Sufficient Conditions for Robustness of RDMA Programs

Guillaume Ambal¹, Ori Lahav², and Azalea Raad¹

Abstract. Remote Direct Memory Access (RDMA) is a modern technology enabling high-performance inter-node communication. Despite its widespread adoption, theoretical understanding of permissible behaviours remains limited, as RDMA follows a very weak memory model. This paper addresses the challenge of establishing sufficient conditions for RDMA robustness. We introduce a set of straightforward criteria that, when met, guarantee sequential consistency and mitigate potential issues arising from weak memory behaviours in RDMA applications. Notably, when restricted to a tree topology, these conditions become even more relaxed, significantly reducing the need for synchronisation primitives. This work provides developers with practical guidelines to ensure the reliability and correctness of their RDMA-based systems.

Keywords: RDMA · Robustness · Weak Memory Models

1 Introduction

Remote Direct Memory Access (RDMA) is a modern technology that enables a machine to have direct read/write access to the memory of another machine over a network, bypassing the operating systems on both ends. This allows such direct memory accesses (reads/writes) to be performed with far fewer CPU cycles, leading to high-throughput, low-latency networking, which is especially useful in massively parallel computer clusters (e.g. data centres). RDMA has achieved widespread adoption as of 2018 [68], thanks to efficient implementations available at comparable cost to traditional infrastructures (e.g. TCP/IP sockets) [31], with several RDMA technologies such as InfiniBand and RDMA over Converged Ethernet (RoCE) readily available.

RDMA networks directly interact with the hardware through read (get) and write (put) operations on remote memory. As a result, programming RDMA systems is conceptually similar to shared memory systems of existing hardware architectures (e.g. Intel-x86 or ARM). A key difference, however, is that on encountering a remote operation, the CPU forwards it onto the *network interface card* (NIC), which subsequently handles the remote operation without further CPU involvement.

The performance gains of RDMA, as well as its wide range of implementations, have led to a surge of RDMA research [4,72,70,26]. RDMA networks exhibit different degrees of concurrency, depending on whether the concurrent threads reside on different nodes (machines) over the network (inter-node concurrency) or on the same node (intra-thread concurrency). To understand the behaviour of RDMA programs and their various notions of concurrency, Ambal et al. [10] recently developed RDMA^{TSO}, a formal semantics of RDMA programs where each node comprises an Intel-x86 CPU and thus intra-node-inter-thread concurrency is governed by the TSO (total store ordering) model [67].

As the real power of RDMA networks is their ability to run parallel programs over different nodes, writing efficient RDMA programs hinges on utilising internode concurrency. However, writing such programs correctly is far from straightforward. A key challenge is that local operations (accessing the local memory of the executing node) are handled by the CPU, while remote operations (accessing remote memory on other nodes) are handled by the NIC independently and in parallel to CPU operations. Hence, operations in the same thread may not be executed in the intended (program) order, leading to surprising outcomes. As Ambal et al. [10] note, this can result in counter-intuitive behaviours even in the case of sequential programs comprising a single thread. This is in stark contrast to all previously existing concurrency models (be they of CPU architectures or programming languages), where sequential programs do behave sequentially.

The permissive nature of RDMA semantics requires developers to carefully consider potential instruction reorderings. Reasoning about concurrent programs and ensuring proper synchronisation between threads is inherently complex, even without instruction reordering. Accounting for instruction reorderings adds another layer of complexity to this challenge.

As such, we should ideally enable reasoning about RDMA programs under a simpler, more intuitive model such as sequential consistency (SC) [42], where no instruction reordering is allowed, and thus instructions in each thread always execute in order. To this end, a common approach to simplify reasoning is to ensure robustness. A program P is robust under a consistency model CM, if its set of possible behaviours under CM coincide with those of its behaviours under SC; i.e. P is robust under CM if it exhibits no non-SC behaviours. If a program is robust under CM, then we can simply reason about it under SC, without considering the complexities of CM.

Contributions. In this paper, we close this gap and simplify reasoning about RDMA program through robustness. To simplify our presentation and not distract the reader from the RDMA complexities by the *orthogonal* intricacies of CPU concurrency, we first present RDMA^{SC}, a simplification of the RDMA^{TSO} model of Ambal et al. [10], where intra-node concurrency follows the simpler SC model [42], while inter-node concurrency is analogous to that of RDMA^{TSO}. We then identify two sets of sufficient constraints that, if satisfied, ensure the robustness of RDMA^{SC} programs. Our proposed constraints are purely *syntactic*, in that they do not require an understanding of the complex RDMA semantics and can be established by simply checking the syntax of the program. The first

set of constraints is restrictive, but can be applied to any RDMA program. The second relaxes the requirements of the first, but requires the RDMA network to follow a tree topology. Our conditions enable a number of useful paradigms for RDMA programs such as the server-client model, which we show can be used for automatically translating existing concurrent algorithms to distributed ones over RDMA, as well as for modelling star network topologies used e.g. in Local Area Networks (LAN). Finally, we adapt our results to the RDMA^{TSO} model and accordingly propose analogous syntactic and topological constraints.

Outline. In §2 we present an intuitive account of the weak RDMA semantics through examples and discuss how we ensure robustness through syntactic constraints. In §3 we present our formal RDMA^{SC} model. In §4 we establish sufficient syntactic conditions that ensure the robustness of RDMA^{SC} programs. In §5 we apply these findings to tree-shaped network topologies, offering a further streamlined set of conditions under RDMA^{SC}. We discuss related work in §6. The appendix contains the proofs of all theorems stated in this paper (§A), as well as the extension of all our results to the RDMA^{TSO} model (§B).

2 Overview

We present an intuitive account of RDMA semantics through several examples, showing the counter-intuitive and unexpected behaviours they can exhibit due to possible *instruction reorderings* (§2.1). We then discuss how we can tame this complexity by introducing *syntactic constraints* that, if fulfilled, prohibit problematic instruction reorderings, pre-empting unexpected behaviours and thus simplifying the task of reasoning about RDMA programs for developers (§2.2).

2.1 RDMA Semantics at a Glance

Consistency (Concurrency) Models and Weak Behaviours. In the literature of shared-memory concurrent (multi-threaded) programming, the set of possible behaviours (i.e. semantics) of a concurrent program is defined via a consistency model (a.k.a. memory model or concurrency model), with a number of such models available in different domains such as hardware architectures (e.g. Intel and ARM) and programming languages (e.g. C/C++ and Java). The most well-known and intuitive consistency model is sequential consistency (SC, a.k.a. interleaving concurrency) [42], where the instructions are interleaved in program order. That is, under SC the instructions in each thread cannot be reordered. While simple, SC is too strong in that it precludes many common hardware/compiler optimisations and thus unduly hinders performance. As such, modern hardware architectures and programming languages adhere to weaker, more lenient models, admitting more behaviours than SC. In this context a program behaviour (outcome) is referred to as weak, if it is not allowed under SC. Such weak behaviours can typically be understood in terms of instruction reorderings within a thread or visibility delays (where the effects of an instruction (e.g. a write) is not observed at the same time by all threads), both of which are disallowed under SC.

4

Conceptual RDMA Model. We model concurrent RDMA programs running over a network of nodes (i.e. computers), where each node hosts zero, one, or more threads, and each thread can directly access remote memory of other nodes through its network interface card (NIC). As we discuss below, RDMA programs exhibit three sources of weak behaviours: 1. CPU weak behaviours, due to the usual interactions (and reordering) of multiple threads on a single node; 2. intra-thread weak behaviours, due to RDMA operations being reordered or delayed; and 3. inter-node weak behaviours, due to multiple nodes executing concurrently. Here we focus on the latter two sources as they are specific to RDMA programs, and discuss how such weak behaviours may be prevented.

CPU Concurrency. RDMA enacts data transfers between nodes via the NIC subsystems of the constituent nodes, which are independent from the CPU subsystems. Consequently, the RDMA technology can be combined with different CPU architectures governed by different memory models (e.g. TSO or ARM). The first validated formal model of RDMA programs, RDMA^{TSO} [10], assumes that CPU concurrency is governed by the TSO model [67]. To simplify our presentation and not distract the reader from the RDMA complexities by the *orthogonal* intricacies of CPU concurrency, we present the simpler RDMA^{SC} model, where CPU concurrency follow the stronger SC model [42]. We later generalise our results to RDMA^{TSO} in the appendix (§B).

Almost all weak behaviours introduced by RDMA stem from the NIC and are independent of CPU concurrency (i.e. CPU and RDMA concurrency can often be decoupled). As such, the distinction between RDMA^{TSO} and RDMA^{SC} is often irrelevant, in which case we write RDMA* to encompass both models. In particular, in this overview section we focus on nodes with *at most one thread each*, i.e. with no CPU concurrency, so all behaviours discussed below hold of both RDMA^{SC} and RDMA^{TSO} (i.e. for RDMA*). Note that this is *merely a presentational choice* we have made in this section, and our formal models, theorems, and examples in subsequent sections also account for CPU concurrency.

Litmus Test Outcome Notation. We frequently present small representative examples (known as litmus tests in the literature). In each example, the outcomes annotated with \checkmark are allowed by the RDMA model under discussion, while those annotated with \checkmark are disallowed.

Remote Direct Memory Access (RDMA). RDMA programs comprise operations that access remote memory, as well as various synchronisation operations. As such, programming RDMA networks is conceptually similar to shared memory systems. To distinguish remote (RDMA) operations from CPU ones, we refer to RDMA reads and writes as get and put operations, respectively. To distinguish local and remote memory locations, we assume nodes do not reuse location names, we write x^n for a location on a remote node n, and write x for a location on the local node. A put operation is of the form $x^n := y$ and consists of reading from a local location y and writing to a remote location x on x. Similarly, a get operation is of the form $x := y^n$ and consists of reading from a remote location y on y and writing to a local location y. We write \overline{n} to identify

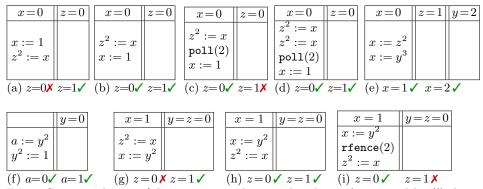


Fig. 1: Sequential RDMA* litmus tests, where each column (separated by ||) denotes a distinct node, the statement on the top line of each column denotes the initial values of locations.

a node other than n. When node n issues a remote operation to be executed on node \overline{n} , we state that the operation is by n towards \overline{n} .

Sequential (Single-Threaded) RDMA* Behaviours. When a thread issues a get or put operation, it is handled by the NIC, in contrast to local reads and writes handled by the CPU. As such, the interaction between CPU and remote operations lead to further behaviours even within a *sequential* (single-threaded) program. We demonstrate this in the examples of Fig. 1, where each column represents a distinct node, numbered from left to right starting from 1. For instance, the example in Fig. 1a comprises a single thread on node 1 (the leftmost column) that writes to the local location x (x := 1) and puts x towards the remote location z on node 2 ($z^2 := x$).

Intuitively, when a thread t on n issues remote operations towards node \overline{n} , one can view these remote operations as if being executed by a thread running in parallel to t. As such, when a remote operation follows a CPU one, the order of the two operations is preserved since the parallel thread is spawned only after the CPU operation is executed. This is illustrated in Fig. 1a. By contrast, when a remote operation precedes a CPU one, the remote operation is performed by a 'separate thread' run in parallel to the later CPU operation in the main thread, and thus may execute before or after the CPU operation, meaning that in the latter case the execution order is not preserved. This is illustrated in Fig. 1b.

Therefore, before using the result of a get or reusing the memory location of a put, it is desirable to avoid such reorderings and to wait for the remote operation to complete. This can be done through a CPU poll operation, poll(n), that blocks until the earliest (in program order) remote operation towards node n has completed. This is shown in Fig. 1c, obtained from Fig. 1b by inserting a poll after the remote operation: poll(2) waits for $z^2 := x$ to complete before proceeding with x := 1, and thus $z^2 := x$ can no longer be reordered after x := 1.

Note that each poll(n) waits for *only one* (the earliest) and *not all* pending remote operations towards n to complete. For instance, in Fig. 1d, poll(2) only blocks until the *first* $z^2 := x$ is complete, and thus z = 1 is once again possible.

| y=0 | x=0 | x=0 | y=0 | x=0 | y=0 | y=0 | x=0 | y=w=0 | x=z=0 |
|--|--|--|--|------------|--------------|--------------|--------------|--------------------------------|--------------|
| | | | | | | | | $x^2 := 1$ | $ y^1 := 1 $ |
| $ _{x^2} \cdot = 1$ | $u^1 \cdot = 1$ | $a := u^2$ | $b = r^1$ | $a := y^2$ | $ b := x^1 $ | $ x^2 := 1 $ | $y^1 := 1$ | $c := z^2$ | $d := w^1$ |
| $\begin{vmatrix} x^2 := 1 \\ a := y \end{vmatrix}$ | $\begin{vmatrix} y & \cdot - 1 \\ h & \cdot - r \end{vmatrix}$ | $\begin{vmatrix} a & -y \\ x & -1 \end{vmatrix}$ | $\begin{vmatrix} o & -x \\ y & -1 \end{vmatrix}$ | poll(2) | poll(1) | pol1(2) | poll(1) | $c := z^2$ $poll(2)$ $poll(2)$ | poll(1) |
| a = g | $0 \cdot - x$ | x .— 1 | $ g \cdot - 1$ | x := 1 | y := 1 | a := y | b := x | pol1(2) | poll(1) |
| | | | | | | | | a := y | |
| (a) $a = b$ | b=0 ✓ | (b) $a = a$ | b=1 | (c) a = | b=1 X | (d) a = | b=0 ✓ | (e) $a =$ | b=0 X |

Fig. 2: Concurrent RDMA* litmus tests.

Two remote operations towards different nodes are independent and can execute in either order, as illustrated in Fig. 1e. The only way to prevent this reordering is to poll the first operation before running the second.

The ordering guarantees on remote operations towards the *same* node are stronger and only certain reorderings are allowed. Recall that a put operation $x^n := y$ comprises two steps: a local read (on y) and a remote write (on x^n). Similarly, a get operation $x := y^n$ comprises two steps: a remote read (on y^n) and a local write (on x). Intuitively, NIC operations follow the *precedence* order: i) local read; ii) remote write; iii) remote read; iv) local write.

If a step with a higher precedence (e.g. a local read) is in program order before one with a lower precedence (e.g. a local write), then their order is preserved and they cannot be reordered. This is illustrated in Fig. 1g. Otherwise the order is not necessarily preserved and these steps can be reordered, as shown in Fig. 1h where an earlier local write on x can occur after the later local read.

As before, the reordering of the two remote operations in Fig. 1h can be prevented by polling the first operation before the second. However, polling is costly as it blocks the current thread, including the submission of remote operations towards any node. Alternatively, we can use a $remote\ fence$, rfence(n), that blocks only the NIC and only towards node n. This in turn ensures that earlier (before the fence) remote operations by the thread towards n are executed before later (after the fence) remote operations towards n. This is illustrated in Fig. 1i, obtained from Fig. 1h by inserting rfence(2) stopping the reordering.

Concurrent (Multi-Threaded) RDMA* Behaviours. The real power of RDMA comes from programs running on different nodes, introducing a wide range of weak behaviours. A network can comprise several nodes, each running several concurrent threads. We limit the examples of Fig. 2 to two nodes, each having a single thread.

As shown in Figs. 2a and 2b, well-known weak behaviours such as store buffering (Fig. 2a) and load buffering (Fig. 2b) are possible. This is because earlier RDMA operations can be delayed after later CPU operations.

As one could expect, most weak behaviours can be prevented by polling the remote operations as needed, as shown for load buffering in Fig. 2c. However, this strategy is not enough to prevent the store buffering weak behaviour, as show in Fig. 2d. This is because the specification of polling offers different guarantees for get and put operations. Polling a get operation $a := x^n$ offers the strong intuitive guarantee that the operation completed, i.e. the value of x^n is fetched

from node n and written to a. By contrast, polling a put operation $x^n := a$ does not guarantee the write on x^n has completed. When sending the value of a towards node n to be put in x^n , the remote NIC merely acknowledges having received the data, but this data may still reside in a buffer (i.e. the PCIe fabric) of the remote node, pending to be written x^n . Polling a put operation only awaits the acknowledgement of the data receipt. As such, it is possible to poll a put operation successfully before the associated remote write has fully completed. In the case of store buffering in Fig. 2d, it is possible for both poll operations to complete before the values of x and y are updated (to 1) in memory.

We also assume NICs are connected to memory though the *Peripheral Component Interconnect Express* (PCIe) fabric, the *de facto* standard for this category of hardware [10]. This ensures that (PCIe) reads cannot overtake (PCIe) writes. As such, a remote read *flushes* (commits) all pending remote writes to memory, and similarly on local memory. This can be used to prevent weak behaviours such as store buffering, as shown in Fig. 2e, obtained from Fig. 2d by adding additional gets and subsequently polling them. Polling a (seemingly unrelated) later get (e.g. $c := z^2$) ensures previous remote writes (e.g. $x^2 := 1$) have been committed to the remote memory.

2.2 Robustness: Taming Weak RDMA* Behaviours

Given the permissive nature of the RDMA* semantics and the numerous weak behaviours it exhibits (even in the case of single-threaded programs), the task of writing *correct* RDMA programs is laborious. Reasoning about concurrent programs is already challenging even in the absence of weak behaviours. Accounting for potential instruction reorderings (which requires experience with RDMA* semantics) introduces yet another layer of complexity for developers.

As such, we should ideally enable reasoning about RDMA programs under a simpler, more intuitive model such as SC (sequential consistency [42]). Specifically, to simplify program reasoning, a common approach is to ensure robustness. A program P is robust under a consistency model CM, if its set of possible behaviours under CM coincide with those of its behaviours under SC; i.e. P is robust if it exhibits no weak behaviours. If a program is robust, then we can reason about it $as\ if$ it were executed under SC, without considering the complexities of RDMA*.

To ensure robustness, we must prevent *observable* reorderings, i.e. those leading to weak behaviours. We can achieve this through *syntactic* requirements (e.g. by inserting sufficient remote fences and poll operations). A naive solution is to wait for each remote operation to fully complete before proceeding further, thereby preventing all reorderings. Unfortunately, this *serialises* these operations, and thus defeats the benefits of RDMA, which is designed to parallelise CPU instructions and data transfers by offloading them to the NIC. Instead, we should account for the RDMA* semantics and only add restrictions when necessary, while allowing non-observable reorderings.

Certain reorderings are observable even when considering a single thread in isolation, as in the examples of Figs. 1b, 1e, 1f, and 1h. Specifically, these exam-

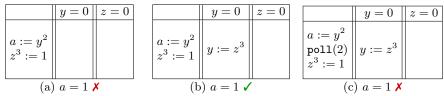


Fig. 3: Examples showing that necessary restrictions depend on other threads.

ples contain data races within a single thread. Beyond robustness, these patterns should be avoided in any sensible program. However, most weak behaviours arise from the interaction of several threads. For instance, in the single-threaded example of Fig. 3a, although the two remote operations $a:=y^2$ and $z^3:=1$ on node 1 may be reordered, this reordering is not observable: it does not lead to additional weak behaviours, and thus no additional constraints are necessary for robustness. By contrast, in the multi-threaded variant of Fig. 3b (with a thread on node 2), nodes 2 and 3 can exchange data and thus we can observe the weak behaviour a=1 due to this reordering. As such, to prohibit this, we must prevent the two operations on node 1 from being reordered, e.g. by polling the first operation, as shown in Fig. 3c.

As seen before, preventing reorderings can be done in different ways. In cases like Fig. 1i, a remote fence is enough. In cases like Fig. 2e, we need dummy get operations. Determining when and how to prevent reorderings is not straightforward. As illustrated in the examples of Fig. 3, it cannot be done *thread-locally*: one must account for the communication between other nodes and thus must take the whole program into account. This raises two questions:

- How do we prevent weak behaviours through simple purely syntactic restrictions? Specifically, how can we ensure that a program has enough constraints (e.g. polls) to prevent weak behaviours, and how do we make sure that waiting for a specific remote operation (as in Fig. 3a) is unnecessary?
- How do we structure RDMA programs to minimise the amount of necessary restrictions in order to maintain efficient implementations?

We set out to answer these questions in the remainder of this paper. Specifically, after defining several formal preliminaries in §3, we present a theorem in §4 stating sufficient syntactic conditions guaranteeing robustness (i.e. the absence of weak behaviours). In §5 we then build on this theorem and present a useful RDMA network topology where fewer limitations are necessary to prevent weak behaviours. Notably, following our prescribed network topology ensures that it is never necessary to poll a remote operation to prevent multi-threaded weak behaviours.

3 RDMA^{SC}: A Declarative Semantics for RDMA Programs

We present the syntax of RDMA programs (taken from [10]) in §3.1. In §3.2 we then present a formal declarative semantics for our RDMA^{sc} model. As we

describe in the appendix (§B), we obtain RDMA^{SC} by *strengthening* the RDMA^{TSO} model of Ambal et al. [10] whereby we make a few simple adjustments to ensure that local (CPU) concurrency follows the SC rather than TSO model.

3.1 RDMA^{SC}: Programming Language

Nodes and Threads. We consider a system with N nodes and M threads in total across all nodes. Let $\mathsf{Node} = \{1, \dots, N\}$ and $\mathsf{Tid} = \{1, \dots, M\}$ denote the sets of *node* and *thread identifiers*, respectively. We use n and t to range over Node and Tid , respectively. Given a node n, we write \overline{n} to range over $\mathsf{Node} \setminus \{n\}$. Each thread $t \in \mathsf{Tid}$ is associated with a node, written n(t).

Memory Locations. Each node n has a set of locations, Loc_n , accessible by all nodes. We define $\mathsf{Loc} \triangleq \biguplus_n \mathsf{Loc}_n$ and $\mathsf{Loc}_{\overline{n}} \triangleq \mathsf{Loc} \setminus \mathsf{Loc}_n$. We use x^n, y^n, z^n, w^n and $x^{\overline{n}}, y^{\overline{n}}, z^{\overline{n}}, w^{\overline{n}}$ to range over Loc_n and $\mathsf{Loc}_{\overline{n}}$, respectively. When the choice of n is clear, we write x for x^n and \overline{x} for $x^{\overline{n}}$. For clarity, we use distinct location names across nodes and write n(x) for the unique $n \in \mathsf{Node}$ where $x \in \mathsf{Loc}_n$. We assume all locations can be accessed by all threads on all nodes. However, for readability, we use a, b, c, and d for (private) locations that are only accessed by a single thread (on a single node).

Values and Expressions. We assume a set of values, Val, with $\mathbb{N} \subseteq \text{Val}$, and use v to range over Val. We assume a language of expressions over Val and Loc, and elide its exact syntax and semantics. We use e to range over expressions, and e^n to range over expressions whose locations are all included in Loc_n .

Sequential Commands and Programs. Sequential programs on node n are described by the C^n grammar below and include primitive commands (c^n) , sequential composition $(C_1^n; C_2^n)$, non-deterministic choice $(C_1^n + C_2^n)$, executing either C_1^n or C_2^n , and non-deterministic loops (C^{n*}, C_2^n) , executing C^n any number of times). A (concurrent) program, P, is a map from thread identifiers to commands, associating each thread $t \in Tid$ with a command on node n(t).

```
\begin{split} \operatorname{\mathsf{Comm}} \ni \mathsf{C}^n &::= \operatorname{\mathsf{skip}} \mid \mathsf{c}^n \mid \mathsf{C}^n_1; \mathsf{C}^n_2 \mid \mathsf{C}^n + \mathsf{C}^n_2 \mid \mathsf{C}^{n*} \quad \operatorname{\mathsf{PComm}} \ni \mathsf{c}^n ::= \mathsf{cc}^n \mid \mathsf{rc}^n \\ \operatorname{\mathsf{CComm}} \ni \mathsf{cc}^n &::= x := e^n \mid \operatorname{\mathsf{assume}}(x = v) \mid \operatorname{\mathsf{assume}}(x \neq v) \\ & \mid x := \operatorname{\mathsf{CAS}}(y, e_1, e_2) \mid \operatorname{\mathsf{poll}}(\overline{n}) \\ \operatorname{\mathsf{RComm}} \ni \mathsf{rc}^n ::= x := \overline{y} \mid \overline{y} := x \mid \operatorname{\mathsf{rfence}}(\overline{n}) \end{split}
```

Primitive commands include CPU (cc^n) and RDMA (rc^n) operations. A CPU operation on n may be a no-op (skip), an assignment to a local location (x := e), an assumption on the value of a local location ($\mathsf{assume}(x = v)$ and $\mathsf{assume}(x \neq v)$), an atomic CAS ('compare-and-set') operation ($x := \mathsf{CAS}(y, e_1, e_2)$), or a 'poll', $\mathsf{poll}(\overline{n})$, that awaits the completion notification of the earliest $\mathsf{put/get}$ that is pending (not yet acknowledged). An RDMA operation may be (i) a 'get', $x := \overline{y}$, reading from remote location \overline{y} and writing the result to local location x; (ii) a 'put', $\overline{y} := x$, reading from local location x and writing the result to remote location \overline{y} ; or (iii) a 'remote fence', $\mathsf{rfence}(\overline{n})$, which ensures that all

later (in program order) RDMA operations towards \overline{n} will await the completion of all earlier RDMA operations towards \overline{n} . poll(\overline{n}) is executed by the CPU and blocks its thread (and prevents the requests of later remote operations), while rfence(\overline{n}) blocks the NIC for the execution of remote operations towards \overline{n} .

3.2 RDMA^{SC}: Declarative Semantics

Events and Executions. In the literature of declarative models, the traces of a program are commonly represented as a set of *executions*, where an execution is a graph comprising: i) a set of *events* (graph nodes); and ii) a number of relations on events (graph edges). Each event is associated with the execution of a primitive command (in PComm) and is a tuple (ι, t, l) , where ι is the (unique) *event identifier*, $t \in \mathsf{Tid}$ identifies the executing thread, and $l \in \mathsf{ELab}$ is the *event label*, defined below.

Definition 1 (Labels and events). An event, $e \in \text{Event}$, is a triple (ι, t, l) , where $\iota \in \mathbb{N}$, $t \in \text{Tid}$ and $l \in \text{ELab}_{n(t)}$. The set of event labels is $\text{ELab} \triangleq \bigcup_n \text{ELab}_n$ for all nodes n. An event label of n, $l \in \text{ELab}_n$, is a tuple of one of the following forms:

```
\begin{array}{lll} - & NIC \; local \; read: \; l = \mathtt{nlR}(x^n, v_r, \overline{n}) & - \; (CPU) \; local \; read: \; l = \mathtt{lR}(x^n, v_r) \\ - & NIC \; remote \; write: \; l = \mathtt{nrR}(y^{\overline{n}}, v_w) & - \; (CPU) \; local \; write: \; l = \mathtt{lW}(x^n, v_w) \\ - & NIC \; local \; write: \; l = \mathtt{nlW}(x^n, v_w, \overline{n}) & - \; (CPU) \; CAS: \; l = \mathtt{CAS}(x^n, v_r, v_w) \\ - & NIC \; fence: \; l = \mathtt{nF}(\overline{n}) & - \; (CPU) \; poll: \; l = \mathtt{P}(\overline{n}) \end{array}
```

Each event label denotes whether the associated primitive command is handled by the NIC (left column, prefixed with n), or the CPU (right column). A poll instruction is handled by the CPU. A put operation $x^{\overline{n}} := y^n$ by node n towards node \overline{n} comprises a NIC local read from y^n and a NIC remote write on $x^{\overline{n}}$ and is thus modelled as two events with labels $\mathtt{nlR}(y^n,v,\overline{n})$ and $\mathtt{nrW}(x^{\overline{n}},v)$, where v denotes the value read from y^n and written to $x^{\overline{n}}$. Similarly, a get $x^n := y^{\overline{n}}$ is modelled as two events with labels of the form $\mathtt{nrR}(y^{\overline{n}},v)$ and $\mathtt{nlW}(x^n,v,\overline{n})$.

CPU operations are modelled by events as expected. A successful operation $x := \mathtt{CAS}(y, v_1, v_2)$ is modelled by two events with labels $\mathtt{CAS}(y, v_1, v_2)$ and $\mathtt{lW}(x, v_1)$. An unsuccessful $x := \mathtt{CAS}(y, v_1, v_2)$ operation is modelled by a CPU read instead: $\mathtt{lR}(y, v)$ and $\mathtt{lW}(x, v)$, with $v \neq v_1$.

We write $\mathsf{type}(l)$, $\mathsf{loc}(l)$, $v_{\mathsf{r}}(l)$, $v_{\mathsf{w}}(l)$, and $\overline{n}(l)$ for the type (e.g. 1R), location, read value, write value, and remote node of l, where applicable; e.g. $\mathsf{loc}(\mathsf{nlR}(x^n,v_{\mathsf{r}},\overline{n})) = x^n$ and $\overline{n}(\mathsf{nlR}(x^n,v_{\mathsf{r}},\overline{n})) = \overline{n}$. We lift these functions to events as expected. We write $\iota(\mathsf{e})$, $\iota(\mathsf{e})$, $\iota(\mathsf{e})$ to project the corresponding components of an event $\mathsf{e} = (\iota,t,l)$, and write $\iota(\mathsf{e})$ for the node $\iota(\mathsf{e})$ of an event.

Queue Pairs. As mentioned in §2 (see Fig. 1e), two remote operations by the same thread towards different remote nodes can be reordered. When using RDMA, each thread establishes a communication channel, called a *queue pair*, towards each remote node. The intuition is that operations on different queue pairs are independent and can always be reordered. Different threads, even on the same node, create different queue pairs to connect to the same remote node.

Notation. Given a relation r and a set A, we write r^+ for the transitive closure of r; 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. $\{(a,a) \mid a \in A\}$. We write $r_1; r_2$ for their relational composition: $\{(a,b) \mid \exists c. (a,c) \in r_1 \land (c,b) \in r_2\}$. When r is a strict partial order, we write $r|_{\text{imm}}$ for the *immediate* edges in r, i.e. $r \setminus (r;r)$. Given a set of events E and a location x, we write E_x for $\{e \in E \mid loc(e) = x\}$. Given a set of events E and a label type X, we write E.X for $\{e \in E \mid type(e) = X\}$, and define its sets of reads as $E.\mathcal{R} \triangleq E.\mathtt{lR} \cup E.\mathtt{CAS} \cup E.\mathtt{nlR} \cup E.\mathtt{nrR}$, writes as $E.\mathcal{W} \triangleq$ $E.1W \cup E.CAS \cup E.n1W \cup E.nrW$, CPU events as $E^{cpu} \triangleq E.1W \cup E.1R \cup E.CAS \cup E.P$, and NIC writes as $E.nW \triangleq E.nlW \cup E.nrW$. We define the 'same-location' relation as $sloc \triangleq \{(e, e') \in Event^2 | loc(e) = loc(e')\};$ the 'same-thread' relation as sthd $\triangleq \{(e, e') \in \text{Event}^2 \mid t(e) = t(e')\};$ and the 'same-queue-pair' relation as $\operatorname{sqp} \triangleq \{(e, e') \in \operatorname{Event}^2 \mid t(e) = t(e') \land \overline{n}(e) = \overline{n}(e')\}.$ We use sqp for events on the same queue pair, i.e. by the same thread and towards the same remote node. Note that $sqp \subseteq sthd$ and that sloc, sthd, and sqp are all symmetric. For a set of events E, we write E.sloc for $\operatorname{sloc}_{|E|}$; similarly for E.sthd and E.sqp.

Definition 2 (Pre-executions). A tuple $\mathcal{G} = \langle E, po, pf \rangle$ is a pre-execution of a program if:

- $-E \subseteq \mathsf{Event}$ is the set of events and includes a set of initialisation events, $E^0 \subseteq E$, comprising a single write with label $\mathtt{lW}(x,0)$ for each $x \in \mathsf{Loc}$.
- po $\subseteq E \times E$ is the 'program order' relation defined as a disjoint union of strict total orders, each ordering the events of one thread, with $E^0 \times (E \setminus E^0) \subseteq$ po, and such that:
 - Each put (resp. get) operation corresponds to two events: a read and a write with the read immediately preceding the write in po: 1. if $r \in G.nlR$ (resp. $r \in G.nrR$), then $(r, w) \in po|_{imm}$ for some $w \in G.nrW$ ($w \in G.nlW$); and 2. if $w \in G.nrW$ (resp. $w \in G.nlW$), then $(r, w) \in po|_{imm}$ for some $r \in G.nlR$ ($r \in G.nrR$).
 - Read and write events of a put (resp. get) have matching values: if $(r,w) \in G.\mathsf{po}|_{imm}$, $\mathsf{type}(r) \in \{\mathsf{nlR},\mathsf{nrR}\}$, and $\mathsf{type}(w) \in \{\mathsf{nlW},\mathsf{nrW}\}$, then $v_r(r) = v_w(w)$.
- pf ⊆ E.nW × E.P is the 'polls-from' relation, relating earlier (in programorder) NIC writes to later poll operations on the same queue pair; i.e. pf ⊆ po ∩ sqp. Moreover, pf is functional on its domain (every NIC write can be be polled at most once), and pf is total and functional on its range (every poll in E.P polls from exactly one NIC write). Also, Poll events poll-from the oldest non-polled remote operation on the same queue pair:
 - if $w_1 \in G.nW$ and $w_1 \xrightarrow{po \cap sqp} w_2 \xrightarrow{pf} p_2$, then there exists p_1 such that $w_1 \xrightarrow{pf} p_2$.

Pre-executions are constructed syntactically by induction on the structure of the corresponding program. This definition is standard and omitted.

Intuitively, a pre-execution can also be seen as a trace of the execution: for each thread t, po restricted to t is a total order, and so $\langle E, po \rangle$ is fundamentally

a sequence of events for each thread. In this view, pf should be considered a well-formedness condition: each prefix of the trace needs to have at least as many remote operations as poll operations. So $\langle E, po, pf \rangle$ can be seen as providing a well-formed trace for each thread. We later define robustness conditions on pre-executions, and as such they can also be considered conditions on traces.

We next extend the notion of a pre-execution to an *execution* by choosing explicitly how the different events interact.

Definition 3 (Executions). $G = \langle E, po, pf, rf, mo, nfo \rangle$ is an execution if:

- $-\langle E, po, pf \rangle$ is a pre-execution.
- rf $\subseteq E.W \times E.R$ is the 'reads-from' relation on events of the same location with matching values; i.e. $(a,b) \in \mathsf{rf} \Rightarrow (a,b) \in \mathsf{sloc} \land v_w(a) = v_r(b)$. Moreover, rf is total and functional on its range: every read in E.R is related to exactly one write in E.W.
- $\text{ mo} \triangleq \bigcup_{x \in \mathsf{Loc}} \mathsf{mo}_x$ is the 'modification-order', where each mo_x is a strict total order on $E.\mathcal{W}_x$ with $E_x^0 \times (E.\mathcal{W}_x \setminus E_x^0) \subseteq \mathsf{mo}_x$ describing the order in which writes on x reach the memory.
- nfo $\subseteq E.$ sqp is the 'NIC flush order', such that for all $(a,b) \in E.$ sqp, if $a \in E.$ nlR, $b \in E.$ nlW, then $(a,b) \in$ nfo \cup nfo⁻¹, and if $a \in E.$ nrR, $b \in E.$ nrW, then $(a,b) \in$ nfo \cup nfo⁻¹.

We define the *reads-before* relation as $\mathsf{rb} \triangleq (\mathsf{rf}^{-1}; \mathsf{mo}) \setminus [E]$, relating each read r to writes that are mo -after the write r reads from. Given a (pre-)execution G (resp. \mathcal{G}), we use the 'G.' prefix to project its various components (e.g. G. rf) and derived relations (e.g. G. rb). When the context is clear, we drop the prefix.

PCIe guarantees that a NIC local read (nlR) propagates all pending NIC local writes (nlW) (processed by the same queue pair) to memory, while a NIC remote read (nrR) propagates all pending NIC remote writes (nrW) (processed by the same queue pair) to memory. We model this total order through the nfo relation, stipulating that all NIC local reads and writes (resp. all NIC remote reads and writes) on the same queue pair be totally ordered.

Issue and Observation Points. In what follows we distinguish between when an instruction is issued and when it is observed. Intuitively, an instruction is issued when it is processed by the CPU or the NIC, and it is observed when its effect is propagated to memory. As such, since NIC writes can be delayed and have an observable effect on memory, the time points at which they are issued and observed may differ. Since we assume CPUs follow the strong SC memory model, CPU writes are issued and observed at the same time. However, the local (resp. remote) write of a get (resp. put) is issued when it is processed by the NIC and sent to the PCIe fabric, and observed when it is propagated to memory. All other events are instantaneous in that either they do not have an observable effect on memory (e.g. reads), or their effect is written to memory immediately (e.g. CAS operations and CPU writes). Given a set of events E, we thus define the set of instantaneous events in E as E.Inst $\triangleq E \setminus (E.nlW \cup E.nrW)$. Intuitively, the effects of NIC local writes and NIC remote writes (labelled nlW and nrW) can be delayed in the PCIe fabric and are thus excluded from the set

| | Later in Program Order | | | | | | | | | |
|---------|------------------------|------------------|------------------|----------------------|----------------|----------------------|-----|----------------------|--|--|
| | ippo | | 1 | 2 | 3 | 4 | 5 | 6 | | |
| | | | E^{cpu} | nlR | \mathtt{nrW} | nrR | nlW | nF | | |
| 0 | $ \mathbf{A} $ | E^{cpu} | √ | > | √ | > | 1 | ✓ | | |
| Ъ | В | nlR | X | sqp | sqp | sqp | sqp | sqp | | |
| in | $ \mathbf{C} $ | nrW | X | X | sqp | sqp | sqp | sqp | | |
| lieı | \mathbf{D} | nrR | X | X | X | X | sqp | sqp | | |
| Earlier | \mathbf{E} | nlW | X | X | X | X | sqp | sqp | | |
| щ | \mathbf{F} | nF | X | sqp | sqp | sqp | sqp | sqp | | |

| Later in Program Order | | | | | | | | | |
|------------------------|--------------|--------------------|------------------|----------------------|----------------------|----------------------|----------------------|----------------------|--|
| onno | | | 1 | 2 | 3 | 4 | 5 | 6 | |
| | oppo | | E^{cpu} | nlR | nrW | nrR | nlW | nF | |
| 0 | A | E^{cpu} | ✓ | \ | √ | > | √ | √ | |
| 1 P | В | nlR | X | sqp | sqp | sqp | sqp | sqp | |
| i. | \mathbf{C} | nrW | | X | sqp | sqp | sqp | X | |
| lier | \mathbf{D} | nrR | X | Х | Х | X | sqp | sqp | |
| Earlier | \mathbf{E} | nlW | X | Х | Х | X | sqp | X | |
| щ | F | nF | X | sqp | sqp | sqp | sqp | sqp | |

Fig. 4: The RDMA^{SC} ordering constraints on ippo (left) and oppo (right), where \checkmark denotes that instructions are ordered (and cannot be reordered), \checkmark denotes they are not ordered (and may be reordered), and sqp denotes they are ordered iff they are on the same queue pair.

of instantaneous events. Note that the observation point either follows the issue point (for NIC writes), or coincides (for instantaneous events).

We next define the 'issue-preserved program order', ippo, as the subset of po edges (ippo \subseteq po) that must be preserved when issuing instructions. That is, if two events are ippo-related, then they must be issued in program order; otherwise they may be processed in either order. The left table of Fig. 4 describes which po edges are included in ippo, where \checkmark denotes that the two instructions are ippo-related (i.e. they must be issued in program order), \checkmark denotes that they are not ippo-related (i.e. they may be issued out of order) and sqp denotes that they are ippo-related iff they are on the same queue pair. For instance, when a CPU instruction is followed by anything, they are issued in order (line A); but when a NIC instruction is followed by a CPU one, they may be reordered (cells B1-F1).

Analogously, we define the 'observation-preserved program order', oppo, as the subset of po edges (oppo \subseteq po) that must be preserved when observing the effects of instructions. I.e., if two events are oppo-related, then they become observable in program order in RDMA^{SC}; otherwise they may become observable in either order. The right table of Fig. 4 describes which po edges are included in oppo. The two tables differ in cells C6 and E6. This is because NIC writes can be delayed, and remote fences do not guarantee propagation to memory.

RDMA^{SC} Consistency. The notion of executions (Def. 3) imposes very few constraints on the po, pf, rf, mo, and nfo relations. Such restrictions and thus the permitted behaviours of a program are determined by defining the set of *consistent* executions, defined below.

Definition 4 (RDMA^{SC}-consistency). An execution $\langle E, po, pf, rf, mo, nfo \rangle$ is RDMA^{SC}-consistent iff ib and ob are irreflexive, where:

The ib (resp. ob) relation is an extension of ippo (resp. oppo), describing the issue (resp. observation) order across the instructions of different threads and nodes. RDMA^{SC}-consistency requires that ib and ob be irreflexive (i.e. yield strict partial orders as they are defined transitively).

The rf (resp. pf) component in ib states that if e reads from (resp. polls from) w, then w must have been issued before e. Recall that nfo totally orders the nlR/nlW and nrR/nrW operations on the same queue pair and is thus in ib. The rf component in ob states that if a read r reads from a write w, then the write has reached memory. This is because reads can only read the main memory and not auxiliary buffers. The [nlW]; pf component states that if p polls from a NIC local write w, then w must have left the PCIe fabric and reached the memory. Note that this is not the case for nrW events: polling an nrW event w might succeeds when w is still in the remote PCIe fabric before reaching the remote memory. The nfo in ob can be justified as in the case of ib. The rb component in ob ensures that a read r on x observes the latest write on x that has reached the memory. As mo describes the order in which the writes on each location reach the memory, it is included in ob. Let (τ_i, τ_o) be the issue and observation points of e and (τ_i', τ_o') be those of e'. The [Inst]; ib in ob ensures that if $e \stackrel{\text{ib}}{\rightarrow} e'$ (i.e. $\tau_i' < \tau_i'$) and e is instantaneous $(\tau_i = \tau_o)$, then $\tau_o = \tau_i < \tau_i' \le \tau_o'$, i.e. $e \stackrel{\text{ob}}{\rightarrow} e'$.

4 Robustness of RDMA^{SC} Programs

In the traditional setting of CPU concurrency (where all threads execute CPU instructions), the most intuitive consistency model is *sequential consistency* (SC) [42] While SC is too strong—in that disallowing *all* reorderings does not enable efficient implementations—it provides an intuitive and commonly understood model, making it easier for developers to reason about their programs.

Although none of the existing well-known consistency models follow SC by default, programmers typically address this difficulty by focusing on *robust* implementations of algorithms. Specifically, a program is robust under a weak consistency model CM if every possible behaviour of the program under CM is also an allowed behaviour under SC. In our model, this is defined as follows.

Definition 5 (SC-consistency and RDMA^{SC}-robustness). Given an execution $\langle E, po, pf, rf, mo, nfo \rangle$, its associated sequential-consistency relation is defined as $sc \triangleq (po \cup rf \cup rb \cup mo)$. An execution G is SC-consistent iff G.sc is acyclic. A pre-execution is robust under RDMA^{SC} iff all of its RDMA^{SC}-consistent executions (Def. 4) are also SC-consistent.

Our aim here is to provide guidelines to ensure the robustness of RDMA^{SC} programs. That is, we identify a number of *syntactic* requirements such that if a program fulfils them, then the behaviours of the program under RDMA^{SC} coincide with its behaviours under SC; i.e. the program does not exhibit any weak behaviours brought about by observable reorderings.

There are two complementary approaches to achieve robustness. The first is to structure the program in a way that limits the very existence of problematic

| | | | | | Late: | r in | ро | | | | | |
|---------|-------------------------|-----|----------------------|-----|-------|------|----------|----------|-----------------|----------|-----|--|
| | | | different queue pair | | | | | | same queue pair | | | |
| | gb | | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | |
| | | | CPU | nlR | nrW | nrR | nlW | nlR | nrW | nrR | nlW | |
| bo | A | CPU | ✓ | 1 | 1 | 1 | / | | N_{i} | /A | | |
| Щ. | В | nlR | Р | Р | Р | Р | Р | \ | \ | ✓ | 1 | |
| er | \mathbf{C} | nrW | GP | GP | GP | GP | GP | GP | ✓ | √ | 1 | |
| Earlier | $\overline{\mathbf{D}}$ | nrR | Р | Р | Р | Р | Р | F | F | F | 1 | |
| 茁 | E | nlW | Р | Р | Р | Р | Р | F | F | Р | 1 | |

Fig. 5: Constraints necessary to guarantee that a pair of po-related events in $\mathcal{R} \cup \mathcal{W}$ will be ob-related for any consistent execution. CPU denotes local events in $1 \text{W} \cup 1 \text{R} \cup \text{CAS}$. The \checkmark denotes that no additional constraint is needed and that the events are already in ob. P denotes that the earlier operation must be polled before executing the later one. F denotes that either the earlier operation must be polled (similar to P) or that a remote fence must be inserted between the two operations. GP denotes that a get operation and its associated poll on the first queue pair must be inserted between the two operations.

cases. The second is to extend the program with enough restrictions (e.g. polls and remote fences) to prohibit reorderings. In the next section (§5) we focus on the former and provide a set of explicit guidelines to avoid most problematic cases by design. In this section we focus on the latter, and describe how to identify problematic cases and how to block them. In what follows, we present the general syntactic restrictions required to forbid the reordering opportunities for specific operations (§4.1). We then propose sufficient syntactic conditions that block observable reorderings, and we prove that these conditions imply robustness (§4.2). Finally, we discuss the limitations of this approach (§4.3).

4.1 A Syntactic Approach to Enforce the Program Order

One of our key results relies on enforcing the program order (i.e. blocking instruction reordering) in potentially observable cases. Recall that given an execution, the observed-before order (ob) describes when an event takes effect before another. That is, for $(e_1, e_2) \in po$, when $e_1 \xrightarrow{ob} e_2$ in an execution G, then they are not reordered in G. Our first aim here is to identify syntactic constraints that ensure that a specific pair of given instructions (of the same thread) are related by ob. However, in order to define syntactic constraints for robustness, we can only rely on the syntax of the program and not components such as rf or mo. Our syntactic constraints can only rely on the pre-execution components po and pf, and we cannot directly use the ob relation derived from a specific execution.

To this end, we first define the *guaranteed-before* relation, $\mathsf{gb} \subseteq \mathsf{po}$, describing when two instructions in the same thread are guaranteed to remain in order (and their reordering is blocked), as shown in Fig. 5. Specifically, if two instructions are related by oppo , then they are guaranteed to be observed in that order and

thus there is no need for additional restrictions; this is denoted by \checkmark in cells A1–A5, B6–B9, C7–C9, D9, and E9 (*cf.* oppo in Fig. 4). For most other cases (noted P or F), polling the earlier instruction enforces the ordering. Recall that polling a NIC remote write does not guarantee its completion, and we need to add a 'dummy' get operation and its corresponding poll to ensure ordering (noted GP).

In most cases, when the two operations are on the *same* queue pair, then a remote fence is sufficient to enforce the ordering (noted F in D6–D8, E6–E7), and is a cheaper alternative to a poll. Perhaps surprisingly, a remote fence is not always sufficient: the two outliers are cells C6 and E8. For C6, consider the program $z^2 := x$; rfence(2); $w^2 := y$: the local value of y might be read before the value of z is changed. This is because rfence (2) (as with poll) only awaits the acknowledgement from the remote side which does not necessarily ensure that the first put has completed. For E8, consider $x := z^2$; rfence(2); $y := w^2$, where w^2 can be read before x is modified: rfence (2) only waits for the NIC local write $(x := v_z)$ to be sent to the local PCIe fabric and thus the put operation $(y := w^2)$ can start earlier than one could expect.

Definition 6 (guaranteed-before). Given a pre-execution $\mathcal{G} = \langle E, \mathsf{po}, \mathsf{pf} \rangle$, its guaranteed-before order, $\mathsf{gb} \subseteq \mathsf{po}$, is defined as $\mathsf{gb} \triangleq \mathsf{gb}_{\mathsf{base}}^+$, with:

```
\begin{array}{lll} {\sf gb_{base}} \triangleq & {\sf oppo} & ({\rm A1-A5,B6-B9,C7-C9,D9,E9 \ in \ Fig. \ 5}) \\ & \cup [{\sf n1R}]; {\sf po|}_{imm}; [{\sf nrW}]; {\sf pf} & ({\rm B1-B5 \ in \ Fig. \ 5}) \\ & \cup [{\sf nrW}]; ({\sf po} \cap {\sf sqp}); [{\sf n1W}]; {\sf pf} & ({\rm C1-C6 \ in \ Fig. \ 5}) \\ & \cup [{\sf nrR}]; {\sf po|}_{imm}; [{\sf n1W}]; {\sf pf} & ({\rm D1-D8 \ in \ Fig. \ 5}) \\ & \cup [{\sf n1W}]; ({\sf po} \cap {\sf sqp}); [{\sf nF}]; ({\sf po} \cap {\sf sqp}) & ({\rm D6-D8 \ in \ Fig. \ 5}) \\ & \cup [{\sf n1W}]; {\sf pf} & ({\rm E1-E8 \ in \ Fig. \ 5}) \\ & \cup [{\sf n1W}]; ({\sf po} \cap {\sf sqp}); [{\sf nF}]; ({\sf po} \cap {\sf sqp}); [{\sf n1R} \cup {\sf nrW}] & ({\rm E6-E7 \ in \ Fig. \ 5}) \end{array}
```

Given an execution $G = \langle E, \mathsf{po}, \mathsf{pf}, \mathsf{rf}, \mathsf{mo}, \mathsf{rb} \rangle$, we write $G.\mathsf{gb}$ for $\langle E, \mathsf{po}, \mathsf{pf} \rangle.\mathsf{gb}$. Finally, we prove that gb implies ob for any RDMA^{SC}-consistent execution (see Theorem 7 for the proof).

Theorem 1 (gb implies ob). Given a pre-execution $\langle E, \mathsf{po}, \mathsf{pf} \rangle$, for all RDMA^{SC}-consistent executions $G = \langle E, \mathsf{po}, \mathsf{pf}, \mathsf{rf}, \mathsf{mo}, \mathsf{nfo} \rangle$ and all $\mathsf{e}_1, \mathsf{e}_2 \in E$, if $(\mathsf{e}_1, \mathsf{e}_2) \in G.\mathsf{gb}$, then $(\mathsf{e}_1, \mathsf{e}_2) \in G.\mathsf{ob}$.

Given Theorem 1 above, we can use gb as a tool to enforce robustness. Specifically, whenever a program order pair $(e_1,e_2) \in po$ may be reordered, we can add the prescribed fences to enforce $(e_1,e_2) \in gb$ and thus block the reordering. The rest of this section describes when we should use this tool.

4.2 Conditions for Robustness under RDMA^{SC}

As mentioned before, blocking all instruction reorderings, i.e. by requiring po = gb, would enforce sequential consistency and thus robustness. However, this is too strict and highly impractical. Instead, we should ideally enforce gb selectively when needed and only prevent observable reorderings.

Two sources of weak behaviours. As presented in §2, RDMA^{SC} programs have two distinct sources of weak behaviours. These come from two different kinds of pairs of events (of the same thread): (1) pairs forming a data race on a certain location, e.g. $a := y^2; y^2 := 1$, as presented in Fig. 1f (copied below-left) and Figs. 1b, 1e, and 1h; and (2) pairs whose reordering can be observed by other threads, e.g. $a := y^2; z^3 := 1$, as in the examples of Fig. 3b (copied below-right).

| | y=0 |
|-----------------------|-----|
| $a := y^2$ $y^2 := 1$ | |
| a = 1 | 1 |

| | y = 0 | z = 0 | | | | |
|--|------------|-------|--|--|--|--|
| $\begin{vmatrix} a := y^2 \\ z^3 := 1 \end{vmatrix}$ | $y := z^3$ | | | | | |
| $a=1$ \checkmark | | | | | | |

As such, stopping these two sources of weak behaviours would be enough to ensure robustness. Data races within a thread are always problematic, no matter the context, and we always need to block the reordering of such pairs (i.e. enforce gb to ensure the pair is ob-ordered in any execution). Pairs of the second kind cannot create weak behaviours by themselves, but they might allow weak behaviours depending on the rest of the program of other threads. In the next section (§5), we show conditions making sure that such pairs can never create weak behaviours by design. In this section, we focus on deciding whether such a pair might lead to a weak behaviour and, if so, how to block the reordering.

To formulate this intuition, we write public(x) to denote that x is a public location accessed by multiple threads, and given a set of events E, we define the set of public events in E as $E^{\text{pub}} \triangleq \{e \in E \mid \text{public}(\log(e))\}$. We further define $E \setminus t \triangleq \{e \in E \mid t(e) \neq t\}$ for the set of events in E that are not by thread t. We can then formulate the two categories of weak behaviours above as two kinds of sc cycles: sc cycles on a single thread (1) and sc cycles on public events across threads (2), as formulated below (see Theorem 8 for the full proof).

Theorem 2 (sc cycle decomposition). Given a RDMA^{SC}-consistent execution $G = \langle E, po, pf, rf, mo, nfo \rangle$, $if \exists e \in E. e \xrightarrow{G.sc}^+ e$ (i.e. a cycle in G.sc), then:

- $\begin{array}{lll} \ \ either \ there \ is \ a \ G.sc \ \ cycle \ on \ a \ single \ thread, \ i.e. \ \exists e \in E. \ e \xrightarrow{G.sc \cap sthd}^+ e; \\ \ \ or \ there \ \ exists \ e_1, e_2 \in E^{pub} \ \ such \ that \ e_1 \xrightarrow{po \setminus G.ob} e_2 \xrightarrow{(G.sc; [E^{pub} \setminus t(e_1)])^+; G.sc} e_1. \end{array}$ That is, there is an sc cycle on public events, with two po-related events on some thread $t(e_1)$ not related in ob, and where the rest of the cycle does not go through the events of $t(e_1)$.

The two kinds of problematic reorderings are tackled separately below, and Theorem 5 confirms the two resulting conditions are sufficient for robustness.

Preventing sc cycles from data races. As shown above, when an allowed reordering is part of a data race, it becomes observable independently from the context. Thus, we should always preclude this kind of reordering. Specifically, in Def. 7 below we present a local data-race freedom property to block data races within each thread and prevent single-threaded weak behaviours.

Definition 7 (Local DRF). Given a pre-execution $\langle E, \mathsf{po}, \mathsf{pf} \rangle$, two events $\mathsf{e}_1, \mathsf{e}_2 \in E$ are locally conflicting iff 1. $(\mathsf{e}_1, \mathsf{e}_2) \in \mathsf{sthd}$; 2. $\mathsf{loc}(\mathsf{e}_1) = \mathsf{loc}(\mathsf{e}_2)$; and 3. at least one of $\mathsf{e}_1, \mathsf{e}_2$ is a write event. A pre-execution $\mathcal G$ is locally datarace free (LDRF), iff for all $\mathsf{e}_1, \mathsf{e}_2 \in \mathcal G.E$, if $\mathsf{e}_1, \mathsf{e}_2$ are locally conflicting, then $(\mathsf{e}_1, \mathsf{e}_2) \in \mathcal G.\mathsf{gb} \cup \mathcal G.\mathsf{gb}^{-1}$. Put differently, given the definition of gb (Fig. 5), a pre-execution $\langle E, \mathsf{po}, \mathsf{pf} \rangle$ is LDRF iff for all locally conflicting accesses $\mathsf{e}_1, \mathsf{e}_2 \in E$, if $(\mathsf{e}_1, \mathsf{e}_2) \in \mathsf{po}$, then the following four conditions hold:

- 1. If $e_1 \in nlW$ and $(e_1, e_2) \notin sqp$, then there exists $e_3 \in P$ such that $(e_1, e_3) \in pf$ and $(e_3, e_2) \in po$ (cells E1, E2, and E5 in Fig. 5).
- 2. If $e_1 \in nlW$, $e_2 \in nlR$, and $(e_1, e_2) \in sqp$, then either there exists $e_3 \in P$ with $(e_1, e_3) \in pf$ and $(e_3, e_2) \in po$; or there exists $e_3 \in nF$ with $(e_1, e_3) \in po$ and $(e_3, e_2) \in po$ (E6).
- 3. If $e_1 \in nlR$, $e_2 \in (nlW \cup lW \cup CAS)$, and $(e_1, e_2) \not\in sqp$, then there exists $e_1' \in nrW$ and $e_3 \in P$ such that $(e_1, e_1') \in po|_{imm}$, $(e_1', e_3) \in pf$, and $(e_3, e_2) \in po$ (cells B1 and B5).
- 4. If $e_1 \in nrR$ and $e_2 \in nrW$, then either there exists $e_3 \in nF$ such that $e_1 \xrightarrow{po \cap sqp} e_3 \xrightarrow{po \cap sqp} e_2$; or there exists $e_1' \in nlW$ and $e_3 \in P$ such that $e_1 \xrightarrow{po|_{imm}} e_1' \xrightarrow{pf} e_3 \xrightarrow{po} e_2$ (cell D7 in Fig. 5).

These cases prohibit all possible races on a location x, i.e. of the form $x := y^n; x := -(E1,E5), x := y^n; -:= x (E2), x := y^n; z^n := x (E6), y^n := x; x := -(B1,B5), or <math>-:= x^n; x^n := -(D7)$. Other entries in Fig. 5 cannot create races as either their ordering is already guaranteed (e.g. \checkmark in E9); or they are on two read events (e.g. B2,D8); or they cannot be on the same location (e.g. D3,E7).

We argue that the constraints in Def. 7 do not restrict RDMA capabilities in that waiting for remote operations to complete before reusing their locations is already considered standard practice when writing RDMA programs.

We next show that LDRF prevents single-threaded weak behaviours (see Theorem 9 for the proof).

Theorem 3. Given a RDMA^{SC}-consistent execution $G = \langle E, \mathsf{po}, \mathsf{pf}, \mathsf{rf}, \mathsf{mo}, \mathsf{nfo} \rangle$, if $\langle E, \mathsf{po}, \mathsf{pf} \rangle$ is locally data-race free, then there is no sc cycle on a single thread; that is, $(G.\mathsf{sc} \cap \mathsf{sthd})$ is acyclic and the first case of Theorem 2 does not arise.

Preventing sc cycles across threads. Unlike data races, pairs of the second kind cannot create weak behaviours by themselves, and their reorderings can only be observed in certain contexts.

The general strategy to prevent observable reorderings is straightforward: for every pair $(e_1, e_2) \in po$ on public locations, either we know for certain that $e_2 \stackrel{sc}{\longrightarrow}^* e_1$ (using other threads) is impossible, or we conservatively block the reordering by enforcing $(e_1, e_2) \in gb$. The challenge is that the relation sc is heavily dependent on the specific execution. So how can we ascertain *syntactically* that a later event e_2 cannot influence an earlier event e_1 ?

One easily accessible syntactic property is the communication pattern between nodes (i.e. when one node performs a remote operation towards another).

Thus, to simplify the task, we over-approximate dependency (i.e. sc) with *communication*. Intuitively, if two nodes do not communicate in the network topology, then they cannot causally influence each other.

We write $n_1 \iff n_2$ (defined below) to denote that nodes n_1 and n_2 communicate via some event in E, in that some thread t on n_1 performs a remote operation $e \in E$ towards n_2 , written hasQP (t, n_2, E) , or vice versa.

$$\begin{split} n_1 & \underset{E}{\longleftrightarrow} n_2 \triangleq \exists t. \; (n(t) = n_1 \land \mathsf{hasQP}(t, n_2, E)) \lor (n(t) = n_2 \land \mathsf{hasQP}(t, n_1, E)) \\ & \mathsf{hasQP}(t, \overline{n}, E) \triangleq \exists \mathsf{e} \in (E.\mathtt{nrW} \cup E.\mathtt{nrR}). \; t(\mathsf{e}) = t \land \overline{n}(\mathsf{e}) = \overline{n} \end{split}$$

We next show that if there is an sc-path from one event e_2 to another e_1 using public events in A, then the corresponding nodes (of the locations) of e_2 and e_1 must communicate via A. This is established in Lem. 1 below, with the proof given in Lem. 2.

Lemma 1. For all
$$A \subseteq E^{\mathsf{pub}}$$
, if $\mathsf{e}_2 \xrightarrow{\mathsf{sc}|_A}^* \mathsf{e}_1$ then $n(\mathsf{loc}(\mathsf{e}_2)) \leftrightsquigarrow^* n(\mathsf{loc}(\mathsf{e}_1))$.

We are interested in the inverse direction of this lemma: a topological connection between the nodes (of the locations) of \mathbf{e}_2 and \mathbf{e}_1 is a necessary condition for an \mathbf{sc} -path from \mathbf{e}_2 to \mathbf{e}_1 . Put differently, if there is no communication between the nodes of \mathbf{e}_2 and \mathbf{e}_1 , then \mathbf{e}_2 cannot influence \mathbf{e}_1 . As such, we can use this to over-approximate safely whether an event can influence another. We conservatively assume that if the two nodes can communicate (outside of the thread) then \mathbf{e}_2 might influence \mathbf{e}_1 . These communications do not depend on a specific execution and can be ascertained syntactically from the pre-execution.

We can then prevent sc cycles across threads using the *fenced* condition below (Def. 8): for all $e_1 \xrightarrow{po} e_2$ on public locations, if e_2 might influence e_1 , then we block the reordering. We subsequently prove that if a pre-execution is fenced, then it does not admit sc cycles across threads.

Definition 8 (fenced). A pre-execution
$$\langle E, \mathsf{po}, \mathsf{pf} \rangle$$
 is fenced iff for all $\mathsf{e}_1, \mathsf{e}_2 \in E^{\mathsf{pub}}$, if $\mathsf{e}_1 \stackrel{\mathsf{po}}{\longrightarrow} \mathsf{e}_2$ and $n(\mathsf{loc}(\mathsf{e}_1)) \stackrel{\longleftrightarrow}{\longleftarrow} {}^* n(\mathsf{loc}(\mathsf{e}_2))$, then $(\mathsf{e}_1, \mathsf{e}_2) \in \mathsf{gb}$.

Theorem 4. Given an RDMA^{SC}-consistent execution $\langle E, \mathsf{po}, \mathsf{pf}, \mathsf{rf}, \mathsf{mo}, \mathsf{nfo} \rangle$, if its associated pre-execution $\langle E, \mathsf{po}, \mathsf{pf} \rangle$ is fenced, then there is no sc cycle of the shape $\mathsf{e}_1 \xrightarrow{\mathsf{po} \backslash \mathsf{ob}} \mathsf{e}_2 \xrightarrow{(\mathsf{sc}; [E^{\mathsf{pub}} \backslash t(\mathsf{e}_1)])^+; \mathsf{sc}} \mathsf{e}_1$ with $\mathsf{e}_1, \mathsf{e}_2 \in E^{\mathsf{pub}}$. That is, the second case of Theorem 2 does not arise.

Robustness. Lastly, we show that LDRF and fenced imply robustness under RDMA^{SC}. Thus, this approach can be used to prevent RDMA weak behaviours.

Theorem 5 (Robustness under RDMA^{SC}). Given a pre-execution $\mathcal{G} = \langle E, po, pf \rangle$, if \mathcal{G} is locally data-race free (Def. 7) and fenced (Def. 8), then \mathcal{G} is also robust under RDMA^{SC} (Def. 5).

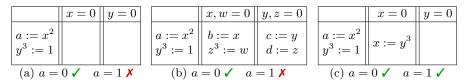


Fig. 6: Examples illustrating the limitation of Theorem 5, where the programs in (a) and (b) are robust (the weak behaviour a=1 is not allowed in either) while that in (c) is not robust (it admits the weak behaviour a=1); while Theorem 5 rightfully identifies (a) as robust (true positive) and (c) as not robust (true negative), it conservatively deems (b) not robust (false negative).

4.3 Usage and Limitations

Local data-race freedom (Def. 7) and fenced (Def. 8) are intuitive properties that can be checked syntactically. Indeed, given a program, it is straightforward to check mechanically whether these properties hold or to provide an explicit counterexample and a suggested fix using the definition of gb (Def. 6). As a result, sufficient constraints can automatically be added to ensure robustness.

However, this simplicity can occasionally be the limitation of our approach. Specifically, as the main theorem does not account for interactions between threads, it takes a conservative approach, which at times can lead to false negatives (where the program is deemed not robust even though no weak behaviours are possible), recommending unnecessary restrictions.

To see this, consider the example in Fig. 6a, where $a := x^2$ and $y^3 := 1$ can be reordered *without* introducing weak behaviours. In this case, Theorem 5 rightfully confirms that no additional restrictions are necessary. By contrast, consider the variant shown in Fig. 6b: although the two extended threads do not introduce any additional weak behaviours, our approach assumes there might be a causal dependency from $y^3 := 1$ to $a := x^2$, as is the case e.g. in Fig. 6c. As such, Theorem 5 cannot determine Fig. 6b as robust, and our approach would recommend inserting a poll operation in the first thread. Note that removing any of the six operations would enable Theorem 5 to ascertain Fig. 6b as robust.

Understanding that the reordering of the instructions in the first thread of Fig. 6b is not problematic would require a more complex static analysis beyond the scope of this paper.

5 Application: Tree Topology

Theorem 5 outlines the conditions under which we can guarantee that a program is robust under RDMA^{SC}. However, while the LDRF property (Def. 7) is reasonable, the fenced property (Def. 8) can lead to excessive restrictions (e.g. as in Fig. 6). Specifically, for every pair of events (e_1, e_2) in program order, we must either verify that e_2 cannot affect e_1 , or ensure that their execution order is preserved. The main issue is that preserving the order of every pair of events can be particularly costly, notably when considering NIC remote write events. In

such cases, the only resort is to introduce a 'dummy' get operation and poll it, which is inefficient. Instead, we propose a strategy whereby we stipulate certain conditions on the network *topology* (i.e. the shape of the RDMA network) so that later events are often unable to influence earlier events.

To this end, we propose a *tree topology* that balances generality (supporting a wide range of programs) with efficiency and restrictiveness (requiring minimal additional constraints to respect the fenced property). In §5.1 we present an overview of our new set of restrictions and illustrate their rationale through examples. In §5.2 we formalises these restrictions and prove that they indeed imply robustness under RDMA^{sc}. Finally, in §5.3 we demonstrate specific applications of the tree topology and how RDMA programs can make use of them.

5.1 Overview of the Restrictions

We describe four different conditions that, if satisfied, ensure the robustness of RDMA programs under RDMA^{SC}, and we justify them through examples.

LDRF. As before, we require that programs satisfy LDRF (Def. 7). As discussed, this is considered standard practice when writing RDMA programs and should not be seen as a limitation.

Private Copies. We require the *local locations* of RDMA operations – e.g. location y in $y:=x^2$ – to be private (i.e. accessed by only one thread, namely that executing the RDMA operation). Intuitively, to maximise the efficiency of RDMA programs, we should ideally allow arbitrary interleaving of RDMA operations and CPU computations. For instance, let us consider the single-threaded program $\mathsf{C} \triangleq y := x^2; \mathsf{c}^{\sharp}_{\mathsf{cpu}}$, where $\mathsf{c}^{\sharp}_{\mathsf{cpu}}$ denotes a block of CPU instructions that does not access location y. If y is private, then although $\mathsf{c}^{\sharp}_{\mathsf{cpu}}$ and the get $y := x^2$ may be reordered, this reordering will not lead to any observable weak behaviours. That is, when we run C concurrently with any RDMA program C' (i.e. as $\mathsf{C}||\mathsf{C}'|$), if y is private, then we do not need to poll $y := x^2$ before proceeding with $\mathsf{c}^{\sharp}_{\mathsf{cpu}}$ (even though they may be reordered), as the reordering cannot be observed by C' .

However, if y is accessible by other threads (on the same node or from a remote node), then the reordering becomes visible, allowing additional, potentially unwanted, weak behaviours. This is illustrated in the example below, where $c_{\text{CDU}}^{\text{y}} \triangleq z := 1$ and y is public (accessed by nodes 1 and 3).

| y, z = 0 | x = 1 | | | | |
|----------------------------|-------|------------|--|--|--|
| $y := x^2$ | | $a := z^1$ | | | |
| z := 1 | | poll(1) | | | |
| | | $b := y^1$ | | | |
| $(a,b) = (1,0) \checkmark$ | | | | | |

More concretely, due to the reordering, the later CPU computation (z := 1) can be observed before the earlier get $(y := x^2)$, leading to the weak outcome (a, b) = (1, 0).

Therefore, to prevent such weak behaviours, we stipulate that local locations of RDMA operations be private. This is not a costly limitation. Specifically, in the case of put operations, the data can easily be copied beforehand to a

one-time-use private location. In the case of get operations, it means the thread running the command needs to acknowledge the data and copy it to make it available to other threads having access to the node.

Get in Order. We stipulate that each get operation be followed by a remote fence. Recall that only certain reorderings are allowed on the operations of the same queue pair. Intuitively, put operations cannot be overtaken, and we do not need to restrict their usage. However, get operations can be overtaken by other get/put operations, as shown in the examples below, where the $a:=x^2$ is overtaken by a later remote operation on the same queue pair, leading to weak behaviours.

| | x = 0 | | x, y = 0 |
|-----------------------|----------------------|-----------------------|----------|
| $a := x^2$ $b := x^2$ | x := 1 | $a := x^2$ $y^2 := 1$ | x := y |
| (a, b) = | $(1,0)$ \checkmark | a = | 1 🗸 |

As such, to prevent non-SC behaviours, we require that each get operation be followed by a remote fence, forcing the queue pair to await the completion of the get before starting the next remote operation. Of course, if the get is polled before another RDMA operation is submitted, the remote fence is not needed. Note that since remote fences do not block CPU computations nor communications with other nodes, they are not very expensive and are a reasonable cost to pay to ensure remote operations towards a specific remote node stay in order.

Tree Topology. Finally, the most important restriction is to constrain the topology of the network over which the program runs. Intuitively, having multiple paths between a set of nodes allows for visible effects to overtake each other (i.e. be reordered) along different paths, leading to weak behaviours. In the extreme case where every thread can communicate directly with every other node, we allow for a large number of visible reorderings, and lose any hope of preventing non-SC behaviours. When such connected topologies are needed to enable more efficient implementations (e.g. consensus algorithms), the developers must carefully account for the possible weak behaviours.

Our proposal is to adhere to a minimal topology where there is (at most) a single communication path between each pair of nodes. In the examples below we show how not adhering to the tree topology can lead to weak behaviours. Note that although we have followed each remote operation with a corresponding (costly) poll, we still cannot prevent the weak behaviours shown.

| y = 0 | x = 0 | | | x = 0 | y, z = 0 | x = | = 0 | y, z = 0 |
|-------------------------------|--|---|--|--------------|-----------------|-------------------------------|-------------|--|
| $x^2 := 1$ $poll(2)$ $a := y$ | $\begin{vmatrix} y^1 := 1 \\ \text{poll}(1) \\ b := x \end{vmatrix}$ | | $\begin{vmatrix} z^3 := 1 \\ \text{poll}(3) \\ x^2 := 1 \end{vmatrix}$ | $y^3 := x$ | a := y $b := z$ | $z^2 := 1$ $poll(2)$ $x := 1$ | $y^2 := x$ | $\begin{vmatrix} a := y \\ b := z \end{vmatrix}$ |
| (a, b) = | (0,0) 🗸 | | (a, b) | (b) = (1, 0) |) ✓ | (a, b) | (0) = (1,0) | ✓ |
| - mi | | • | | | | | /= . | |

The first example shows that queue pairs in both directions (between nodes 1 and 2) can lead to weak behaviours as they can observe the reordering of

operations on the other node. The second example illustrates two paths between node 1 and 3: a direct path from node 1 to 3 (via $z^3 := 1$) and an indirect path through node 2 (from node 1 to 2 via $x^2 := 1$; from node 2 to 3 via $y^3 := x$). As shown, having multiple paths between two nodes allows threads to observe reorderings: $z^3 := 1$ is submitted first, but the effects of $x^2 := 1$, forwarded via $y^3 := x$, is observed first. The third example is a variant of the second, where the middle node is replaced by an additional thread on the left node. As queue pairs from different threads of the same node towards the same remote are still independent, the weak behaviour shown is permitted.

5.2 Tree Robustness

We next formalise the conditions described in §5.1 in Def. 9 below.

Definition 9 (tree-fenced). A pre-execution $\langle E, po, pf \rangle$ is tree-fenced iff:

- 1. Local locations of RDMA operations are private: $E^{\mathsf{pub}}.\mathtt{nlR} = E^{\mathsf{pub}}.\mathtt{nlW} = \emptyset$
- 2. Each get operation is followed by a remote fence (or is polled) before the next remote operation on the same queue pair.
 - That is, for all e_1, e_2 , if $e_1 \in nrR$, $e_2 \in (nrR \cup nrW)$, and $(e_1, e_2) \in (po \cap sqp)$, then: either there exists $f \in nF$ such that $(e_1, f) \in (po \cap sqp)$ and $(f, e_2) \in (po \cap sqp)$; or there exists $e_3 \in nlW$ and $p \in P$ such that $(e_1, e_3) \in po|_{imm}$, $(e_3, p) \in pf$, and $(p, e_2) \in po$.
- 3. There is (at most) a single communication path between any pair of nodes in that the following three properties hold:
 - (a) The network does not have cycles, i.e. for all sets of distinct nodes $\{n_1; \ldots; n_k\}$ with k > 2: $\neg(n_1 \underset{E}{\longleftrightarrow} n_2 \underset{E}{\longleftrightarrow} \ldots \underset{E}{\longleftrightarrow} n_k \underset{E}{\longleftrightarrow} n_1)$
 - (b) No two nodes have queue pairs towards each other: $\neg \exists t_1, t_2$. hasQP $(t_1, n(t_2), E) \land \text{hasQP}(t_2, n(t_1), E)$
 - (c) Each node has at most one queue pair towards each remote node: $\forall t, t', \overline{n}. \ t \neq t' \land \mathsf{hasQP}(t, \overline{n}, E) \land \mathsf{hasQP}(t', \overline{n}, E) \Longrightarrow n(t) \neq n(t')$

Conditions 1 and 2 are purely syntactic and can be straightforwardly checked by examining the RDMA program. Condition 3 pertains to the topology of the RDMA network and can also be checked by examining the RDMA program.

A key advantage of these restrictions is that preventing weak behaviours never requires polling remote operations. This is crucial because the efficiency of RDMA implementations comes from parallelising data transfers and computations. As shown in the overview (§2), polling is very costly as it completely halts local computations and prevents submission of remote operations to any queue pair. With a tree topology, programmers only need to wait for remote operations to use their results (as per LDRF Def. 7), and do not need to sacrifice computation time to prevent reorderings.

We next prove that if a pre-execution is tree-fenced, then it is also fenced. The full proof is given in the appendix (see Theorem 12).

Theorem 6. If a pre-execution is tree-fenced (Def. 9), then it is fenced (Def. 8).

Hence, LDRF and tree-fenced properties imply robustness under RDMA^{SC}.

Corollary 1 (Tree robustness under RDMA^{SC}). If a pre-execution $\mathcal{G} = \langle E, \mathsf{po}, \mathsf{pf} \rangle$ satisfies LDRF (Def. 7) and is tree-fenced (Def. 9), then it is also robust under RDMA^{SC} (Def. 5).

5.3 Specific Applications

The tree-fenced conditions above provide guidelines to ensure programs cannot exhibit weak behaviours. While not all RDMA programs follow the restrictions presented, a tree topology is sufficient for a range of applications. Notably, any setup using RDMA solely for the data transfer capabilities (and not for distributed computations) can easily be configured as a tree.

Star Topology: Single Manager Multiple Workers. The star topology is one of the most typical network configurations, providing simple and reliable communication between nodes, with many common applications such as for implementing local area networks (LAN). The star topology allows a main node to distribute jobs to other nodes and periodically check for progress. As demonstrated in this paper, this setup prevents any network weak behaviour even if communications towards different workers are independent and can be reordered.

Star Topology: Single Server Multiple Clients. The tree-fenced condition (Def. 9) is permissive enough to allow us to translate common concurrent algorithms (comprising loads and stores over shared memory) to distributed ones over RDMA automatically as follows. Specifically, consider a concurrent algorithm P_c using k threads $(t_1, ..., t_k)$. We can translate this to a corresponding RDMA program P_r using k nodes $(n_1, ..., n_k)$, where a designated node (say n_1) is the server and the others $(n_2, ..., n_k)$ are clients, and each node n_i has a single thread simulating t_i . All shared locations and data are located on the server node $(n_1$ running t_1). For each of the remaining nodes n_i , we replace the loads and stores on shared locations with get and put operations, respectively. Moreover, we insert a remote fence after each get operation (to ensure condition (2) of Def. 9) and poll get operations before using their values (to ensure LDRF).

The resulting RDMA program follows a star topology, with n_1 as the central (server) node accessed by multiple clients $(n_2, ..., n_k)$. Client locations are private by definition, ensuring that the tree-fenced condition holds. P_r thus avoids weak behaviours and constitutes a suitable implementation of P_c .

Observe that in this implementation, polling put operations is unnecessary (as long as different local locations are used for copying), and get operations can be optimised by being submitted as early as possible (i.e. after previous RDMA operations and reads on the same location) and before they are needed, allowing them to be interleaved with other computations.

6 Related Work

RDMA Semantics. The first realistic formal model for RDMA programs is RDMA^{TSO} by Ambal et al. [10] (where they assume that CPU concurrency is

governed by TSO) formalised both operationally and declaratively, which they show to be equivalent. They also validate RDMA^{TSO} empirically by running an extensive suite of litmus tests on RDMA hardware. While comprehensive in its formal description of the language, this work does not present strategies for mitigating RDMA weaknesses or optimising the use of this technology by using e.g. minimal poll and fence instructions. The only other work on formal RDMA semantics is that by Dan et al. [26], which as demonstrated by Ambal et al. [10] does not follow the RDMA specification.

Weak Memory Models. Existing literature includes multiple examples of weak consistency models. For hardware, several works have formalised the semantics of the x86, ARMv8 and POWER architectures [67,9,2,62,47,5,58,30,66]. However, none of these works covered the consistency semantics of RDMA programs. For software, there has been a number of formal models for C11 [41,39,11,36,43,52,55,24] with verified compilation schemes [57,56,50], Java [48,14], transactional memory [71,60,59], the Linux kernel [8] and the ext4 filesystem [38]. Additionally, there has been several works on formalising the persistency semantics of programs in the context of non-volatile memory, describing the behaviour of programs in case of crashes [65,64,63,25,37], as well as program logics for verifying such programs [61,16,69].

Robustness. The concept of robustness against weak memory semantics has been extensively studied across various models as a means to simplify programming, reasoning, and verification. Notably, robustness for Total Store Order (TSO) and its Partial Store Order (PSO) variant [35,54,9] has received significant attention, e.g. [22,23,53,19,34,46,17,1,2,18,44,45]. In addition, robustness has been used as a correctness notion in the context of automatic fence insertion for weak hardware memory models [7,28,27,21,6]. More recent work has developed techniques for checking robustness against concurrency semantics in programming languages, particularly the C11 memory model [40,49]. Robustness has also been explored in distributed systems, where Sequential Consistency (SC) is replaced by serialisability [29,15,51,20,13,12]. More recently, [33] addressed the problem of checking robustness in the context of weak persistency models for non-volatile memory.

Some of these works provide sound and complete techniques for verifying robustness, along with complexity bounds for specific models. Others, as with our work on RDMA, focus on practical over-approximations, offering programmers guidelines that, when followed, ensure stronger semantics. The well-known Data-Race-Free (DRF) guarantee [3,32] for multicore hardware and programming language models is a prominent criterion of this type.

Acknowledgements. We thank the anonymous reviewers for their valuable feedback and Viktor Vafeiadis for many fruitful discussions. Ambal is supported by the EPSRC grant EP/X037029/1. Lahav is supported by the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation programme (grant agreement no. 851811) and the Israel Science Foundation (grant no. 814/22). Raad is supported by a UKRI fellowship MR/V024299/1, by the EPSRC grant EP/X037029/1 and by VeTSS.

References

- Abdulla, P.A., Atig, M.F., Lång, M., Ngo, T.P.: Precise and sound automatic fence insertion procedure under PSO. In: NETYS. pp. 32–47. Springer International Publishing, Cham (2015)
- Abdulla, P.A., Atig, M.F., Ngo, T.P.: The best of both worlds: Trading efficiency and optimality in fence insertion for tso. In: Proceedings of the 24th European Symposium on Programming on Programming Languages and Systems - Volume 9032. pp. 308-332. Springer-Verlag New York, Inc., New York, NY, USA (2015). https://doi.org/10.1007/978-3-662-46669-8_13, http://dx.doi.org/ 10.1007/978-3-662-46669-8_13
- Adve, S.V., Hill, M.D.: Weak ordering—a new definition. In: ISCA. pp. 2-14. ACM, New York (1990). https://doi.org/10.1145/325164.325100, http://doi.acm. org/10.1145/325164.325100
- Aguilera, M.K., Ben-David, N., Guerraoui, R., Marathe, V.J., Zablotchi, I.: The impact of RDMA on agreement. In: Robinson, P., Ellen, F. (eds.) Proceedings of the 2019 ACM Symposium on Principles of Distributed Computing, PODC 2019, Toronto, ON, Canada, July 29 - August 2, 2019. pp. 409–418. ACM (2019). https://doi.org/10.1145/3293611.3331601, https://doi.org/10.1145/3293611.3331601
- Alglave, J., Deacon, W., Grisenthwaite, R., Hacquard, A., Maranget, L.: Armed cats: Formal concurrency modelling at arm. ACM Trans. Program. Lang. Syst. 43(2), 8:1–8:54 (2021). https://doi.org/10.1145/3458926, https://doi.org/10.1145/3458926
- Alglave, J., Kroening, D., Nimal, V., Poetzl, D.: Don't sit on the fence: a static analysis approach to automatic fence insertion. ACM Trans. Program. Lang. Syst. 39(2), 6:1-6:38 (May 2017). https://doi.org/10.1145/2994593, http://doi.acm.org/10.1145/2994593
- Alglave, J., Maranget, L.: Stability in weak memory models. In: CAV. pp. 50-66.
 Springer-Verlag, Berlin, Heidelberg (2011), http://dl.acm.org/citation.cfm?
 id=2032305.2032311
- Alglave, J., Maranget, L., McKenney, P.E., Parri, A., Stern, A.: Frightening small children and disconcerting grown-ups: Concurrency in the linux kernel. SIGPLAN Not. 53(2), 405–418 (Mar 2018). https://doi.org/10.1145/3296957.3177156, https://doi.org/10.1145/3296957.3177156
- Alglave, J., Maranget, L., Tautschnig, M.: Herding cats: Modelling, simulation, testing, and data mining for weak memory. ACM Trans. Program. Lang. Syst. 36(2) (Jul 2014). https://doi.org/10.1145/2627752, https://doi.org/10.1145/2627752
- Ambal, G., Dongol, B., Eran, H., Klimis, V., Lahav, O., Raad, A.: Semantics
 of remote direct memory access: Operational and declarative models of rdma on
 tso architectures. Proc. ACM Program. Lang. 8(OOPSLA2) (Oct 2024). https://doi.org/10.1145/3689781
- Batty, M., Owens, S., Sarkar, S., Sewell, P., Weber, T.: Mathematizing c++ concurrency. In: Proceedings of the 38th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages. pp. 55-66. POPL '11, ACM, New York, NY, USA (2011). https://doi.org/10.1145/1926385.1926394, http://doi.acm.org/10.1145/1926385.1926394
- 12. Beillahi, S.M., Bouajjani, A., Enea, C.: Checking robustness against snapshot isolation. In: Computer Aided Verification. pp. 286–304. Springer International Publishing, Cham (2019)

- Beillahi, S.M., Bouajjani, A., Enea, C.: Robustness against transactional causal consistency. In: CONCUR 2019. vol. 140, pp. 30:1–30:18. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, Dagstuhl, Germany (2019). https://doi.org/10.4230/ LIPIcs.CONCUR.2019.30
- Bender, J., Palsberg, J.: A formalization of java's concurrent access modes. Proc. ACM Program. Lang. 3(OOPSLA) (Oct 2019). https://doi.org/10.1145/3360568, https://doi.org/10.1145/3360568
- 15. Bernardi, G., Gotsman, A.: Robustness against consistency models with atomic visibility. In: CONCUR. pp. 7:1-7:15. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, Dagstuhl, Germany (2016). https://doi.org/10.4230/LIPIcs.CONCUR. 2016.7, http://drops.dagstuhl.de/opus/volltexte/2016/6165
- 16. Bila, E.V., Dongol, B., Lahav, O., Raad, A., Wickerson, J.: View-based owickigries reasoning for persistent x86-tso. In: Sergey, I. (ed.) Programming Languages and Systems. pp. 234–261. Springer International Publishing, Cham (2022)
- 17. Bouajjani, A., Derevenetc, E., Meyer, R.: Checking and enforcing robustness against TSO. In: ESOP. pp. 533-553. Springer-Verlag, Berlin, Heidelberg (2013). https://doi.org/10.1007/978-3-642-37036-6_29, http://dx.doi.org/10.1007/978-3-642-37036-6_29
- Bouajjani, A., Enea, C., Mutluergil, S.O., Tasiran, S.: Reasoning about TSO programs using reduction and abstraction. In: CAV. pp. 336–353. Springer, Cham (2018)
- Bouajjani, A., Meyer, R., Möhlmann, E.: Deciding robustness against total store ordering. In: ICALP. pp. 428–440. Springer, Berlin, Heidelberg (2011)
- Brutschy, L., Dimitrov, D., Müller, P., Vechev, M.: Static serializability analysis for causal consistency. In: PLDI. pp. 90-104. ACM, New York (2018). https://doi.org/10.1145/3192366.3192415, http://doi.acm.org/10.1145/3192366. 3192415
- Burckhardt, S., Alur, R., Martin, M.M.K.: CheckFence: Checking consistency of concurrent data types on relaxed memory models. In: PLDI. pp. 12–21. ACM, New York (2007). https://doi.org/10.1145/1250734.1250737, http://doi.acm.org/ 10.1145/1250734.1250737
- Burckhardt, S., Musuvathi, M.: Effective program verification for relaxed memory models. In: CAV. pp. 107–120. Springer-Verlag, Berlin, Heidelberg (2008). https://doi.org/10.1007/978-3-540-70545-1_12, http://dx.doi.org/ 10.1007/978-3-540-70545-1_12
- 23. Burnim, J., Sen, K., Stergiou, C.: Sound and complete monitoring of sequential consistency for relaxed memory models. In: TACAS. pp. 11–25. Springer, Berlin, Heidelberg (2011)
- Chakraborty, S., Vafeiadis, V.: Grounding thin-air reads with event structures. Proc. ACM Program. Lang. 3(POPL) (Jan 2019). https://doi.org/10.1145/3290383, https://doi.org/10.1145/3290383
- 25. Cho, K., Lee, S.H., Raad, A., Kang, J.: Revamping hardware persistency models: View-based and axiomatic persistency models for intel-x86 and armv8. In: Proceedings of the 42nd ACM SIGPLAN International Conference on Programming Language Design and Implementation. p. 16–31. PLDI 2021, Association for Computing Machinery, New York, NY, USA (2021). https://doi.org/10.1145/3453483.3454027, https://doi.org/10.1145/3453483.3454027
- Dan, A.M., Lam, P., Hoefler, T., Vechev, M.: Modeling and analysis of remote memory access programming. SIGPLAN Not. 51(10), 129–144 (oct 2016). https://doi.org/10.1145/3022671.2984033, https://doi.org/10.1145/ 3022671.2984033

- Derevenetc, E.: Robustness against relaxed memory models. Ph.D. thesis, University of Kaiserslautern (2015), http://kluedo.ub.uni-kl.de/frontdoor/index/index/docId/4074
- 28. Derevenetc, E., Meyer, R.: Robustness against Power is PSpace-complete. In: ICALP. pp. 158–170. Springer, Berlin, Heidelberg (2014)
- Fekete, A., Liarokapis, D., O'Neil, E., O'Neil, P., Shasha, D.: Making snapshot isolation serializable. ACM Trans. Database Syst. 30(2), 492–528 (Jun 2005). https://doi.org/10.1145/1071610.1071615, http://doi.acm.org/10.1145/1071610.1071615
- 30. Flur, S., Gray, K.E., Pulte, C., Sarkar, S., Sezgin, A., Maranget, L., Deacon, W., Sewell, P.: Modelling the armv8 architecture, operationally: Concurrency and isa. In: Proceedings of the 43rd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages. p. 608–621. POPL '16, Association for Computing Machinery, New York, NY, USA (2016). https://doi.org/10.1145/2837614. 2837615, https://doi.org/10.1145/2837614.2837615
- 31. Gerstenberger, R., Besta, M., Hoefler, T.: Enabling highly scalable remote memory access programming with mpi-3 one sided. Commun. ACM **61**(10), 106–113 (sep 2018). https://doi.org/10.1145/3264413, https://doi.org/10.1145/3264413
- 32. Gharachorloo, K., Adve, S.V., Gupta, A., Hennessy, J.L., Hill, M.D.: Programming for different memory consistency models. Journal of Parallel and Distributed Computing 15(4), 399 407 (1992). https://doi.org/https://doi.org/10.1016/0743-7315(92)90052-0, http://www.sciencedirect.com/science/article/pii/0743731592900520
- 33. Gorjiara, H., Luo, W., Lee, A., Xu, G.H., Demsky, B.: Checking robustness to weak persistency models. In: Proceedings of the 43rd ACM SIGPLAN International Conference on Programming Language Design and Implementation. p. 490–505. PLDI 2022, Association for Computing Machinery, New York, NY, USA (2022). https://doi.org/10.1145/3519939.3523723, https://doi.org/10.1145/3519939.3523723
- 34. Gotsman, A., Musuvathi, M., Yang, H.: Show no weakness: sequentially consistent specifications of TSO libraries. In: DISC. pp. 31–45. Springer-Verlag, Berlin, Heidelberg (2012). https://doi.org/10.1007/978-3-642-33651-5_3, http://dx.doi.org/10.1007/978-3-642-33651-5_3
- 35. Inc., S.I.: The SPARC architecture manual (version 9). Prentice-Hall, Inc., Upper Saddle River, NJ, USA (1994)
- Kang, J., Hur, C.K., Lahav, O., Vafeiadis, V., Dreyer, D.: A promising semantics for relaxed-memory concurrency. SIGPLAN Not. 52(1), 175–189 (Jan 2017). https://doi.org/10.1145/3093333.3009850, https://doi.org/10.1145/3093333.3009850
- 37. Khyzha, A., Lahav, O.: Taming x86-tso persistency. Proc. ACM Program. Lang. 5(POPL) (Jan 2021). https://doi.org/10.1145/3434328, https://doi.org/10.1145/3434328
- 38. Kokologiannakis, M., Kaysin, I., Raad, A., Vafeiadis, V.: Persevere: Persistency semantics for verification under ext4. Proc. ACM Program. Lang. 5(POPL) (jan 2021). https://doi.org/10.1145/3434324, https://doi.org/10.1145/3434324
- 39. Lahav, O., Giannarakis, N., Vafeiadis, V.: Taming release-acquire consistency. SIGPLAN Not. **51**(1), 649-662 (Jan 2016). https://doi.org/10.1145/2914770. 2837643, https://doi.org/10.1145/2914770.2837643
- 40. Lahav, O., Margalit, R.: Robustness against release/acquire semantics. In: Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design

- and Implementation. p. 126–141. PLDI 2019, Association for Computing Machinery, New York, NY, USA (2019). https://doi.org/10.1145/3314221.3314604, https://doi.org/10.1145/3314221.3314604
- 41. Lahav, O., Vafeiadis, V., Kang, J., Hur, C.K., Dreyer, D.: Repairing sequential consistency in c/c++11. In: Proceedings of the 38th ACM SIGPLAN Conference on Programming Language Design and Implementation. p. 618-632. PLDI 2017, Association for Computing Machinery, New York, NY, USA (2017). https://doi.org/10.1145/3062341.3062352
- 42. Lamport, L.: How to make a multiprocessor computer that correctly executes multiprocess programs. IEEE Trans. Computers 28(9), 690–691 (Sep 1979). https://doi.org/10.1109/TC.1979.1675439, http://dx.doi.org/10.1109/TC.1979.1675439
- 43. Lee, S.H., Cho, M., Podkopaev, A., Chakraborty, S., Hur, C.K., Lahav, O., Vafeiadis, V.: Promising 2.0: Global optimizations in relaxed memory concurrency. In: Proceedings of the 41st ACM SIGPLAN Conference on Programming Language Design and Implementation. p. 362–376. PLDI 2020, Association for Computing Machinery, New York, NY, USA (2020). https://doi.org/10.1145/3385412.3386010, https://doi.org/10.1145/3385412.3386010
- 44. Linden, A., Wolper, P.: A verification-based approach to memory fence insertion in relaxed memory systems. In: SPIN. pp. 144–160. Springer-Verlag, Berlin, Heidelberg (2011), http://dl.acm.org/citation.cfm?id=2032692.2032707
- 45. Linden, A., Wolper, P.: A verification-based approach to memory fence insertion in PSO memory systems. In: TACAS. pp. 339–353. Springer-Verlag, Berlin, Heidelberg (2013). https://doi.org/10.1007/978-3-642-36742-7_24, http://dx.doi.org/10.1007/978-3-642-36742-7_24
- Liu, F., Nedev, N., Prisadnikov, N., Vechev, M., Yahav, E.: Dynamic synthesis for relaxed memory models. In: PLDI. pp. 429-440. ACM, New York (2012). https://doi.org/10.1145/2254064.2254115, http://doi.acm.org/10.1145/2254064.2254115
- 47. Mador-Haim, S., Maranget, L., Sarkar, S., Memarian, K., Alglave, J., Owens, S., Alur, R., Martin, M.M.K., Sewell, P., Williams, D.: An axiomatic memory model for POWER multiprocessors. In: Madhusudan, P., Seshia, S.A. (eds.) Computer Aided Verification 24th International Conference, CAV 2012, Berkeley, CA, USA, July 7-13, 2012 Proceedings. Lecture Notes in Computer Science, vol. 7358, pp. 495–512. Springer (2012). https://doi.org/10.1007/978-3-642-31424-7_36
- 48. Manson, J., Pugh, W., Adve, S.V.: The java memory model. In: Proceedings of the 32nd ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages. p. 378–391. POPL '05, Association for Computing Machinery, New York, NY, USA (2005). https://doi.org/10.1145/1040305.1040336, https://doi.org/10.1145/1040305.1040336
- 49. Margalit, R., Lahav, O.: Verifying observational robustness against a c11-style memory model. Proc. ACM Program. Lang. 5(POPL) (Jan 2021). https://doi.org/10.1145/3434285, https://doi.org/10.1145/3434285
- 50. Moiseenko, E., Podkopaev, A., Lahav, O., Melkonian, O., Vafeiadis, V.: Reconciling Event Structures with Modern Multiprocessors (Artifact). Dagstuhl Artifacts Series 6(2), 4:1–4:3 (2020). https://doi.org/10.4230/DARTS.6.2.4, https://drops.dagstuhl.de/opus/volltexte/2020/13201
- 51. Nagar, K., Jagannathan, S.: Automated detection of serializability violations under weak consistency. In: CONCUR 2018. vol. 118, pp. 41:1–41:18. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik, Dagstuhl,

- Germany (2018). https://doi.org/10.4230/LIPIcs.CONCUR.2018.41, http://drops.dagstuhl.de/opus/volltexte/2018/9579
- 52. Nienhuis, K., Memarian, K., Sewell, P.: An operational semantics for c/c++11 concurrency. In: Proceedings of the 2016 ACM SIGPLAN International Conference on Object-Oriented Programming, Systems, Languages, and Applications. p. 111–128. OOPSLA 2016, Association for Computing Machinery, New York, NY, USA (2016). https://doi.org/10.1145/2983990.2983997, https://doi.org/10.1145/2983990.2983997
- 53. Owens, S.: Reasoning about the implementation of concurrency abstractions on x86-TSO. In: ECOOP. pp. 478–503. Springer-Verlag, Berlin, Heidelberg (2010)
- 54. Owens, S., Sarkar, S., Sewell, P.: A better x86 memory model: x86-TSO. In: TPHOLs. pp. 391-407. Springer, Heidelberg (2009). https://doi.org/10.1007/978-3-642-03359-9_27
- 55. Pichon-Pharabod, J., Sewell, P.: A concurrency semantics for relaxed atomics that permits optimisation and avoids thin-air executions. In: Proceedings of the 43rd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages. p. 622–633. POPL '16, Association for Computing Machinery, New York, NY, USA (2016). https://doi.org/10.1145/2837614.2837616, https://doi.org/10.1145/2837614.2837616
- 56. Podkopaev, A., Lahav, O., Vafeiadis, V.: Promising Compilation to ARMv8 POP. In: Müller, P. (ed.) 31st European Conference on Object-Oriented Programming (ECOOP 2017). Leibniz International Proceedings in Informatics (LIPIcs), vol. 74, pp. 22:1–22:28. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, Dagstuhl, Germany (2017). https://doi.org/10.4230/LIPIcs.ECOOP.2017.22, http://drops.dagstuhl.de/opus/volltexte/2017/7266
- 57. Podkopaev, A., Lahav, O., Vafeiadis, V.: Bridging the gap between programming languages and hardware weak memory models. Proc. ACM Program. Lang. 3(POPL), 69:1-69:31 (Jan 2019). https://doi.org/10.1145/3290382, http://doi.acm.org/10.1145/3290382
- 58. Pulte, C., Flur, S., Deacon, W., French, J., Sarkar, S., Sewell, P.: Simplifying arm concurrency: Multicopy-atomic axiomatic and operational models for armv8. Proc. ACM Program. Lang. 2(POPL), 19:1-19:29 (Dec 2018). https://doi.org/10.1145/3158107, http://doi.acm.org/10.1145/3158107
- 59. Raad, A., Lahav, O., Vafeiadis, V.: On parallel snapshot isolation and release/acquire consistency. In: Ahmed, A. (ed.) Programming Languages and Systems. pp. 940–967. Springer International Publishing, Cham (2018)
- 60. Raad, A., Lahav, O., Vafeiadis, V.: On the semantics of snapshot isolation. In: Enea, C., Piskac, R. (eds.) Verification, Model Checking, and Abstract Interpretation. pp. 1–23. Springer International Publishing, Cham (2019)
- 61. Raad, A., Lahav, O., Vafeiadis, V.: Persistent owicki-gries reasoning: A program logic for reasoning about persistent programs on intel-x86. Proc. ACM Program. Lang. 4(OOPSLA) (nov 2020). https://doi.org/10.1145/3428219, https://doi.org/10.1145/3428219
- 62. Raad, A., Maranget, L., Vafeiadis, V.: Extending intel-x86 consistency and persistency: Formalising the semantics of intel-x86 memory types and non-temporal stores. Proc. ACM Program. Lang. 6(POPL) (jan 2022). https://doi.org/10.1145/3498683, https://doi.org/10.1145/3498683
- Raad, A., Vafeiadis, V.: Persistence semantics for weak memory: Integrating epoch persistency with the tso memory model. Proc. ACM Program. Lang. 2(OOPSLA), 137:1-137:27 (Oct 2018). https://doi.org/10.1145/3276507, http://doi.acm. org/10.1145/3276507

- 64. Raad, A., Wickerson, J., Neiger, G., Vafeiadis, V.: Persistency semantics of the intel-x86 architecture. Proc. ACM Program. Lang. 4(POPL) (Dec 2020). https://doi.org/10.1145/3371079, https://doi.org/10.1145/3371079
- 65. Raad, A., Wickerson, J., Vafeiadis, V.: Weak persistency semantics from the ground up: Formalising the persistency semantics of armv8 and transactional models. Proc. ACM Program. Lang. 3(OOPSLA), 135:1-135:27 (Oct 2019). https://doi.org/10.1145/3360561, http://doi.acm.org/10.1145/3360561
- 66. Sarkar, S., Sewell, P., Alglave, J., Maranget, L., Williams, D.: Understanding power multiprocessors. In: Proceedings of the 32nd ACM SIGPLAN Conference on Programming Language Design and Implementation. p. 175–186. PLDI '11, Association for Computing Machinery, New York, NY, USA (2011). https://doi.org/10.1145/1993498.1993520, https://doi.org/10.1145/1993498.1993520
- Sewell, P., Sarkar, S., Owens, S., Nardelli, F.Z., Myreen, M.O.: X86-TSO: A rigorous and usable programmer's model for x86 multiprocessors. Commun. ACM 53(7), 89–97 (Jul 2010). https://doi.org/10.1145/1785414.1785443, http://doi.acm.org/10.1145/1785414.1785443
- 68. Shpiner, A., Zahavi, E., Dahley, O., Barnea, A., Damsker, R., Yekelis, G., Zus, M., Kuta, E., Baram, D.: Roce rocks without pfc: Detailed evaluation. In: Proceedings of the Workshop on Kernel-Bypass Networks. p. 25–30. KBNets '17, Association for Computing Machinery, New York, NY, USA (2017). https://doi.org/10.1145/3098583.3098588, https://doi.org/10.1145/3098583.3098588
- Vindum, S.F., Birkedal, L.: Spirea: A mechanized concurrent separation logic for weak persistent memory. Proc. ACM Program. Lang. 7(OOPSLA2), 632–657 (2023). https://doi.org/10.1145/3622820, https://doi.org/10.1145/3622820
- Wei, X., Shi, J., Chen, Y., Chen, R., Chen, H.: Fast in-memory transaction processing using rdma and htm. In: Proceedings of the 25th Symposium on Operating Systems Principles. p. 87–104. SOSP '15, Association for Computing Machinery, New York, NY, USA (2015). https://doi.org/10.1145/2815400.2815419, https://doi.org/10.1145/2815400.2815419
- Xiong, S., Cerone, A., Raad, A., Gardner, P.: Data Consistency in Transactional Storage Systems: A Centralised Semantics. In: Hirschfeld, R., Pape, T. (eds.) 34th European Conference on Object-Oriented Programming (ECOOP 2020). Leibniz International Proceedings in Informatics (LIPIcs), vol. 166, pp. 21:1-21:31. Schloss Dagstuhl-Leibniz-Zentrum für Informatik, Dagstuhl, Germany (2020). https://doi.org/10.4230/LIPIcs.ECOOP.2020.21, https://drops.dagstuhl.de/opus/volltexte/2020/13178
- 72. Zhu, Y., Eran, H., Firestone, D., Guo, C., Lipshteyn, M., Liron, Y., Padhye, J., Raindel, S., Yahia, M.H., Zhang, M.: Congestion control for large-scale rdma deployments. In: Proceedings of the 2015 ACM Conference on Special Interest Group on Data Communication. p. 523-536. SIGCOMM '15, Association for Computing Machinery, New York, NY, USA (2015). https://doi.org/10.1145/2785956.2787484, https://doi.org/10.1145/2785956.2787484

A Proofs

In this appendix, we regroup the proofs of the lemmas and theorems presented in this paper.

Theorem 7 (gb implies ob, cf. 1). Given a pre-execution $\langle E, po, pf \rangle$, for all RDMA^{SC}-consistent executions $G = \langle E, po, pf, rf, mo, nfo \rangle$ and all $e_1, e_2 \in E$, if $(e_1, e_2) \in G$.gb, then $(e_1, e_2) \in G$.ob.

Proof. Since ob is transitive, we only need to show the inclusion for each component of gb.

- oppo \subseteq ob by definition.
- $-[nlR]; po|_{imm}; [nrW]; pf \subseteq [Inst]; ippo; pf \subseteq [Inst]; ib \subseteq ob.$
- $-[nrW]; (po \cap sqp); [nlW]; pf \subseteq oppo; ([nlW]; pf) \subseteq ob.$
- $[nrR]; po|_{imm}; [nlW]; pf \subseteq oppo; ([nlW]; pf) \subseteq ob.$
- $-[nrR];(po \cap sqp);[nF];(po \cap sqp) \subseteq [nrR];oppo;[nF];oppo \subseteq ob.$
- [nlW]; $pf \subseteq ob$ by definition.
- Let us assume $e_1 \in nlW$, $e_2 \in nF$, and $e_3 \in nlR$ such that $e_1 \xrightarrow{(po \cap sqp)} e_2 \xrightarrow{(po \cap sqp)} e_3$. We have [nlW]; $(po \cap sqp)$; [nF]; $(po \cap sqp)$; $[nlR] \subseteq ippo$; $ippo \subseteq ib$, so $(e_1, e_3) \in ib$. From Definition 3, we have $(e_1, e_3) \in (nfo \cup nfo^{-1})$. If $(e_3, e_1) \in nfo \subseteq ib$, this create a cycle in ib. Since the execution is consistent, this is not possible and $(e_1, e_3) \in nfo \subseteq ob$. Thus [nlW]; $(po \cap sqp)$; [nF]; $(po \cap sqp)$; $[nlR] \subseteq ob$.
- [nlW]; (po∩sqp); [nF]; (po∩sqp); [nrW] = [nlW]; (po∩sqp); [nF]; (po∩sqp); [nlR]; (po|_{imm} ∩ sqp); [nrW] ⊆ ob; [nlR]; oppo ⊆ ob from above. The nlR event necessarily exists from Definition 2.

Theorem 8 (sc cycle decomposition, *cf.* 2). Given a RDMA^{SC}-consistent execution $G = \langle E, po, pf, rf, mo, nfo \rangle$, if there is a cycle in G.sc (i.e. $\exists e \in E. e \xrightarrow{G.sc}^+$ e), then:

- either there is a G.sc cycle on a single thread, i.e. $\exists e \in E$. $e \xrightarrow{G.sc \cap sthd}^+ e$;
- or there exists $e_1, e_2 \in E^{\mathsf{pub}}$ such that $e_1 \xrightarrow{\mathsf{po} \backslash G.\mathsf{ob}} e_2 \xrightarrow{(G.\mathsf{sc}; [E^{\mathsf{pub}} \backslash t(e_1)])^+; G.\mathsf{sc}} e_1$.

 That is, there is an sc cycle on public events, with two consecutive events on some thread $t(e_1)$ not in ob order and where the rest of the cycle does not go through the events of $t(e_1)$.

Proof. Assuming there is an sc cycle and there is no sc cycle on a single thread (i.e., the first case does not hold), we need to construct an sc cycle corresponding to the second case. By induction on the size of the given sc cycle, we can assume a minimal cycle.

Since ob is acyclic, a cycle contains at least an edge in $(sc \setminus ob)$. Clearly $(sc \setminus po) \subseteq rf \cup rb \cup mo \subseteq ob$, so the edge is in $(po \setminus ob)$. Since sc cycles use at least two threads and thus another event, we can assume a minimum cycle of the form $e_1 \xrightarrow{po \setminus ob} e_2 \xrightarrow{sc} e_2' \xrightarrow{sc}^* e_1' \xrightarrow{sc} e_1$ (note we might have $e_2' = e_1'$).

Is is not possible for this minimal cycle to have three consecutive events on the same thread $(\dots \xrightarrow{sc} e_a \xrightarrow{sc \cap sthd} e_b \xrightarrow{sc \cap sthd} e_c \xrightarrow{sc} \dots)$: if $(e_a, e_b) \in po$ and $(e_b, e_c) \in po$, we can simply drop e_b for a shorter cycle; else we can create a two-event cycle on a single thread. This implies $t(e'_1) \neq t(e_1) = t(e_2) \neq t(e'_2)$.

This also implies that every event e of the cycle is part of an $(sc \setminus sthd)$ edge (either $e \xrightarrow{sc \setminus sthd} e'$ or $e' \xrightarrow{sc \setminus sthd} e$). By definition of rf, rb, and mo, this implies loc(e) = loc(e') and thus public(loc(e)), so every event of the cycle is on a public location.

Finally, we need to assert $e_2' \xrightarrow{\operatorname{sc}|_{E \setminus t(e_1)}}^* e_1'$, i.e. that the rest of the cycle do not reuse the same thread. By contradiction, if it does then it is of the form $e_1 \xrightarrow{po} e_2 \xrightarrow{\operatorname{sc}} e_2' \xrightarrow{\operatorname{sc}}^* e \xrightarrow{\operatorname{sc}}^* e_1' \xrightarrow{\operatorname{sc}} e_1$ with $t(e) = t(e_1)$, and we can create a shorter sc cycle. If $(e, e_1) \in \operatorname{po}$, then $e_2 \xrightarrow{\operatorname{sc}} e_2' \xrightarrow{\operatorname{sc}}^* e \xrightarrow{\operatorname{po}} e_2$ is a shorter cycle. Otherwise, $(e_1, e) \in \operatorname{po}$ and $e_1 \xrightarrow{\operatorname{po}} e \xrightarrow{\operatorname{sc}}^* e_1' \xrightarrow{\operatorname{sc}} e_1$ is a shorter cycle.

Theorem 9 (cf. 3). Given a RDMA^{SC}-consistent execution $G = \langle E, \mathsf{po}, \mathsf{pf}, \mathsf{rf}, \mathsf{mo}, \mathsf{nfo} \rangle$, if $\langle E, \mathsf{po}, \mathsf{pf} \rangle$ is locally data-race free, then there is no sc cycle on a single thread; that is, $(G.\mathsf{sc} \cap \mathsf{sthd})$ is acyclic and thus the first case of Theorem 8 does not arise.

Proof. The LDRF property states that, for every pair $(e_1, e_2) \in (po \cap sloc)$, if at least one of $\{e_1, e_2\}$ is a write event, then $(e_1, e_2) \in gb$. From Theorem 7, this implies $(e_1, e_2) \in ob$. Since ob is acyclic and contains $(rf \cup rb \cup mo)$ relating pairs on the same location with with at least a write, this implies $(rf \cup rb \cup mo) \cap po^{-1} = \emptyset$. Thus $(sc \cap sthd) \subseteq po$ is acyclic.

Lemma 2 (cf. 1). For all $A \subseteq E^{\mathsf{pub}}$, if $\mathsf{e}_1 \xrightarrow{\mathsf{sc}|_A}^* \mathsf{e}_2$ then $n(\mathsf{loc}(\mathsf{e}_1)) \overset{\longleftrightarrow}{\underset{A}{\longleftrightarrow}}^* n(\mathsf{loc}(\mathsf{e}_2))$.

Proof. By induction, it is sufficient to show the result for $e_1 \xrightarrow{\operatorname{sc}|_A} e_2$. If $e_1 \xrightarrow{(\operatorname{sc}\setminus\operatorname{po})|_A} e_2$, then $\operatorname{loc}(e_1) = \operatorname{loc}(e_2)$ (since $\operatorname{rf} \cup \operatorname{rb} \cup \operatorname{mo} \subseteq \operatorname{sloc}$) and the result holds. Else we have $e_1 \xrightarrow{\operatorname{po}|_A} e_2$, and so $n(e_1) = n(e_2)$. If $e_1 \in \{\operatorname{nrR}; \operatorname{nrW}\}$, then e_1 is a witness for $\operatorname{hasQP}(t(e_1), n(\operatorname{loc}(e_1)), A)$, else $e_1 \in E^{\operatorname{pub}}$ so $\operatorname{loc}(e_1)$ is defined, and we simply have $n(e_1) = n(\operatorname{loc}(e_1))$. In both cases, $n(e_1) \longleftrightarrow_A^* n(\operatorname{loc}(e_1))$. Similarly, $n(e_2) \longleftrightarrow_A^* n(\operatorname{loc}(e_2))$ and we have the desired result by transitivity.

Theorem 10 (cf. 4). Given an RDMA^{SC}-consistent execution $\langle E, \mathsf{po}, \mathsf{pf}, \mathsf{rf}, \mathsf{mo}, \mathsf{nfo} \rangle$, if its associated pre-execution $\langle E, \mathsf{po}, \mathsf{pf} \rangle$ is fenced, then there is no sc cycle of the shape $e_1 \xrightarrow{\mathsf{po} \backslash \mathsf{ob}} e_2 \xrightarrow{(\mathsf{sc}; [E^{\mathsf{pub}} \backslash t(e_1)])^+; \mathsf{sc}} e_1$ with $e_1, e_2 \in E^{\mathsf{pub}}$. That is, the second case of Theorem 8 does not arise.

Proof. By contradiction, let us assume such a cycle $e_1 \xrightarrow{po \setminus ob} e_2 \xrightarrow{(sc;[E^{pub} \setminus t(e_1)])^+;sc} e_1$ with $e_1, e_2 \in E^{pub}$. Because $(sc \setminus sthd) \subseteq sloc$, there is e_1' and e_2' such that

 $loc(e_1) = loc(e'_1), loc(e_2) = loc(e'_2), and we can cut the second part of the cycle as <math>e_2 \xrightarrow{sc} e'_2 \xrightarrow{sc|_{E^{pub}\setminus t(e_1)}}^* e'_1 \xrightarrow{sc} e_1$. From Lem. 2, we have $n(loc(e_1)) \xrightarrow[E^{pub}\setminus t(e_1)]{}^* n(loc(e_2))$. Using this result with $e_1 \xrightarrow{po} e_2$, the fenced condition gives us $(e_1, e_2) \in gb$. And from Theorem 7 we have $(e_1, e_2) \in ob$, which contradicts our assumption.

Theorem 11 (Robustness under RDMA^{SC}, cf. 5). Given a pre-execution $\mathcal{G} = \langle E, \mathsf{po}, \mathsf{pf} \rangle$, if \mathcal{G} is locally data-race free (Def. 7) and fenced (Def. 8), then \mathcal{G} is also robust under RDMA^{SC} (Def. 5).

Proof. From Theorem 8, Theorem 9, and Theorem 10.

Theorem 12 (cf. 6). If a pre-execution $\langle E, po, pf \rangle$ is tree-fenced (Def. 9), then it is also fenced (Def. 8).

Proof. Let $\langle E, \mathsf{po}, \mathsf{pf} \rangle$ satisfy the tree-fenced definition (9). Let us assume $\mathsf{e}_1, \mathsf{e}_2 \in E^{\mathsf{pub}}$ such that $\mathsf{e}_1 \xrightarrow{\mathsf{po}} \mathsf{e}_2$ and $n(\mathsf{loc}(\mathsf{e}_1)) \xrightarrow[E^{\mathsf{pub}} \setminus t(\mathsf{e}_1)]{}^* n(\mathsf{loc}(\mathsf{e}_2))$, we then need to prove $(\mathsf{e}_1, \mathsf{e}_2) \in \mathsf{gb}$. Let us name $t = t(\mathsf{e}_1) = t(\mathsf{e}_2)$ and $n = n(\mathsf{e}_1) = n(\mathsf{e}_2)$. From property (1) of tree-fenced, we know that the two events are not NIC local reads nor NIC local writes. If $\mathsf{e}_1 \in E^{\mathsf{cpu}}$, we immediately have $(\mathsf{e}_1, \mathsf{e}_2) \in \mathsf{gb}$ by definition, so we can assume $\mathsf{type}(\mathsf{e}_1) \in \{\mathsf{nrR}, \mathsf{nrW}\}$ and so $n(\mathsf{loc}(\mathsf{e}_1)) = \overline{n}(\mathsf{e}_1)$.

If $\overline{n}(e_1) = n(loc(e_2))$, then e_2 is necessarily a remote event (nrR or nrW) and $(e_1, e_2) \in sqp$. In this case, there is one condition to check to enforce $(e_1, e_2) \in gb$. That is, if $type(e_1) = nrR$, then there is either a remote fence or a poll of e_1 between e_1 and e_2 (cells D7 and D8 of Fig. 5). This condition is exactly property (2) of tree-fenced.

Otherwise, we can assume $\overline{n}(e_1) \neq n(loc(e_2))$, which leads to a contradiction. We can use e_1 as a witness of $hasQP(t,\overline{n}(e_1),E)$. If e_2 is a remote event, we can also use it as a witness of $hasQP(t,n(loc(e_2)),E)$; else e_2 is a CPU event and $n=n(loc(e_2))$. Thus in both cases $\overline{n}(e_1) \underset{E}{\longleftrightarrow} n \underset{E}{\longleftrightarrow} n(loc(e_2))$. Intuitively, our assumption of $\overline{n}(e_1) \underset{E^{pub} \setminus t}{\longleftrightarrow} n(loc(e_2))$ provides a different topological path, which is not allowed. Given property (3a), there is no other node sequence between $\overline{n}(e_1)$ and $n(loc(e_2))$, so our assumption must use a component $\overline{n}(e_1) \underset{E^{pub} \setminus t}{\longleftrightarrow} n$. By property (3b), there is no queue pair from $\overline{n}(e_1)$ towards n. By property (3c), there is no queue pair from other threads of n towards $\overline{n}(e_1)$. So $\overline{n}(e_1) \underset{E^{pub} \setminus t}{\longleftrightarrow} n$ cannot hold.

Corollary 2 (Tree robustness under RDMA^{SC}, cf. 1). If a pre-execution $\mathcal{G} = \langle E, po, pf \rangle$ satisfies LDRF (Def. 7) and is tree-fenced (Def. 9), then it is also robust under RDMA^{SC} (Def. 5).

Proof. Follows directly from Theorem 11 and Theorem 12.

B Extension to RDMA^{TSO}

As discussed in §2, the only realistic, formal semantics for RDMA programs that has been empirically validated is that by Ambal et al. [10], modelled over x86 architecture, i.e. where the *local* concurrency (governing the behaviours of concurrent threads on the same node) is that of TSO [67]. However, local concurrency is decoupled from remote concurrency (describing the behaviours of concurrent threads over different nodes), in that one can understand remote concurrency underpinning RDMA programs separately from its local concurrency. As such, to keep our presentation simple, in §3 we presented the RDMA^{SC} model, where the local concurrency model is that of the strong and intuitive, albeit unrealistic, sequential consistency (SC) [42] model, and we presented our robustness results in §4 and §5 in the context of RDMA^{SC}.

The remote concurrency underpinning RDMA^{SC} is identical to that of RDMA^{TSO}, and the two differ only in their notions of local concurrency (SC versus TSO). Nevertheless, as we show here, we can generalise our results to RDMA^{TSO}. In what follows we present RDMA^{TSO} as formalised by Ambal et al. [10], and reframe our robustness results from §4 and §5 for RDMA^{TSO}. To pinpoint the differences between RDMA^{SC} and RDMA^{TSO}, we visually highlight the changes/extensions, where applicable.

B.1 The RDMA^{TSO} Model

As discussed above, the local concurrency in RDMA^{TSO} is governed by the TSO model. Intuitively, compared to SC, the TSO model allows CPU reads to overtake CPU writes, permitting the well-known weak behaviour known as *store buffering*, illustrated in the example below (comprising two local threads on node 1). Specifically, while the weak outcome (a, b) = (0, 0) is disallowed under SC, it is allowed under TSO as the earlier write in each thread can be reordered after the later read.

```
\begin{array}{l} (a,b)=(0,0) \checkmark \text{ under TSO (and RDMA}^{\text{\tiny TSO}}) \\ (a,b)=(0,0) \checkmark \text{ under SC (and RDMA}^{\text{\tiny SC}}) \end{array}
```

In order to prevent such reordering and prohibit store buffering, programmers can insert a *memory fence*, mfence, between the write and read in question. The same effect can be achieved by separating the write and read with a CAS instruction. The difference between SC and TSO lies in this one weak behaviour (store buffering), and TSO does not admit any other weak behaviours.

Conceptually, one can model this write-read reordering by furnishing each thread with a *store buffer* that is *private* (i.e. only accessible to the thread itself and no other thread). Specifically, executing a write on TSO comprises two steps: 1. when a thread issues a write, the write is only recorded in its store buffer; 2. writes in the buffer are debuffered (in FIFO order) and propagated to the

local memory at a later point. When a thread issues a read from a location x, it first consults its own store buffer. If it contains a write for x, the thread reads the value of the latest such write; otherwise, the thread reads the value of x from its local memory. In other words, one can model the reordering of a write w after a later read r by delaying the debuffering of w until after r has executed. Moreover, executing an mfence or a CAS debuffers all its delayed writes in the store buffer and propagates them to memory (in FIFO order), thus preventing write-read reordering.

Extending the language. In the RDMA^{TSO} language we extend CPU commands (see §3.1) with mfence:

```
\operatorname{cc}^n ::= x := e^n \mid \operatorname{assume}(x = v) \mid \operatorname{assume}(x \neq v) \mid \operatorname{mfence} \mid x := \operatorname{CAS}(y, e_1, e_2) \mid \operatorname{poll}(\overline{n})
```

Updating the declarative semantics. We accordingly extend our labels with F for memory fences. We thus redefine: $E^{\text{cpu}} \triangleq E.1 \text{W} \cup E.1 \text{R} \cup E.\text{CAS} \cup E.\text{P} \cup E.\text{F}$. The translation from commands to events is slightly modified: as an unsuccessful atomic operation $x := \text{CAS}(y, v_1, v_2)$ has the effect of a memory fence, we model it as three events, namely F, 1R(y, v), and 1W(x, v), with $v \neq v_1$.

Note that in RDMA^{TSO} the issue and observation points of CPU writes no longer coincide: once issued, a write can be delayed in the store buffer of the executing thread and is made observable (to other threads and the NIC) once propagated from the buffer to the memory. We thus redefine the set of instantaneous events as $E.Inst \triangleq E \setminus (E.1W \cup E.nlW \cup E.nrW)$, thus excluding 1W events.

We further define the buffered-reads-from relation, rf_{b} , as $\mathsf{rf}_{\mathsf{b}} \triangleq [1\mathtt{W}]$; ($\mathsf{rf} \cap \mathsf{sthd}$); [1R], denoting rf edges between CPU writes and reads by the same thread, i.e. with access to the same store buffer. We define the rf_{b} -complement as $\mathsf{rf}_{\mathsf{b}} \triangleq \mathsf{rf} \setminus \mathsf{rf}_{\mathsf{b}}$, including all other rf edges (i.e. by different threads or involving remote operations). Intuitively, when $w \xrightarrow{\mathsf{rf}_{\mathsf{b}}} r$, then r may read from w before it is observable. Specifically, as CPU writes are delayed in the store buffer and CPU reads first check the buffer, r can read from w either when 1. w is still in the thread's store buffer (i.e. w is not yet observable); or 2. w is in the memory (i.e. w is observable). By contrast, when $w \xrightarrow{\mathsf{rf}_{\mathsf{b}}} r$, then r reads from w only once it is observable (i.e. it has reached the memory). Analogously, we define the buffered-reads-before relation, rb_{b} , as $\mathsf{rb}_{\mathsf{b}} \triangleq [1\mathtt{R}]$; ($\mathsf{rb} \cap \mathsf{sthd}$); [1W].

While ippo (issue-preserved program order) is unchanged from RDMA^{SC} (see Fig. 4), in RDMA^{TSO} we must update oppo (observation-preserved program order) such that 1. a CPU read can overtake a CPU write; 2. a poll does not act as a memory fence; and 3. a memory fence (F) or an atomic operation (CAS) can prevent the reordering of CPU operations. We thus define oppo_{tso} as follows:

$$\mathsf{oppo}_\mathsf{tso} := \mathsf{oppo} \cup \ (E^\mathsf{cpu} \times E) \ \backslash \ (\mathtt{lW} \times (\mathtt{lR} \cup \mathtt{P}))$$

Note that oppo already contains $(E^{\text{cpu}} \times E)$, but we repeat it here as the notion of E^{cpu} has changed.

Lastly, we need to update the definitions of ib and ob to account for the changes discussed above.

Definition 10 (RDMA^{TSO}-consistency, *cf.* **Def. 4).** An execution $\langle E, po, pf, rf, mo, nfo \rangle$ is RDMA^{TSO}-consistent iff ib_{tso} and ob_{tso} are irreflexive, where:

```
\begin{split} & ib_{tso} \triangleq \left( ippo \cup rf \cup pf \cup nfo \ \cup \ rb_b \ \right)^+ \\ & ob_{tso} \triangleq \left( \begin{array}{c} oppo_{tso} \ \cup \ rf_{\overline{b}} \ \cup ([nlW];pf) \cup nfo \cup rb \cup mo \cup ([Inst];ib_{tso}) \end{array} \right)^+ \end{split}
```

The ${\bf rb_b}$ component in ib ensures that a CPU read r by thread t on location x observes all earlier writes on x by t: if r in t reads from w_r ($w_r \stackrel{{\sf rf}}{\to} r$), which is later overwritten by w in t ($w_r \stackrel{{\sf mo}}{\to} w$ and (w, r) \in sthd), i.e. $r \stackrel{{\sf rb_b}}{\to} w$, then w must have been issued after r, as otherwise r should have read from the later w and not w_r . Note that this is not the case for ${\sf rb} \setminus {\sf rb_b}$: r can be issued after a later write w and still not observe it because the effect of w has not yet reached the memory (delayed in the PCIe fabric or a store buffer of another thread).

Note that while ob_{tso} includes $rf_{\overline{b}}$, it no longer includes rf_b since a thread can now read from its store buffer before the write has reached the memory.

B.2 Robustness under RDMA^{TSO}

We first revisit the notion of robustness and redefine it for RDMA^{TSO}.

Definition 11 (RDMA^{TSO}-**robustness, cf. Def. 5).** A pre-execution is robust under RDMA^{TSO} iff all of its RDMA^{TSO}-consistent executions (Def. 10) are also SC-consistent.

As $oppo_{tso}$ differs from oppo, we update gb (Def. 6) under RDMA^{TSO} to use $oppo_{tso}$.

Definition 12 (guaranteed-before under TSO, *cf.* **Def. 6).** Given a preexecution $\mathcal{G} = \langle E, \mathsf{po}, \mathsf{pf} \rangle$, its TSO-guaranteed-before order, $\mathsf{gb}_\mathsf{tso} \subseteq \mathsf{po}$, is defined as $\mathsf{gb}_\mathsf{tso} \triangleq ((\mathsf{gb}_\mathsf{base} \setminus \mathsf{oppo}) \cup \mathsf{oppo}_\mathsf{tso})^+$, or equivalently as:

```
\begin{split} \mathbf{g} \mathbf{b}_{\mathsf{tso}} &\triangleq \big( & & \mathsf{oppo}_{\mathsf{tso}} \big) \\ &\cup [\mathtt{n1R}]; \mathsf{po}|_{\mathit{imm}}; [\mathtt{nrW}]; \mathsf{pf} \\ &\cup [\mathtt{nrW}]; (\mathsf{po} \cap \mathsf{sqp}); [\mathtt{n1W}]; \mathsf{pf} \\ &\cup [\mathtt{nrR}]; \mathsf{po}|_{\mathit{imm}}; [\mathtt{n1W}]; \mathsf{pf} \\ &\cup [\mathtt{nrR}]; (\mathsf{po} \cap \mathsf{sqp}); [\mathtt{nF}]; (\mathsf{po} \cap \mathsf{sqp}) \\ &\cup [\mathtt{n1W}]; \mathsf{pf} \\ &\cup [\mathtt{n1W}]; (\mathsf{po} \cap \mathsf{sqp}); [\mathtt{nF}]; (\mathsf{po} \cap \mathsf{sqp}); [\mathtt{n1R} \cup \mathtt{nrW}] \ \big)^+ \end{split}
```

Note that gb_{tso} no longer includes $1W\times1R$, but does include [1W]; po; [FCAS]; po; [1R] by transitivity.

And gb_{tso} plays the same role for TSO as gb did for SC.

Theorem 13 (gb_{tso} implies ob_{tso}). Given a pre-execution $\langle E, po, pf \rangle$, for all RDMA^{TSO}-consistent executions $G = \langle E, po, pf, rf, mo, nfo \rangle$ and all $e_1, e_2 \in E$, if $(e_1, e_2) \in G$. gb_{tso} , then $(e_1, e_2) \in G$. gb_{tso} (cf. Theorem 1).

We adapt Theorem 2 to refer to RDMA^{TSO}-consistent executions and gb_{tso}.

Theorem 14 (sc cycle decomposition under TSO). Given a RDMA^{TSO}-consistent execution $\langle E, \mathsf{po}, \mathsf{pf}, \mathsf{rf}, \mathsf{mo}, \mathsf{nfo} \rangle$, if there is a cycle in sc (i.e. $\exists \mathsf{e} \in E$. $\mathsf{e} \xrightarrow{\mathsf{sc}^+} \mathsf{e}$), then:

- either there is a sc cycle on a single thread, i.e. $\exists e \in E$. e $\xrightarrow{sc \cap sthd}^+ e$;
- or there exists $e_1, e_2 \in E^{\mathsf{pub}}$ such that $e_1 \xrightarrow{\mathsf{po} \backslash \mathsf{ob}_{\mathsf{tso}}} e_2 \xrightarrow{(\mathsf{sc}; [E^{\mathsf{pub}} \backslash t(e_1)])^+; \mathsf{sc}} e_1$. That is, there is an sc cycle on public events, with two consecutive events on some thread $t(e_1)$ not in $\mathsf{ob}_{\mathsf{tso}}$ order and where the rest of the cycle does not go through the events of $t(e_1)$.

Theorem 3 can be generalised to RDMA^{TSO}-consistent executions at no cost.

Theorem 15. Given a RDMA^{TSO}-consistent execution $\langle E, \mathsf{po}, \mathsf{pf}, \mathsf{rf}, \mathsf{mo}, \mathsf{nfo} \rangle$, if $\langle E, \mathsf{po}, \mathsf{pf} \rangle$ is locally data-race free (Def. 7), then there is no sc cycle on a single thread; that is, (sc \cap sthd) is acyclic and thus the first case of Theorem 14 does not arise (cf. Theorem 3.)

It is important to note that our notion of local data-race freedom (Def. 7) remains unchanged under TSO and does *not* use the new definition of gb_{tso} . This is because LDRF prevents weak behaviours arising from the interaction of the NIC and the CPU for *one thread*, and from the perspective of a single thread, there is no difference between the SC and TSO models. Therefore, a thread running a program such as (x := 1; a := x) is LDRF and is *not* considered racy. That is, even if other threads may observe a reordering, there is no allowed weak behaviour on a single thread.

We then define fenced under TSO by adapting Def. 8 to refer to gb_{tso} , and we used this definition to adapt Theorems 4 and 5.

Definition 13 (fenced under TSO, *cf.* **Def. 8).** A pre-execution $\langle E, \mathsf{po}, \mathsf{pf} \rangle$ is fenced under TSO iff for all $\mathsf{e}_1, \mathsf{e}_2 \in E^{\mathsf{pub}}$, if $\mathsf{e}_1 \xrightarrow{\mathsf{po}} \mathsf{e}_2$ and $n(\mathsf{loc}(\mathsf{e}_1)) \xleftarrow{\mathsf{e}_1} * n(\mathsf{loc}(\mathsf{e}_2))$, then $(\mathsf{e}_1, \mathsf{e}_2) \in \mathsf{gb}_{\mathsf{tso}}$.

Theorem 16. Given an RDMA^{TSO}-consistent execution $\langle E, po, pf, rf, mo, nfo \rangle$, if its associated pre-execution $\langle E, po, pf \rangle$ is fenced under TSO, then there is no

 $\begin{array}{c} \textbf{sc} \ \textit{cycle of the shape} \ \textbf{e}_1 \xrightarrow{\text{po} \backslash \ \text{ob}_{\textbf{tso}}} \textbf{e}_2 \xrightarrow{(\textbf{sc}; [E^{\text{pub}} \backslash t(\textbf{e}_1)])^+; \textbf{sc}} \textbf{e}_1 \ \textit{with} \ \textbf{e}_1, \textbf{e}_2 \in E^{\text{pub}}. \ \textit{That} \\ \textit{is, the second case of Theorem 14 does not arise (cf. Theorem 4).} \end{array}$

Theorem 17 (Robustness under RDMA^{TSO}). Given a pre-execution \mathcal{G} , if \mathcal{G} is LDRF (Def. 7) and fenced under TSO (Def. 13), then \mathcal{G} is also robust under RDMA^{TSO} (Def. 11) – cf. Theorem 5.

B.3 Tree Robustness under RDMA^{TSO}

We next revisit the notion of tree-fenced under TSO. Specifically, we need to extend Def. 9 with an additional condition: if a CPU read on a public location follows a CPU write on a public location, we ask for a memory fence (or an atomic operation) to be introduced between them to prevent them from being reordered, thus pre-empting the well-known store buffering weak behaviour allowed under TSO. Note that the memory fence is not necessary for such CPU operations on private locations.

Definition 14 (tree-fenced under TSO, cf. Def. 9). A pre-execution $\mathcal{G} = \langle E, \mathsf{po}, \mathsf{pf} \rangle$ is tree-fenced under TSO iff it respects the following properties:

- 1. G is tree-fenced (Def. 9); and
- 2. for all $e_1 \in E^{\text{pub}}$.1W, $e_2 \in E^{\text{pub}}$.1R, if $e_1 \xrightarrow{\text{po}} e_2$ then there exists e_3 such that $\mathsf{type}(e_3) \in \{\mathsf{F}, \mathsf{CAS}\}$ and $e_1 \xrightarrow{\mathsf{po}} e_3 \xrightarrow{\mathsf{po}} e_2$.

Lastly, we revisit Theorem 6 and Corollary 1 and re-establish them for TSO.

Theorem 18. If a pre-execution $\langle E, po, pf \rangle$ is tree-fenced under TSO (Def. 14), then it is also fenced under TSO (Def. 13) – cf. Theorem 6.

Corollary 3 (Tree robustness under RDMA^{TSO}). Given a pre-execution \mathcal{G} , if \mathcal{G} is LDRF (Def. 7) and is tree-fenced under TSO (Def. 14), then it is also robust under RDMA^{TSO}. (cf. Corollary 1).

B.4 Proof Updates for RDMA^{TSO}

Most of our proofs carry over to RDMA^{TSO} without significant changes.

Some proofs of this paper (notably Theorem 8) rely on the fact that, given a consistent execution $\langle E, \mathsf{po}, \mathsf{pf}, \mathsf{rf}, \mathsf{mo}, \mathsf{nfo} \rangle$, we necessarily have $(\mathsf{rf} \setminus \mathsf{po}) \subseteq \mathsf{ob}$. This was previously trivial since $\mathsf{rf} \subseteq \mathsf{ob}$ by definition. Now we need to show that $(\mathsf{rf} \setminus \mathsf{po}) \subseteq \mathsf{ob}_\mathsf{tso}$. This comes from the face that, if the execution is RDMA^{TSO}-consistent, $(([1R];\mathsf{po};[1W]) \cup \mathsf{rf})^+ \subseteq (\mathsf{ippo} \cup \mathsf{rf})^+ \subseteq \mathsf{ib}_\mathsf{tso}$ is acyclic. This implies $[1W]; (\mathsf{rf} \cap \mathsf{po}^{-1}); [1R] = \emptyset$ and $\mathsf{rf}_b \triangleq [1W]; (\mathsf{rf} \cap \mathsf{sthd}); [1R] \subseteq \mathsf{po}$ (i.e., a CPU read cannot read from a later CPU write). Thus $(\mathsf{rf} \setminus \mathsf{po}) \subseteq (\mathsf{rf} \setminus \mathsf{rf}_b) = \mathsf{rf}_b \subseteq \mathsf{ob}_\mathsf{tso}$.