Put}, \text{Get}, \text{Wait}\}$ are methods of the $\text{RDMA}^{\text{WAIT}}$ library (see §3.3).

A read/write on a thread t accesses the location of its node $n(t)$. A broadcast executes multiple Put operations. Each of them reads the location of its node and overwrites the location of a designated

node. A wait operation works similarly to $\text{RDMA}^{\text{WAIT}}$. Finally, a global fence executes a Get operation towards each node requiring fencing, and waits for the completion of all the Get operations. As mentioned in the overview, this ensures that all previous NIC operations towards these nodes are completely finished.

We can easily see that I_{SV} is well defined, as it cannot return a break number greater than zero, and every (succeeding) implementation generates at least one event.

Using these definitions, we arrive at a notion of a sound implementation, which holds whenever the implementation is a refinement of the library specification.

Definition 3.10. We say that I is a *sound implementation* of L using Λ if, for any program \tilde{p} such that $\text{loc}(I) \cap \text{loc}(\tilde{p}) = \emptyset$, we have that $\text{outcome}_{\Lambda}(\llbracket \tilde{p} \rrbracket_I) \subseteq \text{outcome}_{\Lambda \uplus \{L\}}(\tilde{p})$.

For a concurrent program \tilde{p} using methods of $(\Lambda \uplus \{L\})$, $\llbracket \tilde{p} \rrbracket_I$ only uses methods of Λ . The implementation I is sound if the translation does not introduce any new outcomes. We can assume I and \tilde{p} use disjoint locations to avoid capture of location names.

3.6 Abstractions and Locality

We now work towards the modular proof technique for verifying soundness of an implementation against a library in MowGLI. As is common in proofs of refinement, we use an *abstraction function* [Abadi and Lamport 1991] mapping the concrete implementation to its abstract library specification. For $f : A \rightarrow B$ and $r \subseteq A \times A$, we note $f(r) \triangleq \{\langle f(x), f(y) \rangle \mid \langle x, y \rangle \in r\}$.

Definition 3.11. Suppose I is a well-defined implementation of a library L using Λ , and that $G = \langle E, \text{po} \rangle$ and $G' = \langle E', \text{po}' \rangle$ are plain executions using methods of Λ and L respectively. We say that a surjective function $f : E \rightarrow E'$ abstracts G to G' , denoted $\text{abs}_{I,L}^f(G, G')$, iff

- $E|_L = \emptyset$ (i.e. G contains no calls to the abstract library L) and $E'|_L = E'$ (i.e. G' only contains calls to the abstract library L);
- $f(\text{po}) \subseteq (\text{po}')^*$ and $\forall e_1, e_2, \langle f(e_1), f(e_2) \rangle \in \text{po}' \implies \langle e_1, e_2 \rangle \in \text{po}$; and
- if $e' = \langle t, \iota, \langle m, \tilde{v}, v' \rangle \rangle \in E'$ then $\langle \langle v', 0 \rangle, G|_{f^{-1}(e')} \rangle \in \llbracket I(t, m, \tilde{v}) \rrbracket_t$

Intuitively, $\text{abs}_{I,L}^f(G, G')$ means there is some abstract concurrent program \tilde{p} on library L such that $\langle _, G' \rangle \in \llbracket \tilde{p} \rrbracket$ is a plain execution of the abstract program, $\langle _, G \rangle \in \llbracket \llbracket \tilde{p} \rrbracket_I \rrbracket$ is a plain execution of its implementation, and G and G' behave similarly. The abstraction function f maps every event of the implementation to the abstract method call it was created for. The second requirement states that the program order is preserved in both directions. The last requirement states that, for each abstract event e' , its implementation $G|_{f^{-1}(e')}$ behaves properly. We ask that this subgraph be a valid plain execution of the implementation with the same output value.

LEMMA 3.12. Given \tilde{p} on library L and a well-defined implementation I of L , if $\langle \tilde{v}, G \rangle \in \llbracket \llbracket \tilde{p} \rrbracket_I \rrbracket$ then there exists $\langle \tilde{v}, G' \rangle \in \llbracket \tilde{p} \rrbracket$ and f such that $\text{abs}_{I,L}^f(G, G')$.

Finally, we can define a notion of local soundness for an implementation.

Definition 3.13. We say that a well defined implementation I of a library L is *locally sound* iff, whenever we have a Λ -consistent execution $\mathcal{G} = \langle E, \text{po}, \text{stmp}, \text{so}, \text{hb} \rangle$ and $\text{abs}_{I,L}^f(\langle E, \text{po} \rangle, \langle E', \text{po}' \rangle)$, then there exists stmp', so' , and a concretisation function $g : \langle E', \text{po}', \text{stmp}' \rangle.\text{SEvent} \rightarrow \mathcal{G}.\text{SEvent}$ such that:

- $g(\langle e', a' \rangle) = \langle e, a \rangle$ implies $f(e) = e'$ and
 - For all a_0 such that $\langle a_0, a' \rangle \in \text{to}$, there exists $\langle e_1, a_1 \rangle \in \mathcal{G}.\text{SEvent}$ such that $f(e_1) = e'$, $\langle a_0, a_1 \rangle \in \text{to}$, and $\langle \langle e_1, a_1 \rangle, \langle e, a \rangle \rangle \in \text{hb}^*$;

- For all a_0 such that $\langle a', a_0 \rangle \in \text{to}$, there exists $\langle e_2, a_2 \rangle \in \mathcal{G}.\text{SEvent}$ such that $f(e_2) = e'$, $\langle a_2, a_0 \rangle \in \text{to}$, and $\langle \langle e, a \rangle, \langle e_2, a_2 \rangle \rangle \in \text{hb}^*$.
- $g(\text{so}') \subseteq \text{hb}$;
- For all hb' transitive such that $(\text{ppo}' \cup \text{so}')^+ \subseteq \text{hb}'$ and $g(\text{hb}') \subseteq \text{hb}$, we have $\langle E', \text{po}', \text{stmp}', \text{so}', \text{hb}' \rangle \in L.C$, where $\text{ppo}' \triangleq \langle E', \text{po}', \text{stmp}' \rangle.\text{ppo}$.

Unlike the notion of soundness (cf. Def. 3.10) expressed using an arbitrary program, local soundness is expressed using an arbitrary abstraction. It states that whenever we have an abstraction from $\langle E, \text{po} \rangle$ to $\langle E', \text{po}' \rangle$ and we know the implementation $\langle E, \text{po} \rangle$ has a Λ -consistent execution \mathcal{G} , then the abstract plain execution $\langle E', \text{po}' \rangle$ also has an L -consistent execution (third point) and the implementation respects the synchronisation promises made by the abstract library L (first and second point).

To translate the synchronisation promises, we require a *concretisation function* g that maps every subevent of the abstraction to a subevent in their implementation. The library L makes two kinds of synchronisation promises: **to** (via stamps) and **so'**. If we have $\langle s'_1, s'_2 \rangle \in \text{so}'$ in the abstraction, then we require that the concretisation of s'_1 synchronises with the concretisation of s'_2 , i.e. we ask that $g(\text{so}') \subseteq \text{hb}$.

Whenever the abstraction contains a subevent of the form $\langle e', a' \rangle$, the usage of the stamp a' carries an obligation. The subevent promises to synchronise with *any* earlier or later subevent, not necessarily from library L , according to the **to** relation (cf. Fig. 10 for RDMA). In most cases, the concretisation uses the same stamp, i.e. $g(\langle e', a' \rangle) = \langle e, a \rangle$ with $a' = a$, and the property is trivially respected by the implementation with $\langle e_1, a_1 \rangle = \langle e_2, a_2 \rangle = \langle e, a \rangle$. Otherwise we have $a' \neq a$, and so for any earlier (resp. later) stamp a_0 that a' should synchronise with, we need to justify this synchronisation happens in the implementation, i.e. that we have $\langle e_1, a_1 \rangle \xrightarrow{\text{hb}^*} \langle e, a \rangle$, where a_1 can perform the expected stamp synchronisation $\langle a_0, a_1 \rangle \in \text{to}$.

An important point to note is that **hb** is potentially bigger than $(\text{ppo} \cup \text{so})^+$. In which case, we need to prove the result for any reasonable **hb'** bigger than $(\text{ppo}' \cup \text{so}')^+$. Thus local soundness states that if the implementation has a Λ -consistent execution *with additional constraints*, then the abstraction similarly has an L -consistent execution *with these additional constraints*. This is required for the implementation to work in any context, i.e. for programs using L in conjunction to other libraries, as expressed by the following theorem.

THEOREM 3.14. *If a well-defined implementation is locally sound, then it is sound.*

PROOF. See the extended version [EV]. □

In the case of the shared variable library, we can use this proof technique to verify the implementation I_{SV} .

THEOREM 3.15. *I_{SV} is locally sound, and hence I_{SV} is sound.*

PROOF. See the extended version [EV]. □

4 Barrier Library

As discussed informally in §2.2, LOCO implements a barrier library (BAL), which supports synchronisation of threads across multiple threads. Note that each barrier corresponds to a set of threads, which we refer to as the “participating threads” of a barrier. Each participating thread must wait for *all* operations towards *all* participating threads (including its own) that are po-before each barrier to be completed. We first present a generic specification for barriers with participating nodes in §4.1, and the LOCO barrier and its correctness proof in §4.2. In §4.3 we discuss an issue with such a barrier that only synchronises participating nodes and a possible fix.

4.1 Generic Barrier Specification

The barrier library (BAL) only has the single method $\text{BAR}_{\text{BAL}} : \text{Loc} \rightarrow ()$, taking a location as an input and producing no output. Thus, we have $\text{loc}(\text{BAR}_{\text{BAL}}(x)) = \{x\}$. The input location x defines the set of threads that synchronise via $\text{BAR}_{\text{BAL}}(x)$. In our model, we assume a function $b : \text{Loc} \rightarrow \mathcal{P}(\text{Tid})$ associating each location x with a set of threads that perform a barrier synchronisation on x .

While the LOCO barrier implementation (see §4.2) supports synchronisation across nodes connected by RDMA, our specification is more general and abstracts away the notion of nodes. Instead, our library defines synchronisation between *threads*, providing freedom to implement different synchronisation mechanisms depending on whether the threads are on the same or on different nodes.

Since MOWGLI allows libraries to be defined in isolation, we only consider E containing barrier calls. Let $E_x \triangleq \{e \in E \mid \text{loc}(e) = \{x\}\}$ denote the set of barrier calls on the location x .

Definition 4.1 (BAL-consistency). We say that $\mathcal{G} = \langle E, \text{po}, \text{stmp}, \text{so}, \text{hb} \rangle$ is BAL-consistent iff:

- $\text{stmp} = \text{stmp}_{\text{BAL}}$, defined as $\text{stmp}_{\text{BAL}}(\langle _, _, \langle \text{BAR}_{\text{BAL}}, (x), () \rangle \rangle) = \bigcup_{t \in b(x)} \{\text{aGF}_{n(t)}\} \cup \{\text{aCR}\}$;
- for all x and $e \in E_x$, $t(e) \in b(x)$; i.e. non-participating threads do not participate;
- for all $x \in \text{Loc}$, there is an integer c_x such that for all thread $t \in b(x)$ we have $\#(E_x|_t) = c_x$; i.e. each participating thread makes exactly c_x calls to the barrier on x ;
- there is an ordering function $o : E \rightarrow \mathbb{N}$ such that for all location x :
 - if $e \in E_x$ then $1 \leq o(e) \leq c_x$;
 - if $e_1, e_2 \in E_x$ and $\langle e_1, e_2 \rangle \in \text{po}$ then $o(e_1) < o(e_2)$; and
- $\text{so} = \bigcup_{x \in \text{Loc}} \bigcup_{1 \leq i \leq c_x} \{ \langle \langle e_1, \text{aGF}_n \rangle, \langle e_2, \text{aCR} \rangle \rangle \mid e_1, e_2 \in (E_x \cap o^{-1}(i)) \}$

This predicate clearly respects monotonicity (since **hb** is unrestricted) and decomposability (since each location is treated independently).

The function o associates each barrier call to the number of times the location has been used by this thread, in program order. We say that e_1 and e_2 *synchronise together* iff $\text{loc}(e_1) = \text{loc}(e_2)$ and $o(e_1) = o(e_2)$. The stamps of the form aGF correspond to the *entry points* of the barrier calls, waiting for previous operations to finish before the synchronisation. The stamp aCR represents the *exit point* of the barrier, after the synchronisation. The synchronisation is then an **so** ordering between aGF and aCR for barrier calls that synchronise together.

4.2 LOCO Implementation

Given $b : \text{Loc} \rightarrow \mathcal{P}(\text{Tid})$, for each location x with $b(x) = \{t_1, \dots, t_k\}$ synchronising k threads, we create a set of k shared variables (i.e. sv locations) $\{x_{t_1}, \dots, x_{t_k}\}$. Each shared variable x_t is used as a counter indicating how many times thread t has executed a barrier on x . The LOCO implementation decomposes the barrier into three steps: (1) wait for previous operations to finish; (2) increase your counter; (3) wait for the counters of other threads to increase. We define the implementation I_{BAL}^b in Fig. 11. Clearly, the implementation is well defined: it cannot return a break number greater than zero, since all break commands have a break number of 1 and are inside loops; and every succeeding implementation generates at least one event.

For $t \notin b(x)$: $I_{\text{BAL}}^b(t, \text{BAR}_{\text{BAL}}(x)) \triangleq \text{loop } \{ () \}$

For $t \in b(x) = \{t_1, \dots, t_k\}$:

```

 $I_{\text{BAL}}^b(t, \text{BAR}_{\text{BAL}}(x)) \triangleq$ 
  let  $s_n = \{n(t_i) \mid t_i \in b(x)\}$  in
  GFsv( $s_n$ );
  let  $v = \text{Read}_{\text{sv}}(x_t)$  in
  Writesv( $x_t, v + 1$ );
  Bcastsv( $x_t, \_, (s_n \setminus \{n(t)\})$ );
  loop {
    let  $v' = \text{Read}_{\text{sv}}(x_{t_1})$  in
    if  $v' > v$  then break1() else () ;
    ...
  }
  loop {
    let  $v' = \text{Read}_{\text{sv}}(x_{t_k})$  in
    if  $v' > v$  then break1() else () ;
  }

```

Fig. 11. I_{BAL}^b implementation

If a method call is made by a non-participating thread, the call is invalid and we implement it using a non-terminating loop. This is necessary for soundness, as the outcomes of the implementation must be valid, and in this situation the BAL specification does not allow any valid outcomes.

If a method call is made by a participating thread t , the implementation starts with a global fence ensuring any previous operation towards any relevant node is fully finished. Then, it increments its counter x_t to indicate to other threads that the barrier has been reached and executed. The value of x_t is immediately available to other threads on the same node, and is made available to other participating nodes using a broadcast. Note that the broadcast does not perform a loopback (i.e. we exclude $n(t)$ from the targets), as asking the NIC to overwrite x_t with itself might cause the new value of a later barrier call to be reverted to the current value. Then, we repeatedly read the (local) values of the other counters x_{t_i} and wait for each of them to indicate other threads have reached their matching barrier call. Note that there is no reason to wait for the broadcast to finish: the implementation on t might go ahead before other threads are aware that t reached the barrier, but that does not break the guarantees provided by the barrier.

THEOREM 4.2. *The implementation I_{BAL}^b is locally sound.*

PROOF. See the extended version [EV]. □

4.3 Supporting Transitivity

The barrier semantics in §4.1 only performs a global fence on nodes with participating threads. While this appears intuitive and reduces assumptions about other nodes, barrier synchronisation using such a library is *not* transitive. For example, consider the program in Fig. 12. Since $\bar{x} := 1$ is an operation towards node 3, the barrier $\text{BAR}_{\text{BAL}}(b_1)$ does not wait for it to finish, allowing $a = 0$.

Such a transitive barrier can straightforwardly be obtained by synchronising across *all* nodes, instead of just “participating” threads. For the specification, we define $\text{stmp}_{\text{BAL}}(\langle _, _, \langle \text{BAR}_{\text{BAL}}, (x), () \rangle \rangle) = \bigcup_{n \in \text{Node}} \{\text{aGF}_n\} \cup \{\text{aCR}\}$ and for the implementation, we define $I_{\text{BAL}}^b(t, \text{BAR}_{\text{BAL}}, (x)) \triangleq \text{let } s_n = \text{Node in } \dots$. This stronger version is the one implemented in LOCO (see Fig. 7).

		$x = 0$
$\bar{x} := 1$	$\text{BAR}_{\text{BAL}}(b_1)$	$\text{BAR}_{\text{BAL}}(b_2)$
	$\text{BAR}_{\text{BAL}}(b_1)$	$a := x$

$a = 0 \checkmark$

Fig. 12. Allowed weak barrier behaviour

5 Ring Buffer Library

The ring buffer library (RBL) provides methods for a single-writer-multiple-reader FiFo queue for messages of any size, where each message is duplicated as necessary and can be read once by each reader. Here, we present its specification (§5.1), and an implementation and correctness proof (§5.2).

5.1 Ring Buffer Specification

The ring buffer library has two methods $\text{Submit}^{\text{RBL}} : \text{Loc} \times \text{Val}^* \rightarrow \mathbb{B}$ and $\text{Receive}^{\text{RBL}} : \text{Loc} \rightarrow \text{Val}^* \uplus \{\perp\}$, with $\text{loc}(\text{Submit}^{\text{RBL}}(x, _)) = \text{loc}(\text{Receive}^{\text{RBL}}(x)) = \{x\}$. $\text{Submit}^{\text{RBL}}(x, \tilde{v})$ tries to add a new message \tilde{v} to the ring buffer x . It can either fail if the ring buffer is full, returning false, or succeed returning true. $\text{Receive}^{\text{RBL}}(x)$ tries to read a message from the ring buffer x . It can either succeed if there is at least one pending message, returning the next message, or fail if there is no pending messages, returning \perp .

In our model, we assume two functions $\text{wthd} : \text{Loc} \rightarrow \text{Tid}$ and $\text{rthd} : \text{Loc} \rightarrow \mathcal{P}(\text{Tid})$ associating each location x with a writing thread $\text{wthd}(x)$ and a set of reader threads $\text{rthd}(x)$. For subevents, we define the stamping function stmp_{RBL} as follows:

$$\begin{aligned} \text{stmp}_{\text{RBL}}(\langle t, _, \langle \text{Submit}^{\text{RBL}}, (x, _), \text{true} \rangle \rangle) &\triangleq \{\text{aNRW}_{n(t')} \mid t' \in \text{rthd}(x) \wedge n(t') \neq n(t)\} \cup \{\text{aCW}\} \\ \text{stmp}_{\text{RBL}}(\langle t, _, \langle \text{Submit}^{\text{RBL}}, (x, _), \text{false} \rangle \rangle) &\triangleq \{\text{aWT}\} \end{aligned}$$

$$\text{stmp}_{\text{RBL}}(\langle _, _, \langle \text{Receive}^{\text{RBL}}, (x), \widetilde{v} \rangle \rangle) \triangleq \{\text{aCR}\}$$

$$\text{stmp}_{\text{RBL}}(\langle _, _, \langle \text{Receive}^{\text{RBL}}, (x), \perp \rangle \rangle) \triangleq \{\text{aWT}\}$$

A successful call to $\text{Submit}^{\text{RBL}}$ (with return value true) is denoted by a write stamp for each relevant node: the stamp aCW is used by the writer node, and the stamps $\text{aNRW}_{n(t')}$ are used by the corresponding remote nodes. Failing calls (with return value false or \perp) are depicted by the stamp aWT . Finally, a succeeding $\text{Receive}^{\text{RBL}}$ call uses the reading stamp aCR .

We note different sets corresponding to calls to $\text{Submit}^{\text{RBL}}$ succeeding (\mathcal{W}) and calls to $\text{Receive}^{\text{RBL}}$ failing (\mathcal{F}) or succeeding (\mathcal{R}). Calls to $\text{Submit}^{\text{RBL}}$ failing are ignored by the specification.

$$\begin{aligned} \mathcal{W}_x^n &\triangleq \{ \langle e, \text{aNRW}_n \rangle \mid e = \langle t, _, \langle \text{Submit}^{\text{RBL}}, (x), _ \rangle, \text{true} \rangle \in E \wedge \text{aNRW}_n \in \text{stmp}_{\text{RBL}}(e) \} \\ &\cup \{ \langle e, \text{aCW} \rangle \mid e = \langle t, _, \langle \text{Submit}^{\text{RBL}}, (x), _ \rangle, \text{true} \rangle \in E \wedge n(t) = n \} \\ \mathcal{F}_x^n &\triangleq \{ \langle e, \text{aWT} \rangle \mid e = \langle t, _, \langle \text{Receive}^{\text{RBL}}, (x), \perp \rangle \rangle \in E \wedge n(t) = n \} \\ \mathcal{R}_x^n &\triangleq \{ \langle e, \text{aCR} \rangle \mid e = \langle t, _, \langle \text{Receive}^{\text{RBL}}, (x), \widetilde{v} \rangle \rangle \in E \wedge n(t) = n \} \end{aligned}$$

We then define the reads-from relation rf matching successful $\text{Submit}^{\text{RBL}}$ and $\text{Receive}^{\text{RBL}}$ events.

Definition 5.1. Given $\mathcal{G} = \langle E, \text{po}, \text{stmp}_{\text{RBL}}, _, _ \rangle$, we say that rf is *well-formed* iff each of the following holds:

- (1) $\text{rf} = \bigcup_{n,x} \text{rf}_x^n$ with $\text{rf}_x^n \subseteq \mathcal{W}_x^n \times \mathcal{R}_x^n$
- (2) rf_x^n is total and functional on its range, i.e. each read subevent in \mathcal{R}_x^n is related to exactly one write subevent in \mathcal{W}_x^n .
- (3) If $(\langle _, _, \langle \text{Submit}^{\text{RBL}}, (x), \widetilde{v} \rangle, \text{true} \rangle, a) \xrightarrow{\text{rf}} (\langle _, _, \langle \text{Receive}^{\text{RBL}}, (x), \widetilde{v}' \rangle, a' \rangle)$ then $\widetilde{v} = \widetilde{v}'$, i.e. related events write and read the same tuple of values.
- (4) If $\langle s_1, s_2 \rangle \in \text{rf}$, $\langle s_1, s_3 \rangle \in \text{rf}$, and $s_2 \neq s_3$, then $t(s_2) \neq t(s_3)$, i.e. each thread can read each message at most once.
- (5) If $s_1, s_2 \in \mathcal{W}_x^n$, $\langle s_1, s_2 \rangle \in \text{po}$, and $\langle s_2, s_4 \rangle \in \text{rf}$, then there is s_3 such that $\langle s_1, s_3 \rangle \in \text{rf}$, and $\langle s_3, s_4 \rangle \in \text{po}$, i.e. threads cannot jump a message.

We define the *fails-before* relation fb expressing that a failing $\text{Receive}^{\text{RBL}}$ occurs before a succeeding $\text{Submit}^{\text{RBL}}$ as follows:

$$\text{fb} \triangleq \bigcup_{n,x} (\mathcal{F}_x^n \times \mathcal{W}_x^n \setminus (\text{po}^{-1}; \text{rf}^{-1}))$$

If $s_1 \in \mathcal{W}_x^n$ and $s_3 \in \mathcal{F}_x^n$, then the contents written by s_1 is not available when s_3 is executed. Either there is s_2 such that $\langle s_1, s_2 \rangle \in \text{rf}$ and $\langle s_2, s_3 \rangle \in \text{po}$, in which case the message has been read; or there is no such s_2 and we have $\langle s_3, s_1 \rangle \in \text{fb}$ to express that the message was not yet written.

Definition 5.2 (RBL-consistency). We say that an execution $\mathcal{G} = \langle E, \text{po}, \text{stmp}_{\text{RBL}}, \text{so}, \text{hb} \rangle$ is **BAL**-consistent iff:

- if $\langle t, _, \langle \text{Submit}^{\text{RBL}}, (x), _ \rangle, _ \rangle \in E$ then $t = \text{wthd}(x)$; and if $\langle t, _, \langle \text{Receive}^{\text{RBL}}, (x), _ \rangle \rangle \in E$ then $t \in \text{rthd}(x)$; and
- there exists a well-formed rf such that $\text{so} = \text{rf} \cup \text{fb}$.

Note that this definition allows the writer thread to also be a reader, and nodes to have multiple reading threads. Moreover, the consistency predicate does not tell us anything about failing writes; they may fail spuriously.

```

For  $wthd(x) = t \wedge rthd(x) = \{t_1; \dots; t_k\}$  :
 $I_{S,RBL}^{wthd,rthd}(t, Submit^{RBL}, (x, \tilde{v} = (v_1, \dots, v_V)) \triangleq$ 
  let  $s_n = \{n(t_i) \mid t_i \in rthd(x)\} \setminus \{n(t)\}$  in
  let  $V = \text{len}(\tilde{v})$  in
  let  $H = \text{Read}_{sv}(h^x)$  in
  let  $H_1 = \text{Read}_{sv}(h_{t_1}^x)$  in
  ...
  let  $H_k = \text{Read}_{sv}(h_{t_k}^x)$  in
  let  $M = \min(\{H_1, \dots, H_k\})$  in
  if  $(H - M) + (V + 1) > S$  then false else {
    Writesv( $x_{H \% S}, V$ ); Bcastsv( $x_{H \% S}, \_, s_n$ );
    Writesv( $x_{(H+1) \% S}, v_1$ ); Bcastsv( $x_{(H+1) \% S}, \_, s_n$ );
    ...
    Writesv( $x_{(H+V) \% S}, v_V$ ); Bcastsv( $x_{(H+V) \% S}, \_, s_n$ );
    Waitsv( $d_x$ );
    Writesv( $h^x, H + V + 1$ ); Bcastsv( $h^x, d_x, s_n$ );
    true };

For  $t \notin rthd(x)$ :
 $I_{S,RBL}^{wthd,rthd}(t, Receive^{RBL}, (x)) \triangleq \text{loop } \{()\}$ 

For  $t \in rthd(x)$ :
 $I_{S,RBL}^{wthd,rthd}(t, Receive^{RBL}, (x)) \triangleq$ 
  let  $H = \text{Read}_{sv}(h_t^x)$  in
  let  $H' = \text{Read}_{sv}(h^x)$  in
  if  $H \geq H'$  then  $\perp$  else {
    let  $V = \text{Read}_{sv}(x_{H \% S})$  in
    let  $v_1 = \text{Read}_{sv}(x_{(H+1) \% S})$  in
    ...
    let  $v_V = \text{Read}_{sv}(x_{(H+V) \% S})$  in
    Writesv( $h_t^x, H + V + 1$ );
    if  $n(wthd(x)) = n(t)$  then  $()$  else
      { Bcastsv( $h_t^x, \_, \{n(wthd(x))\}$ ) };
    ( $v_1, \dots, v_V$ ) };

```

Fig. 14. Implementation $I_{S,RBL}^{wthd,rthd}$ of the ring buffer library into sv

Alternative weaker semantics. Instead of requiring $\text{so} = \text{rf} \cup \text{fb}$, we could give an alternative specification with $\text{so} = \text{rf}$ and $\text{hb}^{-1} \cap \text{fb} = \emptyset$. The latter says that you still cannot ignore (fb) a write that you know (hb) has finished; but if you do ignore a write, you do not have to export the guarantee (so) that the write has not finished. For instance, take the litmus test in Fig. 13. With the semantics in Def. 5.2, at least one of the two Receive^{RBL} has to succeed. With the weaker semantics, they are allowed to both fail, even when both Submit^{RBL} calls succeed.

$a := \text{Submit}^{RBL}(x, 1)$	$b := \text{Submit}^{RBL}(y, 1)$
$\text{GF}_{sv}(\{n_2\})$	$\text{GF}_{sv}(\{n_1\})$
$c := \text{Receive}^{RBL}(y)$	$d := \text{Receive}^{RBL}(x)$
$(a, b, c, d) = (\text{true}, \text{true}, \perp, \perp) \times$	

Fig. 13. Alternative ring buffer semantics

5.2 LOCO Implementation

As before, we assume given the functions $wthd : \text{Loc} \rightarrow \text{Tid}$ and $rthd : \text{Loc} \rightarrow \mathcal{P}(\text{Tid})$. We also assume an integer S representing the size of the ring buffer. We implement the ring buffer library (RBL) using the shared variable library (sv). For each location x with $rthd(x) = \{t_1, \dots, t_k\}$ we create the shared variable (i.e. sv locations) x_0, \dots, x_{S-1} for the content of the buffer, as well as shared variables h^x for the writer and $h_{t_1}^x; \dots; h_{t_k}^x$ for the readers. We also use a work identifier d_x .

Events that do not respect $rthd$ or $wthd$ are implemented using an infinite loop (i.e. $\text{loop } \{()\}$), similarly to other implementations. Otherwise, we use the implementation $I_{S,RBL}^{wthd,rthd}$ given in Fig. 14, where $\%$ represents the modulo operation.

The value of h^x represents the next place to write for the writing thread. The value of $h_{t_i}^x$ represents the next place thread t_i needs to read. If $h^x = h_{t_i}^x$ then thread t_i is up-to-date and needs to wait for the writer to send additional data. If the difference between h^x and $h_{t_i}^x$ gets close to S , then the buffer is full and the writer cannot send any more data.

In the implementation of Submit^{RBL} , the value M represents the minimum of all $h_{t_i}^x$. As such, $(H - M)$ represents the amount of space currently in use. Since $(V + 1)$ represents the number of cells necessary to submit a new message (the size V itself is also submitted), we can proceed if $H - M + V + 1 \leq S$, i.e. if there is enough free space.

Since, for a specific remote node, the broadcasts complete in order, when a reader sees the new value of h^x it means the written data is available. We need to take care that the broadcast of h^x must read from the write of the *same* function call, and not from the write of a later call to `SubmitRBL`. Otherwise, the value of h^x for the second submit might be available to readers before the data of the second submit. For this, we simply need to wait for the broadcast of previous function calls, using `Waitsv(dx)`, before modifying h^x .

When thread t_i wants to receive, it only proceeds if $h^x > h_{t_i}^x$, otherwise t_i is up-to-date and returns \perp . After reading a message, the reader updates $h_{t_i}^x$ to signal to the writer the space of the message is no longer in use. If the reader is on the same node as the writer, there is no need for a broadcast, otherwise the reader broadcasts to the node of the writer.

With this implementation, each participating node possesses only one copy of the data, and potentially multiple readers per node can read from the same memory locations.

THEOREM 5.3. *The implementation $I_{S,RBL}^{\text{wthd},\text{rthd}}$ is locally sound.*

PROOF. See the extended version [EV]. □

6 Evaluation

In this section, we explore the performance of our LOCO primitives, then use them to build a high performance key-value store. Further applications can be found in the extended version [EV].

All results were collected using c6525 – 25g nodes on the Cloudlab platform [clo [n. d.]]. These machines each have a 16-core AMD 7302P CPU, running Ubuntu 22.04. Nodes communicate over a 25 Gbps Ethernet fabric using Mellanox ConnectX-5 NICs.

6.1 LOCO Primitives

First, we compare the performance of the verified barrier (BAL) and ring buffer (RBL) primitives to equivalent operations in OpenMPI [Gabriel et al. 2004], a message-passing library commonly used to build distributed applications. We compare against OpenMPI 5.0.5, using the PML/UCX backend for RoCE support. Results are shown in Figure 15.

For the barrier experiments, we compare to the `MPI_Barrier` operation, varying both thread count per node and node count. The MPI barrier does not actually provide synchronization, expecting the user to instead appropriately track and fence operations before using the primitive. We compare the barrier to our LOCO barrier, both with and without the synchronization fence, and show that the LOCO barrier with equivalent semantics (no fence) performs as well or better than the MPI barrier. Note the MPI barrier dynamically switches between several internal algorithms adjusting to load leading to non-smooth performance across the test domain.

For the ring buffer experiments, we compare a ring buffer broadcast to the `MPI_Ibcast` (non-blocking broadcast) operation. We measure across different node counts and amounts of “network load”, that is, the number n of outstanding broadcast operations in the network, along with total node count. A single node acts as the sender: it starts by sending n broadcasts, then sends a new one every time a prior message completes. All other nodes wait to receive and acknowledge messages. Messages have a fixed size of 64 bytes. Here, we find that the formally verified LOCO ring buffer provides better broadcast performance than MPI in most configurations, with MPI performance falling drastically as the number of outstanding messages rises.

6.2 Example Application: A Key-Value Store

Beyond our microbenchmarks, we describe an example LOCO application: a key-value store, built using composable LOCO primitives.

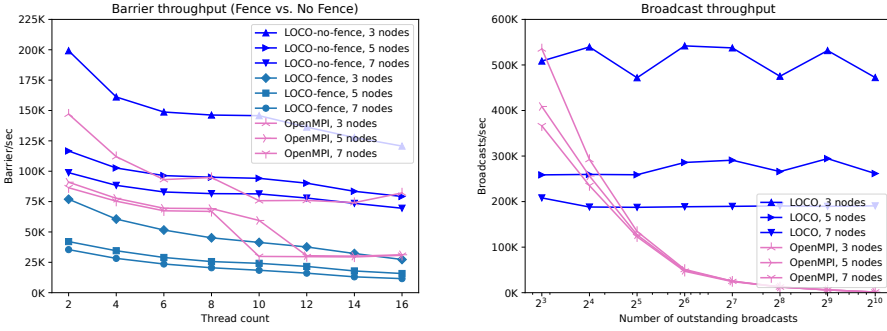


Fig. 15. Comparison of barrier and broadcast operations for LOCO and OpenMPI.

Our kvstore object is a distributed key-value store with a lookup operation that takes no locks, and insertion, deletion, and update operations protected by locks. Lookup and update are depicted in Fig. 16. Each node allocates a remotely-accessible memory region that is used to store values and consistency metadata (a checksum for atomicity, a counter for garbage collection, and a valid bit).

Each node also maintains a local index (a C++ unordered_map), protected by a local reader-writer lock, which records the locations of all keys in the kvstore as (node_id, array_index) pairs, along with a counter matching the one stored with the data. The kvstore is linearisable, with a proof given in the extended version [EV] – our proof is simplified by leveraging the compositional properties of LOCO. Note that RDMA^{TSO} does not have a semantics for locks or RDMA read-modify-write operations, which means that this proof currently does not use Mowgli. We consider an extension of RDMA^{TSO} with synchronisation operations (and hence a full proof of kvstore) to be future work. Almost all RDMA maps [Barthels et al. 2015; Kalia et al. 2014; Li et al. 2023; Lu et al. 2024; Wang et al. 2022] lack any formal safety specification (we are only aware of two [Dragojević et al. 2014], [Alquraan et al. 2024]), likely due to difficulties in encapsulation, which the LOCO philosophy solves.

We compared our key-value store design against Sherman [she [n. d.]; Wang et al. 2022] and the MicroDB from Scythe [scy [n. d.]; Lu et al. 2024], two state-of-the-art RDMA key-value stores. We also compare against Redis-cluster [Ltd. 2021] as a non-RDMA baseline. Results are shown in Figure 17. We measured throughput on read-only, mixed read-write, and write-only operation distributions, across both uniform and Zipfian ($\theta = 0.99$) key distributions, and across different node counts and per-node thread counts. Each data point is the geometric mean of 5 runs with a 20 second duration, not including prefill.

All benchmarks use a 10MB keyspace, filled to 80% capacity with 64-bit keys and values. All benchmarks use the CityHash64 key hashing function [Pike and Alakuijala [n. d.]], and the YCSB-C implementation of a Zipfian distribution [ycs [n. d.]].

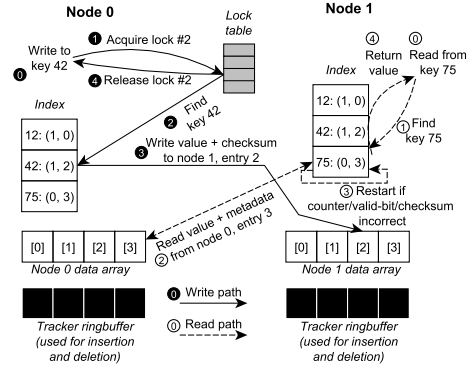


Fig. 16. kvstore read and write operations

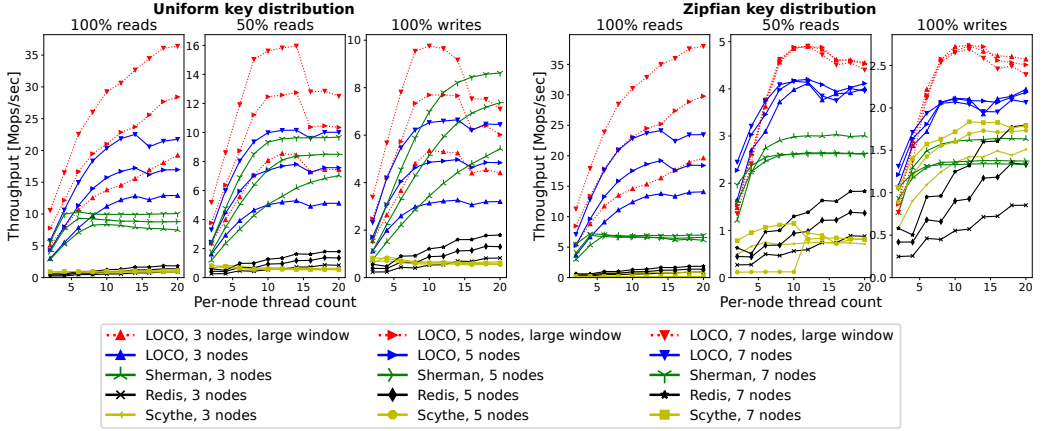


Fig. 17. Throughput comparison of key-value stores.

We modified Sherman to issue a fence (GF_{sv}) between lock-protected writes and lock releases to solve a bug related to consistency issues. Our kvstore also issues a fence for the same reason. For both, this fence incurs a 15% overhead.

For LOCO, Sherman, and Redis, write operations are updates. For Scythe, we found that stressing update operations led to program instability and very low throughput, so we use the performance of insertion operations as an upper bound on write performance. For Redis, we configure a cluster with no replication or persistence. Since each Redis server instance uses 4 threads, we create $\text{ceil}(\text{num_threads}/4)$ server instances for a given thread count. We use Memtier [Ltd. 2024] as a benchmark client. Each node runs a single Memtier instance with threads equal to the thread count, and 128 clients per thread (matching the LOCO large window size).

In addition, all systems expose a parameter we call the *window size*, which specifies the maximum number of outstanding operations per application thread (note this is not a batch size – each operation is started and completed individually). Increasing LOCO’s window size to 128 yielded significant improvement (the “large window” series). However, increasing Sherman’s and Scythe’s window sizes appeared to cause internal errors, so the main results for all systems except Redis (see above) use a window size of 3 for accurate comparison.

LOCO outperforms Sherman on read-only configurations. We believe this is because Sherman reads whole sections of the tree from remote memory, while the LOCO design looks up the location locally and only remotely reads the value. On the other hand, LOCO’s advantage over Sherman for Zipfian writes likely comes from the better performance under contention.

Sherman outperforms LOCO (with a window size of 3) on mixed read-write and write-only distributions on uniform keys, while the reverse is true for Zipfian keys. Sherman’s advantage here is likely due to the fact that, unlike LOCO, Sherman colocates locks with data, allowing them to issue lock releases in a batch with writes.

7 Related and Future Work

Although the formal semantics of RDMA has only recently been established [Ambal et al. 2024], our work is able to take advantage of earlier results in weak memory hardware [Alglave et al. 2014; Flur et al. 2016] and programming languages [Batty et al. 2011; Lahav et al. 2017]. We do not provide their details here since they are rather expansive.

RDMA Semantics. Prior works on RDMA semantics include coreRMA [Dan et al. 2016] (which formalises RDMA over the SC memory model) and RDMA^{TSO} [Ambal et al. 2024], a more realistic formal model that is very close to the Verbs library [linux-rdma 2018], describing the behaviour of RDMA over TSO. These semantics are however low-level and are difficult for programmers to use directly, as illustrated by examples such as those in Fig. 2.

RDMA Libraries. Much prior work in RDMA focuses on *upper-level primitives*, e.g. consensus protocols [Aguilera et al. 2019, 2020; Izraelevitz et al. 2023; Jha et al. 2019; Poke and Hoefler 2015], distributed maps or databases [Alquraan et al. 2024; Barthels et al. 2015; Dragojević et al. 2014, 2015; Gavrielatos et al. 2020; Kalia et al. 2014; Li et al. 2023; Wang et al. 2022], graph processing [Wang et al. 2023a], distributed learning [Ren et al. 2017; Xue et al. 2019], stand-alone data structures [Brock et al. 2019; Devarajan et al. 2020], disaggregated scheduling [Ruan et al. 2023a,b] or file systems [Yang et al. 2019, 2020]. These works focus on the final application, rather than considering the programming model as its own, partitionable problem. As a result, the intermediate library between RDMA and the exported primitive is usually ad-hoc and tightly coupled to the application, or effectively non-existent. In general, these applied, specific, projects manage raw memory explicitly statically allocated to particular nodes, use ad-hoc atomicity and consistency mechanisms, and do not consider the possibility of primitive reuse. This design is not a fundamentally flawed approach, but it does raise the possibility of a better mechanism, which likely could underlie all the above solutions.

Some works have considered this intermediate layer explicitly, however, the general approach for this intermediate layer has been to encapsulate local and remote memory as *distributed shared memory*, that is, a flat, uniform, coherent, and consistent address space hiding the relaxed consistency and non-uniform performance of the underlying RDMA network. These works generally focus on transparently (or mostly-transparently [Ruan et al. 2020; Zhang et al. 2022]) porting existing shared memory applications. We argue that this technique, either with purely software-based virtualisation [Cai et al. 2018; Gouk et al. 2022; Ruan et al. 2020; Wang et al. 2020; Zhang et al. 2022], or by extending hardware [Calciu et al. 2021], is unlikely to gain traction because the performance will always be worse than an approach which takes into account the underlying memory network.

Other programming models have simply used RDMA to implement existing distributed system abstractions. For example, both MPI [Message Passing Interface Forum 2023] and NCCL [NVIDIA Corporation 2020] can use RDMA for inter-node communication. However, fundamentally, these are *message passing programming models* with explicit send and receive primitives. While MPI does support some remote memory accesses, this support is best seen as a zero-copy send/receive mechanism where synchronisation is either coarse-grained and inflexible, or simply nonexistent. While message-passing is well-suited for dataflow applications (e.g. machine learning and signal processing) and highly parallel scale-out workloads (e.g. physical simulation), it is less useful for workloads that exhibit data-dependent communication [Liu et al. 2021], such as transaction processing or graph computations. In these applications, cross-node synchronisation is unavoidable and unpredictable, so the ideal performance strategy shifts from simply avoiding synchronisation to minimising contention, accelerating synchronisation use, and reducing data movement.

Compared to prior art, LOCO aims to build composable, reusable, and performant primitives for complicated memory networks, suitable for irregular workloads. No such option currently exists in the literature.

Verification. Our proofs have followed the declarative style [Raad et al. 2019; Stefanescu et al. 2024] enabling modular verification. RDMA^{TSO} [Ambal et al. 2024] also includes an operational model, which could form a basis for a program logic (e.g., [Bila et al. 2022; Lahav et al. 2023]), ultimately enabling operational abstractions and proofs of refinement [Dalvandi and Dongol 2022]. Other modular approaches include modular proofs through separation logics [Jung et al. 2018],

but this additionally requires a separation logic encoding of the $\text{RDMA}^{\text{WAIT}}$ memory model (and an associated proof of soundness) before it can be applied to verify libraries such as LOCO. We consider operational proofs and those involving separation logic as a topic for future work.

Nagasamudram et al. [2024] have verified, in Rocq, key properties of a coordination service known as Derecho [Jha et al. 2018], which can be configured to run over RDMA. However, their proofs start with a very high-level model called a *shared-state table*, which is an array of shared variables (cf. Fig. 7). Unlike our work, these assumed shared state table semantics have not been connected to any formal RDMA semantics. In future work, it would be interesting to connect our work to middleware such as Derecho, ultimately leading to a fully verified RDMA application stack.

There is a rich literature of work around model checking under weak and persistent memory [Abdulla et al. 2023; Kokologiannakis and Vafeiadis 2021] including recent works that tackle refinement and linearisability [Golovin et al. 2025; Raad et al. 2024]. It would be interesting to know whether these techniques can be extended to support RDMA^{TSO} (and by extension $\text{RDMA}^{\text{WAIT}}$).

8 Conclusion

In this paper, we describe LOCO, a verified library for building composable and reusable objects in network memory and its associated proof system MOWGLI. Our results show that LOCO can expose the full performance of underlying network memory to applications, while simultaneously easing proof burden.

Acknowledgments

This work was partially funded by industry partner Genuen, which provides hardware validation services using the harness described in the extended version [EV]. Izraelevitz also privately contracted with this company to assist with commercialization efforts of this harness. Ambal is supported by the EPSRC grant EP/X037029/1 and Raad is supported by a UKRI fellowship MR/V024299/1, by the EPSRC grant EP/X037029/1, and by VeTSS. Dongol and Chockler are supported by EPSRC grants EP/Y036425/1, EP/X037142/1 and EP/X015149/1 and Royal Society grant IES\R1\221226. Dongol is additionally supported by EPSRC grant EP/V038915/1 and VeTSS. Vafeiadis is supported by ERC Consolidator Grant for the project “PERSIST” (grant agreement No. 101003349).

References

- [n. d.]. The CloudLab Manual: Hardware. ([n. d.]). <http://docs.cloudlab.us/hardware.html>.
- [n. d.]. Scythe. ([n. d.]). <https://github.com/PDS-Lab/Scythe>.
- [n. d.]. Sherman: A Write-Optimized Distributed B+Tree Index on Disaggregated Memory. ([n. d.]). <https://github.com/thustorage/Sherman>.
- [n. d.]. Yahoo! Cloud Serving Benchmark in C++. ([n. d.]). <https://github.com/basicthinker/YCSB-C>.
- Martín Abadi and Leslie Lamport. 1991. The Existence of Refinement Mappings. *Theor. Comput. Sci.* 82, 2 (1991), 253–284. [https://doi.org/10.1016/0304-3975\(91\)90224-P](https://doi.org/10.1016/0304-3975(91)90224-P)
- Parosh Aziz Abdulla, Mohamed Faouzi Atig, S. Krishna, Ashutosh Gupta, and Omkar Tuppe. 2023. Optimal Stateless Model Checking for Causal Consistency. In *Tools and Algorithms for the Construction and Analysis of Systems - 29th International Conference, TACAS 2023, Held as Part of the European Joint Conferences on Theory and Practice of Software, ETAPS 2022, Paris, France, April 22-27, 2023, Proceedings, Part I (Lecture Notes in Computer Science)*, Sriram Sankaranarayanan and Natasha Sharygina (Eds.), Vol. 13993. Springer, 105–125. https://doi.org/10.1007/978-3-031-30823-9_6
- Marcos K. Aguilera, Naama Ben-David, Rachid Guerraoui, Virendra Marathe, and Igor Zablotchi. 2019. The Impact of RDMA on Agreement. In *Proceedings of the 2019 ACM Symposium on Principles of Distributed Computing (PODC '19)*. Association for Computing Machinery, New York, NY, USA, 409–418. <https://doi.org/10.1145/3293611.3331601>
- Marcos K. Aguilera, Naama Ben-David, Rachid Guerraoui, Virendra J. Marathe, Athanasios Xygkis, and Igor Zablotchi. 2020. Microsecond Consensus for Microsecond Applications. In *14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20)*. USENIX Association, 599–616. <https://www.usenix.org/conference/osdi20/presentation/aguilera>

- Jade Alglave, Luc Maranget, and Michael Tautschnig. 2014. Herding Cats: Modelling, Simulation, Testing, and Data Mining for Weak Memory. *ACM Trans. Program. Lang. Syst.* 36, 2 (2014), 7:1–7:74. <https://doi.org/10.1145/2627752>
- Ahmed Alquraan, Sreeharsha Udayashankar, Virendra Marathe, Bernardo Wong, and Samer Al-Kiswani. 2024. LoLKV: the logless, line the logless, linearizable, RDMA-based key-value storage system arizable, RDMA-based key-value storage system. In *Proceedings of the 21st USENIX Symposium on Networked Systems Design and Implementation (NSDI'24)*. USENIX Association, USA, Article 3, 14 pages.
- Guillaume Ambal, Brijesh Dongol, Haggai Eran, Vasileios Klimis, Ori Lahav, and Azalea Raad. 2024. Semantics of Remote Direct Memory Access: Operational and Declarative Models of RDMA on TSO Architectures. *Proc. ACM Program. Lang.* 8, OOPSLA2 (2024), 1982–2009. <https://doi.org/10.1145/3689781>
- Guillaume Ambal, George Hodgkins, Mark Madler, Gregory Chockler, Brijesh Dongol, Joseph Izraelevitz, Azalea Raad, and Viktor Vafeiadis. 2025. A Verified High-Performance Composable Object Library for Remote Direct Memory Access (Extended Version). (2025). [arXiv:cs.PL/2510.10531](https://arxiv.org/abs/2510.10531) <https://arxiv.org/abs/2510.10531>
- Andrew W. Appel and Sandrine Blazy. 2007. Separation Logic for Small-Step cminor. In *Theorem Proving in Higher Order Logics, 20th International Conference, TPHOLs 2007, Kaiserslautern, Germany, September 10-13, 2007, Proceedings (Lecture Notes in Computer Science)*, Klaus Schneider and Jens Brandt (Eds.), Vol. 4732. Springer, 5–21. https://doi.org/10.1007/978-3-540-74591-4_3
- Claude Barthels, Simon Loesing, Gustavo Alonso, and Donald Kossmann. 2015. Rack-Scale In-Memory Join Processing using RDMA. In *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data (SIGMOD '15)*. Association for Computing Machinery, New York, NY, USA, 1463–1475. <https://doi.org/10.1145/2723372.2750547>
- Mark Batty, Scott Owens, Susmit Sarkar, Peter Sewell, and Tjark Weber. 2011. Mathematizing C++ concurrency. In *Proceedings of the 38th ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL 2011, Austin, TX, USA, January 26-28, 2011*, Thomas Ball and Mooly Sagiv (Eds.). ACM, 55–66. <https://doi.org/10.1145/1926385.1926394>
- Eleni Vafeiadi Bila, Brijesh Dongol, Ori Lahav, Azalea Raad, and John Wickerson. 2022. View-Based Owicki–Gries Reasoning for Persistent x86-TSO. In *Programming Languages and Systems, Ilya Sergey (Ed.)*. Springer International Publishing, Cham, 234–261.
- Benjamin Brock, Aydın Buluç, and Katherine Yelick. 2019. BCL: A Cross-Platform Distributed Data Structures Library. In *Proceedings of the 48th International Conference on Parallel Processing (ICPP '19)*. Association for Computing Machinery, New York, NY, USA, Article 102, 10 pages. <https://doi.org/10.1145/3337821.3337912>
- Qingchao Cai, Wentian Guo, Hao Zhang, Divyakant Agrawal, Gang Chen, Beng Chin Ooi, Kian-Lee Tan, Yong Meng Teo, and Sheng Wang. 2018. Efficient distributed memory management with RDMA and caching. *Proc. VLDB Endow.* 11, 11 (jul 2018), 1604–1617. <https://doi.org/10.14778/3236187.3236209>
- Irina Calciu, M. Talha Imran, Ivan Puddu, Sanidhya Kashyap, Hasan Al Maruf, Onur Mutlu, and Aasheesh Kolli. 2021. Rethinking software runtimes for disaggregated memory. In *Proceedings of the 26th ACM International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS '21)*. Association for Computing Machinery, New York, NY, USA, 79–92. <https://doi.org/10.1145/3445814.3446713>
- Sadegh Dalvandi and Brijesh Dongol. 2022. Implementing and verifying release-acquire transactional memory in C11. *Proc. ACM Program. Lang.* 6, OOPSLA2 (2022), 1817–1844. <https://doi.org/10.1145/3563352>
- Andrei Marian Dan, Patrick Lam, Torsten Hoefer, and Martin Vechev. 2016. Modeling and Analysis of Remote Memory Access Programming. *SIGPLAN Not.* 51, 10 (oct 2016), 129–144. <https://doi.org/10.1145/3022671.2984033>
- Hariharan Devarajan, Anthony Kougkas, Keith Bateman, and Xian-He Sun. 2020. HCL: Distributing Parallel Data Structures in Extreme Scales. In *2020 IEEE International Conference on Cluster Computing (CLUSTER)*. 248–258. <https://doi.org/10.1109/CLUSTER49012.2020.00035>
- Brijesh Dongol, Radha Jagadeesan, James Riely, and Alasdair Armstrong. 2018. On abstraction and compositionality for weak-memory linearisability. In *Verification, Model Checking, and Abstract Interpretation - 19th International Conference, VMCAI 2018, Los Angeles, CA, USA, January 7-9, 2018, Proceedings (Lecture Notes in Computer Science)*, Isil Dillig and Jens Palsberg (Eds.), Vol. 10747. Springer, 183–204. https://doi.org/10.1007/978-3-319-73721-8_9
- Aleksandar Dragojević, Dushyanth Narayanan, Orion Hodson, and Miguel Castro. 2014. FaRM: Fast Remote Memory. In *Proceedings of the 11th USENIX Conference on Networked Systems Design and Implementation (NSDI'14)*. USENIX Association, Berkeley, CA, USA, 401–414. <http://dl.acm.org/citation.cfm?id=2616448.2616486>
- Aleksandar Dragojević, Dushyanth Narayanan, Edmund B. Nightingale, Matthew Renzelmann, Alex Shamis, Anirudh Badam, and Miguel Castro. 2015. No Compromises: Distributed Transactions with Consistency, Availability, and Performance. In *Proceedings of the 25th Symposium on Operating Systems Principles (SOSP '15)*. ACM, New York, NY, USA, 54–70. <https://doi.org/10.1145/2815400.2815425>
- Shaked Flur, Kathryn E. Gray, Christopher Pulte, Susmit Sarkar, Ali Sezgin, Luc Maranget, Will Deacon, and Peter Sewell. 2016. Modelling the ARMv8 architecture, operationally: concurrency and ISA. In *Proceedings of the 43rd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages, POPL 2016, St. Petersburg, FL, USA, January 20 - 22, 2016*, Rastislav Bodik and Rupak Majumdar (Eds.). ACM, 608–621. <https://doi.org/10.1145/2837614.2837615>

- Edgar Gabriel, Graham E. Fagg, George Bosilca, Thara Angskun, Jack J. Dongarra, Jeffrey M. Squyres, Vishal Sahay, Prabhanjan Kambadur, Brian Barrett, Andrew Lumsdaine, Ralph H. Castain, David J. Daniel, Richard L. Graham, and Timothy S. Woodall. 2004. Open MPI: Goals, Concept, and Design of a Next Generation MPI Implementation. In *Proceedings, 11th European PVM/MPI Users' Group Meeting*. Budapest, Hungary, 97–104.
- Adithya Gangidi, Rui Miao, Shengbao Zheng, Sai Jayesh Bondu, Guilherme Goes, Hany Morsy, Rohit Puri, Mohammad Riftadi, Ashmitha Jeevaraj Shetty, Jingyi Yang, et al. 2024. Rdma over ethernet for distributed training at meta scale. In *Proceedings of the ACM SIGCOMM 2024 Conference*. 57–70.
- Vasilis Gavrielatos, Antonios Katsarakis, Vijay Nagarajan, Boris Grot, and Arpit Joshi. 2020. Kite: efficient and available release consistency for the datacenter. In *Proceedings of the 25th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP '20)*. Association for Computing Machinery, New York, NY, USA, 1–16. <https://doi.org/10.1145/3332466.3374516>
- Pavel Golovin, Michalis Kokologiannakis, and Viktor Vafeiadis. 2025. RELINCHE: Automatically Checking Linearizability under Relaxed Memory Consistency. *Proc. ACM Program. Lang.* 9, POPL (2025), 2090–2117. <https://doi.org/10.1145/3704906>
- Donghyun Gouk, Sangwon Lee, Miryeong Kwon, and Myoungsoo Jung. 2022. Direct Access, High-Performance Memory Disaggregation with DirectCXL. In *2022 USENIX Annual Technical Conference (USENIX ATC 22)*. USENIX Association, Carlsbad, CA, 287–294. <https://www.usenix.org/conference/atc22/presentation/gouk>
- Richard L Graham, Timothy S Woodall, and Jeffrey M Squyres. 2006. Open MPI: A flexible high performance MPI. In *Parallel Processing and Applied Mathematics: 6th International Conference, PPAM 2005, Poznań, Poland, September 11-14, 2005, Revised Selected Papers 6*. Springer, 228–239.
- R. Gupta, V. Tipparaju, J. Nieplocha, and D. Panda. 2002. Efficient barrier using remote memory operations on VIA-based clusters. In *Proceedings. IEEE International Conference on Cluster Computing*. 83–90. <https://doi.org/10.1109/CLUSTER.2002.1137732>
- Danny Hendler, Itai Incze, Nir Shavit, and Moran Tzafrir. 2010. Flat Combining and the Synchronization-Parallelism Tradeoff. In *Proceedings of the 22nd ACM Symposium on Parallelism in Algorithms and Architectures (SPAA '10)*. Santorini, Greece, 355–364.
- Maurice Herlihy and Jeannette M. Wing. 1990. Linearizability: A Correctness Condition for Concurrent Objects. *ACM Trans. Program. Lang. Syst.* 12, 3 (1990), 463–492. <https://doi.org/10.1145/78969.78972>
- George Hodgkins, Mark Madler, and Joseph Izraelevitz. 2025. LOCO: Rethinking Objects for Network Memory. (2025). arXiv:cs.DC/2503.19270 <https://arxiv.org/abs/2503.19270>
- Joseph Izraelevitz, Gaukas Wang, Rhett Hanscom, Kayli Silvers, Tamara Silbergleit Lehman, Gregory Chockler, and Alexey Gotsman. 2023. Acuerdo: Fast Atomic Broadcast over RDMA. In *Proceedings of the 51st International Conference on Parallel Processing (ICPP '22)*. Association for Computing Machinery, New York, NY, USA, Article 59, 11 pages. <https://doi.org/10.1145/3545008.3545041>
- Sagar Jha, Jonathan Behrens, Theo Gkountouvas, Matthew Milano, Weijia Song, Edward Tremel, Robbert Van Renesse, Sydney Zink, and Kenneth P. Birman. 2019. Derecho: Fast State Machine Replication for Cloud Services. *ACM Trans. Comput. Syst.* 36, 2, Article 4 (April 2019), 49 pages. <https://doi.org/10.1145/3302258>
- Sagar Jha, Jonathan Behrens, Theo Gkountouvas, Mae Milano, Weijia Song, Edward Tremel, Robbert van Renesse, Sydney Zink, and Kenneth P. Birman. 2018. Derecho: Fast State Machine Replication for Cloud Services. *ACM Trans. Comput. Syst.* 36, 2 (2018), 4:1–4:49. <https://doi.org/10.1145/3302258>
- Sagar Jha, Jonathan Behrens, Theo Gkountouvas, Matthew Milano, Weijia Song, Edward Tremel, Sydney Zink, Ken Birman, and Robbert Van Renesse. 2017. Building Smart Memories and High-speed Cloud Services for the Internet of Things with Derecho. In *Proceedings of the 2017 Symposium on Cloud Computing (SoCC '17)*. ACM, New York, NY, USA, 632–632. <https://doi.org/10.1145/3127479.3134597> Extended version available from www.cs.cornell.edu/ken/derecho-tocs.pdf.
- Ralf Jung, Robbert Krebbers, Jacques-Henri Jourdan, Ales Bizjak, Lars Birkedal, and Derek Dreyer. 2018. Iris from the ground up: A modular foundation for higher-order concurrent separation logic. *J. Funct. Program.* 28 (2018), e20. <https://doi.org/10.1017/S0956796818000151>
- Anuj Kalia, Michael Kaminsky, and David G. Andersen. 2014. Using RDMA Efficiently for Key-value Services. In *Proceedings of the 2014 ACM Conference on SIGCOMM (SIGCOMM '14)*. ACM, New York, NY, USA, 295–306. <https://doi.org/10.1145/2619239.2626299>
- Stefanos Kaxiras, David Klaftenegger, Magnus Norgren, Alberto Ros, and Konstantinos Sagonas. 2015. Turning Centralized Coherence and Distributed Critical-Section Execution on their Head: A New Approach for Scalable Distributed Shared Memory. In *Proceedings of the 24th International Symposium on High-Performance Parallel and Distributed Computing (HPDC '15)*. Association for Computing Machinery, New York, NY, USA, 3–14. <https://doi.org/10.1145/2749246.2749250>
- Pete Keleher, Alan L. Cox, Sandhya Dwarkadas, and Willy Zwaenepoel. 1994. TreadMarks: distributed shared memory on standard workstations and operating systems. In *Proceedings of the USENIX Winter 1994 Technical Conference on USENIX Winter 1994 Technical Conference (WTEC '94)*. USENIX Association, USA, 10.

- Michalis Kokologiannakis and Viktor Vafeiadis. 2021. GenMC: A Model Checker for Weak Memory Models. In *Computer Aided Verification - 33rd International Conference, CAV 2021, Virtual Event, July 20-23, 2021, Proceedings, Part I (Lecture Notes in Computer Science)*, Alexandra Silva and K. Rustan M. Leino (Eds.), Vol. 12759. Springer, 427–440. https://doi.org/10.1007/978-3-030-81685-8_20
- Ori Lahav, Brijesh Dongol, and Heike Wehrheim. 2023. Rely-Guarantee Reasoning for Causally Consistent Shared Memory. In *Computer Aided Verification - 35th International Conference, CAV 2023, Paris, France, July 17-22, 2023, Proceedings, Part I (Lecture Notes in Computer Science)*, Constantin Enea and Akash Lal (Eds.), Vol. 13964. Springer, 206–229. https://doi.org/10.1007/978-3-031-37706-8_11
- Ori Lahav, Viktor Vafeiadis, Jeehoon Kang, Chung-Kil Hur, and Derek Dreyer. 2017. Repairing sequential consistency in C/C++11. In *Proceedings of the 38th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2017, Barcelona, Spain, June 18-23, 2017*, Albert Cohen and Martin T. Vechev (Eds.). ACM, 618–632. <https://doi.org/10.1145/3062341.3062352>
- Pengfei Li, Yu Hua, Pengfei Zuo, Zhangyu Chen, and Jiajie Sheng. 2023. ROLEX: A Scalable RDMA-oriented Learned Key-Value Store for Disaggregated Memory Systems. In *21st USENIX Conference on File and Storage Technologies (FAST 23)*. USENIX Association, Santa Clara, CA, 99–114. <https://www.usenix.org/conference/fast23/presentation/li-pengfei>
- linux-rdma. 2018. RDMA core. (2018). <https://github.com/linux-rdma/rdma-core/> (Accessed: Jul. 2025).
- Feilong Liu, Claude Barthels, Spyros Blanas, Hideaki Kimura, and Garret Swart. 2021. Beyond MPI: New Communication Interfaces for Database Systems and Data-Intensive Applications. *SIGMOD Rec.* 49, 4 (March 2021), 12–17. <https://doi.org/10.1145/3456859.3456862>
- Xu Liu and John Mellor-Crummey. 2014. A tool to analyze the performance of multithreaded programs on NUMA architectures. In *Proceedings of the 19th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP '14)*. Association for Computing Machinery, New York, NY, USA, 259–272. <https://doi.org/10.1145/2555243.2555271>
- Redis Ltd. 2021. Redis v6.0.16. (2021). <https://github.com/redis/redis/releases/tag/6.0.16>.
- Redis Ltd. 2024. Memtier v2.1.2. (2024). https://github.com/RedisLabs/memtier_benchmark/releases/tag/2.1.2.
- Kai Lu, Siqi Zhao, Haikang Shan, Qiang Wei, Guokuan Li, Jiguang Wan, Ting Yao, Huatao Wu, and Daohui Wang. 2024. Scythe: A Low-latency RDMA-enabled Distributed Transaction System for Disaggregated Memory. *ACM Trans. Archit. Code Optim.* 21, 3, Article 57 (Sept. 2024), 26 pages. <https://doi.org/10.1145/3666004>
- Yuanwei Lu, Guo Chen, Bojie Li, Kun Tan, Yongqiang Xiong, Peng Cheng, Jiansong Zhang, Enhong Chen, and Thomas Moscibroda. 2018. {Multi-Path} transport for {RDMA} in datacenters. In *15th USENIX symposium on networked systems design and implementation (NSDI 18)*. 357–371.
- Zoltan Majo and Thomas R. Gross. 2017. A Library for Portable and Composable Data Locality Optimizations for NUMA Systems. *ACM Trans. Parallel Comput.* 3, 4, Article 20 (mar 2017), 32 pages. <https://doi.org/10.1145/3040222>
- Message Passing Interface Forum. 2023. MPI: A Message-Passing Interface Standard Version 4.1. <https://www.mpi-forum.org/docs/mpi-4.1/mpi41-report.pdf>
- Adam Morrison and Yehuda Afek. 2013. Fast Concurrent Queues for x86 Processors. In *Proceedings of the 18th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP '13)*. ACM, New York, NY, USA, 103–112. <https://doi.org/10.1145/2442516.2442527>
- Ramana Nagasamudram, Lennart Beringer, Ken Birman, Mae Milano, and David A. Naumann. 2024. Verifying a C Implementation of Derecho's Coordination Mechanism Using VST and Coq. In *NASA Formal Methods - 16th International Symposium, NFM 2024, Moffett Field, CA, USA, June 4-6, 2024, Proceedings (Lecture Notes in Computer Science)*, Nathaniel Benz, Divya Gopinath, and Nija Shi (Eds.), Vol. 14627. Springer, 99–117. https://doi.org/10.1007/978-3-031-60698-4_6
- J. Nieplocha, R.J. Harrison, and R.J. Littlefield. 1994. Global Arrays: a portable "shared-memory" programming model for distributed memory computers. In *Supercomputing '94: Proceedings of the 1994 ACM/IEEE Conference on Supercomputing*. 340–349. <https://doi.org/10.1109/SUPERC.1994.344297>
- NVIDIA Corporation. 2020. NVIDIA Collective Communication Library (NCCL) Documentation. (2020). <https://docs.nvidia.com/deeplearning/nccl/user-guide/docs/index.html>
- Scott Owens, Susmit Sarkar, and Peter Sewell. 2009. A Better x86 Memory Model: x86-TSO. In *Theorem Proving in Higher Order Logics, 22nd International Conference, TPHOLs 2009, Munich, Germany, August 17-20, 2009. Proceedings (Lecture Notes in Computer Science)*, Stefan Berghofer, Tobias Nipkow, Christian Urban, and Makarius Wenzel (Eds.), Vol. 5674. Springer, 391–407. https://doi.org/10.1007/978-3-642-03359-9_27
- Geoff Pike and Jyrki Alakuijala. [n. d.]. Introducing CityHash. ([n. d.]). <https://opensource.googleblog.com/2011/04/introducing-cityhash.html>.
- Marius Poke and Torsten Hoefler. 2015. DARE: High-Performance State Machine Replication on RDMA Networks. In *Proceedings of the 24th International Symposium on High-Performance Parallel and Distributed Computing (HPDC '15)*. ACM, New York, NY, USA, 107–118. <https://doi.org/10.1145/2749246.2749267>
- Azalea Raad, Marko Doko, Lovro Rozic, Ori Lahav, and Viktor Vafeiadis. 2019. On library correctness under weak memory consistency: specifying and verifying concurrent libraries under declarative consistency models. *Proc. ACM Program.*

- Lang. 3, POPL (2019), 68:1–68:31. <https://doi.org/10.1145/3290381>
- Azalea Raad, Ori Lahav, John Wickerson, Piotr Balcer, and Brijesh Dongol. 2024. Intel PMDK Transactions: Specification, Validation and Concurrency. In *Programming Languages and Systems - 33rd European Symposium on Programming, ESOP 2024, Held as Part of the European Joint Conferences on Theory and Practice of Software, ETAPS 2024, Luxembourg City, Luxembourg, April 6-11, 2024, Proceedings, Part II (Lecture Notes in Computer Science)*, Stephanie Weirich (Ed.), Vol. 14577. Springer, 150–179. https://doi.org/10.1007/978-3-031-57267-8_6
- Yufei Ren, Xingbo Wu, Li Zhang, Yandong Wang, Wei Zhang, Zijun Wang, Michel Hack, and Song Jiang. 2017. iRDMA: Efficient Use of RDMA in Distributed Deep Learning Systems. In *2017 IEEE 19th International Conference on High Performance Computing and Communications; IEEE 15th International Conference on Smart City; IEEE 3rd International Conference on Data Science and Systems (HPCC/SmartCity/DSS)*. 231–238. <https://doi.org/10.1109/HPCC-SmartCity-DSS.2017.30>
- Zhenyuan Ruan, Shihang Li, Kaiyan Fan, Marcos K. Aguilera, Adam Belay, Seo Jin Park, and Malte Schwarzkopf. 2023a. Unleashing True Utility Computing with Quicksand. In *Proceedings of the 19th Workshop on Hot Topics in Operating Systems (HOTOS '23)*. Association for Computing Machinery, New York, NY, USA, 196–205. <https://doi.org/10.1145/3593856.3595893>
- Zhenyuan Ruan, Seo Jin Park, Marcos K. Aguilera, Adam Belay, and Malte Schwarzkopf. 2023b. Nu: Achieving Microsecond-Scale Resource Fungibility with Logical Processes. In *20th USENIX Symposium on Networked Systems Design and Implementation (NSDI 23)*. USENIX Association, Boston, MA, 1409–1427. <https://www.usenix.org/conference/nsdi23/presentation/ruan>
- Zhenyuan Ruan, Malte Schwarzkopf, Marcos K. Aguilera, and Adam Belay. 2020. AIFM: High-Performance, Application-Integrated Far Memory. In *14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20)*. USENIX Association, 315–332. <https://www.usenix.org/conference/osdi20/presentation/ruan>
- Léo Stefanescu, Azalea Raad, and Viktor Vafeiadis. 2024. Specifying and Verifying Persistent Libraries. In *Programming Languages and Systems - 33rd European Symposium on Programming, ESOP 2024, Held as Part of the European Joint Conferences on Theory and Practice of Software, ETAPS 2024, Luxembourg City, Luxembourg, April 6-11, 2024, Proceedings, Part II (Lecture Notes in Computer Science)*, Stephanie Weirich (Ed.), Vol. 14577. Springer, 185–211. https://doi.org/10.1007/978-3-031-57267-8_8
- Lingjia Tang, Jason Mars, Xiao Zhang, Robert Hagmann, Robert Hundt, and Eric Tune. 2013. Optimizing Google’s warehouse scale computers: The NUMA experience. In *2013 IEEE 19th International Symposium on High Performance Computer Architecture (HPCA)*. 188–197. <https://doi.org/10.1109/HPCA.2013.6522318>
- Chenxi Wang, Haoran Ma, Shi Liu, Yuanqi Li, Zhenyuan Ruan, Khanh Nguyen, Michael D. Bond, Ravi Netravali, Miryung Kim, and Guoqing Harry Xu. 2020. Semeru: A Memory-Disaggregated Managed Runtime. In *14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20)*. USENIX Association, 261–280. <https://www.usenix.org/conference/osdi20/presentation/wang>
- Jing Wang, Chao Li, Yibo Liu, Taolei Wang, Junyi Mei, Lu Zhang, Pengyu Wang, and Minyi Guo. 2023a. Fargraph+: Excavating the parallelism of graph processing workload on RDMA-based far memory system. *J. Parallel and Distrib. Comput.* 177 (2023), 144–159. <https://doi.org/10.1016/j.jpdc.2023.02.015>
- Qing Wang, Youyou Lu, and Jiwu Shu. 2022. Sherman: A Write-Optimized Distributed B+Tree Index on Disaggregated Memory. In *Proceedings of the 2022 International Conference on Management of Data (SIGMOD '22)*. Association for Computing Machinery, New York, NY, USA, 1033–1048. <https://doi.org/10.1145/3514221.3517824>
- Zilong Wang, Layong Luo, Qingsong Ning, Chaoliang Zeng, Wenxue Li, Xinchun Wan, Peng Xie, Tao Feng, Ke Cheng, Xiongfei Geng, et al. 2023b. {SRNIC}: A scalable architecture for {RDMA} {NICs}. In *20th USENIX Symposium on Networked Systems Design and Implementation (NSDI 23)*. 1–14.
- Jilong Xue, Youshan Miao, Cheng Chen, Ming Wu, Lintao Zhang, and Lidong Zhou. 2019. Fast Distributed Deep Learning over RDMA. In *Proceedings of the Fourteenth EuroSys Conference 2019 (EuroSys '19)*. Association for Computing Machinery, New York, NY, USA, Article 44, 14 pages. <https://doi.org/10.1145/3302424.3303975>
- Jian Yang, Joseph Izraelevitz, and Steven Swanson. 2019. Orion: A Distributed File System for Non-Volatile Main Memories and RDMA-Capable Networks. In *17th USENIX Conference on File and Storage Technologies (FAST '19)*. USENIX Association.
- Jian Yang, Joseph Izraelevitz, and Steven Swanson. 2020. FileMR: Rethinking RDMA Networking for Scalable Persistent Memory. In *Proceedings of the 17th USENIX Conference on Networked Systems Design and Implementation (NSDI'20)*. USENIX Association.
- Qizhen Zhang, Xinyi Chen, Sidharth Sankhe, Zhilei Zheng, Ke Zhong, Sebastian Angel, Ang Chen, Vincent Liu, and Boon Thau Loo. 2022. Optimizing Data-intensive Systems in Disaggregated Data Centers with TELEPORT. In *Proceedings of the 2022 International Conference on Management of Data (SIGMOD '22)*. Association for Computing Machinery, New York, NY, USA, 1345–1359. <https://doi.org/10.1145/3514221.3517856>
- Yili Zheng, Amir Kamil, Michael B. Driscoll, Hongzhang Shan, and Katherine Yelick. 2014. UPC++: A PGAS Extension for C++. In *2014 IEEE 28th International Parallel and Distributed Processing Symposium*. 1105–1114. <https://doi.org/10.1109/>

[IPDPS.2014.115](#)

Yibo Zhu, Hagga Eran, Daniel Firestone, Chuanxiong Guo, Marina Lipshteyn, Yehonatan Liron, Jitendra Padhye, Shachar Raindel, Mohamad Haj Yahia, and Ming Zhang. 2015. Congestion control for large-scale RDMA deployments. *ACM SIGCOMM Computer Communication Review* 45, 4 (2015), 523–536.

Received 2025-07-10; accepted 2025-11-06