## Specifying and Verifying RDMA Synchronisation

Abstract. Remote direct memory access (RDMA) allows a machine to directly read from and write to the memory of remote machine, enabling high-throughput, low-latency data transfer. Ensuring correctness of RDMA programs has only recently become possible with the formalisation of RDMA semantics (describing the behaviour of RDMA networking over a TSO CPU). However, this semantics currently lacks a formalisation of remote synchronisation, meaning that the implementations of common abstractions such as locks cannot be verified. In this paper, we close this gap by presenting RDMARMW, the first semantics for remote 'read-modify-write' (RMW) instructions over TSO. It turns out that remote RMW operations are weak and only ensure atomicity against other remote RMWs. We therefore build a set of composable synchronisation abstractions starting with the RDMA<sub>RMW</sub> library. Underpinned by RDMA<sub>RMW</sub>, we then specify, implement and verify three classes of remote locks that are suitable for different scenarios. Additionally, we develop the notion of a strong RDMA model, RDMA<sub>RMW</sub>, which is akin to sequential consistency in shared memory architectures. Our libraries are built to be compatible with an existing set of high-performance libraries called LOCO, which ensures compositionality and verifiability.

## 1 Introduction

Remote Direct Memory Access (RDMA), as implemented by RoCE and Infiniband, is a high-performance networking technology that enables low-latency wire-speed data transmission. Specifically, an RDMA device can directly read and write from the memory of a remote (network) node (machine), bypassing the remote CPU and operating system. RDMA technology has been used in high-performance computing applications (including supercomputers) since the early 2000s, and is being branched out to support a much wider range of applications, ranging from production-grade data centres [25, 32, 34] to distributed AI training [17]. Thus, there is currently a push towards developing programmer-friendly libraries to improve the reliability and robustness of such applications.

To enable rigorous development and verification, there is ongoing work aimed at formalising the semantics of RDMA architectures, primarily the RDMA memory model. Dan et al [13] proposed an early model, called coreRMA, which was used to formalise the behaviours of remote read/write operations, assuming a sequentially consistent CPU. Ambal et al [3] have presented a more realistic RDMA<sup>TSO</sup> specification, which assumes a total-store-order (TSO) CPU (e.g. as implemented by Intel processors) that (unlike coreRMA) has been validated against real RoCE and Infiniband hardware. RDMA<sup>TSO</sup> precisely describes the interaction between the CPU and NIC (Network Interface Card) and the reorderings that they allow. Their formalisation comprises both declarative and operational models (which are proved equivalent). However, RDMA<sup>TSO</sup> only covers a subset

of RDMA instructions. In particular it only covers local (i.e. CPU-level) 'read-modify-write' (RMW) synchronisation, relegating remote (i.e. RDMA) RMWs to future work. This means that RDMA<sup>TSO</sup> cannot be used to specify and verify locks and other related high-level mechanisms that require synchronisation at the network level.

In this work, we address this gap and extend the existing efforts with a notion of remote (RDMA) synchronisation. Specifically, we develop the RDMA<sup>TSO</sup> model by extending RDMA<sup>TSO</sup> to account for remote RMWs. To ensure the fidelity of our extension, we developed RDMA<sup>TSO</sup><sub>RMW</sub> by careful inspection of the Infiniband technical manual [21] and in close consultation with engineers at NVIDIA, the largest manufacturer of RDMA products worldwide (after acquiring Mellanox in 2019). We then build a series of synchronisation libraries and prove them correct (as we discuss below). An overview of our development is given in Fig. 1.

Remote RMW instructions are surprisingly weak in that they only guarantee a weak form of isolation: remote RMWs are atomic only with respect to other remote RMWs and not CPU accesses or remote read and write accesses operations (cf. weak transactional isolation [8,10,16,28,29]). We provide a set of litmus tests that exemplify these behaviours in two- and three-node configurations. A second challenge is that (like RDMA<sup>TSO</sup>) RDMA<sup>TSO</sup><sub>RMW</sub> is not compositional: the semantics of a certain remote operation, Poll, directly depends on the exact number of remote operations in the program up to that point! As such, one cannot specify the behaviour of Poll modularly (in isolation).

To address both issues, we build on the *Library of Composable Objects* (LOCO) framework [4, 20], which is a modular set of objects for constructing RDMA libraries. We start at the lowest level of LOCO, called RDMA<sup>WAIT</sup>, which is a compositional analogue of RDMA<sup>TSO</sup> (i.e. also does not support remote RMWs). As shown in Fig. 1, the RDMA<sup>WAIT</sup> library itself is implemented using RDMA<sup>TSO</sup>.

Importantly, RDMAWAIT abstracts RDMATSO

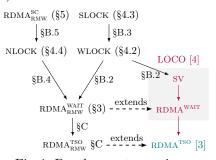


Fig. 1: Development overview

by replacing its non-modular operation (Poll) with a modular analogue (Wait, see §2.1). As such, unlike RDMA<sup>TSO</sup>, RDMA<sup>WAIT</sup> is *modular* and can be *composed* with other LOCO libraries (thanks to its Wait operation). Accordingly, we develop RDMA<sup>WAIT</sup><sub>RMW</sub> by extending RDMA<sup>WAIT</sup> with RMW operations. Specifically, in RDMA<sup>WAIT</sup><sub>RMW</sub> we specify two remote RMWs: RCAS (remote compare-and-swap) and RFAA (remote fetch-and-add). In doing so, we also ensure that our extensions are compatible with RDMA<sup>WAIT</sup> and the modular design of LOCO, thus guaranteeing that RDMA<sup>WAIT</sup><sub>RMW</sub> is also modular and can be composed with other LOCO libraries.

We next use RDMA<sup>WAIT</sup><sub>RMW</sub> to develop several RDMA libraries (Fig. 1). First, we combine RDMA<sup>WAIT</sup><sub>RMW</sub> with the *shared variable* (SV) library (that provides a mechanism for broadcasting to many nodes) of LOCO to develop three lock libraries with varying synchronisation guarantees (§4), each offering a different trade-

off between intuitive behaviours and efficiency. Second, we develop an RDMA library with strong sequential consistency (SC) [23] semantics (§5).

Our first lock library is a weak lock, WLOCK, that provides mutual exclusion across multiple threads over the network, but does not provide any ordering guarantees on RDMA instructions enclosed within critical sections. Nevertheless, it is possible to recover such strong ordering guarantees on RDMA operations within a WLOCK critical section by inserting a global fence immediately before the lock is released. To capture this, we thus develop a strong lock, SLOCK, that guarantees the desired strong guarantees by executing a global fence before releasing the lock. The most novel aspect of our library is the notion of a node lock, NLOCK, that takes a node n as a parameter, and only guarantees synchronisation on RDMA operations specific to n, while operations within a critical section acting on a different node  $n' \neq n$  are left unsynchronised.

Interestingly, we show that it is possible to build a novel, strong model for RDMA using NLOCK. Specifically, we develop the RDMA<sup>SC</sup><sub>RMW</sub> library, which, unlike RDMA<sup>WAIT</sup>, provides support for strong isolation of remote RMW instructions, with strong synchronisation akin to SC.<sup>1</sup>

For each library L in our development (Fig. 1), we 1) formally specify L; 2) develop a reference implementation of L using lower-level libraries; and 3) prove our implementation is correct against its specification. For (1) and (3), we use MOWGLI [4], a declarative framework previously used to verify a subset of LOCO (those without RMWs). MOWGLI is a compositional framework for specification and verification of very weak libraries where program order is not preserved (e.g. RDMA programs). However, previous definitions [4] are not sufficient to specify remote RMWs out of the box, and we extend them with the features needed (§3).

Contributions. Our core contributions are as follows. (1) We develop the first formal semantics of remote RMWs through the RDMA<sub>RMW</sub><sup>TSO</sup> and RDMA<sub>RMW</sub><sup>WAIT</sup> models by carefully inspecting the (informal) technical specification [21]. Our models have further been validated by NVIDIA engineers. (2) We generalise MOWGLI to add the intricate features needed for specifying and verifying very weak libraries. We then use LOCO and our extension of MOWGLI to develop several programmer-friendly and composable RDMA libraries. Specifically, (3) we specify, implement and verify three lock libraries offering varying degrees of synchronisation guarantees and efficiency; and (4) we develop a novel, strong RDMA model, RDMA<sub>RMW</sub><sup>SC</sup>, ensuring strong isolation of RDMA instructions with strong synchronisation guarantees of SC.

Outline. The remainder of this article is organised as follows. In §2 we discuss the necessary background and present an intuitive overview of our contributions. In §3 we describe how we extend the MOWGLI framework and present our RDMA<sub>RMW</sub> model. In §4 we present our three lock libraries (including their specification, implementation, and verification), which we build on top of RDMA<sub>RMW</sub>. In §5 we specify, implement, and verify our RDMA<sub>RMW</sub> library (simulating SC in RDMA programs). Finally, we discuss related work in §6.

<sup>&</sup>lt;sup>1</sup> In related work, Ambal et al. [5] write RDMA<sup>SC</sup> for an RDMA model where the underlying CPU is SC (instead of TSO). This is unrelated to RDMA<sup>SC</sup><sub>RMW</sub>.

## 2 Background and Overview

We present an intuitive account of our contributions via a series of litmus tests. We begin with a summary of necessary background (§2.1 and §2.2). We discuss the behaviour of remote RMW ('read-modify-write') synchronisation, culminating in our formal RDMA<sup>WAIT</sup><sub>RMW</sub> model (§2.3). We then describe our RDMA libraries (§2.4), including locks and a library for sequential consistency we build from it.

**Terminology and Litmus Test Notation.** Throughout this article, we present small examples (litmus tests) to highlight particular behaviours. A single vertical bar (e.g. in Fig. 8a) separates threads on the *same* (network) node, while a double vertical bar (e.g. in Fig. 2a) separates distinct nodes. For each annotated outcome,  $\checkmark$  denotes that the outcome is allowed by the semantics, while  $\nearrow$  states that the outcome is disallowed. To distinguish local and remote (memory) locations, we write  $x^n$  for a location on a remote node n, and write x for a location on the current local node. We number nodes from left to right, starting at 1. The statement on the top line of each column denotes where locations reside as well as their initial values; e.g. x=0 and z=0 on top of Fig. 2a denote that x and z respectively reside on nodes 1 and 2 with initial value 0. When a thread on local node n issues a remote operation to be executed on remote node n', we denote this by stating that the operation is by n towards n'.

## 2.1 Background: RDMA<sup>TSO</sup>, RDMA<sup>WAIT</sup>, and LOCO

The RDMA<sup>TSO</sup> Model. Ambal et al. [3] developed RDMA<sup>TSO</sup>, the first formal model of RDMA programs where the underlying CPUs are assumed to follow the x86-TSO memory model [26]. RDMA<sup>TSO</sup> formalises the semantics of RDMA Writes (referred to as puts), RDMA Reads (referred to as gets) and polling instructions, executed by the network interface card (NIC). A put operation towards n, written  $x^n := y$ , reads from local location y (referred to as a NIC) local read) and writes to remote location x on node n (a NIC remote write). Similarly, a get operation towards n, written  $y := x^n$  reads from remote location x (a NIC remote read) and writes to local location y (a NIC local write). The RDMA<sup>TSO</sup> semantics is unintuitive as remote operations are executed by NIC independently from later CPU operations, as if run in parallel to them. For instance, the program  $z^2 := x$ ; x := 1 (comprising a put towards node 2, followed by a standard CPU store) can result in z containing value 1 as follows: 1) CPU offloads the put instruction to the NIC; 2) CPU executes x := 1; 3) NIC executes the put, fetching the new value 1 of x and updating the remote location z in node 2 to this new value. To prevent this weak behaviour, a programmer can poll the remote instruction (towards node 2) by executing Poll(2), as shown in Fig. 2a: this blocks the CPU until the NIC confirms that the put has been executed, thereby preventing the above scenario.

The polling system on RDMA hardware (and thus RDMA<sup>TSO</sup>) is highly brittle in that it synchronises with the *earliest* (in program order) unpolled remote operation. For instance, in Fig. 2b the single poll only acknowledges the first

	$ \begin{array}{c c} x = 0 & z = 0 \\ \hline z^2 := x \\ z^2 := x \\ \text{Poll}(2) \\ x := 1 \\ \end{array} $	$ \begin{array}{ c c c c } \hline x=0 & z=0 \\ \hline z^2:=x & \\ z^2:=x \\ \text{Poll}(2) \\ \text{Poll}(2) \\ x:=1 \\ \hline \end{array} $		$ \begin{array}{ c c c }\hline x = 0 & z = 0\\\hline z^2 := ^e x\\z^2 := ^d x\\ \text{Wait}(d)\\x := 1\\\hline \end{array} $
(a) $z = 0$ $\checkmark$ $z = 1$ $\checkmark$	(b) $z = 0 \checkmark$ $z = 1 \checkmark$	(c) $z = 0$ $\checkmark$ $z = 1$ $\checkmark$	(a) $z = 0$ $\checkmark$ $z = 1$ $\checkmark$	(b) $z = 0 \checkmark$ $z = 1 \checkmark$

Fig. 2: Polling on RDMA<sup>TSO</sup>

Fig. 3: Waiting on RDMAWAIT

put, and the second put can be arbitrarily delayed, once again enabling the outcome z=1. Preventing unintended weak behaviours therefore often relies on counting remote operations and polling them accordingly; e.g. in this case we must use two polls to prevent the weak outcome, as in Fig. 2c.

The RDMA WAIT Model. The non-local semantics of polls does not lend itself to compositional programming and verification. That is, the polling semantics depends on the exact number of earlier remote operations towards the same node. To address this, recent work developed LOCO [4] as an RDMA library for composable objects with a more abstract completion system that ensures modularity and compositionality through a waiting instruction that is analogous to polling but is compositional. Specifically, in LOCO each remote operation is associated with a work identifier,  $d \in \text{Wid}$ , and the wait operation Wait(d) ensures the completion of all previous operations with this identifier (multiple remote operations may have the same identifier). This is illustrated in Figs. 3a and 3b (obtained from Figs. 2a and 2b by replacing polls with waits), where  $z^2 := d x$  denotes a put (as before) with work id d. Note that unlike in Fig. 2b, adding an earlier put in Fig. 3b (with different work id e) does not alter the behaviour of Wait(d) and the weak outcome z=1 remains prohibited.

From a reordering perspective, RDMA<sup>WAIT</sup> is still quite permissive. For example, because a remote NIC sends an acknowledgement for a put as soon as it is received (but before the put takes effect in memory), RDMA<sup>WAIT</sup> permits the store-buffering behaviour in Fig. 4. Therefore, using RDMA<sup>WAIT</sup>, LOCO additionally implements a global-fence operation towards a node n, written GFence( $\{n\}$ ), that blocks until all previous remote operations towards n are fully completed (see §4.1). Replacing Wait(d) and Wait(e) in Fig. 4 respectively with fences GFence( $\{2\}$ ) and GFence( $\{1\}$ ) would prevent the store-buffering behaviour.

#### 2.2 Background: MOWGLI

To support compositional specification and verification, Ambal et al. have developed the MOWGLI framework [4]. They have specified the RDMA<sup>WAIT</sup> formal model (obtained from RDMA<sup>TSO</sup> by replacing the poll instruction with Wait) in MOWGLI and subsequently used it as a foundation for developing and verifying a suite of RDMA libraries. The principal one is a *shared variable* (sv) library (see §4.1), where each node possesses a local copy of each variable x. The methods include store ( $x :=_{SV} v$ ) and load ( $a :=_{SV} x$ ) operations to access the local

y = 0	x = 0
$x^2 :=^d 1$ $\text{Wait}(d)$ $a := y$	$ \begin{vmatrix} y^1 :=^e 1 \\ \mathtt{Wait}(e) \\ b := x \end{vmatrix} $
(a,b) =	(0,0) ✓

SVar x	=0
	z = 0
$z^2 := 1$ $x :=_{ ext{SV}} 1$ $ ext{Bcast}_{ ext{SV}}(x)$	$a :=_{SV} x$ $b := z$
(a) $(a, b) =$	(1,0) X

SVar x = 0					
	y, z = 0, 0				
$z^2 := 1 \ x :=_{ ext{SV}} 1 \  ext{Bcast}_{ ext{SV}}(x)$	a := y $b := z$	$c :=_{sv} x$ $y^2 := 1$			

(b) (a, b, c) = (1, 0, 1)

Fig. 4: Store buffering

Fig. 5: Shared variable examples

copy, as well as a broadcast  $(\mathtt{Bcast}_{\mathtt{SV}}(x))$  operation to forward the local value to other nodes.

**Specification.** MOWGLI [4] is a declarative framework for modularly specifying and verifying libraries in the context of (very) weak concurrency models. Unlike other declarative frameworks in the literature [27,31], MOWGLI can handle the behaviours allowed by RDMA programs. The key novelty in MOWGLI enabling this is the use of a fixed set of stamps,  $Stamp = \{a_1, ...\}$ , and the stamp-order relation,  $sto \subseteq Stamp \times Stamp$ , defined as a subset of the program order that is preserved. This then allows one to define weak libraries where the program order is not fully preserved, as is the case in RDMA.

We present the stamps and their ordering in Fig. 9 (assuming that the underlying CPUs follow the TSO model). Intuitively, each stamp denotes a behaviour category, such as a CPU write (aCW), a CPU read (aCR), a NIC remote read (aNRR<sub>n</sub>) or write (aNRW<sub>n</sub>) towards n, or a NIC local read (aNLR<sub>n</sub>) or write (aNLW<sub>n</sub>) towards n. Compared to [4], we also introduce a new stamp aNAR<sub>n</sub> to represent the ordering guarantees of remote RMWs (see §2.3).

This stamp mechanism addresses two problems. The first is the reordering of methods of different libraries. As libraries are defined independently, the exact interaction between pairs of methods of different libraries cannot be explicit. Instead, libraries can associate their method calls with generic behaviour categories (stamps), so that their interactions can be implicitly deduced. For instance, in Fig. 5a,  $z^2 := 1$  and b := z are part of RDMA<sub>RMW</sub>, while  $\mathtt{Bcast}_{\mathrm{SV}}(x)$  and  $a :=_{\mathrm{SV}} x$ are part of the SV library. To determine if the outcome (a,b)=(1,0) is allowed, we need to check if  $z^2 := 1$  and  $B_{cast_{SV}}(x)$  can be reordered on node 1, and if  $a :=_{SV} x$  and b := z can be reordered on node 2. The semantics of the two libraries (§3 and §A) ensure that  $z^2 := 1$  and  $B_{cast_{sv}}(x)$  behave as remote writes towards node 2 (stamp aNRW<sub>2</sub>) and that  $a :=_{SV} x$  and b := z behave as CPU reads (stamp aCR). This enforces their respective program orders as  $\langle z^2 := 1, aNRW_2 \rangle$  $\xrightarrow{\text{ppo}} \langle \text{Bcast}_{\text{SV}}(x), \text{aNRW}_2 \rangle \text{ and } \langle a :=_{\text{SV}} x, \text{aCR} \rangle \xrightarrow{\text{ppo}} \langle b := z, \text{aCR} \rangle, \text{ where ppo is the } \langle b := z, \text{aCR} \rangle$ preserved program order, i.e. they cannot be reordered. Moreover, if a=1, then we have the happens-before (hb) relation  $\langle Bcast_{SV}(x), aNRW_2 \rangle \xrightarrow{hb} \langle a :=_{SV} x, aCR \rangle$ , and as  $ppo \subseteq hb$ , by transitivity we have  $\langle z^2 := 1, aNRW_2 \rangle \xrightarrow{hb} \langle b := z, aCR \rangle$ , i.e. the weak outcome (a, b) = (1, 0) is prohibited.

The second problem stamps address is the *partial* execution of methods. A method call may have multiple visible effects, and observing one does not necessarily imply that others are also observed. In Fig. 5b the shared variable

x is read by the third node, which then sends a message to node 2 (through  $y^2:=1$ ). As such, when (a,c)=(1,1), we have a  $\frac{\mathsf{hb}}{\mathsf{chain}}$  chain from  $\mathsf{Bcast}_{\mathsf{SV}}(x)$  to b:=z and may naturally expect c=1. However, this is not case. Specifically, as per the semantics of  $\mathsf{SV}$ ,  $\mathsf{Bcast}_{\mathsf{SV}}(x)$  is associated with (at least) two stamps,  $\mathsf{aNRW}_2$  (remote write towards node 2) and  $\mathsf{aNRW}_3$  (remote write towards node 3), where the latter is observed but is not ordered with the earlier  $z^2:=1$  operation (as they are toward different nodes). That is, we have the  $\mathsf{hb}$  orders  $\langle z^2:=1,\mathsf{aNRW}_2\rangle$   $\xrightarrow{\mathsf{ppo}\subseteq\mathsf{hb}}$   $\langle \mathsf{Bcast}_{\mathsf{SV}}(x),\mathsf{aNRW}_2\rangle$  (as in example Fig. 5a) and  $\langle \mathsf{Bcast}_{\mathsf{SV}}(x),\mathsf{aNRW}_3\rangle \xrightarrow{\mathsf{hb}} \langle b:=z,\mathsf{aCR}\rangle$ , and when put together they do not imply  $\langle z^2:=1,\mathsf{aNRW}_2\rangle \xrightarrow{\mathsf{hb}} \langle b:=z,\mathsf{aCR}\rangle$ , allowing the weak outcome b=0. In other words,  $z^2:=1$  and  $\mathsf{Bcast}_{\mathsf{SV}}(x)$  can be partially reordered: although their respective updates (on z and x) towards node 2 stay ordered, the update on x towards other nodes (i.e. node 3) may take place before  $z^2:=1$  is executed. Associating a method call with multiple stamps allows us to express such nuances.

Implementation and Soundness. Within the MOWGLI framework, Ambal et al. [4] also formalise the notion of a library implementation and what it means for an implementation I to be sound against its specification, i.e. that the behaviours of the implementation are contained in those of its specification. To enable proving implementation soundness compositionally, they establish a local soundness theorem. Specifically, to show that an implementation I of library Lis correct, one must show that for all client programs P with calls to L (where P may in general contain calls to libraries other than L), replacing the calls to L with their corresponding (inlined) implementation yields the same outcome. Intuitively, as the only calls being replaced (inlined) are those of L, the calls to libraries other than L should not affect the outcome. That is, it should be sufficient to show that the implementation is locally sound by considering client programs that only constitute calls to L. Ambal et al. then prove that local soundness implies soundness: if I is a locally sound implementation of L(i.e. for all client programs that only comprise calls to L), then I is a sound implementation of L (i.e. for all client programs).

As we discuss below, we use MOWGLI to specify several RDMA libraries and verify their implementations, as shown in Fig. 1.

#### 2.3 Remote Read-Modify-Write Operations

**CPU Read-Modify-Writes.** Read-modify-writes (RMW) are a category of synchronisation operations that simultaneously read the value v of a location and update (modify-write) it in place. Examples of common RMWs include the compare-and-swap,  $CAS(x, v_1, v_2)$ , instruction (it reads the current value v of x and updates it to  $v_2$  if  $v = v_1$  and otherwise leaves it unchanged); and the fetch-and-add, FAA(x, v), instruction (it increments the value of x by v unconditionally). Both operations return the old value of x. These operations are useful for ensuring inter-thread synchronisation and are often used to implement strong synchronisation mechanisms such as locks (mutexes).

	x = 0			x = 0			x = 0
$\begin{array}{c} \text{RCAS} \\ (a, x^2, 0, 2) \end{array}$	x := 1	$\begin{array}{c} {\tt RCAS} \\ (a,x^3,0,2) \end{array}$	$x^3 := 1$		$\begin{array}{c} \mathtt{RCAS} \\ (a, x^3, 0, 2) \end{array}$	$\begin{bmatrix} \mathtt{RFAA} \\ (b, x^3, 1) \end{bmatrix}$	
(a) $x =$	2 🗸	(b)	x = 2		(c) x	$=2$ $\chi$	

Fig. 6: Examples showcasing the limited atomicity of remote RMW operations

CPU RMWs behave *atomically*: their 'read' and 'modify-write' phases cannot be interleaved by concurrent instructions. As such, RMWs are commonly referred to as 'atomic operations'. This is illustrated

x = 0	
$a:=\mathtt{CAS}(x,0,2)$	x := 1
x=2 X	

in the example across where the outcome x=2 is disallowed. If the right thread executes first, then x is updated to 1 and subsequently the CAS fails. If the left thread executes first, then the right thread overwrites x to 1.

**Remote RMWs.** The RDMA hardware specification [21] optionally supports two remote RMW instructions, referred to as 'atomics<sup>2</sup>': RCAS $(a, x, v_1, v_2)$ , analogous to  $a := CAS(x, v_1, v_2)$  on CPUs, and RFAA(a, x, v), analogous to a := FAA(x, v) on CPUs, where x is a remote location in both cases.

Unlike CPU RMWs, remote RMWs do *not* always behave atomically: their 'read' and 'modify-write' phases may be interleaved by other CPU or (remote) put/get instructions. This is illustrated in the examples of Figs. 6a and 6b, executing a remote CAS in parallel with a CPU store (Fig. 6a) and a put (Fig. 6b), where the remote CAS can first read 0 from x, be interleaved with the concurrent CPU store/put writing 1 to x, and then update x to 2.

This weakness is due to an inherent hardware limitation. Atomicity is possible on CPUs because a CPU core can: 1) request exclusive access to a cache line; 2) read the cache line; 3) write to the cache line; 4) release the cache line. During periods of exclusive access, other components (e.g. other CPU cores or the NIC) cannot access the cache line in-between the 'read' and 'modify-write'. This, however, is not feasible over RDMA since NICs cannot lock a cache line; they can only submit read and write operations to their PCIe root complex. As such, it is not possible to block accesses by other components (e.g. the CPU) interleaving between the NIC's 'read' and the 'modify-write'.

Nevertheless, remote RMWs do behave atomically with respect to other remote RMWs. For instance, as shown in Fig. 6c, a remote FAA cannot interleave between the 'read' and 'modify-write' phases of a remote CAS.

In practice, one can ensure atomicity of accesses to a location x by ensuring x is *only* ever accessed through remote RMWs. As such, it is common for RDMA programs to access *local* locations (i.e. those residing on their node) through remote RMWs (via loop-back). To provide atomicity between remote RMWs and other operations, we require software solutions, as supported by RDMA $_{\rm RMW}^{\rm SC}$ .

<sup>&</sup>lt;sup>2</sup> Although RMWs are commonly referred to as 'atomics' in the RDMA specification, they do *not* always behave atomically.

				y = 0	x = 0
	$ \begin{array}{ c c } \hline \text{RCAS}(a, x^2, 8, 9) \\ y^2 := 1 \\ \hline \end{array} $	$\begin{vmatrix} x, y = 0, 0 \\ b := y \\ x := 1 \end{vmatrix}$	Po	` '	$ \begin{vmatrix} \mathtt{RFAA}(_{-}, y^1 \\ \mathtt{Poll}(1) \\ b := x \end{vmatrix} $
(a) $(a,b) = (1,1)$	(b) $(a, b) =$	$(1,1)$ $\times$		(c) (a, b)	=(0,0)

Fig. 7: Examples of remote RMW behaviours and how they compare from Puts.

 $RFAA(_{-},y^1,1)$ Poll(1)

Extending RDMATSO with RMWs. The RDMATSO and RDMAWAIT models do not include the semantics of remote RMWs; we close this gap in this work. Specifically, starting from  $RDMA^{TSO}$  [3], we formulate  $RDMA^{TSO}_{RMW}$  both declaratively and operationally and prove the two characterisations are equivalent (see §D).

Our main reference for modelling the semantics of remote RMWs is the Infiniband technical specification [21]. However, as the specification is often ambiguous, we developed our model in close collaboration with NVIDIA experts specialising in RDMA hardware who confirmed the expected behaviours of RMWs and that our model captures them faithfully.

Compared to  $RDMA^{TSO}_{RMW}$ , our  $RDMA^{TSO}_{RMW}$  declarative model brings two important changes. The first is a new relation, the 'remote-atomic-order' rao, capturing the mutual exclusion of remote RMWs. We require a total order  $rao_n$  on remote RMWs towards each node n, such that  $rao_n \subseteq hb$  (i.e. it induces synchronisation). The second is a new stamp,  $aNAR_n$  (the 'NIC atomic read'), encoding the new ordering guarantees of the read phase of a remote RMW (see Fig. 9). Recall that a Get performs a NIC remote read (stamp  $aNRR_n$ ) followed by a NIC local write  $(aNLW_n)$ , while a Put performs a NIC local read  $(aNLR_n)$  followed by a NIC remote write (aNRW<sub>n</sub>). Analogously, a remote RMW, e.g. RFAA(x, y, v), performs (up to) three NIC accesses: 1) it remotely reads y (aNAR<sub>n</sub>); 2) remotely updates y (aNRW<sub>n</sub>); and 3) locally writes (the return value) to x (aNLW<sub>n</sub>).

Note that the stamp  $aNAR_n$  is required because a remote read stamp  $(aNRR_n)$ is insufficient for modelling the stronger ordering guarantees of the 'read' phase of an RMW. We show an example of this in Figs. 7a and 7b. The (remote) read phase of a Get (aNRR<sub>n</sub>) may be delayed (reordered) after a later remote write (aNRW<sub>n</sub> of a Put or RMW). As such, the weak 'load-buffering' behaviour in Fig. 7a is allowed. By contrast, the read phase of a remote RMW (aNAR<sub>n</sub>) cannot be delayed, and thus the analogous behaviour is prohibited in Fig. 7b.

Finally, the example in Fig. 7c shows that the remote write ('modify') phase  $(aNRW_n)$  of an RMW behaves similarly to that of a Put. In particular, a poll does not enforce the full completion of the remote write and thus the weak 'store-buffering' behaviour presented is allowed, similarly to Fig. 4.

Supporting Modularity with RDMA $_{\rm RMW}^{\rm WAIT}$ . As RDMA $_{\rm RMW}^{\rm TSO}$  is not modular, we develop RDMA $_{\rm RMW}^{\rm WAIT}$  by adapting RDMA $_{\rm RMW}^{\rm WAIT}$  [4, 20]. We then implement RDMA $_{\rm RMW}^{\rm WAIT}$  using RDMA $_{\rm RMW}^{\rm TSO}$  and prove that it is correct (§C) against its specification.

#### 2.4 Modular RDMA Synchronisation Libraries

We now implement RDMA libraries modularly and specify and verify them using MOWGLI. A key use case of our remote RMWs is for implementing network-wide locks that ensure mutual exclusion of critical sections. A lock library provides two main operations, Acq(l) and Rel(l), for acquiring and releasing a lock l, respectively. When specifying such a network lock, there are several choices for defining its semantics as there are trade-offs between the guarantees (strength) of a lock and the efficiency of its implementation.

Fig. 8 presents several variants of an example where two threads use a lock l to access locations x and y in a critical section (the first thread writing to x and y and the last thread reading from x and y). As the locks are expected to ensure atomicity of the critical sections (enclosed within the lock acquisition and release blocks), the expected outcomes are either a=b=0 or a=b=1, i.e. not  $a\neq b$ . However, ensuring this strong guarantee for locks is not straightforward over RDMA. Specifically, while ensuring mutual exclusion is necessary for prohibiting the weak  $a\neq b$  behaviour, it is not sufficient. We must additionally ensure that the operations enclosed in a critical section are *completed* before the end of the critical section (and hence are not *reordered* past the lock release). However, as we demonstrated above, meeting these latter constraints are not always straightforward due to the weak ordering guarantees on remote operations.

Weak Lock Library. The weakest network lock that we consider ensures mutual exclusion only, but does not prohibit the enclosed operations from being reordered. As shown in Fig. 8a, when the operations enclosed in a critical section are CPU loads and stores, the weak outcome  $a \neq b$  is prohibited. By contrast, as shown in Fig. 8b, when the enclosed operations are over RDMA (two Puts in Fig. 8b), then a weak lock is insufficient and we may observe  $a \neq b$ . This is because the remote operations may not complete before the critical section ends.

Thus, we require an operation akin to a global fence (see §2.1) to ensure that these remote operations are completed. Note that as shown in Fig. 8c, the global fence in isolation (without the protection provided by a weak lock) is also insufficient for prohibiting the weak behaviour as it only provides intra-thread synchronisation (and does not ensure mutual exclusion). However, as shown in Fig. 8d, if we combine a weak lock with a global fence, we can attain the desired strong guarantees and prohibit  $a \neq b$ .

Strong Lock Library. The weak lock library discussed above is efficient and gives programmers full control over synchronisation. However, if not used correctly, without the relevant global fences, its guarantees are not as strong as one may expect. That is, in designing the weak lock library, we opted for better performance over the strength of guarantees. We next develop a *strong* lock library that achieves the desired strong guarantees (without the need for additional synchronisation via fences). Specifically, on releasing a strong lock *all* earlier operations are guaranteed to have *fully* completed.

This is illustrated in Fig. 8e, where the outcome  $a \neq b$  is once again prohibited. However, the strong guarantees of strong locks come at the cost of their implementation efficiency. Intuitively, an implementation of a strong lock release

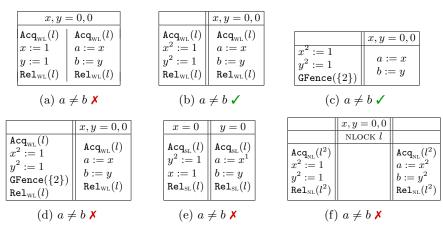


Fig. 8: Examples of weak, strong, and node lock behaviours

issues a global fence towards *every* node on the network to ensure that there are no pending remote operations. This is in contrast to a weak lock release implementation that issues no fence by default, and developers have full control over fencing the relevant nodes.

**Node Lock Library.** As a midway between the efficient weak locks, requiring manual synchronisation, and the inefficient strong locks, we develop the concept of (more fine-grained) node locks. Intuitively, as shown in Fig. 8f, a node lock l is associated with a specific node n (node 2 in Fig. 8f) and provides strong guarantees only for locations on n.

A node lock is stronger than a weak lock: as shown in Fig. 8f the weak behaviour  $a \neq b$  is prohibited without the need for additional synchronisation. Moreover, a node lock is weaker than a strong lock in two ways. First, it only provides guarantees for operations towards one node. For instance, consider a variant of the example in Fig. 8f where the lock l is associated with node 1 (instead of 2); the outcome  $a \neq b$  would once again be allowed. Second, it does not provide any intra-thread ordering guarantees in that releasing a node lock does not guarantee that previous operations (even towards the associated node) have completed. For instance, in the program  $\text{Acq}_{\text{NL}}(l^2); z^2 := x; \text{Rel}_{\text{NL}}(l^2); x := 1$  the outcome z = 1 is allowed:  $z^2 := x$  and the lock release may not have fully completed before the CPU runs the subsequent x := 1 store; i.e., while  $z^2 := x$  cannot be reordered past  $\text{Rel}_{\text{NL}}(l^2)$ , the x := 1 can be reordered before both of them. More concretely, our implementation of  $\text{Rel}_{\text{NL}}$  (§4.4) comprises RDMA operations that can be delayed after later CPU operations.

Nevertheless, as a common usage of a lock is to protect a specific object that is likely to reside on a single node, this level of guarantee is sufficient for many applications, while enabling efficient implementations.

The RDMA<sub>RMW</sub> Library. Lastly, to simplify RDMA programming, we specify and implement the RDMA<sub>RMW</sub> library that fully abstracts away the notion of nodes and provides strong *sequentially consistent* (SC) [23] semantics via four (perlocation) instructions, Write<sub>SC</sub>, Read<sub>SC</sub>, CAS<sub>SC</sub>, and FAA<sub>SC</sub> analogous to stores,

loads, and RMWs on CPUs with strong SC semantics. Our implementation uses node locks to wrap RDMA operations and ensure they become visible in the order they are submitted. Indeed, as we discuss later in §5, we can use the same approach to implement any concurrent data structure over RDMA, and show that it is correct in that it is linearisable [19].

# ${\bf 3}\quad {\bf Extending}\,\,{\rm RDMA}^{\rm WAIT}\,\,{\bf to}\,\,{\rm RDMA}^{\rm WAIT}_{\rm RMW}$

We present RDMA<sup>WAIT</sup> model, an extension of RDMA<sup>WAIT</sup> [4] with remote RMW instructions. Our definitions naturally extend those of RDMA<sup>WAIT</sup>. To underline the distinction between the two, we have highlighted our extensions from RDMA<sup>WAIT</sup> to RDMA<sup>WAIT</sup>. We specify RDMA<sup>WAIT</sup> in MOWGLI [4], yielding a modular semantics that enables compositional reasoning. In particular, as we show below, since LOCO libraries can be freely composed together, this allows us to use the locality result of MOWGLI to verify each library modularly (in isolation). We proceed with an account of MOWGLI preliminaries (§3.1) and present RDMA<sup>WAIT</sup><sub>RMW</sub> in §3.2.

#### 3.1 The MOWGLI Framework Preliminaries

**Libraries.** Intuitively, a library L specification identifies its associated methods as well as the semantics of these methods. A method call is of the form  $m(\widetilde{v})$ , where m denotes the method name and  $\widetilde{v}$  denote its arguments. Ambal et al. capture the method semantics in MOWGLI by identifying the set of executions that are L-consistent in that they uphold the guarantees promised by L. To this end, they associate L with a set  $\mathcal{C}$  of L-consistent executions. A library is then formally defined as a triple  $L = \langle M, \mathsf{loc}, \mathcal{C} \rangle$ , where M is its set of method names (e.g. Write or Put); loc associates each method call with its set of accessed locations (within the method call arguments); and  $\mathcal{C}$  is its set of L-consistent executions. (MOWGLI further requires  $\mathcal{C}$  to adhere to some basic properties to ensure modularity [4], which we elide here.) We use the prefix 'L.' to project the components of a library L, e.g. L.M.

Later (in Program Order) Stamp															
					single				families						
		st	0	1	2	3	4	5	6	7	8	9	10	11	12
ďυ				aCR	aCW	aCAS	aMF	aWT	$\mathtt{aNLR}_n$	$\mathtt{aNRW}_n$	$\mathtt{aNAR}_n$	$\mathtt{aNRR}_n$	$\mathtt{aNLW}_n$	$\mathtt{aRF}_n$	$\mathtt{aGF}_n$
Stamp		A	aCR	1	/	1	1	1	✓	1	✓	✓	1	<b>/</b>	1
S	le.	В	aCW	Х	1	1	1	Х	✓	1	1	✓	1	<b>√</b>	<b>✓</b>
Order)	single	$\mathbf{C}$	aCAS	1	1	1	1	1	✓	1	1	✓	1	<b>√</b>	<b>✓</b>
)rc	ß	D	aMF	1	1	1	1	1	✓	1	1	1	1	<b>√</b>	<b>✓</b>
B		E	aWT	<b>✓</b>	/	1	1	<b>✓</b>	✓	✓	✓	✓	1	<b>✓</b>	<b>✓</b>
Program		F	$\mathtt{aNLR}_n$	Х	Х	Х	Х	Х	SN	SN	SN	SN	SN	SN	SN
rog		$\mathbf{G}$	$\mathtt{aNRW}_n$	Х	Х	Х	Х	Х	X	SN	SN	SN	SN	Х	SN
ı P	lies	H	$\mathtt{aNAR}_n$	Х	Х	Х	Х	Х	X	SN	SN	SN	SN	SN	SN
(in	families	Ι	$\mathtt{aNRR}_n$	Х	Х	Х	Х	Х	X	X	X	X	SN	SN	SN
ier	$_{\mathrm{fa}}$	J	$\mathtt{aNLW}_n$	Х	Х	Х	Х	X	Х	X	Х	X	SN	X	SN
Earlier		K	$\mathtt{aRF}_n$	Х	Х	Х	Х	Х	SN	SN	SN	SN	SN	SN	SN
뙤		L	$\mathtt{aGF}_n$	✓	<b>√</b>	1	<b>√</b>	<b>✓</b>	<b>√</b>	✓	<b>✓</b>	<b>√</b>	<b>√</b>	<b>√</b>	<b>✓</b>

Fig. 9: The sto order in RDMA<sup>WAIT</sup> and RDMA<sup>WAIT</sup>, where highlighted cells denote our extensions from RDMA<sup>WAIT</sup> to RDMA<sup>WAIT</sup>. The  $\checkmark$  denotes that the (program-order-related) stamps are *ordered*; the  $\checkmark$  denotes that the stamps are *not ordered*; the SN denotes the stamps are ordered iff they are associated with the *same node*.

Events and Executions. In the literature of declarative models, traces of a program are represented as a set of *executions*. An execution is a graph comprising: 1) a set of *events* (graph nodes), where each event is associated with the execution of a method call; and 2) a number of relations on events (graph edges). For instance, if thread t executes a  $\operatorname{Read}(x)$  and reads value v, the corresponding event is of the form  $\langle t, \iota, \langle \operatorname{Read}, (x), v \rangle \rangle$ , where  $\iota$  denotes its (unique) event identifier. Identifiers serve to distinguish calls to the same method (with same arguments and output) by the same thread in an execution. For an event e, we write t(e) and m(e) to extract its thread and method name, respectively.

**Definition 1 (Events and Executions).** An event is a tuple  $\langle t, \iota, \langle m, \widetilde{v}, v' \rangle \rangle$ , where  $t \in \mathsf{Tid}$  denotes the executing thread,  $\iota$  is the (unique) event identifier, m denotes the method being executed,  $\widetilde{v} \in \mathsf{Val}^*$  is the method input (arguments) and  $v' \in \mathsf{Val}$  is its output (return value, which may be unit ()). An execution  $\mathcal{G}$  is a tuple  $\langle E, \mathsf{po}, \mathsf{stmp}, \mathsf{so}, \mathsf{hb} \rangle$  where:

- E is the set of events;
- po  $\subseteq E \times E$  is the (strict) program order, total for each thread;
- stmp:  $E \to \mathcal{P}(\mathsf{Stamp})$  associates each event with a non-empty set of stamps and induces a set of subevents,  $\mathsf{SEvent} \triangleq \{ \langle \mathsf{e}, a \rangle \mid \mathsf{e} \in E \land a \in \mathsf{stmp}(\mathsf{e}) \};$
- so ⊆ SEvent × SEvent is the synchronisation order, representing the intralibrary dependencies exported by each library;
- hb ⊆ SEvent × SEvent is the happens-before order, a strict partial order such that so ∪ ppo ⊆ hb, where ppo ⊆ SEvent × SEvent denotes the preserved program order capturing inter-library dependencies and is defined as follows:

$$\mathsf{ppo} \triangleq \{ \langle \langle \mathsf{e}_1, a_1 \rangle, \langle \mathsf{e}_2, a_2 \rangle \rangle \mid \langle \mathsf{e}_1, \mathsf{e}_2 \rangle \in \mathsf{po} \land a_i \in \mathsf{stmp}(\mathsf{e}_i) \land \langle a_1, a_2 \rangle \in \mathsf{sto} \}$$

**Notations.** Given a set A and a relation  $\mathbf{r} \subseteq A \times A$ , we write  $\mathbf{r}^+$  for the transitive closure of  $\mathbf{r}$ ;  $\mathbf{r}^*$  for its reflexive transitive closure;  $\mathbf{r}^{-1}$  for the inverse of  $\mathbf{r}$ ; and [A] for the identity relation on A, i.e.  $\{\langle a,a\rangle \mid a\in A\}$ . We write  $\mathbf{r}_1; \mathbf{r}_2$  for the relational composition of  $\mathbf{r}_1$  and  $\mathbf{r}_2$ :  $\{\langle a,b\rangle \mid \exists c. \langle a,c\rangle \in \mathbf{r}_1 \wedge \langle c,b\rangle \in \mathbf{r}_2\}$ . We write  $A|_c$  to restrict A with condition c. For instance, given a set of events E, we define  $E|_L \triangleq \{\mathbf{e} \in E \mid \mathbf{m}(\mathbf{e}) \in L.M\}$ ,  $E|_t \triangleq \{\mathbf{e} \in E \mid \mathbf{t}(\mathbf{e}) = t\}$ , and we write  $E|_d$  for the set of events in E with work identifier d. We define  $E_x \triangleq \{\mathbf{e} \in E \mid x \in \mathbf{loc}(\mathbf{e})\}$ . Similarly, we define  $\mathbf{r}|_c \triangleq [A|_c]; \mathbf{r}; [A|_c]$  (e.g.  $\mathbf{po}|_t$ ) and  $\mathbf{r}_x \triangleq [E_x]; \mathbf{r}; [E_x]$  (e.g.  $\mathbf{po}_x$ ). Given a subset  $A' \subseteq A$ , we define  $\mathbf{r}|_{A'} \triangleq [A']; \mathbf{r}; [A']$ . When  $\mathbf{r}$  is a strict partial order, we write  $\mathbf{r}|_{\text{imm}}$  for its immediate edges, i.e.  $\mathbf{r} \setminus (\mathbf{r}; \mathbf{r})$ .

Given execution  $\mathcal{G} = \langle E, \mathsf{po}, \mathsf{stmp}, \mathsf{so}, \mathsf{hb} \rangle$ , we write  $\mathcal{G}|_L$  for  $\langle E|_L, \mathsf{po}|_L, \mathsf{stmp}|_L$ ,  $\mathsf{so}|_L, \mathsf{hb}|_L \rangle$ , where  $\mathsf{stmp}|_L$  denotes the function obtained by restricting the domain of  $\mathsf{stmp}$  (i.e. E) to  $E|_L$ . We use the prefix ' $\mathcal{G}$ .' to project the components of  $\mathcal{G}$  (e.g.  $\mathcal{G}.\mathsf{po}$ ), including its derived ones (e.g.  $\mathcal{G}.\mathsf{SEvent}$ ). Given a stamp a, we write  $\mathcal{G}.a$  for  $\{\mathsf{s} \in \mathcal{G}.\mathsf{SEvent} \mid \mathsf{s} = \langle \_, a \rangle\}$ ; analogously for a stamp family, e.g.  $\mathcal{G}.\mathsf{aNRR}$ . We define the set of read subevents as  $\mathcal{G}.\mathcal{R} \triangleq \mathcal{G}.\mathsf{aCR} \cup \mathcal{G}.\mathsf{aCAS} \cup \mathcal{G}.\mathsf{aNLR} \cup \mathcal{G}.\mathsf{aNAR} \cup \mathcal{G}.\mathsf{aNRR}$ , and write subevents as  $\mathcal{G}.\mathcal{W} \triangleq \mathcal{G}.\mathsf{aCW} \cup \mathcal{G}.\mathsf{aCAS} \cup \mathcal{G}.\mathsf{aNLW} \cup \mathcal{G}.\mathsf{aNRW}$ . Given a set of subevents A, we define  $A_x \triangleq \{\mathsf{s} \in A \mid \mathsf{loc}(\mathsf{s}) = \{x\}\}$ ; e.g.  $\mathcal{G}.\mathcal{W}_x$  is the set of write subevents on x. When the choice of  $\mathcal{G}$  is clear, we omit ' $\mathcal{G}.$ ', e.g. we simply write  $\mathcal{W}$  for  $\mathcal{G}.\mathcal{W}$  and  $[\mathsf{aCW}]$  for  $[\mathcal{G}.\mathsf{aCW}]$ .

Consistency. An execution is *consistent* against a set of libraries  $\Lambda$  iff 1)  $\mathcal{G}|_L$  is L-consistent for each  $L \in \Lambda$ ; 2) its events and their synchronisation are those of the libraries in  $\Lambda$ ; and 3) its happens-before relation is irreflexive. Note that the first condition ensures *modularity* as each library can specify independently the visible behaviours of its functions (stamps), its allowed outcomes (consistency) and the synchronisation (guarantees) it offers (so).

**Definition 2 (Consistency).** Let  $\Lambda$  be a set of libraries where  $L_1.M \cap L_2.M = \emptyset$  for distinct  $L_1, L_2$ . An execution  $\mathcal{G} = \langle E, \mathsf{po}, \mathsf{stmp}, \mathsf{so}, \mathsf{hb} \rangle$  is  $\Lambda$ -consistent iff:

- For all  $L \in \Lambda$ :  $\mathcal{G}|_{L} \in L.\mathcal{C}$  (i.e.  $\mathcal{G}$  is L-consistent for each  $L \in \Lambda$ );
- $E = \bigcup_{L \in \Lambda} E|_L$  and so  $= \bigcup_{L \in \Lambda} \operatorname{so}|_L$ ; and
- hb is irreflexive (i.e. hb is a strict partial order).

## 3.2 The Declarative $RDMA_{RMW}^{WAIT}$ Model

We present RDMA $_{\rm RMW}^{\rm WAIT}$  as an extension of RDMA $^{\rm WAIT}$  [4] with remote RMWs. Our definitions naturally extend those of RDMA $^{\rm WAIT}$ . To underline the distinction between the two, we have highlighted our extensions from RDMA $^{\rm WAIT}$  to RDMA $^{\rm WAIT}$ 

The RDMA $_{RMW}^{WAIT}$  Methods. RDMA $_{RMW}^{WAIT}$  methods extend those of RDMA $_{RMW}^{WAIT}$  with remote RMWs as defined by the following grammar, where RDMA $_{RMW}^{WAIT}$  methods comprise local (CPU) operations on TSO machines and remote operations.

```
\begin{split} m(\widetilde{v}) &::= \mathtt{Write}(x,v) \mid \mathtt{Read}(x) \mid \mathtt{CAS}(x,v_1,v_2) \mid \mathtt{Mfence}() \quad /\!/ \; \mathtt{RDMA}^{\mathtt{WAIT}} \colon \mathtt{Local} \\ &\mid \mathtt{Get}(x,y,d) \mid \mathtt{Put}(x,y,d) \mid \mathtt{Wait}(d) \mid \mathtt{Rfence}(n) \quad /\!/ \; \mathtt{RDMA}^{\mathtt{WAIT}} \colon \mathtt{Remote} \\ &\mid \mathtt{RCAS}(x,y,v_1,v_2,d) \mid \mathtt{RFAA}(x,y,v,d) \qquad \qquad /\!/ \; \mathtt{Remote} \; \mathtt{RMWs} \end{split}
```

The remote operations comprise  $\operatorname{Get}$ ,  $\operatorname{Put}$ , Wait (as described in §2.1), and Rfence instructions. Note that for readability in our examples we write  $x:=^dy^n$  (resp.  $x^n:=^dy$ ) for  $\operatorname{Get}(x,y,d)$  (resp.  $\operatorname{Put}(x,y,d)$ ). Similarly, we write x:=v (resp. a:=x) for Write(x,v) (resp.  $\det a = \operatorname{Read}(x)$  in . . .). The Rfence(n) denotes a remote fence that strongly orders all operations towards n without blocking the (local) CPU. That is, given a (sequential) program of the form C; Rfence(n); C', all remote operations towards n in C are ordered before those in C'. The RCAS $(x,y,v_1,v_2,d)$  is the remote analogue of writing  $\det v = \operatorname{CAS}(y,v_1,v_2)$  in Write(x,v) with work identifier d, where the RMW is run on remote location y and the result is written to local location x. Similarly, RFAA(x,y,v,d) increments (remote) y by v and writes its old value to x.

Well-addressed RDMA\*\* Executions. We assume each location x is associated with exactly one node denoted by  $\mathbf{n}(x)$ . We write  $\mathbf{n}(t)$  to denote the node on which t is run. An execution  $\mathcal{G}$  is well-addressed iff it comprises method calls (in  $\mathcal{G}.E$ ) with appropriate local locations when expected; e.g. for each  $\mathsf{Write}(x, \_)$  or  $\mathsf{Put}(\_,x,\_)$  call by thread t in  $\mathcal{G}$ ,  $\mathbf{n}(x) = \mathbf{n}(t)$ . We define  $\mathsf{loc}$  for  $\mathsf{RDMA}^{\mathsf{WAIT}}_{\mathsf{RMW}}$  as expected; e.g.  $\mathsf{loc}(\mathsf{Write}(x,\_)) = \{x\}$ ,  $\mathsf{loc}(\mathsf{Put}(x,y,\_)) = \{x,y\}$  and  $\mathsf{loc}(\mathsf{Mfence}) = \emptyset$ . Well-stamped  $\mathsf{RDMA}^{\mathsf{WAIT}}_{\mathsf{RMW}}$  Executions. An execution  $\mathcal{G}$  is well-stamped if for all  $\mathsf{e} = \langle \_, \_, \langle m, (\widetilde{v}), v' \rangle \rangle \in \mathcal{G}.E$ :  $\mathcal{G}.\mathsf{stmp}(\mathsf{e}) \in \mathsf{stmp}_{\mathsf{RW}}(m(\widetilde{v}), v')$ , with  $\mathsf{stmp}_{\mathsf{RW}}$  defined as follows. Note that depending on whether RCAS calls succeed, they may have multiple valid sets of stamps; as such, the  $\mathsf{stmp}_{\mathsf{RW}}$  function returns a set of stamp sets (set of set of stamps), though in all cases but for RCAS this set is a singleton.

```
\begin{split} & \operatorname{stmp}_{\operatorname{RW}}(\operatorname{CAS}(x,v_1,{}_{\scriptscriptstyle{-}}),v_2) \triangleq \begin{cases} \left\{\{\operatorname{aMF},\operatorname{aCR}\right\} \right\} & \operatorname{if}\ v_1 \neq v_2 \\ \left\{\{\operatorname{aCAS}\right\} & \operatorname{if}\ v_1 = v_2 \end{cases} & \operatorname{stmp}_{\operatorname{RW}}(\operatorname{Read}(x),{}_{\scriptscriptstyle{-}}) \triangleq \left\{\{\operatorname{aCW}\right\} \\ & \operatorname{stmp}_{\operatorname{RW}}(\operatorname{Get}(x,y,{}_{\scriptscriptstyle{-}}),{}_{\scriptscriptstyle{-}}) \triangleq \left\{\left\{\operatorname{aNRR}_{\operatorname{n}(y)},\operatorname{aNLW}_{\operatorname{n}(y)}\right\}\right\} & \operatorname{stmp}_{\operatorname{RW}}(\operatorname{Mfence}(),{}_{\scriptscriptstyle{-}}) \triangleq \left\{\{\operatorname{aMF}\right\} \\ & \operatorname{stmp}_{\operatorname{RW}}(\operatorname{Mid}(x),{}_{\scriptscriptstyle{-}}) \triangleq \left\{\{\operatorname{aMF}\right\} \\ & \operatorname{stmp}_{\operatorname{RW}}(\operatorname{Mid}(x),{}_{\scriptscriptstyle{-}}) \triangleq \left\{\{\operatorname{aMT}\right\} \\ & \operatorname{stmp}_{\operatorname{RW}}(\operatorname{Mid}(x),{}_{\scriptscriptstyle{-}}) \triangleq \left\{\{\operatorname{aWT}\right\} \\ & \operatorname{stmp}_{\operatorname{RW}}(\operatorname{Read}(x),{}_{\scriptscriptstyle{-}}) \triangleq \left\{\{\operatorname{aWT}\right\} \\ & \operatorname{stmp}_{\operatorname{RW}}(\operatorname{Mid}(x),{}_{\scriptscriptstyle{-}}) \triangleq \left\{\{\operatorname{aWT}\right\} \\ & \operatorname{stmp}_{\operatorname{RW}}(\operatorname{Read}(x),{}_{\scriptscriptstyle{-}}) \triangleq \left\{\{\operatorname{aWT}\right\} \\ & \operatorname{stmp}_{\operatorname{RW}}(\operatorname{Mid}(x),{}_{\scriptscriptstyle{-}}) \triangleq \left\{\{\operatorname{AWT}\right\} \\ & \operatorname{Mid}(x),{}_{\scriptscriptstyle{-}} \triangleq \left\{\operatorname{AWT}\right\} \\ & \operatorname{Mid}(x),{}_{\scriptscriptstyle{-}} \triangleq \left\{\operatorname{AWT}\right\} \\ & \operatorname{Mid}(x),{}_{\scriptscriptstyle{-}} \triangleq \left
```

A successful remote RMW has three stamps for reading the remote location, modifying it, and writing it to the local location, while a failed RCAS does not modify the remote location. Recall that the remote read of a remote RMW yields stamp  $aNAR_n$ , which offers more guarantees than the stamp  $aNRR_n$  of Gets.

We extend the location function (loc, defined above for RDMA\_RMW] to subevents. For method calls corresponding to *local operations* (with one or zero locations) their subevents have the same locations. The subevents of Get, Put, RCAS, and RFAA are associated with the relevant location as expected. For instance, if  $e = \langle -, -, \langle \text{Get}, (x, y, d), - \rangle \rangle$  (with subevents aNRR<sub>n</sub> and aNLW<sub>n</sub>), then loc( $\langle e, \text{aNRR}_n \rangle$ ) =  $\{y\}$  and loc( $\langle e, \text{aNLW}_n \rangle$ ) =  $\{x\}$ ; whereas if  $e = \langle -, -, \langle \text{RFAA}, (x, y, v, d), - \rangle \rangle$ , then loc( $\langle e, \text{aNAR}_n \rangle$ )= $\{y\}$ , loc( $\langle e, \text{aNRW}_n \rangle$ )= $\{y\}$  and loc( $\langle e, \text{aNLW}_n \rangle$ )= $\{x\}$ .

Well-formed RDMA\_{RMW}^{WAIT} Executions. We shortly define the notion of RDMA\_{RMW}^{WAIT}-consistency for an execution  $\mathcal{G}$ . To do this, we need a few auxiliary functions and relations as follows. We assume functions  $v_R:\mathcal{G}.\mathcal{R}\to Val$  and  $v_W:\mathcal{G}.\mathcal{W}\to Val$ , which associate each read (resp. write) subevent with the value returned (resp. written). We define the 'reads-from' relation,  $rf\subseteq\mathcal{G}.\mathcal{W}\times\mathcal{G}.\mathcal{R}$ , on subevents of

the same location with matching values (formalised below); the 'modification-order' relation,  $mo \subseteq \mathcal{G}.\mathcal{W} \times \mathcal{G}.\mathcal{W}$ , describing a (total) order in which writes reach the memory; and the 'NIC flush order', nfo, capturing the PCIe guarantee that NIC reads flush previous NIC writes. For remote RMWs, we define the 'remote-atomic-order', rao, describing the (total) order in which remote read parts of remote RMWs towards each node is executed. A tuple  $\langle v_R, v_W, rf, mo, nfo, rao \rangle$  is well-formed if the following hold for all  $e, v, v', v_1, v_2, s_1, s_2, n, x, y$ .

- If e is of the form  $\langle \text{Read}, -, v \rangle$  or  $\langle \text{CAS}, -, v \rangle$ , then  $v_{R}(e) = v$ .
- If e is of the form  $\langle \text{Write}, (\_, v), \_ \rangle$  or  $\langle \text{CAS}, (\_, v', v), v' \rangle$ , then  $v_{\text{W}}(e) = v$
- If  $s_1 = \langle e, aNLR_n \rangle \land s_2 = \langle e, aNRW_n \rangle$ , then  $v_R(s_1) = v_W(s_2)$ ; mutatis mutandis for  $s_1 = \langle e, aNRR_n \rangle$ ,  $s_2 = \langle e, aNLW_n \rangle$  and  $s_1 = \langle e, aNAR_n \rangle$ ,  $s_2 = \langle e, aNLW_n \rangle$ .
- $\langle s_1, s_2 \rangle \in \mathsf{rf} \Rightarrow \mathsf{loc}(s_1) = \mathsf{loc}(s_2) \wedge \mathsf{v}_{\mathtt{W}}(s_1) = \mathsf{v}_{\mathtt{R}}(s_2).$
- rf<sup>-1</sup> is a function, i.e. every read is related to at most one write. If a read is not related to a write, it returns zero:  $s_2 \notin img(rf) \Rightarrow v_R(s_2) = 0$ .
- mo  $\triangleq \bigcup_{x \in \mathsf{Loc}} \mathsf{mo}_x$ , where each  $\mathsf{mo}_x$  is a strict total order on  $\mathcal{G}.\mathcal{W}_x$ .
- if  $\langle s_1, s_2 \rangle \in (\mathtt{aNLR}_n \times \mathtt{aNLW}_n) \cup ((\mathtt{aNRR}_n \cup \mathtt{aNAR}_n) \times \mathtt{aNRW}_n)$  and  $\mathtt{t}(s_1) = \mathtt{t}(s_2)$  then  $\langle s_1, s_2 \rangle \in \mathsf{nfo} \cup \mathsf{nfo}^{-1}$ .
- RCAS succeeds iff it reads the expected value, in which case it overwrites with the given value. That is, given  $\mathbf{e} = \langle \_, \_, \langle \mathtt{RCAS}, (x,y,v_1,v_2,\_),\_ \rangle \rangle$ : if  $\mathtt{stmp}(\mathbf{e}) = \{\mathtt{aNAR}_{\mathtt{n}(y)}, \mathtt{aNLW}_{\mathtt{n}(y)}\}$ , then  $\mathtt{v}_\mathtt{R}(\langle \mathtt{e}, \mathtt{aNAR}_{\mathtt{n}(y)} \rangle) \neq v_1$ ; and if  $\mathtt{stmp}(\mathbf{e}) = \{\mathtt{aNAR}_{\mathtt{n}(y)}, \mathtt{aNRW}_{\mathtt{n}(y)}, \mathtt{aNLW}_{\mathtt{n}(y)}\}$ , then  $\mathtt{v}_\mathtt{R}(\langle \mathtt{e}, \mathtt{aNAR}_{\mathtt{n}(y)} \rangle) = v_1$  and  $\mathtt{v}_\mathtt{W}(\langle \mathtt{e}, \mathtt{aNRW}_{\mathtt{n}(y)} \rangle) = v_2$ .
- If  $e = \langle \_, \_, \langle \mathtt{RFAA}, (x, y, v, \_), \_ \rangle \rangle$ , then  $v_{\mathtt{W}}(\langle \mathtt{e}, \mathtt{aNRW}_{\mathtt{n}(y)} \rangle) = v_{\mathtt{R}}(\langle \mathtt{e}, \mathtt{aNAR}_{\mathtt{n}(y)} \rangle) + v$ .
- rao  $\triangleq \bigcup_{n \in \mathsf{Node}} \mathsf{rao}_n$ , where rao<sub>n</sub> is a strict total order on the set of subevents  $\{\langle \mathsf{e}, \mathtt{aNAR}_n \rangle \mid \mathsf{e} = \langle \_, \_, \langle m, (x, y, \ldots), \_ \rangle \rangle \land m \in \{\mathtt{RFAA}, \mathtt{RCAS}\} \land \mathtt{n}(y) = n\}$

We distinguish the point subevents start executing (point of 'issue') from when they complete. We define the issued-before relation, ib, to record dependencies between the starts of subevents, while so records dependencies between their ends. Note that ib and so are incomparable:  $\langle s_1, s_2 \rangle \in ib$  does not imply  $\langle s_1, s_2 \rangle \in so$  and vice versa. We define instantaneous subevents,  $\mathcal{G}.Inst \triangleq \mathcal{G}.SEvent \setminus (\mathcal{G}.aCW \cup \mathcal{G}.aNLW \cup \mathcal{G}.aNRW)$ , as those that start and end at the same time.

Given an execution  $\mathcal{G}$  and well-formed  $\langle v_R, v_W, rf, mo, nfo, rao \rangle$ , we further define the following relations that will help us define ib and so for RDMA $_{\rm RMW}^{\rm WAIT}$ :

```
\begin{split} \operatorname{rb} &\triangleq \left\{ \langle r, w \rangle \in \mathcal{G}.\mathcal{R} \times \mathcal{G}.\mathcal{W} \;\middle|\; \left( \langle r, w \rangle \in (\operatorname{rf}^{-1}; \operatorname{mo}) \vee r \not\in \operatorname{img}(\operatorname{rf}) \right) \right\} \setminus [\mathcal{G}.\mathsf{SEvent}] \\ \operatorname{rb}_i &\triangleq [\operatorname{aCR}]; ((\operatorname{po} \cup \operatorname{po}^{-1}) \cap \operatorname{rb}); [\operatorname{aCW}] & \operatorname{pfg} \triangleq \left\{ \left\langle \left\langle e_1, \operatorname{aNLW}_n \right\rangle, \left\langle e_2, \operatorname{aWT} \right\rangle \right\rangle \mid \exists d. \left\langle e_1, e_2 \right\rangle \in \operatorname{po} \middle|_d \right\} \\ \operatorname{rf}_i &\triangleq [\operatorname{aCW}]; (\operatorname{po} \cap \operatorname{rf}); [\operatorname{aCR}] & \operatorname{rf}_e \triangleq \operatorname{rf} \setminus \operatorname{rf}_i & \operatorname{pfp} \triangleq \left\{ \left\langle \left\langle e_1, \operatorname{aNRW}_n \right\rangle, \left\langle e_2, \operatorname{aWT} \right\rangle \right\rangle \mid \exists d. \left\langle e_1, e_2 \right\rangle \in \operatorname{po} \middle|_d \right\} \\ \operatorname{iso} &\triangleq & \left\{ \left\langle \left\langle e, \operatorname{aMF} \right\rangle, \left\langle e, \operatorname{aCR} \right\rangle \right\rangle \mid \operatorname{m}(e) = \operatorname{CAS} \right\} \\ &\cup \left\{ \left\langle \left\langle e, \operatorname{aNRR}_n \right\rangle, \left\langle e, \operatorname{aNLW}_n \right\rangle \right\rangle \mid \operatorname{m}(e) = \operatorname{Get} \right\} \cup \left\{ \left\langle \left\langle e, \operatorname{aNLR}_n \right\rangle, \left\langle e, \operatorname{aNRW}_n \right\rangle \right\rangle \mid \operatorname{m}(e) = \operatorname{Put} \right\} \\ &\cup \left\{ \left\langle \left\langle e, \operatorname{aNAR}_n \right\rangle, \left\langle e, \operatorname{aNRW}_n \right\rangle \right\rangle \mid \operatorname{m}(e) \in \left\{ \operatorname{RCAS}, \operatorname{RFAA} \right\} \wedge \operatorname{aNRW}_n \in \operatorname{stmp}(e) \right\} \end{split}
```

The rb denotes the 'reads-before' relation: given a read r that reads from a write  $w_r$ , i.e.  $\langle w_r, r \rangle \in \text{rf}$ , then rb relates r to all writes w (on the same location) that are mo-later than  $w_r$ . The internal rb relation, rb, restricts rb to CPU reads

and writes on the same thread; similarly for rf<sub>i</sub> (internal rf). The external rf, rf<sub>e</sub>, is defined as rf edges that are not internal. The pfg (resp. pfp) relation captures the synchronisation between the local write subevent of a Get or remote RMW (resp. remote write subevent of a Put or remote RMW) and a later Wait with the same work identifier. As we describe shortly, while both are included in ib, only pfg is included in so as waiting for a Put (or remote RMW) does not guarantee that the NIC remote write has completed. The 'internal synchronisation order', iso, captures ordering between subevents of the same event and ensure that a failing CPU CAS performs a memory fence before reading; RDMA operations (Get, Put, and remote RMW) read before copying the value; and a successful remote RMW reads before updating the remote value.

Finally, we define ib as follows and it includes a superset ippo of ppo. Specifically, while a later CPU read might finish before an earlier CPU write or wait (cells B1 and B5, in Fig. 9), they start (are issued) in order; and while a remote fence does not guarantee previous NIC writes have completed (cells G11 and J11, in Fig. 9), it guarantees they have at least started.

```
\begin{array}{c} \mathsf{ib} \triangleq (\mathsf{ippo} \cup \mathsf{iso} \cup \mathsf{rf} \cup \mathsf{pfg} \cup \mathsf{pfp} \cup \mathsf{nfo} \cup \mathsf{rb_i})^+ \\ \mathsf{with} \ \mathsf{ippo} \triangleq \mathsf{ppo} \cup ([\mathsf{aCW}]; \mathsf{po}; [\mathsf{aCR} \cup \mathsf{aWT}]) \cup \ \bigcup_{n \in \mathsf{Node}} ([\mathsf{aNRW}_n \cup \mathsf{aNLW}_n]; \mathsf{po}; [\mathsf{aRF}_n]) \end{array}
```

We next define *consistency* for RDMA<sub>RMW</sub>. We require that ib and so be irreflexive (the latter is implied by irreflexivity of hb in Def. 2 as so  $\subseteq$  hb (Def. 1)).

**Definition 3** (RDMA<sup>WAIT</sup><sub>RMW</sub>-consistency). An execution  $\mathcal{G} = \langle E, \mathsf{po}, \mathsf{stmp}, \mathsf{so}, \mathsf{hb} \rangle$  is RDMA<sup>WAIT</sup><sub>RMW</sub>-consistent iff it is well-addressed, well-stamped, and there exists a well-formed tuple  $\langle \mathsf{v}_\mathsf{R}, \mathsf{v}_\mathsf{W}, \mathsf{rf}, \mathsf{mo}, \mathsf{nfo}, \mathsf{rao} \rangle$  such that:

```
1) ib is irreflexive; and 2) so = iso \cup rf<sub>e</sub> \cup pfg \cup nfo \cup rb \cup mo \cup rao \cup ([aNRW]; iso<sup>-1</sup>; rao) \cup ([Inst]; ib).
```

As described above, rao captures the order in which remote read parts of remote RMWs towards a node is executed. The extension ([aNRW]; iso<sup>-1</sup>; rao) ensures that remote RMWs towards the same node do not overlap: if a remote RMW succeeds, then its remote write completes before the next RMW can read.

#### 4 Specifying and Verifying RDMA Lock Libraries

We use the RDMA<sub>RMW</sub> library to *specify, implement, and verify* three RDMA lock libraries. As discussed in §2.4, designing an RDMA lock presents a trade-off between strong, intuitive behaviours and efficient implementations. As such, after introducing the required preliminaries (§4.1), we develop a weak (WLOCK), strong (SLOCK), and node (NLOCK) lock library.

#### 4.1 Preliminaries

Well-formed Locks. A lock library typically provides two methods Acq(x) and Rel(x) for acquiring and releasing a (network-shared) lock x, ensuring mutual

exclusion; i.e. no two thread can hold the lock on x simultaneously. We assume the existence of a location function loc such that  $loc(Acq(x)) = loc(Rel(x)) = \{x\}$ . We further assume that locks are used in a well-formed fashion: a thread only acquires (resp. releases) lock x if it has not (resp. has) already acquired x. We formalise this in Def. 4 below, requiring that each loc Acq(x) (resp. loc Rel(x)) is followed (resp. preceded) by loc Rel(x) (resp. loc Acq(x)) in program order.

**Definition 4.** An execution  $\langle E, po, -, -, - \rangle$  is lock-well-formed iff for all x:

```
1) for all e_a \in E_x there exists an e_r \in E_x such that \langle e_a, e_r \rangle \in po_x|_{imm}; and 2) for all e_r \in E_x there exists an e_a \in E_x such that \langle e_a, e_r \rangle \in po_x|_{imm}
```

where  $e_a, e_r$  are acquire and release events:  $m(e_a) = Acq$  and  $m(e_r) = Rel$ .

Library guarantees only hold for programs that adhere to this well-formedness requirement. For those that do not, *any* behaviour is allowed.

Background: SV Library. Ambal et al. [4] use RDMA<sup>WAIT</sup> to define higher-level libraries such as a shared-variable library (SV) where each node maintains its own copy for each location x. A thread then accesses (reads/writes) its own local copies, and can broadcast its local value to other nodes. The SV library comprises these methods:  $M = \{\text{Write}_{\text{SV}}, \text{Read}_{\text{SV}}, \text{Bcast}_{\text{SV}}, \text{Wait}_{\text{SV}}, \text{GFence}\}$ . The Write<sub>SV</sub>(x, v) (resp. Read<sub>SV</sub>(x)) writes (resp. reads) value v to the local copy of x on the current node. The Bcast<sub>SV</sub> $(x, d, \{n_1, \ldots, n_k\})$  broadcasts the local value of x and overwrites x on nodes  $n_1, \ldots, n_k$ , which may include the local node itself (where d is the work id). The Wait<sub>SV</sub>(d) waits for previous broadcasts of the thread associated with work id  $d \in \text{Wid}$ . Finally, the global fence GFence( $\{n_1, \ldots, n_k\}$ ) ensures every previous operation of the thread towards nodes  $n_1, \ldots, n_k$  is fully completed. We repeat the formal semantics of SV in §A. In the remainder of this article we use SV to implement several libraries.

#### 4.2 The Weak Lock Library

We present our WLOCK library, which only guarantees *mutual exclusion*, without any guarantees on the completion order of submitted RDMA operations.

The WLOCK Specification. The stamps for WLOCK are defined through the  $\mathsf{stmp}_{\mathtt{WL}}$  function as follows. That is, acquiring a weak lock behaves as a memory fence (stamp  $\mathtt{aMF}$ ) on TSO, while releasing it behaves merely as a write ( $\mathtt{aCW}$ ).

$$\mathtt{stmp}_{\mathtt{WL}}(\langle t, \_, \langle \mathtt{Acq}_{\mathtt{WL}}, (x), () \rangle \rangle) \triangleq \{\mathtt{aMF}\} \qquad \mathtt{stmp}_{\mathtt{WL}}(\langle t, \_, \langle \mathtt{Rel}_{\mathtt{WL}}, (x), () \rangle \rangle) \triangleq \{\mathtt{aCW}\}$$

As we formulate in Def. 5 below (the second condition), WLOCK provides synchronisation between lock releases and acquisitions of each lock.

**Definition 5** (WLOCK-consistency). A lock-well-formed execution  $\mathcal{G} = \langle E, po, stmp, so, hb \rangle$  is WLOCK-consistent iff:

```
1)  \begin{split} \text{stmp} &= \text{stmp}_{\text{WL}} \ (\textit{where} \ \text{stmp}_{\text{WL}} \ \textit{is as defined above}); \ \textit{and} \\ 2) \ &= \bigcup_x \left\{ \langle \langle \mathsf{e}_1, \mathsf{aCW} \rangle, \langle \mathsf{e}_2, \mathsf{aMF} \rangle \rangle \ \middle| \ \langle \mathsf{e}_1, \mathsf{e}_2 \rangle \in (\mathsf{po}_x|_{imm})^{-1}; \mathsf{lo}_x \right\}, \ \textit{where} \ \mathsf{lo}_x \ \textit{is a total order on acquisition events on} \ x, \ \textit{i.e. on} \ \{ \mathsf{e} \in E_x \ \middle| \ \mathsf{m}(\mathsf{e}) = \mathsf{Acq}_{\mathrm{WL}} \}. \end{split}
```

```
\begin{split} I_{\text{WL}}(t, \texttt{Acq}_{\text{WL}}, (x)) &\triangleq \\ & \texttt{RFAA}(p_x^t, x_a, 1, d); \; \texttt{Wait}(d); & I_{\text{WL}}(t, \texttt{Rel}_{\text{WL}}, (x)) \triangleq \\ & \texttt{let} \; v = \texttt{Read}(p_x^t) \; \texttt{in} & \texttt{let} \; v = \texttt{Read}(p_x^t) \; \texttt{in} \\ & \texttt{loop} \; \{ \texttt{if} \; \texttt{Read}_{\text{SV}}(x_1) = v \; \texttt{then break else} & & \texttt{Write}_{\text{SV}}(x_t, v + 1); \\ & \dots & \texttt{Bcast}_{\text{SV}}(x_t, -, \mathsf{Node} \setminus \{ \texttt{n}(t) \}) \end{split}
```

Fig. 10: The WLOCK implementation using RDMA<sub>RMW</sub> and SV libraries.

Given a release event  $e_1$  on x (in a lock-well-formed execution), the  $(po_x|_{imm})^{-1}$  component identifies an acquire event  $e_3$  that is the latest corresponding acquire event on x preceding  $e_1$  (in po). As such, so induces synchronisation between  $e_1$  and all later (in  $lo_x$ ) acquisition events  $e_2$ . Note that  $lo_x$  is also indirectly included in hb, since the acquire and release operations stay in order.

The release stamp (aCW) does not synchronise with previous RDMA-specific stamps (bottom-left part of Fig. 9). As such, reacquiring a lock does not guarantee that previous RDMA operations submitted with the lock are completed.

The WLOCK (**Distributed**) Implementation. We present our WLOCK implementation in Fig. 10 (via the  $I_{\rm WL}$  function), inspired by the well-known ticket lock implementation. For each lock location x, we create a ticket dispenser  $x_a$  (on some arbitrary node) that records the value of the next available ticket, thread-local locations ( $p_x^t$  for each  $t \in \mathsf{Tid} = \{1, \ldots, T\}$ ) to track the ticket allocated to t (i.e. its turn), and shared variables  $x_t$  (for each  $t \in \mathsf{Tid}$ ) to signal releasing the lock.

To release the lock on x, thread t writes the next turn, i.e. v+1 when t holds ticket v (obtained by reading  $p_x^t$ ), to its release location  $x_t$  and subsequently broadcasts it to all nodes other than itself  $(\mathbf{n}(t))$ . To acquire the lock on x, thread t calls a fetch-and-add on  $x_a$  to fetch the next available ticket (i.e. its turn) in  $p_x^t$  and increments  $x_a$ . It then records its turn in v and repeatedly examines the release location  $x_{t'}$  of each thread  $t' \in \{1, \ldots, T\}$  until one has value v, indicating that its turn has come and thus t holds the lock. Note that t' may be t itself, i.e. t = t', if it was the last thread to release the lock.

At the cost of more network messages (through broadcasts), our implementation provides lower latency than centralised systems (e.g. in Fig. 13) as messages are transmitted directly from the thread releasing the lock to the next thread acquiring the lock. We next prove (Theorem 1) that our implementation is correct against the WLOCK specification with the full proof given in §B.2.

**Theorem 1.** The implementation  $I_{WL}$  is sound.

#### 4.3 The Strong Lock Library

We present our strong lock library SLOCK that, as well as ensuring mutual exclusion of critical sections, additionally guarantees that *all* earlier operations have *fully* completed on releasing a strong lock. We present several examples of the

x, y = 0, 0			x, y =		
$\mathtt{Acq}_{\scriptscriptstyle \mathrm{WL}}(l)$	$\mathtt{Acq}_{\scriptscriptstyle \mathrm{WL}}(l)$		Acq,		
x := 1	$a := x^1$		x :=		
y := 1	$b := y^1$		y :=		
$\mathtt{Rel}_{\scriptscriptstyle \mathrm{WL}}(l)$	$\mathtt{Rel}_{\scriptscriptstyle \mathrm{WL}}(l)$		Rel		
. ,	. ,				
(a) $a \neq b \checkmark$					

x, y = 0, 0			
$\begin{aligned} & \mathtt{Acq}_{\scriptscriptstyle\mathrm{WL}}(l) \\ & x := 1 \\ & y := 1 \\ & \mathtt{Rel}_{\scriptscriptstyle\mathrm{WL}}(l) \end{aligned}$	$\begin{array}{l} \texttt{Acq}_{\scriptscriptstyle{\mathrm{WL}}}(l) \\ a :=^d x^1 \\ b :=^d y^1 \\ \texttt{Wait}(d) \\ \texttt{Rel}_{\scriptscriptstyle{\mathrm{WL}}}(l) \end{array}$		
(b) <i>a</i> ≠ <i>b</i> ×			

x, y = 0, 0			
$\begin{aligned} & \mathtt{Acq}_{\scriptscriptstyle{\mathrm{SL}}}(l) \\ & x := 1 \\ & y := 1 \\ & \mathtt{Rel}_{\scriptscriptstyle{\mathrm{SL}}}(l) \end{aligned}$	$\begin{array}{c} \mathtt{Acq}_{\scriptscriptstyle{\mathrm{SL}}}(l) \\ a := x^1 \\ b := y^1 \\ \mathtt{Rel}_{\scriptscriptstyle{\mathrm{SL}}}(l) \end{array}$		
(c) $a \neq b \times$			

Fig. 11: Weak versus strong locks when interacting with Get instructions.

'message-passing' behaviour in Fig. 11 contrasting the behaviour of weak and strong locks when interacting with Gets and whether the weak outcome  $a \neq b$  is allowed. In particular, we may observe  $a \neq b$  when using a weak lock (Fig. 11a) and this can be prohibited by explicitly waiting (using Wait(d)) on the completion of the Gets before releasing the weak lock (Fig. 11b). By contrast, when using a strong lock we no longer need to wait for their completion as this is guaranteed by the strong lock release (Fig. 11c).

The SLOCK Specification. The SLOCK stamps are defined (via  $stmp_{SL}$ ) as:

$$\mathtt{stmp}_{\mathtt{SL}}(\langle t, \_, \langle \mathtt{Acq}_{\mathtt{SL}}, (x), () \rangle \rangle) \triangleq \{\mathtt{aMF}\} \quad \mathtt{stmp}_{\mathtt{SL}}(\langle t, \_, \langle \mathtt{Rel}_{\mathtt{SL}}, (x), () \rangle \rangle) \triangleq \bigcup_{n \in \mathsf{Node}} \{\mathtt{aGF}_n\}$$

As with WLOCK, acquiring a strong lock behaves as a memory fence (aMF), while releasing it behaves as a global fence (aGF), ensuring that all previous remote operations are completed.

**Definition 6** (SLOCK-consistency). A lock-well-formed execution  $\mathcal{G} = \langle E, po, stmp, so, hb \rangle$  is SLOCK-consistent iff:

```
1)  \begin{split} \text{stmp} &= \text{stmp}_{\text{SL}} \ (\textit{where} \ \text{stmp}_{\text{SL}} \ \textit{is} \ \textit{defined} \ \textit{above}); \ \textit{and} \\ 2) \ &= \bigcup_{x \in \mathsf{Loc}, n \in \mathsf{Node}} \left\{ \langle \langle \mathsf{e}_1, \mathsf{aGF}_n \rangle, \langle \mathsf{e}_2, \mathsf{aMF} \rangle \rangle \ \big| \ \langle \mathsf{e}_1, \mathsf{e}_2 \rangle \in (\mathsf{po}_x|_{imm})^{-1}; \mathsf{lo}_x \right\}, \ \textit{where} \\ &\mathsf{lo}_x \ \textit{is} \ \textit{a} \ \textit{total} \ \textit{order} \ \textit{on} \ \big\{ \mathsf{e} \in E_x \ \big| \ \mathsf{m}(\mathsf{e}) = \mathsf{Acq}_{\mathsf{SL}} \big\}.  \end{split}
```

**Strong Lock Implementation.** We implement SLOCK (via  $I_{SL}$ ) simply by combining the weak locks and global fences (from the SV library) as follows:

$$I_{\mathrm{SL}}(t, \mathrm{Acq}_{\mathrm{SL}}, (x)) \triangleq \mathrm{Acq}_{\mathrm{WL}}(x) \qquad I_{\mathrm{SL}}(t, \mathrm{Rel}_{\mathrm{SL}}, (x)) \triangleq \mathrm{GFence}(\mathsf{Node}); \mathrm{Rel}_{\mathrm{WL}}(x)$$

Finally, we prove (Theorem 2) that our implementation is sound against the SLOCK specification with the full proof given in §B.3.

**Theorem 2.** The implementation  $I_{SL}$  is sound.

#### 4.4 The Node Lock Library

A common use case of locks is to protect an object (set of locations) on a specific node. In such cases, neither weak nor strong locks are suitable as they either incur a high programmer burden (weak locks) or a high performance overhead (strong

				x = 0	y = 0	z = 0
	x = 0	y = 0		NLOCK $l$		
x=0   y=0	NLOCK l		$\mathtt{Acq}_{\scriptscriptstyle{\mathrm{NL}}}(l^2)$		a := y	
$\begin{bmatrix} \mathtt{Acq}_{\mathtt{SL}}(l) \\ x^2 := 1 \end{bmatrix} \qquad \boxed{a := y}$	$\begin{bmatrix} \mathtt{Acq}_{\scriptscriptstyle{\mathrm{NL}}}(l^2) \\ x^2 := 1 \end{bmatrix}$	a :- u	$x^2 := 1$ $z^4 := 1$		$\begin{aligned} a &:= y \\ \mathtt{Acq}_{\scriptscriptstyle{\mathrm{NL}}}(l^2) \\ b &:= x^2 \\ c &:= z^4 \end{aligned}$	
$\begin{vmatrix} x^2 := 1 \\ \operatorname{Rel}_{\operatorname{SL}}(l) \end{vmatrix} \qquad \begin{vmatrix} a := y \\ b := x^2 \end{vmatrix}$	$\left\  \frac{x}{\text{Rel}_{\text{NL}}(l^2)} \right\ $	a := y $b := x^2$	z = 1		$\begin{vmatrix} o & -x \\ c & -x^4 \end{vmatrix}$	
$\begin{bmatrix} y^3 := 1 \end{bmatrix} \qquad \begin{bmatrix} b := x \\ \end{bmatrix}$	$y^3 := 1$	0 = x	$\begin{array}{l} \mathtt{Rel}_{\scriptscriptstyle \mathrm{NL}}(l^2) \\ y^3 := 1 \end{array}$		$\left  egin{array}{ccc} {\sf Rel}_{\scriptscriptstyle  m NL}(l^2) \end{array}  ight $	
(a) $(a,b) = (1,0) \times$	(b) $(a, b) = (1,$	0) 🗸	(c)	(a, b, c) =	=(1,0,_) ×	•
				(a,b,c)=	(1,, 0)	

Fig. 12: Strong (left) versus node (middle and right) locks examples.

locks). To address this, we develop  $node\ locks$ , NLOCK, a novel lock library that provides synchronisation on a specific node. Given a node lock x on node n, we write  $\mathbf{n}(x)$  for n. A node lock x ensures that on re-acquiring it all previous remote operations (within a critical section of x) towards n are observable.

The NLOCK Specification. The NLOCK stamps are defined (via stmp<sub>M</sub>) as:

$$\mathtt{stmp}_{\mathtt{NL}}(\langle t, \_, \langle \mathtt{Acq}_{\mathtt{NL}}, (x), () \rangle \rangle) \triangleq \{\mathtt{aMF}\} \qquad \mathtt{stmp}_{\mathtt{NL}}(\langle t, \_, \langle \mathtt{Rel}_{\mathtt{NL}}, (x), () \rangle \rangle) \triangleq \{\mathtt{aRF}_{\mathtt{n}(x)}, \mathtt{aNRW}_{\mathtt{n}(x)}\}$$

Note that unlike WLOCK, the NLOCK releases use  $aRF_n$  and  $aNRW_n$  stamps to synchronise with previous remote operations towards n (i.e. those with stamps  $aNAR_n$ ,  $aNRR_n$ , and  $aNRW_n$ ). Importantly, note that unlike in SLOCK, the release should not include a global fence stamp  $(aGF_n)$  as that would be too strong. By using  $aRF_n$  and  $aNRW_n$ , we ensure that previous operations towards n are completed only when the lock is later re-acquired, and they may not have yet completed on release. This means that, when appropriate, using a node lock is more efficient than combining a weak lock with a global fence.

To understand the difference between strong and node locks, consider the examples in Figs. 12a and 12b, where the  $x^2 := 1$  Put by node 1 is enclosed within a lock, while the  $b := x^2$  Get by node 3 is without a lock. In Fig. 12a,  $\text{Rel}_{\text{SL}}(l)$  ensures that the earlier  $x^2 := 1$  has completed. As such, a = 1 implies that  $y^3 := 1$  has been executed and that x (in node 2) has been modified, ensuring b = 1. By contrast, the  $\text{Rel}_{\text{NL}}(l^2)$  in Fig. 12b does not wait for  $x^2 := 1$  to complete, i.e.  $x^2 := 1$  may complete after  $y^3 := 1$ . We can prevent this by enclosing  $b := x^2$  within the node lock, as shown in Fig. 12c. Specifically,  $y^3 := 1$  in Fig. 12c may still complete before earlier remote operations. However, a = 1 implies that  $y^3 := 1$  is executed, and that thread 1 has at least acquired the lock. As such, when thread 3 acquires the lock via  $\text{Acq}_{\text{NL}}(l)$ , it synchronises with  $\text{Rel}_{\text{NL}}(l)$  in thread 1 and ensures that  $x^2 := 1$  is completed on lock acquisition.

Note that the node lock l protects the accesses towards locations on node 2 only. Thus, in Fig. 12c, it only guarantees that  $x^2 := 1$  is completed but not necessarily  $z^4 := 1$  (towards node 4), and thus (a,c)=(1,0) is an allowed outcome. By contrast, Fig. 8f in the overview showcases the lock guarantees: as x and y both reside on node 2, the lock ensures that their accesses by threads 1 and 3 are mutually exclusive, i.e.  $a \neq b$  is disallowed. Specifically, if thread 1 acquires l first, x and y are modified before being read by thread 3, i.e. a = b = 1.

Conversely, if thread 3 acquires l first, x and y are read before being modified by thread 1, i.e. a=b=0.

**Definition 7** (NLOCK-consistency). A lock-well-formed execution  $\mathcal{G} = \langle E, po, stmp, so, hb \rangle$  is NLOCK-consistent iff:

```
1)  \begin{split} \text{stmp} &= \text{stmp}_{\text{NL}} \ (\textit{where} \ \text{stmp}_{\text{NL}} \ \textit{is as defined above}); \ \textit{and} \\ 2) \ &\text{so} &= \left\{ \left\langle \left\langle \mathsf{e}, \mathsf{aRF}_{\mathsf{n(loc(e))}} \right\rangle, \left\langle \mathsf{e}, \mathsf{aNRW}_{\mathsf{n(loc(e))}} \right\rangle \right\rangle \, \middle| \ \mathsf{m}(\mathsf{e}) = \mathsf{Rel}_{\text{NL}} \right\} \\ &\qquad \bigcup_{x \in \mathsf{Loc}} \left\{ \left\langle \left\langle \mathsf{e}_1, \mathsf{aNRW}_{\mathsf{n(loc(e_1))}} \right\rangle, \left\langle \mathsf{e}_2, \mathsf{aMF} \right\rangle \right\rangle \, \middle| \ \left\langle \mathsf{e}_1, \mathsf{e}_2 \right\rangle \in (\mathsf{po}_x|_{imm})^{-1}; \mathsf{lo}_x \right\} \\ \textit{where} \ &\mathsf{lo}_x \ \textit{is a total order on} \ \left\{ \mathsf{e} \in E_x \mid \mathsf{m}(\mathsf{e}) = \mathsf{Acq}_{\text{NL}} \right\}. \end{split}
```

The NLOCK Implementation. We implement NLOCK as a centralised ticket lock using remote RMWs. For each (node) lock x associated with node  $\mathbf{n}(x)$ , we create two remote locations  $x_a$  and  $x_r$  on  $\mathbf{n}(x)$ . As before,  $x_a$  is the ticket dispenser and records the next available ticket. The  $x_r$  tracks the release counter and indicates which ticket currently holds the lock. Each thread also uses a local location  $p_x^t$  to hold the result of remote operations.

Acquiring the lock on x calls a fetchand-add on  $x_a$  to fetch the next available ticket in  $p_x^t$  and increments  $x_a$ . It then

```
\begin{split} I_{\mathrm{NL}}(t, \mathrm{Acq_{_{\mathrm{NL}}}}, (x)) &\triangleq \\ &\mathrm{RFAA}(p_x^t, x_a, 1, d); \mathrm{Wait}(d); \\ &\mathrm{let} \ v = \mathrm{Read}(p_x^t) \ \mathrm{in} \\ &\mathrm{loop} \ \{ \\ &\mathrm{Get}(p_x^t, x_r, d); \mathrm{Wait}(d); \\ &\mathrm{if} \ \mathrm{Read}(p_x^t) = v \ \mathrm{then} \ \mathrm{break} \ \}; \\ &\mathrm{Write}(p_x^t, v + 1) \\ I_{\mathrm{NL}}(t, \mathrm{Rel_{_{\mathrm{NL}}}}, (x)) &\triangleq \\ &\mathrm{Rfence}(\mathbf{n}(x)); \\ &\mathrm{Put}(x_r, p_x^t, -) \end{split}
```

Fig. 13: Node lock implementation  $(I_{NL})$  using RDMA $_{\text{RMW}}^{\text{WAIT}}$ 

records the ticket value in v and repeatedly examines  $x_r$  until it has value v, indicating that its turn has come and thus t holds the lock. Finally, it increments its ticket value in  $p_x^t$  in preparation for later releasing the lock; i.e.  $p_x^t$  now records the ticket whose turn is next. As such, releasing the lock simply updates  $x_r$  to  $p_x^t$  using a Put rather than an RMW; this is because only the lock holder can write to  $x_r$ . Note that the preceding Rfence ensures that earlier Get operations towards  $\mathbf{n}(x)$  have completed before the lock is release. We prove (Theorem 3) that our implementation is correct against the NLOCK specification with the full proof given in §B.4.

**Theorem 3.** The implementation  $I_{NL}$  is sound.

# ${f 5}$ The RDMA $_{ m RMW}^{ m SC}$ Library

We specify (§5.1), implement, and verify (§5.2) the RDMA<sup>SC</sup><sub>RMW</sub> library that provides intuitive read, write, and RMW operations with the strong semantics of sequential consistency (SC) [23]. That is, as with SC, the instructions in each thread under RDMA<sup>SC</sup><sub>RMW</sub> are always observed in (program) order. Moreover, unlike in RDMA<sup>MAT</sup><sub>RMW</sub> or the lock libraries in §4, the users do not need to specify whether a location is local or remote and which node it resides on. For instance, a user can simply call Write<sub>SC</sub>(x, v) to write (with SC semantics) to location x,

regardless of whether x is local (on the current node) or remote. As such, we use the typewriter font and write x to denote an abstract RDMA<sup>SC</sup><sub>RMW</sub> location whose underlying memory address may be local (i.e. x=x) or on a remote node n (i.e.  $x=x^n$ ).

## 5.1 The RDMARMW Specification

The RDMA<sup>SC</sup><sub>RMW</sub> Methods. The RDMA<sup>SC</sup><sub>RMW</sub> library has four methods: Read<sub>SC</sub>(x), to read from x; Write<sub>SC</sub>(x, v) to write v to x; CAS<sub>SC</sub>(x,  $v_1, v_2$ ), a compare-and-swap on x; and FAA<sub>SC</sub>(x, v), a fetch-and-add on x. We define loc as expected, i.e. loc(Write<sub>SC</sub>(x, v)) = loc(Read<sub>SC</sub>(x)) = loc(CAS<sub>SC</sub>(x,  $v_1, v_2$ )) = loc(FAA<sub>SC</sub>(x, v)) = {x}. We extend po and loc to subevents as expected.

**Well-formedness.** Given an RDMA<sup>SC</sup><sub>RMW</sub> execution  $\mathcal{G}$ , we define the sets of read subevents ( $\mathcal{R}$ ) to comprise all subevents except writes and the set of write subevents ( $\mathcal{W}$ ) to include all subevents except reads and failed RMWs.

$$\mathcal{R} \triangleq \{ \langle \mathsf{e}, \mathsf{aMF} \rangle \mid \mathsf{e} \in \mathcal{G}.E \setminus \{ \langle \_, \_, \langle \mathsf{Write}_{\mathsf{SC}}, \_, \_ \rangle \rangle \} \}$$

$$\mathcal{W} \triangleq \{ \langle \mathsf{e}, \mathsf{aMF} \rangle \mid \mathsf{e} \in \mathcal{G}.E \setminus \{ \langle \_, \_, \langle \mathsf{Read}_{\mathsf{SC}}, \_, \_ \rangle \rangle \} \setminus \{ \langle \_, \_, \langle \mathsf{CAS}_{\mathsf{SC}}, (\_, v, \_), v' \rangle \rangle \mid v \neq v' \} \}$$

As before, a tuple  $\langle v_R, v_W, rf, mo \rangle$  is well-formed if the following holds:

•  $v_R/v_W$  map each read/write subevent to the value read/written:

$$\begin{split} \mathbf{v}_{\mathbf{R}}(\langle\langle -, -, \langle -, -, v \rangle \rangle, - \rangle) &\triangleq v & \quad \mathbf{v}_{\mathbf{W}}(\langle\langle -, -, \langle \mathtt{CAS}_{\mathrm{SC}}, (-, v_{1}, v_{2}), v_{1} \rangle\rangle, - \rangle) \triangleq v_{2} \\ \mathbf{v}_{\mathbf{W}}(\langle\langle -, -, \langle \mathtt{Write}_{\mathrm{SC}}, (-, v), - \rangle\rangle, - \rangle) &\triangleq v & \quad \mathbf{v}_{\mathbf{W}}(\langle\langle -, -, \langle \mathtt{FAA}_{\mathrm{SC}}, (-, v), v' \rangle\rangle, - \rangle) \triangleq v + v' \end{split}$$

 $\bullet$  rf and mo satisfy the same constraints as well-formedness of RDMA  $_{\rm RMW}^{\rm WAIT}$  (§3.2).

We next define RDMA<sup>SC</sup><sub>RMW</sub>-consistency, which requires that 1) each event be associated with (single) stamp aMF; and 2) so = po  $\cup$  rf  $\cup$  mo  $\cup$  rb. The former ensures that RDMA<sup>SC</sup><sub>RMW</sub> calls remain ordered with respect to other non-RDMA operations. The latter captures the standard notion of happens-before in SC [27].

**Definition 8** (RDMA<sub>RMW</sub>-consistency). Execution  $\mathcal{G}$  is RDMA<sub>RMW</sub>-consistent if:

- 1)  $\forall e \in E. \ \text{stmp}(e) = \{aMF\}, \ and$
- 2) there exists a well-formed  $\langle v_R, v_W, rf, mo \rangle$  such that  $\mathcal{G}.so = \mathcal{G}.po \cup rf \cup mo \cup rb$ , where rb is defined as in §3.2.

## 5.2 The $RDMA_{RMW}^{SC}$ Implementation

We implement RDMA<sup>SC</sup><sub>RMW</sub> using node locks and RDMA<sup>WAIT</sup><sub>RMW</sub> operations, as shown in Fig. 14. For each RDMA<sup>SC</sup><sub>RMW</sub> location  $\mathbf{x}$ , we create an RDMA<sup>WAIT</sup><sub>RMW</sub> location x on some arbitrary node. We assume each thread t has access to a private location  $r_t$  for recording the remote data it reads, and a private location  $p_x^t$  for recording

```
I_{\mathtt{SC}}(t, \mathtt{Write}_{\mathtt{SC}}, (\mathtt{x}, v)) \triangleq
                                                                        I_{\mathtt{SC}}(t,\mathtt{Read}_{\mathtt{SC}},(\mathtt{x})) \triangleq
                                                                                                                                       I_{SC}(t, CAS_{SC}, (x, v_1, v_2)) \triangleq
                                                                                                                                                                                                                   I_{SC}(t, FAA_{SC}, (x, v)) \triangleq
     Acq_{NL}(l_x);
                                                                              Acq_{NL}(l_x);
                                                                                                                                             Acq_{NL}(l_x);
                                                                                                                                                                                                                          Acq_{_{\mathrm{NL}}}(l_{\mathtt{x}});
                                                                              Get(r_t, x, d);
                                                                                                                                            RCAS(r_t, x, v_1, v_2, d);
                                                                                                                                                                                                                          RFAA(r_t, x, v, d);
     Write(p_x^t, v);
                                                                              \mathtt{Rel}_{\scriptscriptstyle \mathrm{NL}}(l_{\mathtt{x}});
                                                                                                                                            \mathtt{Rel}_{\scriptscriptstyle \mathrm{NL}}(l_{\mathtt{x}});
                                                                                                                                                                                                                          \mathtt{Rel}_{\scriptscriptstyle \mathrm{NL}}(l_{\mathtt{x}});
    \operatorname{Put}(x,p_{\mathtt{x}}^{t},\underline{\ \ });
                                                                              Wait(d):
                                                                                                                                            Wait(d):
                                                                                                                                                                                                                          Wait(d):
    \mathtt{Rel}_{\scriptscriptstyle \mathrm{NL}}(l_{\mathtt{x}})
                                                                              Read(r_t)
                                                                                                                                            Read(r_t)
                                                                                                                                                                                                                          \mathtt{Read}(r_t)
```

Fig. 14: The implementation of RDMA $_{\rm RMW}^{\rm SC}$  (through the  $I_{\rm SC}$  function)

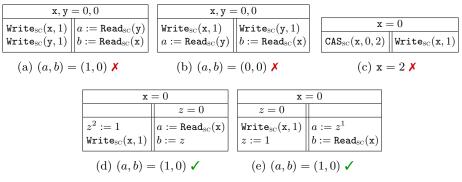


Fig. 15: RDMA<sub>RMW</sub> examples

the value to be put to a remote location (i.e. the second argument of a Put)<sup>3</sup>. Moreover, each location x is associated with a node lock  $l_x$  hosted on the same node as x. We implement RDMA<sup>SC</sup><sub>RMW</sub> writes, reads, and RMWs respectively using Put, Get, and remote RMWs of RDMA<sup>WAIT</sup><sub>RMW</sub> while holding the  $l_x$  lock.

Note that the  $\mathtt{Write}_{\operatorname{SC}}(\mathtt{x},v)$  implementation does not wait for  $\mathtt{Put}(x,p_{\mathtt{x}}^t,\lrcorner)$  to complete. As such, when running  $\mathtt{Write}_{\operatorname{SC}}(\mathtt{x},1); \mathtt{Write}_{\operatorname{SC}}(\mathtt{y},1)$  in Fig. 15a with  $\mathtt{n}(x) \neq \mathtt{n}(y)$ , location y may be modified before x. However, this out-of-order completion is not observable, i.e. the (non-SC) outcome (a,b)=(1,0) is disallowed, because re-acquiring a node lock makes all previous operations towards its node visible (see §4.4). Specifically, a=1 implies that  $\mathtt{Put}(x,p_{\mathtt{x}}^t,\lrcorner)$  has been issued. As the implementation of  $\mathtt{Read}_{\operatorname{SC}}(\mathtt{x})$  acquires  $l_{\mathtt{x}}$ , this enforces  $\mathtt{Put}(x,p_{\mathtt{x}}^t,\lrcorner)$  to become visible; i.e.  $\mathtt{Read}_{\operatorname{SC}}(\mathtt{x})$  reads 1 and (a,b)=(1,0) is disallowed.

In contrast to  $\mathtt{Write}_{\mathtt{SC}}(\mathtt{x}, v)$ , the implementations of the other three operations must wait (via  $\mathtt{Wait}(d)$ ) for their remote operations to complete prior to reading the result via  $\mathtt{Read}(r_t)$  to ensure they observe the correct value. For instance, were we to remove  $\mathtt{Wait}(d)$  in the implementation of  $\mathtt{Read}_{\mathtt{SC}}(\mathtt{x})$ , the  $\mathtt{Read}(r_t)$  could read a stale value from  $r_t$  before  $\mathtt{Get}(r_t, x, d)$  completes and updates  $r_t$ . Nevertheless, it is sufficient to wait for the remote operation to complete after releasing the lock. That is, it is possible for another thread to acquire  $l_{\mathtt{x}}$  (and modify x) before  $\mathtt{Get}(r_t, x, d)$  completes; however, the semantics of NLOCK ensures that  $\mathtt{Get}(r_t, x, d)$  reads the old value into  $r_t$ .

The RDMA $_{\rm RMW}^{\rm SC}$  library, when used in isolation (without calls to e.g. RDMA $_{\rm RMW}^{\rm WAIT}$ ), ensures SC behaviour. As such, the weak behaviours of 'message-passing' in

<sup>&</sup>lt;sup>3</sup> In practice, we can use a Put with 'inlined data' and forgo temporary location  $p_{\mathbf{x}}^t$ .

Fig. 15a and 'store-buffering' in Fig. 15b are disallowed. Moreover, RDMA\_{RMW}^{SC} RMW operations are strongly isolated with RDMA\_{RMW}^{SC} reads and writes; e.g. outcome  $\mathbf{x}=2$  is disallowed in Fig. 15c. This is in contrast to remote RMWs of RDMA\_{RMW}^{WAIT}, where outcome x=2 is allowed in Fig. 6b. However, RDMA\_{RMW}^{SC} operations does not ensure that earlier remote operations by *other libraries* are completed, and thus outcome (a,b)=(1,0) is allowed in both Figs. 15d and 15e.

More generally, we can use this strategy to linearise [19] accesses to any sequential data structure D by wrapping each call to D inside a node lock. This allows us to port existing sequential data structures to RDMA settings with minimal effort. Finally, we prove (Theorem 4) that our implementation is correct against the RDMA $_{\rm RMW}^{\rm SC}$  specification with the full proof given in §B.5.

**Theorem 4.** The implementation  $I_{SC}$  is sound.

#### 6 Related Work

RDMA Semantics. The coreRMA model [13] is an early attempt at formalising remote memory accesses, but this semantics does not match the RDMA technical specification. This gap is addressed by RDMA<sup>TSO</sup> [3], which formalises the actual RDMA semantics over TSO, but the formalisation did not cover remote RMWs. A later model, RDMA<sup>SC</sup> [5], explored the semantics from RDMA<sup>TSO</sup> [3] but over an SC CPU alongside programming strategies to efficiently prevent weak behaviours. RDMA<sup>SC</sup> is unrelated to our work, including RDMA<sup>SC</sup><sub>RMW</sub>.

**RDMA-Based Distributed Systems.** Besides LOCO [4, 20], prior work has covered a range of distributed systems, e.g. consensus protocols [1], databases [2, 24], stand-alone data structures [9, 14]. However, unlike LOCO (and our work), these are bespoke systems rather than a programming methodology or library.

**Verification.** Our proofs for the soundness of library implementations have followed the declarative style [4, 27, 31]. For RDMA<sup>TSO</sup><sub>RMW</sub> (like RDMA<sup>TSO</sup>), we also provide an operational model (§D) which could ultimately form a basis for a program logic (e.g., [7, 22]), ultimately enabling operational abstractions and proofs of refinement [12, 30]. We consider such extensions to be future work.

RDMA Locks. There are several implementations of network locks using RDMA operations, including centralised lock managers [11], decentralised algorithms [33], asymmetric implementations to favour local accesses [6], and technology-agnostic designs that are more general than RDMA [15]. Other stated objectives of these implementations can include fairness, starvation freedom, low latency, load balancing, scalability, contention mitigation, fault tolerance [18], etc.

However, none of these existing implementations have been formally verified (since RDMA<sub>RMW</sub> is the first formal semantics of remote RMWs). These works, at most, have offered intuitive explanations to support the correctness of their approach. Moreover, these implementations lack an explicit description of the interaction guarantees between locks and other RDMA operations, which as we have seen can be subtle. In most cases, programmers are made responsible to ensure relevant operations are completed before releasing the lock, thus aligning with the weak lock semantics that we have presented.

#### References

- Aguilera, M.K., Ben-David, N., Guerraoui, R., Marathe, V.J., Xygkis, A., Zablotchi, I.: Microsecond consensus for microsecond applications. In: 14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20). pp. 599-616. USENIX Association (Nov 2020), https://www.usenix.org/ conference/osdi20/presentation/aguilera
- Alquraan, A., Udayashankar, S., Marathe, V., Wong, B., Al-Kiswany, S.: Lolkv: the logless, line the logless, linearizable, rdma-based key-value storage system arizable, rdma-based key-value storage system. In: Proceedings of the 21st USENIX Symposium on Networked Systems Design and Implementation. NSDI'24, USENIX Association, USA (2024)
- 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), 1982–2009 (2024). https://doi.org/10.1145/3689781, https://doi.org/10.1145/3689781
- Ambal, G., Hodgkins, G., Madler, M., Chockler, G., Dongol, B., Izraelevitz, J., Raad, A., Vafeiadis, V.: A verified high-performance composable object library for remote direct memory access (extended version) (2025), https://arxiv.org/abs/ 2510.10531
- Ambal, G., Lahav, O., Raad, A.: Sufficient conditions for robustness of RDMA programs. In: Vafeiadis, V. (ed.) Programming Languages and Systems 34th European Symposium on Programming, ESOP 2025, Held as Part of the International Joint Conferences on Theory and Practice of Software, ETAPS 2025, Hamilton, ON, Canada, May 3-8, 2025, Proceedings, Part I. Lecture Notes in Computer Science, vol. 15694, pp. 56-87. Springer (2025). https://doi.org/10.1007/978-3-031-91118-7\_3, https://doi.org/10.1007/978-3-031-91118-7\_3
- Baran, A., Nelson-Slivon, J., Tseng, L., Palmieri, R.: Alock: Asymmetric lock primitive for rdma systems. In: Proceedings of the 36th ACM Symposium on Parallelism in Algorithms and Architectures. p. 15–26. SPAA '24, Association for Computing Machinery, New York, NY, USA (2024). https://doi.org/10.1145/3626183.3659977, https://doi.org/10.1145/3626183.3659977
- 7. 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)
- 8. Blundell, C., Lewis, E.C., Martin, M.M.: Subtleties of transactional memory atomicity semantics. IEEE Computer Architecture Letters 5(2), 17–17 (2006). https://doi.org/10.1109/L-CA.2006.18
- Brock, B., Buluç, A., Yelick, K.: Bcl: A cross-platform distributed data structures library. In: Proceedings of the 48th International Conference on Parallel Processing. ICPP '19, Association for Computing Machinery, New York, NY, USA (2019). https://doi.org/10.1145/3337821.3337912, https://doi.org/10.1145/3337821.3337912
- 10. Chong, N., Sorensen, T., Wickerson, J.: The semantics of transactions and weak memory in x86, power, arm, and C++. In: Foster, J.S., Grossman, D. (eds.) Proceedings of the 39th ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI 2018, Philadelphia, PA, USA, June 18-22, 2018. pp. 211–225. ACM (2018). https://doi.org/10.1145/3192366.3192373, https://doi.org/10.1145/3192366.3192373

- 11. Chung, Y., Zamanian, E.: Using RDMA for lock management. CoRR abs/1507.03274 (2015), http://arxiv.org/abs/1507.03274
- 12. Dalvandi, S., Dongol, B.: Implementing and verifying release-acquire transactional memory in C11. Proc. ACM Program. Lang. 6(OOPSLA2), 1817–1844 (2022). https://doi.org/10.1145/3563352, https://doi.org/10.1145/3563352
- 13. 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
- Devarajan, H., Kougkas, A., Bateman, K., Sun, X.H.: Hcl: Distributing parallel data structures in extreme scales. In: 2020 IEEE International Conference on Cluster Computing (CLUSTER). pp. 248–258 (2020). https://doi.org/10.1109/CLUSTER49012.2020.00035
- Devulapalli, A., Wyckoff, P.: Distributed queue-based locking using advanced network features. In: 34th International Conference on Parallel Processing (ICPP 2005), 14-17 June 2005, Oslo, Norway. pp. 408-415. IEEE Computer Society (2005). https://doi.org/10.1109/ICPP.2005.34, https://doi.org/10. 1109/ICPP.2005.34
- Dongol, B., Jagadeesan, R., Riely, J.: Transactions in relaxed memory architectures. Proc. ACM Program. Lang. 2(POPL), 18:1–18:29 (2018). https://doi.org/10.1145/3158106, https://doi.org/10.1145/3158106
- Gangidi, A., Miao, R., Zheng, S., Bondu, S.J., Goes, G., Morsy, H., Puri, R., Riftadi, M., Shetty, A.J., Yang, J., et al.: Rdma over ethernet for distributed training at meta scale. In: Proceedings of the ACM SIGCOMM 2024 Conference. pp. 57–70 (2024)
- 18. Gao, J., Wang, Q., Shu, J.: Shiftlock: Mitigate one-sided RDMA lock contention via handover. In: Gunawi, H.S., Tarasov, V. (eds.) 23rd USENIX Conference on File and Storage Technologies, FAST 2025, Santa Clara, CA, February 25-27, 2025. pp. 355-372. USENIX Association (2025), https://www.usenix.org/conference/fast25/presentation/gao
- 19. Herlihy, M., Wing, J.M.: Linearizability: A correctness condition for concurrent objects. ACM Trans. Program. Lang. Syst. 12(3), 463–492 (1990). https://doi.org/10.1145/78969.78972, https://doi.org/10.1145/78969.78972
- 20. Hodgkins, G., Madler, M., Izraelevitz, J.: Loco: Rethinking objects for network memory (2025), https://arxiv.org/abs/2503.19270
- 21. IBTA: Infiniband architecture specification volume 1 release 1.6. https://www.infinibandta.org/ibta-specification/ (2022)
- 22. Lahav, O., Dongol, B., Wehrheim, H.: Rely-guarantee reasoning for causally consistent shared memory. In: Enea, C., Lal, A. (eds.) Computer Aided Verification 35th International Conference, CAV 2023, Paris, France, July 17-22, 2023, Proceedings, Part I. Lecture Notes in Computer Science, vol. 13964, pp. 206–229. Springer (2023). https://doi.org/10.1007/978-3-031-37706-8\_11, https://doi.org/10.1007/978-3-031-37706-8\_11
- 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
- Li, P., Hua, Y., Zuo, P., Chen, Z., Sheng, J.: ROLEX: A scalable RDMA-oriented learned Key-Value store for disaggregated memory systems. In: 21st USENIX Conference on File and Storage Technologies (FAST 23). pp. 99–114. USENIX Associa-

- tion, Santa Clara, CA (Feb 2023), https://www.usenix.org/conference/fast23/presentation/li-pengfei
- 25. Lu, Y., Chen, G., Li, B., Tan, K., Xiong, Y., Cheng, P., Zhang, J., Chen, E., Moscibroda, T.: {Multi-Path} transport for {RDMA} in datacenters. In: 15th USENIX symposium on networked systems design and implementation (NSDI 18). pp. 357–371 (2018)
- Owens, S., Sarkar, S., Sewell, P.: A better x86 memory model: x86-tso. In: Berghofer, S., Nipkow, T., Urban, C., Wenzel, M. (eds.) Theorem Proving in Higher Order Logics, 22nd International Conference, TPHOLs 2009, Munich, Germany, August 17-20, 2009. Proceedings. Lecture Notes in Computer Science, vol. 5674, pp. 391-407. Springer (2009). https://doi.org/10.1007/978-3-642-03359-9\_ 27, https://doi.org/10.1007/978-3-642-03359-9\_27
- Raad, A., Doko, M., Rozic, L., Lahav, O., Vafeiadis, V.: On library correctness under weak memory consistency: specifying and verifying concurrent libraries under declarative consistency models. Proc. ACM Program. Lang. 3(POPL), 68:1–68:31 (2019). https://doi.org/10.1145/3290381, https://doi.org/10.1145/3290381
- 28. 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)
- 29. 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)
- Singh, A.K., Lahav, O.: An operational approach to library abstraction under relaxed memory concurrency. Proc. ACM Program. Lang. 7(POPL), 1542–1572 (2023). https://doi.org/10.1145/3571246, https://doi.org/10.1145/3571246
- 31. Stefanesco, L., Raad, A., Vafeiadis, V.: Specifying and verifying persistent libraries. In: Weirich, S. (ed.) Programming Languages and Systems 33rd European Symposium on Programming, ESOP 2024, Held as Part of the European Joint Conferences on Theory and Practice of Software, ETAPS 2024, Luxembourg City, Luxembourg, April 6-11, 2024, Proceedings, Part II. Lecture Notes in Computer Science, vol. 14577, pp. 185–211. Springer (2024). https://doi.org/10.1007/978-3-031-57267-8\_8, https://doi.org/10.1007/978-3-031-57267-8\_8
- 32. Wang, Z., Luo, L., Ning, Q., Zeng, C., Li, W., Wan, X., Xie, P., Feng, T., Cheng, K., Geng, X., et al.: {SRNIC}: A scalable architecture for {RDMA}{NICs}. In: 20th USENIX Symposium on Networked Systems Design and Implementation (NSDI 23). pp. 1–14 (2023)
- 33. Yoon, D.Y., Chowdhury, M., Mozafari, B.: Distributed lock management with RDMA: decentralization without starvation. In: Das, G., Jermaine, C.M., Bernstein, P.A. (eds.) Proceedings of the 2018 International Conference on Management of Data, SIGMOD Conference 2018, Houston, TX, USA, June 10-15, 2018. pp. 1571–1586. ACM (2018). https://doi.org/10.1145/3183713.3196890, https://doi.org/10.1145/3183713.3196890
- 34. 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. ACM SIGCOMM Computer Communication Review 45(4), 523–536 (2015)

### A SV Library Semantics

As per other RDMA libraries, we assume a set of nodes, Node, of fixed size. Each thread t is associated to a node  $\mathbf{n}(t)$ . The sv library uses the following 5 methods:

```
\begin{split} m(\widetilde{v}) &::= \text{ Write}_{\text{SV}}(x,v) \mid \text{Read}_{\text{SV}}(x) \mid \text{Bcast}_{\text{SV}}(x,d,\{n_1,\ldots,n_k\}) \\ &\mid \text{Wait}_{\text{SV}}(d) \mid \text{GFence}(\{n_1,\ldots,n_k\}) \end{split}
```

```
 \begin{array}{ll} \bullet \  \, \mathsf{Write}_{\mathrm{SV}} : \mathsf{Loc} \times \mathsf{Val} \to () \\ \bullet \  \, \mathsf{Read}_{\mathrm{SV}} : \mathsf{Loc} \to \mathsf{Val} \\ \end{array} \quad \begin{array}{ll} \bullet \  \, \mathsf{Wait}_{\mathrm{SV}} : \mathsf{Wid} \to () \\ \bullet \  \, \mathsf{GFence} : \mathcal{P}(\mathsf{Node}) \to () \\ \end{array}
```

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

We then require the stamping function stmp<sub>sv</sub>:

```
\begin{split} \operatorname{stmp}_{\operatorname{SV}}(\langle {\scriptscriptstyle{-}},{\scriptscriptstyle{-}},\langle\operatorname{Write}_{\operatorname{SV}},{\scriptscriptstyle{-}},{\scriptscriptstyle{-}}\rangle\rangle) &= \{\operatorname{aCW}\} \\ \operatorname{stmp}_{\operatorname{SV}}(\langle {\scriptscriptstyle{-}},{\scriptscriptstyle{-}},\langle\operatorname{Read}_{\operatorname{SV}},{\scriptscriptstyle{-}},{\scriptscriptstyle{-}}\rangle\rangle) &= \{\operatorname{aCR}\} \\ \operatorname{stmp}_{\operatorname{SV}}(\langle {\scriptscriptstyle{-}},{\scriptscriptstyle{-}},\langle\operatorname{Wait}_{\operatorname{SV}},{\scriptscriptstyle{-}},{\scriptscriptstyle{-}}\rangle\rangle) &= \{\operatorname{aWT}\} \\ \operatorname{stmp}_{\operatorname{SV}}(\langle {\scriptscriptstyle{-}},{\scriptscriptstyle{-}},\langle\operatorname{GFence},(\{n_1,\ldots,n_k\}),{\scriptscriptstyle{-}}\rangle\rangle) &= \{\operatorname{aGF}_{n_1},\ldots,\operatorname{aGF}_{n_k}\} \\ \operatorname{stmp}_{\operatorname{SV}}(\langle {\scriptscriptstyle{-}},{\scriptscriptstyle{-}},\langle\operatorname{Bcast}_{\operatorname{SV}},({\scriptscriptstyle{-}},{\scriptscriptstyle{-}},\{n_1,\ldots,n_k\}),{\scriptscriptstyle{-}}\rangle\rangle) &= \{\operatorname{aNLR}_{n_1},\operatorname{aNRW}_{n_1},\ldots,\operatorname{aNLR}_{n_k},\operatorname{aNRW}_{n_k}\} \end{split}
```

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

With this, the stamp order is enough to enforce the behaviour of the global fence. If we have a program  $\mathtt{Bcast}_{\mathtt{SV}}(x,d,\{\ldots,n,\ldots\}); \mathtt{GFence}(\{\ldots,n,\ldots\}),$  the plain execution has two events  $\mathtt{e_{BR}}$  and  $\mathtt{e_{GF}}$ , and the definitions of  $\mathtt{stmp_{SV}}$  and  $\mathtt{sto}$  (cell G12 in Fig. 9) imply  $\langle \mathtt{e_{BR}},\mathtt{aNRW}_n \rangle \xrightarrow{\mathtt{ppo}} \langle \mathtt{e_{GF}},\mathtt{aGF}_n \rangle$ .

Recall that, for an execution  $\mathcal{G}$ ,  $\mathcal{G}.\mathcal{W}$  represents write subevents (stamps aCW, aCAS, aNLW, and aNRW), while  $\mathcal{G}.\mathcal{R}$  represents read subevents (stamps aCR, aCAS, aNLR, aNAR, and aNRR). Recall also that we note e.g.  $\mathcal{G}.\mathcal{W}_x \triangleq \{\mathbf{s} \in \mathcal{G}.\mathcal{W} \mid \mathtt{loc}(\mathbf{s}) = \{x\}\}$  to constrain a set to subevents on a specific location x. For the SV library, we additionally define  $\mathcal{G}.\mathcal{W}^n \triangleq \{\langle \mathbf{e}, \mathbf{aCW} \rangle \mid \mathbf{n}(\mathbf{t}(\mathbf{e})) = n\} \cup \mathcal{G}.\mathbf{aNRW}_n$  as the set of write subevents occurring on node n. This includes CPU writes on the node, as well as broadcast writes towards n from all threads. We also note  $\mathcal{G}.\mathcal{W}_x^n \triangleq \mathcal{G}.\mathcal{W}_x \cap \mathcal{G}.\mathcal{W}^n$ 

as expected. Similarly,  $\mathcal{G}.\mathcal{R}^n \triangleq \{ s \mid s \in \mathcal{G}.\mathcal{R} \land n(t(s)) = n \}$  covers reads occurring on n, either by a CPU read or as part of a broadcast. We now work towards a definition of consistency for shared variables.

**Definition 9.** For an execution  $\mathcal{G} = \langle E, po, stmp_{sy}, -, - \rangle$ , we define the following:

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

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

We further define the *reads-from-internal* relation as  $rf_i \triangleq [aCW]$ ;  $(po\cap rf)$ ; [aCR] (which corresponds to CPU reads and writes using the same TSO store buffer), and the *reads-from-external* relation as  $rf_e \triangleq rf \setminus rf_i$ . Moreover, given an execution  $\mathcal G$  and well-formed  $rf_e$  and  $rf_e$  defined and  $rf_e$  and  $rf_e$  defined  $rf_e$  defined  $rf_e$  and  $rf_e$  defined  $rf_e$ 

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

We can then define the consistency predicate SV.C as follows.

**Definition 10** (SV-consistency).  $\langle E, po, stmp, so, hb \rangle$  is SV-consistent if:

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

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

#### **B** Correctness Proofs

Correctness proofs of the MOWGLI framework can be found in [4]. We recall the main definitions and results in Appendix B.1 before proving the soundness of the weak lock (§B.2), strong lock (§B.3), node lock (§B.4), and RDMA<sup>SC</sup><sub>RMW</sub> libraries (§B.5).

#### B.1 Background: MOWGLI Definitions and Results

MOWGLI assumes a type Val of values, a type  $Loc \subseteq Val$  of locations, and a type Method of methods. The syntax of sequential programs is given by the following grammar:

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

MOWGLI assumes top-level concurrency, i.e. there is a fixed set of threads  $\mathsf{Tid} \triangleq \{1, 2, \dots, T\}$ , and a concurrent program is given by a tuple  $\widetilde{\mathsf{p}} = \langle \mathsf{p}_1, \dots, \mathsf{p}_T \rangle$ , where each thread t corresponds to a program  $\mathsf{p}_t \in \mathsf{SeqProg}$ .

The semantics of a program is given by an execution, which is a graph over events. Recall that events are defined in Definition 17. The first two components  $\langle E, po \rangle$  of an execution form a *plain execution*:

**Definition 11.** We say that  $\langle E, po \rangle$  is a plain execution iff  $E \subseteq \text{Event}$ ,  $po \subseteq E \times E$ , and  $po = \bigcup_{t \in \text{Tid}} po|_t$  where every  $po|_t$  (i.e. po restricted to the events of thread t) is a total order.

We write  $\emptyset_G \triangleq \langle \emptyset, \emptyset \rangle$  for the empty execution and  $\{e\}_G \triangleq \langle \{e\}, \emptyset \rangle$  for the execution with a single event e. Given two executions,  $G_1 = \langle E_1, \mathsf{po}_1 \rangle$  and  $G_2 = \langle E_2, \mathsf{po}_2 \rangle$ , with disjoint sets of events (i.e.  $E_1 \cap E_2 = \emptyset$ ), we define their sequential composition  $G_1$ ;  $G_2$  and parallel composition  $G_1 \parallel G_2$  as follows:

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

The plain semantics of a program p executed by a thread t is given by  $[\![p]\!]_t$ , which is a set of pairs of the form  $\langle r, G \rangle$ , where r is the output and G is a plain execution. This set represents all conceivable unfoldings of the program into method calls, even those that will be rejected by the semantics of the corresponding libraries. Each output is a pair  $\langle v, k \rangle$ , where v is a value and k a break

number, indicating the program terminates by requesting to exit k nested loops and returning the value v.

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

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

$$[\![\tilde{\mathbf{p}}]\!] \triangleq \{ \langle \langle v_1, \dots, v_T \rangle, \|_{t \in \mathsf{Tid}} \ G_t \rangle \ \big| \ \forall t \in \mathsf{Tid}. \langle \langle v_t, 0 \rangle, G_t \rangle \in [\![\mathbf{p}_t]\!]_t \}$$

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

Then, we can define executions (Def. 1), libraries (§3.1), and consistent executions (Def. 2).

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

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

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

**Definition 12.** We say that I is well defined for a library L using  $\Lambda$  iff for all  $t \in \mathsf{Tid}$ ,  $m \in L.M$  and  $\widetilde{v} \in \mathsf{Val}^*$ , we have:

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

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

$$\|v\|_{t,I} \triangleq v \qquad \qquad \|m(v_1,\ldots,v_k)\|_{t,I} \triangleq \begin{cases} I(t,m,\langle v_1,\ldots,v_k\rangle) & \text{if } m \in L.M \\ m(v_1,\ldots,v_k) & \text{otherwise} \end{cases}$$
 
$$\|\log p\|_{t,I} \triangleq \log \|p\|_{t,I} \qquad \| \|p\|_{t,I} \triangleq \|p\|_{t,I} (\lambda v.\|f v\|_{t,I})$$
 
$$\|\log p\|_{t,I} \triangleq \log p\|_{t,I} \triangleq \|p\|_{t,I} (\lambda v.\|f v\|_{t,I})$$
 
$$\|\log p\|_{t,I} \triangleq \|p\|_{t,I} \|\log p\|_{t,I}$$

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

**Definition 13.** We say that I is a sound implementation of L using  $\Lambda$  if, for any program  $\widetilde{p}$  such that  $loc(I) \cap loc(\widetilde{p}) = \emptyset$ , we have that  $outcome_{\Lambda}(\llbracket \widetilde{p} \rrbracket_I) \subseteq outcome_{\Lambda \uplus \{L\}}(\widetilde{p})$ .

As soundness is difficult to prove directly, MOWGLI develops a modular proof technique using an abstraction function mapping the implementation to its abstract library specification. For  $f:A\to B$  and  $r\subseteq A\times A$ , we note  $f(r)\triangleq\{\langle f(x),f(y)\rangle\mid \langle x,y\rangle\in r\}.$ 

**Definition 14.** Suppose I is a well-defined implementation of a library L using  $\Lambda$ , and that  $G = \langle E, po \rangle$  and  $G' = \langle E', po' \rangle$  are plain executions using methods of  $\Lambda$  and L respectively. We say that a surjective function  $f: E \to E'$  abstracts G to G', denoted  $\mathtt{abs}_{I,L}^f(G,G')$ , iff

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

**Lemma 1.** Given  $\widetilde{\mathbf{p}}$  on library L and a well-defined implementation I of L, if  $\langle \widetilde{v}, G \rangle \in \llbracket \| \widetilde{\mathbf{p}} \|_{I} \rrbracket$  then there exists  $\langle \widetilde{v}, G' \rangle \in \llbracket \widetilde{\mathbf{p}} \rrbracket$  and f such that  $\mathtt{abs}_{LL}^f(G, G')$ .

**Definition 15.** We say that a well defined implementation I of a library L is locally sound iff, whenever we have a  $\Lambda$ -consistent execution  $\mathcal{G} = \langle E, \mathsf{po}, \mathsf{stmp}, \mathsf{so}, \mathsf{hb} \rangle$  and  $\mathsf{abs}_{I,L}^f(\langle E, \mathsf{po} \rangle, \langle E', \mathsf{po'} \rangle)$ , then there exists  $\mathsf{stmp'}, \mathsf{so'}$ , and a concretisation function  $g : \langle E', \mathsf{po'}, \mathsf{stmp'} \rangle$ . SEvent  $\to \mathcal{G}$ . SEvent such that:

- $g(\langle e', a' \rangle) = \langle e, a \rangle$  implies f(e) = e' and
  - For all  $a_0$  such that  $\langle a_0, a' \rangle \in \mathsf{sto}$ , there exists  $\langle \mathsf{e}_1, a_1 \rangle \in \mathcal{G}$ . SEvent such that  $f(\mathsf{e}_1) = \mathsf{e}', \langle a_0, a_1 \rangle \in \mathsf{sto}, \ and \ \langle \langle \mathsf{e}_1, a_1 \rangle, \langle \mathsf{e}, a \rangle \rangle \in \mathsf{hb}^*;$
  - For all  $a_0$  such that  $\langle a', a_0 \rangle \in \mathsf{sto}$ , there exists  $\langle \mathsf{e}_2, a_2 \rangle \in \mathcal{G}$ . SEvent such that  $f(\mathsf{e}_2) = \mathsf{e}', \langle a_2, a_0 \rangle \in \mathsf{sto}$ , and  $\langle \langle \mathsf{e}, a \rangle, \langle \mathsf{e}_2, a_2 \rangle \rangle \in \mathsf{hb}^*$ .
- $g(so') \subseteq hb$ ;
- For all hb' transitive such that  $(ppo' \cup so')^+ \subseteq hb'$  and  $g(hb') \subseteq hb$ , we have  $\langle E', po', stmp', so', hb' \rangle \in L.C$ , where  $ppo' \triangleq \langle E', po', stmp' \rangle$ .ppo.

**Theorem 5.** If a well-defined implementation is locally sound, then it is sound.

We show the local soundness of the different implementations given in this paper. Thus, from the theorem above, these implementations are sound.

#### B.2 WLOCK Library

Theorem 1. The implementation  $I_{\mathtt{WL}}$  is sound.

Proof. We assume an  $\{\text{RDMA}_{\text{RMW}}^{\text{WAIT}}, \text{SV}\}$ -consistent execution  $\mathcal{G} = \langle E, \text{po}, \text{stmp}, \text{so}, \text{hb} \rangle$  which is abstracted via f to  $\langle E', \text{po'} \rangle$  that uses (only) the WLOCK library, i.e.  $\text{abs}_{I_{\text{WLOCK}}}^f(\langle E, \text{po} \rangle, \langle E', \text{po'} \rangle)$  holds. We need to provide stmp', so', and g:

 $\langle E', \mathsf{po'}, \mathsf{stmp'} \rangle$ . SEvent  $\to \mathcal{G}$ . SEvent respecting some conditions. From  $\langle E', \mathsf{po'} \rangle$ , we simply take  $\mathsf{stmp'} = \mathsf{stmp}_{\mathsf{WL}}$ .

Since  $\mathcal{G}$  is  $\{RDMA_{RMW}^{WAIT}, SV\}$ -consistent, it means  $(ppo \cup so|_{RDMA_{RMW}^{WAIT}} \cup so|_{SV}) \subseteq hb$ , hb is transitive and irreflexive, and the two restrictions of  $\mathcal{G}$  are respectively  $RDMA_{RMW}^{WAIT}$ -consistent and SV-consistent.

 $\begin{aligned} & \text{RDMA}_{\text{RMW}}^{\text{WAIT}}\text{-consistency implies there is some well-formed } v_{\text{R}}, \ v_{\text{W}}, \ \text{rf}, \ \underset{\text{no}}{\text{mo}}, \ \text{nfo}, \\ & \text{and rao such that ib is irreflexive}, \forall \text{e.stmp}|_{\text{RDMA}_{\text{RMW}}^{\text{WAIT}}}(\textbf{e}) \in \text{stmp}_{\text{RW}}(\textbf{e}), \ \text{and} \ \underset{\text{so}}{\text{so}}|_{\text{RDMA}_{\text{RMW}}^{\text{WAIT}}} = \\ & \text{iso} \cup \text{rf}_{\textbf{e}} \cup \text{pfg} \cup \text{nfo} \cup \text{rb} \cup \text{mo} \cup \text{rao} \cup ([\texttt{aNRW}_n]; \texttt{iso}^{-1}; \texttt{rao}) \cup ([\texttt{Inst}]; \texttt{ib}). \end{aligned}$ 

SV-consistency implies there is some well-formed  $\mathbf{v}_{R}'', \mathbf{v}_{W}'', \mathbf{r}f'',$ and  $\mathbf{mo}'',$ such that  $\mathtt{stmp}|_{\mathrm{SV}} = \mathtt{stmp}_{\mathrm{SV}}, \ [\mathtt{aCR}]; (\mathtt{po}|_{\mathrm{SV}}^{-1} \cap \mathtt{rb}''); \ [\mathtt{aCW}] = \emptyset,$ and  $\mathtt{so}|_{\mathrm{SV}} = \mathtt{iso}'' \cup \mathtt{rf}_{\mathrm{e}}'' \cup \mathtt{pf}'' \cup \mathtt{rb}'' \cup \mathtt{mo}''.$ (We will use double apostrophes for references to the SV library.) We define g as follows.

- For an event  $e' = (t, -, (Acq_{WL}, (x), ()))$ , we choose  $g(e', aMF) = (e_r, aCR)$  with  $e_r = (t, -, (Read_{SV}, (x_{t'}), (v))) \in f^{-1}(e')$  the last event of the implementation (reading a shared variable owned by some thread t').
- For an event  $\mathbf{e}' = (t, \neg, (\mathtt{Rel}_{\mathrm{WL}}, (x), ()))$ , we choose  $g(\mathbf{e}', \mathtt{aCW}) = (\mathbf{e}_w, \mathtt{aCW})$  with  $\mathbf{e}_w = (t, \neg, (\mathtt{Write}_{\mathrm{SV}}, (x_t, v+1), ())) \in f^{-1}(\mathbf{e}')$  the second event of the implementation.

First, let us show that g preserves  $\operatorname{sto}$  (first property of local soundness). For  $\operatorname{Rel}_{\operatorname{WL}}$  this is trivial using the identity function. For  $\operatorname{Acq}_{\operatorname{WL}}$ , the stamp  $\operatorname{aCR}$  is similar to  $\operatorname{aMF}$  w.r.t. later stamps, so  $(\mathsf{e}_2,a_2)=(\mathsf{e}_r,\operatorname{aCR})$  is enough. For an earlier stamp  $a_0$  such that  $(a_0,\operatorname{aMF})\in\operatorname{sto}$ , we take  $(\mathsf{e}_1,a_1)=((t,\_,(\operatorname{RFAA},(\ldots,d),())),\operatorname{aNLW}_n)$  the first event of the implementation, and with  $\mathsf{e}_{wt}=(t,\_,(\operatorname{Wait},(d),()))$  the second event we have  $(\mathsf{e}_1,a_1) \xrightarrow{\operatorname{pfg}} (\mathsf{e}_{wt},\operatorname{aWT}) \xrightarrow{\operatorname{ppo}} (\mathsf{e}_r,\operatorname{aCR})$  (thus included in  $\operatorname{hb}$ ) with  $(a_0,\operatorname{aNLW}_n)\in\operatorname{sto}$ .

Now we need to pick a suitable so' such that  $g(\mathbf{so'}) \subseteq \mathsf{hb}$  and  $\langle E', \mathsf{po'}, \mathsf{stmp'}, \mathsf{so'}, \_\rangle$  is WLOCK-consistent. We can assume that  $\langle E', \mathsf{po'} \rangle$  respects locks, as otherwise  $\mathsf{so'} = \emptyset$  is enough. Thus, for each location x we need to define a total order  $\mathsf{lo'}_x$  on  $A'_x \triangleq \{\mathsf{e'} \mid \mathsf{e'} \in E'_x \land \mathsf{m}(\mathsf{e'}) = \mathsf{Acq}_{\mathsf{WL}} \}$ . Each event  $\mathsf{e'} \in A'_x$  can be associated to its first subevent of the form  $((t', \_, (\mathsf{RFAA}, (p_x^{t'}, x_a, 1, d), ())), \mathsf{aNAR}_n)$ , with  $n = \mathsf{n}(x)$ . From  $\mathsf{RDMA}^{\mathsf{WAIT}}_{\mathsf{RMW}}$ -consistency, rao induces a total ordering on these subevents, and we simply keep the same ordering for  $A'_x$ . As such, we define  $\mathsf{so'} = \bigcup_x \big\{ \langle \mathsf{e'}_1, \mathsf{aCW} \rangle, \langle \mathsf{e'}_2, \mathsf{aMF} \rangle \ \big| \ (\mathsf{e'}_1, \mathsf{e'}_2) \in (\mathsf{po'}_x|_{\mathsf{imm}})^{-1}; \mathsf{lo'}_x \big\}$  as expected, and we have that  $\langle E', \mathsf{po'}, \mathsf{stmp'}, \mathsf{so'}, \_ \rangle$  is WLOCK-consistent.

Thus, the rest of the proof is to show that  $g(\mathbf{so'}) \subseteq \mathsf{hb}$ , i.e. that the synchronisations promised by the WLOCK library are enforced in the implementation. We can assume  $(\mathsf{e'}_0, \mathsf{aMF}) \xrightarrow{\mathsf{lo'}_x} (\mathsf{e'}_2, \mathsf{aMF})$  and  $(\mathsf{e'}_0, \mathsf{aMF}) \xrightarrow{\mathsf{po'}_x|_{\mathrm{imm}}} (\mathsf{e'}_1, \mathsf{aCW})$ , with  $\mathsf{e'}_0$  running  $\mathsf{Acq}_{\mathsf{WL}}(x)$  by thread  $t_1$ ,  $\mathsf{e'}_1$  running  $\mathsf{Rel}_{\mathsf{WL}}(x)$  by thread  $t_1$ , and  $\mathsf{e'}_2$  running  $\mathsf{Acq}_{\mathsf{WL}}(x)$  by thread  $t_2$ . We also note  $(\mathsf{e}_1, \mathsf{aCW}) = g(\mathsf{e'}_1, \mathsf{aCW})$  and  $(\mathsf{e}_2, \mathsf{aCR}) = g(\mathsf{e'}_2, \mathsf{aMF})$ . Our goal is then to show  $(\mathsf{e}_1, \mathsf{aCW}) \xrightarrow{\mathsf{hb}} (\mathsf{e}_2, \mathsf{aCR})$ .

We proceed by induction on the ordering  $|o'_x|$ . The base case is for  $(e'_0, aMF) \xrightarrow{|o'_x|_{\text{limm}}} (e'_2, aMF)$ . This base case trivially implies the general case by transitivity, since

the program respects locks (i.e. intermediate acquires are being released) and  $(aCR, aCW) \in sto$ .

Let  $\mathbf{e}_0^{faa} = (t_1, \neg, (\mathsf{RFAA}, (p_x^{t_1}, x_a, 1, d), ()))$  be the FAA in the implementation of  $\mathbf{e}_0'$  and  $\mathbf{e}_2^{faa} = (t_2, \neg, (\mathsf{RFAA}, (p_x^{t_2}, x_a, 1, d), ()))$  in the implementation of  $\mathbf{e}_2'$ . By definition we have  $(\mathbf{e}_0^{faa}, \mathsf{aNAR}_n) \xrightarrow{(\mathsf{rao}|_{E_{x_a}})|_{\mathrm{imm}}} (\mathbf{e}_2^{faa}, \mathsf{aNAR}_n)$ , since any remote RMW in  $E_{x_a}$  is from an implementation of some  $\mathsf{Acq}_{\mathsf{WL}}(x)$  event. From the semantics of  $\mathsf{RDMA}_{\mathsf{RMW}}^{\mathsf{WAIT}}$  we have  $(\mathbf{e}_0^{faa}, \mathsf{aNRW}_n) \xrightarrow{\mathsf{hb}} (\mathbf{e}_2^{faa}, \mathsf{aNAR}_n)$  (from the  $([\mathsf{aNRW}]; \mathsf{iso}^{-1}; \mathsf{rao})$  component), and thus we necessarily have  $(\mathbf{e}_0^{faa}, \mathsf{aNRW}_n) \xrightarrow{\mathsf{rf}} (\mathbf{e}_2^{faa}, \mathsf{aNAR}_n)$ , i.e. the second FAA reads the modified value of the first. This is because  $\mathbf{e}_2^{faa}$  cannot read from an earlier write (or the initial value of 0) as that would imply an  $\mathsf{rb}$  dependency and an  $\mathsf{hb}$  cycle; and cannot read  $(\mathsf{rf}_{\mathsf{e}} \subseteq \mathsf{hb})$  from a later write, as any later write is  $\mathsf{hb}$  after  $\mathbf{e}_2^{faa}$  (via rao and  $\mathsf{ppo}$ ).

There is some value  $v_0 = \mathsf{v}_\mathsf{R}((\mathsf{e}_0^{faa}, \mathsf{aNAR}_n))$  read by the first FAA operation. By well-formedness of  $\mathsf{v}_\mathsf{R}$ ,  $\mathsf{v}_\mathsf{W}$ , and rf, we have  $\mathsf{v}_\mathsf{R}((\mathsf{e}_2^{faa}, \mathsf{aNAR}_n)) = \mathsf{v}_\mathsf{W}((\mathsf{e}_0^{faa}, \mathsf{aNAR}_n)) = v_0 + 1$ , i.e. the following  $\mathsf{Acq}_\mathsf{WL}(x)$  gets the next ticket. More generally, it is clear every  $\mathsf{Acq}_\mathsf{WL}(x)$  gets a different ticket. We also have  $\mathsf{v}_\mathsf{W}((\mathsf{e}_0^{faa}, \mathsf{aNLW}_n)) = v_0$ , i.e.  $p_x^{t_1}$  is modified to contain  $v_0$ . Respectively  $p_x^{t_2}$  is modified to contain  $v_0 + 1$ .

Let  $\mathbf{e}_0^r$  be the third event of the implementation of  $\mathbf{e}_0'$  reading  $p_x^{t_1}$ . We necessarily have  $(\mathbf{e}_0^{faa}, \mathtt{aNLW}_n) \stackrel{\mathsf{rf}}{\to} (\mathbf{e}_0^r, \mathtt{aCR})$ . This is because  $\mathbf{e}_0^r$  cannot read from the future (it would create an  $\mathsf{rf}; \mathsf{ippo}$  cycle in  $\mathsf{ib}$ ) and the second event  $\mathsf{Wait}(d)$  makes sure all previous modifications of  $p_x^{t_1}$  are available (ignoring the last one would be an  $\mathsf{rb}; \mathsf{hb}$  cycle). Thus, in the implementation of  $\mathsf{e}_0'$ , the meta-variable  $\mathsf{v}$  corresponds to the value  $v_0$ . More generally, in any implementation of  $\mathsf{Acq}_{\mathsf{WL}}(x)$ ,  $\mathsf{v}$  corresponds to the ticket obtained (e.g.  $v_0 + 1$  for  $\mathsf{e}_2'$ ).

The implementation of  $\mathbf{e}_1'$  (running  $\mathtt{Rel}_{\mathtt{WL}}(x)$ ) also reads  $p_x^{t_1}$ . For the same reason,  $\mathbf{v}$  corresponds to the ticket of the previous  $\mathtt{Acq}_{\mathtt{WL}}(x)$ , i.e.  $v_0$  in our case. Since the program respects locks, every  $\mathtt{Rel}_{\mathtt{WL}}(x)$  handles a different ticket, and  $\mathbf{e}_1'$  is the only one handling ticket  $v_0$  for x.

The second event in the implementation of  $\mathbf{e}_1'$  is  $\mathbf{e}_1 = (t_1, \neg, (\mathtt{Write}_{\mathtt{SV}}, (x_{t_1}, v_0 + 1), ()))$  modifying  $x_{t_1}$ . (There is also a broadcast propagating the new value across the network.) By well-formedness we have  $\mathbf{v}_{\mathtt{W}}''((\mathbf{e}_1, \mathtt{aCW})) = v_0 + 1$ . The last event in the implementation of  $\mathbf{e}_2'$  is of the form  $\mathbf{e}_2 = (t_2, \neg, (\mathtt{Read}_{\mathtt{SV}}, (x_{t_2}), (v_0 + 1)))$  returning a value of  $v_0 + 1$ , and by well-formedness  $\mathbf{v}_{\mathtt{R}}''(\mathbf{e}_2, \mathtt{aCR}) = v_0 + 1$ . The read is necessarily on  $x_{t_2}$  as other  $x_t$  shared variables are never modified to contain  $v_0 + 1$ . Now, by well-formedness of rf" we can create a dependency between  $\mathbf{e}_1$  and  $\mathbf{e}_2$ .

If  $\mathbf{n}(t_1) = \mathbf{n}(t_2)$  (i.e. the two commands are on the same node, perhaps even the same thread), then we have  $(\mathbf{e}_1, \mathbf{aCW}) \xrightarrow{\mathsf{rf}''} (\mathbf{e}_2, \mathbf{aCR})$  as  $(\mathbf{e}_1, \mathbf{aCW})$  is the only element of  $\mathcal{G}.\mathcal{W}^{\mathbf{n}(t_1)}$  writing  $v_0 + 1$ . If they are different threads or  $\mathbf{e}_2 \xrightarrow{\mathsf{po}} \mathbf{e}_1$ , then  $\mathsf{rf}''_{\mathsf{e}} \subseteq \mathsf{hb}$  is enough. Otherwise  $t_1 = t_2$  with  $\mathbf{e}_1 \xrightarrow{\mathsf{po}} \mathbf{e}_2$  and the RFAA/Wait in-between  $\mathbf{e}_1$  and  $\mathbf{e}_2$  forces a sequence of dependencies  $\mathsf{ppo}$ ;  $\mathsf{pfg}$ ;  $\mathsf{ppo} \subseteq \mathsf{hb}$ .

Else  $\mathbf{n}(t_1) \neq \mathbf{n}(t_2)$  and  $\mathbf{e}_2$  reads from an element of  $\mathcal{G.W}^{\mathbf{n}(t_2)}$ , which is a subevent of a broadcast reading from  $\mathbf{e}_1$ . (Technically, this could be from

a delayed broadcast of a previous  $\mathtt{Rel}_{\mathtt{WL}}(x)$  by thread  $t_1$ , not necessarily the broadcast immediately after  $e_1$ .) Thus we similarly have  $((e_1, \mathtt{aCW}), (e_2, \mathtt{aCR})) \in \mathsf{rf}_e''; \mathsf{iso}''; \mathsf{rf}_e'' \subseteq \mathsf{hb}$ .

#### B.3 SLOCK Library

```
I_{SL}(t, Acq_{SL}, (x)) \triangleq Acq_{WL}(x) I_{SL}(t, Rel_{SL}, (x)) \triangleq GFence(Node); Rel_{WL}(x)
```

**Theorem 2.** The implementation  $I_{SL}$  is sound.

*Proof.* This is very straightforward from the semantics of the different libraries. If an execution is lock-well-formed (Definition 4) with respect to strong locks, the implementation is clearly lock-well-formed with respect to weak locks.

A strong acquire  $\mathtt{Acq}_{\mathtt{SL}}(x)$  should behave as stamp  $\mathtt{aMF}$ , which is the case of the implementation  $\mathtt{Acq}_{\mathtt{WL}}(x)$ . A strong release  $\mathtt{Rel}_{\mathtt{SL}}(x)$  should behave as a global fence (stamps  $\mathtt{aGF}_n$ ) and synchronise with later acquires. In the implementation, the first call executes a global fence (stamps  $\mathtt{aGF}_n$ , see Appendix A), while the latter call is a weak release that synchronises with later acquires (Definition 5). The two components execute in order according to  $\mathtt{ppo} \subseteq \mathtt{hb}$  (cell L2 in Fig. 9).

#### B.4 NLOCK Library

**Theorem 3.** The implementation  $I_{NL}$  is sound.

Proof. We assume an {RDMA}\_{RMW}^{WAIT}-consistent execution  $\mathcal{G} = \langle E, \mathsf{po}, \mathsf{stmp}, \mathsf{so}, \mathsf{hb} \rangle$  which is abstracted via f to  $\langle E', \mathsf{po'} \rangle$  that uses (only) the NLOCK library, i.e.  $\mathsf{abs}_{I_{\mathsf{NL}},\mathsf{NLOCK}}^f(\langle E, \mathsf{po} \rangle, \langle E', \mathsf{po'} \rangle)$  holds. We need to provide  $\mathsf{stmp'}, \mathsf{so'}, \mathsf{and} g : \langle E', \mathsf{po'}, \mathsf{stmp'} \rangle$ . SEvent  $\to \mathcal{G}$ . SEvent respecting some conditions. From  $\langle E', \mathsf{po'} \rangle$ , we simply take  $\mathsf{stmp'} = \mathsf{stmp}_{\mathsf{NL}}$ .

Since  $\mathcal{G}$  is  $\{RDMA_{RMW}^{WAIT}\}$ -consistent, it means  $(ppo \cup so) \subseteq hb$ , hb is transitive and irreflexive, and  $\mathcal{G}$  is  $RDMA_{RMW}^{WAIT}$ -consistent. Thus there is some well-formed  $v_R$ ,  $v_W$ , rf, mo, nfo, and rao such that ib is irreflexive,  $\forall e.stmp(e) \in stmp_{RW}(e)$ , and rao so  $iso \cup rf_e \cup pfg \cup nfo \cup rb \cup mo \cup rao \cup ([aNRW]; iso^{-1}; rao) \cup ([Inst]; ib).$ 

As an intermediate result: for each thread t we have  $\mathsf{mo}_{p_x^t} \subseteq \mathsf{po}$ , i.e. the modifications of the temporary location  $p_x^t$  happen in program order. Since  $\mathsf{hb}$  containing  $\mathsf{mo}$  is acyclic, it is enough to show that whenever  $\mathsf{s}_1, \mathsf{s}_2 \in \mathcal{G}.\mathcal{W}_{p_x^t}$  and  $(\mathsf{s}_1, \mathsf{s}_2) \in \mathsf{po}$  then we have  $(\mathsf{s}_1, \mathsf{s}_2) \in \mathsf{hb}$ . If  $\mathsf{s}_1$  has a stamp aCW we immediately have  $(\mathsf{s}_1, \mathsf{s}_2) \in \mathsf{ppo} \subseteq \mathsf{hb}$ . From the implementation, in the other cases  $\mathsf{s}_1$  has a stamp aNLW<sub>n</sub> from either a RFAA or Get operation. In each case,  $\mathsf{s}_1$  is immediately followed by some  $(\mathsf{e}, \mathsf{aWT})$  forcing the write to finish. Thus we have  $(\mathsf{s}_1, \mathsf{s}_2) \in \mathsf{pfg}$ ;  $\mathsf{ppo} \subseteq \mathsf{hb}$ .

We now define g as follows.

• For an event  $e' = (t, \_, (Acq_{NL}, (x), ()))$ , we choose  $g(e', aMF) = (e_r, aCR)$  with  $e_r = (t, \_, (Read, (p_x^t), (v))) \in f^{-1}(e')$  the last read event before breaking the loop, and penultimate event of the implementation.

• For an event  $\mathbf{e}' = (t, \_, (\mathtt{Rel}_{\mathtt{NL}}, (x), ()))$ , we choose:  $g(\mathbf{e}', \mathtt{aRF}_n) = (\mathbf{e}_{rf}, \mathtt{aRF}_n)$  with  $\mathbf{e}_{rf} = (t, \_, (\mathtt{Rfence}, (\mathbf{n}(x)), ())) \in f^{-1}(\mathbf{e}')$  the first event of the implementation; and  $g(\mathbf{e}', \mathtt{aNRW}_n) = (\mathbf{e}_{put}, \mathtt{aNRW}_n)$  with  $\mathbf{e}_{put} = (t, \_, (\mathtt{Put}, (x_r, p_x^t, \_), ())) \in f^{-1}(\mathbf{e}')$  the second event of the implementation.

First, let us show that g preserves  $\operatorname{sto}$  (first property of local soundness). For  $\operatorname{Rel}_{\operatorname{NL}}$  this is trivial as g maps to the same stamps. For  $\operatorname{Acq}_{\operatorname{NL}}$ , the stamp  $\operatorname{aCR}$  is similar to  $\operatorname{aMF}$  w.r.t. later stamps, so  $(\mathsf{e}_2,a_2)=(\mathsf{e}_r,\operatorname{aCR})$  is enough. For an earlier stamp  $a_0$  such that  $(a_0,\operatorname{aMF})\in\operatorname{sto}$ , we take  $(\mathsf{e}_1,a_1)=((t,{}_{-},(\operatorname{RFAA},(\ldots,d),())),\operatorname{aNLW}_n)$  the first event of the implementation, and with  $\mathsf{e}_{wt}=(t,{}_{-},(\operatorname{Wait},(d),()))$  the second event we have  $(\mathsf{e}_1,a_1) \xrightarrow{\operatorname{pfg}} (\mathsf{e}_{wt},\operatorname{aWT}) \xrightarrow{\operatorname{ppo}} (\mathsf{e}_r,\operatorname{aCR})$  (thus included in  $\operatorname{hb}$ ) with  $(a_0,\operatorname{aNLW}_n)\in\operatorname{sto}$ .

Now we need to pick a suitable so' such that  $g(\mathbf{so'}) \subseteq \mathsf{hb}$  and  $\langle E', \mathsf{po'}, \mathsf{stmp'}, \mathsf{so'}, \_\rangle$  is NLOCK-consistent. We can assume that  $\langle E', \mathsf{po'} \rangle$  respects locks, as otherwise  $\mathsf{so'} = \emptyset$  is enough. Thus, for each location x we need to define a total order  $\mathsf{lo'}_x$  on  $A'_x \triangleq \{\mathsf{e'} \mid \mathsf{e'} \in E'_x \land \mathsf{m}(\mathsf{e'}) = \mathsf{Acq}_\mathsf{NL}\}$ . Each event  $\mathsf{e'} \in A'_x$  can be associated to its first subevent of the form  $((t', \_, (\mathsf{RFAA}, (p_x^{t'}, x_a, 1, d), ())), \mathsf{aNAR}_n)$ , with  $n = \mathsf{n}(x)$ . From  $\mathsf{RDMA}^\mathsf{WAIT}_\mathsf{RMW}$ -consistency, rao induces a total ordering on these subevents, and we simply keep the same ordering for  $A'_x$ . As such, we define

$$\begin{split} \mathbf{so}' &= \left\{ \langle \mathbf{e}', \mathtt{aRF}_{\mathtt{n}(\mathtt{loc}(\mathbf{e}'))} \rangle, \langle \mathbf{e}', \mathtt{aNRW}_{\mathtt{n}(\mathtt{loc}(\mathbf{e}'))} \rangle \; \middle| \; \mathtt{m}(\mathbf{e}') = \mathtt{Rel}_{\mathtt{NL}} \right\} \\ &\qquad \bigcup_{x \in \mathtt{Loc}} \left\{ \langle \mathbf{e}'_1, \mathtt{aNRW}_{\mathtt{n}(\mathtt{loc}(\mathbf{e}'_1))} \rangle, \langle \mathbf{e}'_2, \mathtt{aMF} \rangle \; \middle| \; (\mathbf{e}'_1, \mathbf{e}'_2) \in (\mathtt{po}'_x|_{\mathrm{imm}})^{-1}; \mathsf{lo}'_x \right\} \end{split}$$

as expected, and we have that  $\langle E', po', stmp', so', ... \rangle$  is NLOCK-consistent.

Thus, the rest of the proof is to show that  $g(so') \subseteq hb$ , i.e. that the synchronisations promised by the NLOCK library are enforced in the implementation. The easy case is for the internal synchronisation. For  $(\langle e', aRF_n \rangle, \langle e', aNRW_n \rangle) \in so'$ , we clearly have  $(g(\langle e', aRF_n \rangle), g(\langle e', aNRW_n \rangle)) \in ppo \subseteq hb$ .

For the main case, we can assume  $(e'_0, aMF) \xrightarrow{lo'_x} (e'_2, aMF)$  and  $(e'_0, aMF) \xrightarrow{po'_x|_{\text{imm}}} (e'_1, aNRW_n)$ , with  $\mathbf{n}(x) = n$ ,  $e'_0$  running  $\mathtt{Acq}_{\mathrm{NL}}(x)$  by thread  $t_1$ ,  $e'_1$  running  $\mathtt{Rel}_{\mathrm{NL}}(x)$  by thread  $t_1$ , and  $e'_2$  running  $\mathtt{Acq}_{\mathrm{NL}}(x)$  by thread  $t_2$ . We also note  $(e_1, aNRW_n) = g(e'_1, aNRW_n)$  and  $(e_2, aCR) = g(e'_2, aMF)$ . Our goal is then to show  $(e_1, aNRW_n) \xrightarrow{\mathrm{hb}} (e_2, aCR)$ .

We proceed by induction on the ordering  $|o'_x|$ . The base case is for  $(e'_0, aMF) \xrightarrow{|o'_x|_{\text{imm}}} (e'_2, aMF)$ . This base case trivially implies the general case by transitivity, since the program respects locks (i.e. intermediate acquires are being released) and  $(aCR, aNRW_n) \in sto$ .

Let  $\mathbf{e}_0^{faa} = (t_1, \neg, (\mathsf{RFAA}, (p_x^{t_1}, x_a, 1, d), ()))$  be the FAA in the implementation of  $\mathbf{e}_0'$  and  $\mathbf{e}_2^{faa} = (t_2, \neg, (\mathsf{RFAA}, (p_x^{t_2}, x_a, 1, d), ()))$  in the implementation of  $\mathbf{e}_2'$ . By definition we have  $(\mathbf{e}_0^{faa}, \mathsf{aNAR}_n) \xrightarrow{(\mathsf{rao}|_{E_{x_a}})|_{\mathsf{imm}}} (\mathbf{e}_2^{faa}, \mathsf{aNAR}_n)$ , since any remote RMW in  $E_{x_a}$  is from an implementation of some  $\mathsf{Acq}_{\mathsf{WL}}(x)$  event. From the semantics of  $\mathsf{RDMA}_{\mathsf{RMW}}^{\mathsf{WAIT}}$  we have  $(\mathbf{e}_0^{faa}, \mathsf{aNRW}_n) \xrightarrow{\mathsf{hb}} (\mathbf{e}_2^{faa}, \mathsf{aNAR}_n)$  (from the  $([\mathsf{aNRW}]; \mathsf{iso}^{-1}; \mathsf{rao})$  component), and thus we necessarily have  $(\mathbf{e}_0^{faa}, \mathsf{aNRW}_n) \xrightarrow{\mathsf{rf}}$ 

 $(e_2^{faa}, aNAR_n)$ , i.e. the second FAA reads the modified value of the first. This is because  $e_2^{faa}$  cannot read from an earlier write (or the initial value of 0) as that would imply an rb dependency and an hb cycle; and cannot read  $(rf_e \subseteq hb)$  from a later write, as any later write is hb after  $e_2^{faa}$  (via rao and ppo).

There is some value  $v_0 = v_R((e_0^{faa}, aNAR_n))$  read by the first FAA operation. By well-formedness of  $v_R$ ,  $v_W$ , and rf, we have  $v_R((e_2^{faa}, aNAR_n)) = v_W((e_0^{faa}, aNAR_n)) = v_0 + 1$ , i.e. the following  $Acq_{WL}(x)$  gets the next ticket. More generally, it is clear every  $Acq_{WL}(x)$  gets a different ticket. We also have  $v_W((e_0^{faa}, aNLW_n)) = v_0$ , i.e.  $p_x^{t_1}$  is modified to contain  $v_0$ . Respectively  $p_x^{t_2}$  is modified to contain  $v_0 + 1$ .

Let  $\mathbf{e}_0^r$  be the third event of the implementation of  $\mathbf{e}_0'$  reading  $p_x^{t_1}$ . We necessarily have  $(\mathbf{e}_0^{faa}, \mathtt{aNLW}_n) \stackrel{\mathrm{rf}}{\to} (\mathbf{e}_0^r, \mathtt{aCR})$ . This is because  $\mathbf{e}_0^r$  cannot read from the future (it would create an rf; ippo cycle in ib) and the second event  $\mathtt{Wait}(d)$  makes sure all previous modifications of  $p_x^{t_1}$  are available (ignoring the last one would be an rb; hb cycle since  $\mathbf{mo}_{p_x^{t_1}} \subseteq \mathsf{po}$ ). Thus, in the implementation of  $\mathbf{e}_0'$ , the meta-variable v corresponds to the value  $v_0$ . More generally, in any implementation of  $\mathtt{Acq}_{\mathtt{WL}}(x)$ , v corresponds to the ticket obtained (e.g.  $v_0+1$  for  $\mathbf{e}_2'$ ). So the last event  $\mathbf{e}_0^w$  of the implementation of  $\mathbf{e}_0'$  modifies  $p_x^{t_1}$  to  $v_0+1$ .

The implementation of  $\mathbf{e}_1'$  (running  $\mathtt{Rel}_{\mathtt{WL}}(x)$ ) has an operation  $\mathtt{Put}$  (event  $\mathbf{e}_1$ ) reading  $p_x^{t_1}$  to send to  $x_r$ . We necessarily have  $(\mathbf{e}_0^w,\mathtt{aCW}) \xrightarrow{\mathrm{rf}} (\mathbf{e}_1,\mathtt{aNLR}_n)$ , since the write is available  $((\mathtt{aCW},\mathtt{aNLR}_n) \in \mathsf{sto})$  and later write on  $p_x^{t_1}$  from later RFAA are not finished  $((\mathtt{aNLR}_n,\mathtt{aNLW}_n) \in \mathsf{sto}$  and  $\mathtt{n}(x_r) = \mathtt{n}(x_a)$ ). Thus  $\mathtt{v}_{\mathtt{W}}((\mathbf{e}_1,\mathtt{aNRW}_n)) = \mathtt{v}_{\mathtt{R}}((\mathbf{e}_1,\mathtt{aNLR}_n)) = v_0 + 1$ . More generally, each  $\mathtt{Rel}_{\mathtt{WL}}(x)$  modifies  $x_r$  to contain the next value after the ticket obtained by the previous  $\mathtt{Acq}_{\mathtt{WL}}(x)$  operation. Since each  $\mathtt{Acq}_{\mathtt{WL}}(x)$  handles a different ticket, this is the only modification of  $x_r$  to contain  $v_0 + 1$ .

The penultimate event in the implementation of  $\mathbf{e}_2'$  (causing the loop break) is of the form  $\mathbf{e}_2 = (t_2, \neg, (\mathtt{Read}, (p_x^{t_2}), (v_0+1)))$  returning a value of  $v_0+1$ , and by well-formedness  $\mathbf{v}_{\mathtt{R}}((\mathbf{e}_2, \mathtt{aCR})) = v_0+1$ . If we note  $\mathbf{e}_2^{get} = (t_2, \neg, (\mathtt{Get}, (p_x^{t_2}, x_r, d), ()))$  the last  $\mathtt{Get}$  event preceding  $\mathbf{e}_2$  in the implementation, we clearly have  $(\mathbf{e}_2^{get}, \mathtt{aNLW}_n) \xrightarrow{\mathsf{rf}} (\mathbf{e}_2, \mathtt{aCR})$  (as previously, the intermediate  $\mathtt{Wait}$  makes the write available), and  $\mathbf{v}_{\mathtt{R}}((\mathbf{e}_2^{get}, \mathtt{aNRR}_n)) = \mathbf{v}_{\mathtt{W}}((\mathbf{e}_2^{get}, \mathtt{aNLW}_n)) = \mathbf{v}_{\mathtt{R}}((\mathbf{e}_2, \mathtt{aCR})) = v_0 + 1$ .

By well-formedness of rf, we also have  $(e_1, aNRW_n) \xrightarrow{rf} (e_2^{get}, aNRR_n)$  from the only write of  $v_0 + 1$  on  $x_r$ . Finally, we have  $((e_1, aNRW_n), (e_2, aCR)) \in rf_e$ ; iso;  $rf_e \subseteq hb$ .

# B.5 RDMA<sub>RMW</sub> Library

**Theorem 4.** The implementation  $I_{SC}$  is sound.

Proof. We assume an {RDMA\_{\rm RMW}^{\rm WAIT}}, NLOCK}-consistent execution  $\mathcal{G}=\langle E, {\rm po,stmp,so,hb}\rangle$  which is abstracted via f to  $\langle E', {\rm po'}\rangle$  that uses (only) the RDMA\_{\rm RMW}^{\rm SC} library, i.e. abs $_{I_{\rm SC},{\rm RDMA_{\rm RMW}^{\rm SC}}}^{\rm SC}(\langle E, {\rm po}\rangle, \langle E', {\rm po'}\rangle)$  holds. We need to provide stmp', so', and g:  $\langle E', {\rm po'}, {\rm stmp'}\rangle$ . SEvent  $\to \mathcal{G}$ . SEvent respecting some conditions. From  $\langle E', {\rm po'}\rangle$ , we simply take stmp' = stmp\_{SC}.

Since  $\mathcal{G}$  is  $\{RDMA_{RMW}^{WAIT}, NLOCK\}$ -consistent, it means  $(ppo \cup so|_{RDMA_{RMW}^{WAIT}} \cup so|_{NLOCK}) \subseteq hb$ , hb is transitive and irreflexive, and the two restrictions of  $\mathcal{G}$  are respectively  $RDMA_{RMW}^{WAIT}$ -consistent and NLOCK-consistent.

 $\begin{aligned} & \text{RDMA}^{\text{WAIT}}_{\text{RMW}}\text{-consistency implies there is some well-formed } v_{\text{R}}, \ v_{\text{W}}, \ \text{rf}, \ \text{mo}, \ \text{nfo}, \\ & \text{and rao such that ib is irreflexive}, \ \forall \text{e.stmp}|_{\text{RDMA}^{\text{WAIT}}_{\text{RMW}}}(\textbf{e}) \in \text{stmp}_{\text{RW}}(\textbf{e}), \ \text{and} \ \text{so}|_{\text{RDMA}^{\text{WAIT}}_{\text{RMW}}} = \\ & \text{iso} \cup \text{rf}_{\textbf{e}} \cup \text{pfg} \cup \text{nfo} \cup \text{rb} \cup \text{mo} \cup \text{rao} \cup ([\textbf{aNRW}]; \textbf{iso}^{-1}; \text{rao}) \cup ([\textbf{Inst}]; \textbf{ib}). \end{aligned}$ 

 $I_{\rm SC}$  respects locks, as every operation is implemented to contain an  ${\sf Acq}_{\rm NL}$  (first) and a  ${\sf Rel}_{\rm NL}$  operation (later) on the same lock location. As such  $\langle E|_{\rm NLOCK}, {\sf po}|_{\rm NLOCK} \rangle$  respects locks. So NLOCK-consistency implies  ${\sf stmp}|_{\rm NLOCK} = {\sf stmp}_{\rm NL}$  and for each lock location l there is a total order  $|{\sf o}_l|$  on  $\{{\sf e} \mid {\sf e} \in E_l \land {\sf m}({\sf e}) = {\sf Acq}_{\rm NL}\}$  for the acquiring of location l such that:

$$\begin{split} & \text{SO}|_{\text{NLOCK}} = \big\{ \langle e, \text{aRF}_{\text{n(loc(e))}} \rangle, \langle e, \text{aNRW}_{\text{n(loc(e))}} \rangle \bigm| \text{m(e)} = \text{Rel}_{\text{NL}} \big\} \\ & \qquad \bigcup_{l \in l \text{ oc}} \big\{ \langle e_1, \text{aNRW}_{\text{n(loc(e_1))}} \rangle, \langle e_2, \text{aMF} \rangle \bigm| (e_1, e_2) \in (\text{po}_l|_{\text{imm}})^{-1}; \text{lo}_l \big\} \end{split}$$

We define g to map to the first subevent of the implementation. For an event  $e' \in E'$ , we choose g(e', aMF) = (e, aMF) with  $e = (t, \_, (Acq_{NL}, \_, \_)) \in f^{-1}(e')$  the first event of the implementation. This g clearly preserves sto (first property of local soundness), as it maps subevents to subevents using the same stamp.

Now we need to pick a suitable so' such that  $g(so') \subseteq hb$  and  $\langle E', po', stmp', so', \_\rangle$  is RDMA<sup>SC</sup><sub>RMW</sub>-consistent. I.e., we need well-formed  $v'_R$ ,  $v'_W$ , rf', and mo' such that g(po'), g(rf'), g(mo'), and g(rb') are all included in hb. We immediately have  $g(po') \in ppo \subseteq hb$  since  $(aMF, aMF) \in sto$ . For the other relations, we can consider each location x independently. Let us note  $n = n(x) = n(l_x)$ . All the relevant operations acquire the lock  $l_x$ , as such we can use  $lo_{l_x}$  to order them.

We define  $mo'_x$  and  $rf'_x$  as follows:

$$\mathbf{mo}_x' \triangleq \{(\mathbf{s}_1', \mathbf{s}_2') \mid \mathbf{s}_1', \mathbf{s}_2' \in \mathcal{G}'.\mathcal{W} \ \land \ (g(\mathbf{s}_1'), g(\mathbf{s}_2')) \in \mathsf{lo}_{l_x}\}$$

$$\mathsf{rf}_x' \triangleq \left\{ (\mathsf{s}_1', \mathsf{s}_2') \;\middle|\; \begin{aligned} \mathsf{s}_1' \in \mathcal{G}'.\mathcal{W} \; \wedge \; \mathsf{s}_2' \in \mathcal{G}'.\mathcal{R} \; \wedge \; (g(\mathsf{s}_1'), g(\mathsf{s}_2')) \in \mathsf{lo}_{l_x} \; \wedge \\ \forall \mathsf{s}_0'.(\mathsf{s}_1', \mathsf{s}_0') \in \mathsf{mo}_x' \implies (g(\mathsf{s}_0'), g(\mathsf{s}_2')) \not \in \mathsf{lo}_{l_x} \end{aligned} \right\}$$

with the slight abuse of notation of writing  $((e_1, a_1), (e_2, a_2)) \in lo_l$  to mean  $(e_1, e_2) \in lo_l$ . I.e., the location x is modified in the order of the acquires, and reads read from the latest previous write.

We define  $\mathbf{v}_{\mathtt{R}}'$  and  $\mathbf{v}_{\mathtt{W}}'$  from the values of  $\mathbf{v}_{\mathtt{R}}$  and  $\mathbf{v}_{\mathtt{W}}$  on the RDMA subevent (on x) of the implementation. E.g., for  $\mathbf{e}'$  running  $\mathtt{FAA}_{\mathtt{SC}}(x,v)$ , there is an event  $\mathbf{e} = (\_,\_,(\mathtt{RFAA},(\_,x,v,\_),())) \in f^{-1}(\mathbf{e}')$  and we note  $\mathbf{v}_{\mathtt{R}}'((\mathbf{e}',\mathtt{aMF})) = \mathbf{v}_{\mathtt{R}}((\mathtt{e},\mathtt{aNAR}_n))$  and  $\mathbf{v}_{\mathtt{W}}'((\mathtt{e}',\mathtt{aMF})) = \mathbf{v}_{\mathtt{W}}((\mathtt{e},\mathtt{aNRW}_n))$ .

We can easily see that  $g(\mathsf{rf}')$ ,  $g(\mathsf{mo}')$ , and  $g(\mathsf{rb}')$  are all included in hb by design. This comes from the fact that  $(\mathsf{e}_1, \mathsf{e}_2) \in \mathsf{lo}_{l_x}$  implies  $((\mathsf{e}_1, \mathsf{aMF}), (\mathsf{e}_2, \mathsf{aMF})) \in \mathsf{ppo}; \mathsf{so} \subseteq \mathsf{hb}$  (since E respects nodes and the release operation exists) and for  $\mathsf{rb}'$  because  $\mathsf{lo}_{l_x}$  is total on the acquiring of the lock  $l_x$  (Thus if  $(g(\mathsf{s}'_0), g(\mathsf{s}'_2)) \not\in \mathsf{lo}_{l_x}$  and  $\mathsf{s}'_0 \neq \mathsf{s}'_2$  then  $(g(\mathsf{s}'_2), g(\mathsf{s}'_0)) \in \mathsf{lo}_{l_x}$ ).

The remaining part of the proof is to show that  $v_R'$ ,  $v_W'$ , and rf' are well-formed.

Firstly, let's consider  $v'_w$ . For RMW operations, the value is correct from the well-formedness of  $v_w$ . For an event e' running  $\operatorname{Write}_{\operatorname{SC}}(x,v)$ , the implementation contains  $\operatorname{Write}(p^t_x,v); \operatorname{Put}(x,p^t_x,\cdot)$  (let's call them  $e_1$  and  $e_2$ ), and we need to show  $v'_w((e',\operatorname{aMF})) = v$ . By definition  $v'_w((e',\operatorname{aMF})) = v_w((e_2,\operatorname{aNRW}_n)) = v_w((e_2,\operatorname{aNLR}_n))$  and  $v_w((e_1,\operatorname{aCW})) = v$ . To conclude, it is enough to show  $((e_1,\operatorname{aCW}),(e_2,\operatorname{aNLR}_n)) \in r$ . Clearly  $e_1$  is finished when we run  $e_2$  (i.e.,  $((e_2,\operatorname{aNLR}_n),(e_1,\operatorname{aCW})) \in r$ b would create an hb cycle). It is less obvious that  $e_2$  cannot read from a later  $\operatorname{Write}(p^t_x,v')$  (let's call it event  $e_3$ ) of a later operation  $\operatorname{Write}_{\operatorname{SC}}(x,v')$  by the same thread. This is because this later operation would need to acquire the lock  $l_x$ . By the semantics of NLOCK, this creates a synchronisation (since  $e_1$  is also towards node n), and we have  $((e_2,\operatorname{aNLR}_n),(e_3,\operatorname{aCW})) \in \operatorname{ppo}; \operatorname{sol}_{\operatorname{NLOCK}};\operatorname{ppo} \subseteq \operatorname{hb}$ . As such, reading from  $e_3$  would create an hb cycle and is not possible.

Secondly, for  $\mathbf{v}_{R}'$  and non-Write<sub>SC</sub> operations, we need to show that the value returned (i.e. by  $\mathbf{e}_{r}$  running  $\mathtt{Read}(r_{t})$ ) is the value read by the RDMA operation  $\mathbf{e}$  running  $m(r_{t},x,\ldots,d)$  with  $m\in\{\mathtt{Read}_{SC},\mathtt{CAS}_{SC},\mathtt{FAA}_{SC}\}$ . For this, we simply show  $((\mathtt{e},\mathtt{aNLW}_{n}),(\mathtt{e}_{r},\mathtt{aCR}))\in\mathtt{rf}$ . From the in-between Wait operation, we have  $((\mathtt{e},\mathtt{aNLW}_{n}),(\mathtt{e}_{r},\mathtt{aCR}))\in\mathtt{pfg}$ ;  $\mathtt{ppo}\subseteq\mathtt{hb}$ . Thus,  $\mathtt{e}_{r}$  cannot ignore  $\mathtt{e}$  (i.e.  $\mathtt{rb}$  would create an  $\mathtt{hb}$  cycle), and cannot read from a later operations (it would create an  $\mathtt{ib}$  cycle).

Finally, we are left with checking that rf' is well-formed. We need to show that whenever  $(s'_1, s'_2) \in \text{rf}'$ , with  $s'_i = (e'_i, \text{aMF})$ , we have  $v'_w(s'_1) = v'_R(s'_2)$ . (Technically, also that  $(., s'_2) \notin \text{rf}'$  implies  $v'_R(s'_2) = 0$ , but this follows from a similar reasoning.) Let  $e_i$  be the RDMA operation in the implementation of  $e'_i$ , by definition we have  $v'_w(s'_1) = v_w((e_1, \text{aNRW}_n))$  and  $v'_R(s'_2) = v_R((e_2, a_2))$  (with  $a_2 \in \{\text{aNAR}_n, \text{aNRR}_n\}$  depending on the case). Our sufficient goal is then to show that we necessarily have  $((e_1, \text{aNRW}_n), (e_2, a_2)) \in \text{rf}$ . By definition of rf', we have  $(g(s'_1), g(s'_2)) \in \text{lo}_{l_x}$  as well as  $\forall s'_0 \in \mathcal{G}'.\mathcal{W}$ .  $((s'_1, s'_0) \in \text{mo}'_x \implies (g(s'_0), g(s'_2)) \notin \text{lo}_{l_x})$ .

The first point implies  $((e_1, aNRW_n), (e_2, a_2)) \in ppo; so|_{NLOCK}; ppo \subseteq hb$  by the semantics of locks. This makes  $((e_2, a_2), (e_1, aNRW_n)) \in rb$  impossible (hb cycle), and  $e_2$  reads from either  $e_1$  or a later write: there exists an RDMA operation  $e_3$  (in the implementation of some  $e_3'$ ) such that  $((e_3, aNRW_n), (e_2, a_2)) \in rf$  with  $(e_1, aNRW_n) \xrightarrow{mo^*} (e_3, aNRW_n)$ . Note that  $e_2 \neq e_3$  or it would create a  $rf_e$ ; iso cycle in hb; i.e. an event cannot read from itself. Thus we need to show  $e_1 = e_3$ , and by contradiction let us assume  $(e_1, aNRW_n) \xrightarrow{mo} (e_3, aNRW_n)$ .

We then show  $(\mathbf{e}_1', \mathbf{aMF}) \xrightarrow{\mathbf{mo}'} (\mathbf{e}_3', \mathbf{aMF})$ . Since  $\mathsf{lo}_{l_x}$  is a total order, we have either  $(g(\mathsf{s}_1'), g(\mathsf{s}_3')) \in \mathsf{lo}_{l_x}$  or  $(g(\mathsf{s}_3'), g(\mathsf{s}_1')) \in \mathsf{lo}_{l_x}$ . To show the first, we assume the second by contradiction, i.e. that  $\mathbf{e}_3'$  acquires first. Given the implementation  $I_{\mathsf{SC}}$ , there is a  $\mathsf{Rel}_{\mathsf{NL}}(x)$  event  $\mathbf{e}_3^r$  such that  $g(\mathsf{e}_3') \xrightarrow{\mathsf{po}} \mathsf{e}_3 \xrightarrow{\mathsf{po}} \mathsf{e}_3^r$ . Thus from the semantics of NLOCK we have an hb cycle  $(\mathsf{e}_3, \mathsf{aNRW}_n) \xrightarrow{\mathsf{popo}} (\mathsf{e}_3^r, \mathsf{aNRW}_n) \xrightarrow{\mathsf{so}|_{\mathsf{NLOCK}}} g(\mathsf{e}_1') \xrightarrow{\mathsf{ppo}} (\mathsf{e}_1, \mathsf{aNRW}_n) \xrightarrow{\mathsf{mo}'} (\mathsf{e}_3, \mathsf{aNRW}_n)$  providing a contradiction, and we necessarily have  $(\mathsf{e}_1', \mathsf{aMF}) \xrightarrow{\mathsf{mo}'} (\mathsf{e}_3', \mathsf{aMF})$ .

Now, from the definition of rf', the fact  $|o_{l_x}|$  is a total order, and that  $e_3 \neq e_2$ , we have  $(g(s_2'), g(s_3')) \in |o_{l_x}|$ . Similarly to previously, this implies  $((e_2, a_2), (e_3, aNRW_n)) \in ppo; ppo; so|_{NLOCK}; ppo \subseteq hb$  by the semantics of locks (using both the aRF<sub>n</sub> and

aNLW<sub>n</sub> stamps of the release). This contradicts  $((e_3, aNRW_n), (e_2, a_2)) \in rf_e \subseteq hb$ . Thus  $e_1 = e_3$ ,  $((e_1, aNRW_n), (e_2, a_2)) \in rf$ , and rf' is well-formed.

#### Declarative Semantics of RDMARMW à la MOWGLI $\mathbf{C}$

In this appendix, we first (§C.1) present the declarative semantics of  $RDMA_{RMW}^{TSO}$ in a format similar to that of RDMATSO in [4], but extended with remote RMW operations similarly to the semantics of  $RDMA_{RMW}^{WAIT}$  given in §3. It is slightly different from the one in §D, as we use the stamps and subevents system of MOWGLI.

We then (§C.2) provide a definition of the implementation of RDMA<sub>RMW</sub> into RDMA<sub>RMW</sub>. Finally (§C.3), we give a proof of the soundness of this implementation, similarly to [4].

#### C.1Semantics

Our definition of  $RDMA_{RMW}^{TSO}$  is closer to an independent language than a library. We do not need a relation hb to represent the potential rest of the program, as a program cannot combine instructions from  $RDMA_{RMW}^{TSO}$  and other libraries presented in this paper.

We use the following 13 methods:

```
\mid \mathtt{Get}^{\mathrm{TSO}}(x,y) \mid \mathtt{Put}^{\mathrm{TSO}}(x,y) \mid \mathtt{Poll}(n) \mid \mathtt{Rfence}^{\mathrm{TSO}}(n)
           | \operatorname{RCAS}^{\mathrm{TSO}}(x, y, v_1, v_2) | \operatorname{RFAA}^{\mathrm{TSO}}(x, y, v)
           \mid \mathtt{SetAdd}(x,v) \mid \mathtt{SetRemove}(x,v) \mid \mathtt{SetIsEmpty}(x)
```

- $$\begin{split} \bullet \ \, & \mbox{Write}^{\rm TSO} : \mbox{Loc} \times \mbox{Val} \rightarrow () \\ \bullet \ \, & \mbox{Read}^{\rm TSO} : \mbox{Loc} \rightarrow \mbox{Val} \\ \bullet \ \, & \mbox{CAS}^{\rm TSO} : \mbox{Loc} \times \mbox{Val} \times \mbox{Val} \rightarrow \mbox{Val} \\ \end{split}$$
- Mfence $^{\mathrm{TSO}}:() \rightarrow ()$
- $\mathsf{Get^{TSO}}:\mathsf{Loc}\times\mathsf{Loc}\to\mathsf{Val}$   $\mathsf{Put^{TSO}}:\mathsf{Loc}\times\mathsf{Loc}\to\mathsf{Val}$
- Poll : Node  $\rightarrow$  Val

- Rfence $^{TSO}$ : Node  $\rightarrow$  ()
- $\mathtt{RCAS}^{\mathrm{TSO}}: \mathsf{Loc} \times \mathsf{Loc} \times \mathsf{Val}^2 \to \mathsf{Val}$
- RFAA $^{\mathrm{TSO}}:\mathsf{Loc}\times\mathsf{Loc}\times\mathsf{Val}\to\mathsf{Val}$
- SetAdd : Loc  $\times$  Val  $\rightarrow$  ()
- SetRemove : Loc  $\times$  Val  $\rightarrow$  ()
- SetIsEmpty : Loc  $\rightarrow \mathbb{B}$

This version is based on [4] extended with remote RMW. Compared to RDMA<sup>TSO</sup> from [3], we slightly extend the language so that RDMA operations return an arbitrary unique identifier, and polling also returns the same identifier of the operation being polled. In addition, we also assume basic set operations SetAdd, SetRemove, and SetIsEmpty to store these new identifiers, where the locations used for sets do not overlap with locations used for other operations.

Consistency predicate. An execution of an RDMARMW program is of the form  $\mathcal{G} = \langle E, \mathsf{po}, \mathsf{stmp}, \mathsf{so} \rangle$ . Note that  $\mathsf{hb} = (\mathsf{ppo} \cup \mathsf{so})^+$  does not have the flexibility of containing additional external constraints.

We say that a stamping function  $stmp_{TSO}$  is valid if:

- Polls have stamp aWT:  $stmp_{TSO}((-, -, (Poll, -, -))) = \{aWT\}.$
- Auxiliary set operations have stamp aMF:  $stmp_{TSO}((\_,\_,(SetAdd,\_,\_))) = stmp_{TSO}((\_,\_,(SetIsEmpty,\_,\_))) = stmp_{TSO}((\_,\_,(SetIsEmpty,\_,\_))) = {aMF}.$  Other events follow the validity constraints of RDMA\_RMW\_(cf. Section 3). E.g.,
- Other events follow the validity constraints of RDMA<sup>WAIT</sup><sub>RMW</sub> (cf. Section 3). E.g., events calling Write<sup>TSO</sup> have stamp aCW, while events calling Get<sup>TSO</sup> towards node n have stamps aNRR<sub>n</sub> and aNLW<sub>n</sub>. We also define loc on subevents similarly to RDMA<sup>WAIT</sup><sub>RMW</sub>.

We mark set operations with aMF to simplify the consistency conditions, as we do not want to explicitly integrate them in the read  $(\mathcal{R})$  and write  $(\mathcal{W})$  subevents.

Given  $\mathcal{G} = \langle E, po, stmp_{TSO}, so \rangle$ , we say that  $v_R$ ,  $v_W$ , rf, mo, nfo, pf, and rao are well-formed if:

- $\bullet~v_R,\,v_W,\,rf,\, {\color{red}mo},\, nfo,\, {\rm and}~rao~{\rm are~well\text{-}formed},\, {\rm as~in~RDMA_{RMW}^{WAIT}};$
- Let  $P_n \triangleq \{(\mathsf{e}, \mathsf{aWT}) \mid \mathsf{e} = (\_, \_, (\mathsf{Poll}, (n), \_)) \in E\}$  be the set of poll (sub)events towards node n. Let  $C_n \triangleq \{(\mathsf{e}, \mathsf{aNLW}_n) \mid \mathsf{m}(\mathsf{e}) \in \{\mathsf{Get}^{\mathsf{TSO}}, \mathsf{RCAS}^{\mathsf{TSO}}, \mathsf{RFAA}^{\mathsf{TSO}}\}\} \cup \{(\mathsf{e}, \mathsf{aNRW}_n) \mid \mathsf{m}(\mathsf{e}) = \mathsf{Put}^{\mathsf{TSO}}\}$  be the set of (final writes of) remote operations towards node n that need polling. Note that, for remote RMW operations, polling only synchronises with the  $\mathsf{aNLW}_n$  part and not with the (potential)  $\mathsf{aNRW}_n$  part.

Then  $\operatorname{pf} \subseteq \bigcup_{n \in \operatorname{Node}} C_n \times P_n$  is the *polls-from* relation, relating earlier NIC writes to later polls. Moreover:

- $pf \subseteq po$  (we can only poll previous operations of the same thread);
- pf is functional on its domain (every NIC write can be polled at most once);
- pf is total and functional on its range (every Poll polls from exactly one NIC write);
- Poll events poll-from the oldest non-polled remote operation towards the given node:
  - for each node n, if  $w_1, w_2 \in C_n$  and  $w_1 \xrightarrow{po} w_2 \xrightarrow{pf} p_2$ , then there exists  $p_1$  such that  $w_1 \xrightarrow{pf} p_1 \xrightarrow{po} p_2$ ;
- and a Poll returns the unique identifier of the polled operation: if  $((-, -, (-, -, v_1)), -) \xrightarrow{pf} ((-, -, (Poll, -, v_2)), aWT)$  then  $v_1 = v_2$ .

We use the derived relations rb,  $rb_i$ ,  $rf_e$ ,  $rf_i$ , ippo, and iso as defined for RDMA $_{\rm RMW}^{\rm WAIT}$ . We can then define ib as follows:

```
\mathsf{ib} \triangleq (\mathsf{ippo} \cup \mathsf{iso} \cup \mathsf{rf} \cup \mathsf{pf} \cup \mathsf{nfo} \cup \mathsf{rb}_{\mathsf{i}} \cup \{(\mathsf{e}, \mathsf{aNRW}_n), (\mathsf{e}, \mathsf{aNLW}_n)\})^+
```

The last new component states that, for remote RMW operations, the remote write part *starts* before the local write part. As mentioned previously, this does not imply that they finish in order, and this component is not included in so. Since it does not prevent any behaviour, we remove this component in the definition RDMA<sub>RMW</sub> (§3), but we keep it here as it simplifies the soundness proof (Theorem 6) and the equivalence proof with the operational semantics (Appendix D.5).

**Definition 16** (RDMA<sub>RMW</sub>-consistency).  $\mathcal{G} = \langle E, po, stmp, so \rangle$  is RDMA<sub>RMW</sub>-consistent if:

- $(ppo \cup so)^+$  is irreflexive;
- $\langle E, po \rangle$  respects nodes (as in RDMA<sub>RMW</sub>);
- stmp is valid;
- there exists well-formed  $v_R$ ,  $v_W$ , rf, mo, nfo, pf, and rao such that ib is irreflexive and
  - $\mathsf{so} = \mathsf{iso} \cup \mathsf{rf}_\mathsf{e} \cup [\mathsf{aNLW}]; \mathsf{pf} \cup \mathsf{nfo} \cup \mathsf{rb} \cup \mathsf{mo} \cup \mathsf{rao} \cup ([\mathsf{aNRW}]; \mathsf{iso}^{-1}; \mathsf{rao}) \cup ([\mathsf{Inst}]; \mathsf{ib});$
- identifiers for RDMA operations are unique: if  $e_1$  and  $e_2$  are both of the form (-,-,(m,-,v)) with  $m \in \{\text{Put}^{\mathrm{TSO}}, \text{Get}^{\mathrm{TSO}}, \text{RCAS}^{\mathrm{TSO}}, \text{RFAA}^{\mathrm{TSO}}\}$  then  $e_1 = e_2$ ;
- and the set operations are (per-thread) sound: if SetIsEmpty returns true, then every value added to the set was subsequently removed. I.e., if  $e_1 = (t, \_, (SetAdd, (x, v), \_), e_3 = (t, \_, (SetIsEmpty, (x), true)), and e_1 \xrightarrow{po} e_3$ , then there exists  $e_2 = (t, \_, (SetRemove, (x, v), \_)$  such that  $e_1 \xrightarrow{po} e_2 \xrightarrow{po} e_3$ .

### C.2 Implementation Function

In Fig. 16 we define the implementation  $I_{\mathbb{W}}$  from a full program using only the RDMA<sub>RMW</sub> library into a program using only RDMA<sub>RMW</sub>. We assume threads use disjoint work identifiers  $d \in \mathsf{Wid}$ , otherwise it is straightforward to rename them.

For each location x of RDMA<sub>RMW</sub>, we also use a location x for RDMA<sub>RMW</sub>. For each work identifier d of RDMA<sub>RMW</sub>, we use new RDMA<sub>RMW</sub> locations  $\{d^1, \ldots, d^N\}$  where  $N \triangleq \#(\mathsf{Node})$  is the number of nodes. Each location  $d^n$  is used as a set containing the identifiers of ongoing operations towards node n.

Most RDMA<sub>RMW</sub> operations (Write, Read, CAS, Mfence, and Rfence) are directly translated into their RDMA<sub>RMW</sub> counterparts. An operation  $\mathtt{Get}(x,y,d)$  towards node n is translated into a similar  $\mathtt{Get}^{\mathrm{TSO}}(x,y)$  whose output is added to the set  $d^n$ ; We proceed similarly for other RDMA operations. Finally, a  $\mathtt{Wait}(d)$  operation needs to poll until all relevant operations are finished, i.e. the sets  $\{d^1,\ldots,d^N\}$  are all empty. Whenever we poll, we obtain the identifier of a finished operation, and we remove it from all sets where it might be held. We remove it from  $d^n$  but also from any other set  $d^n_k$  tracking a different group of operations, as otherwise a later call to  $\mathtt{Wait}(d_k)$  would hang and never return.

#### C.3 Soundness

We do not prove that the implementation above is locally sound as it does not apply for this case. Instead, we assume a full program using only the  $RDMA_{RMW}^{WAIT}$  library and compile it into  $RDMA_{RMW}^{TSO}$ .

**Theorem 6.** Let  $\widetilde{p}$  be a program using only the RDMA<sub>RMW</sub> library. Then we have  $\operatorname{outcome}_{\operatorname{RDMA}_{\operatorname{BDMY}}^{\operatorname{MSD}}}(\|\widetilde{p}\|_{I_{\overline{w}}}) \subseteq \operatorname{outcome}_{\{\operatorname{RDMA}_{\operatorname{BDMY}}^{\operatorname{WAIT}}\}}(\widetilde{p})$ , where:

```
\begin{aligned} & \mathtt{outcome}_{\{\mathtt{RDMA}_{\mathtt{RMW}}^{\mathtt{MAIT}}\}}(\widetilde{\mathsf{p}}) = \{\widetilde{v} \mid \exists \langle E, \mathtt{po}, \mathtt{stmp}, \mathtt{so}, \mathtt{hb} \rangle \ \{\mathtt{RDMA}_{\mathtt{RMW}}^{\mathtt{MAIT}}\} \text{-} consistent. \ \langle \widetilde{v}, \langle E, \mathtt{po} \rangle \rangle \in \llbracket \widetilde{\mathsf{p}} \rrbracket \} \\ & \mathtt{outcome}_{\mathtt{RDMA}_{\mathtt{RMW}}^{\mathtt{TSO}}}(\lVert \widetilde{\mathsf{p}} \rVert_{I_{\mathtt{k}}}) = \{\widetilde{v} \mid \exists \langle E, \mathtt{po}, \mathtt{stmp}, \mathtt{so} \rangle \ \mathtt{RDMA}_{\mathtt{RMW}}^{\mathtt{TSO}} \text{-} consistent. \ \langle \widetilde{v}, \langle E, \mathtt{po} \rangle \rangle \in \llbracket \lVert \widetilde{\mathsf{p}} \rVert_{I_{\mathtt{k}}} \rrbracket \} \end{aligned}
```

For a thread t using work identifiers  $\{d_1, \ldots, d_K\}$ :

```
I_{\mathsf{W}}(t, \mathsf{Wait}, (d)) \triangleq
I_{\mathsf{W}}(t,\mathsf{Write},(x,v)) \triangleq \mathsf{Write}^{\mathrm{TSO}}(x,v)
                                                                                                For n in 1, \ldots, N do \{
I_{\mathtt{W}}(t,\mathtt{Read},(x)) \triangleq \mathtt{Read}^{\mathrm{TSO}}(x)
                                                                                                    While (SetIsEmpty(d^n) \neq true) do {
I_{\mathsf{W}}(t, \mathsf{CAS}, (x, v_1, v_2)) \triangleq \mathsf{CAS}^{\mathsf{TSO}}(x, v_1, v_2)
                                                                                                         let v = Poll(n) in
I_{\mathsf{W}}(t, \mathsf{Mfence}, ()) \triangleq \mathsf{Mfence}^{\mathsf{TSO}}()
                                                                                                         For k in 1, \ldots, K do \{
I_{\mathsf{W}}(t, \mathsf{Rfence}, (n)) \triangleq \mathsf{Rfence}^{\mathsf{TSO}}(n)
                                                                                                              SetRemove(d_k^n, v) \} \} 
I_{\mathtt{W}}(t,\mathtt{Get},(x,y,d)) \triangleq
                                                                                               I_{\mathbb{W}}(t, \mathtt{RCAS}, (x, y, v_1, v_2, d)) \triangleq
\mathtt{let}\,v = \mathtt{Get}^{\mathrm{TSO}}(x,y)\,\mathtt{in}\,\mathtt{SetAdd}(d^{\mathtt{n}(y)},v)
                                                                                               let v = \mathtt{RCAS}^{\mathrm{TSO}}(x, y, v_1, v_2) in \mathtt{SetAdd}(d^{\mathtt{n}(y)}, v)
I_{\mathbb{W}}(t, \mathtt{Put}, (x, y, d)) \triangleq
                                                                                               I_{\mathtt{W}}(t,\mathtt{RFAA},(x,y,v',d)) \triangleq
let v = Put^{TSO}(x, y) in SetAdd(d^{n(x)}, v)
                                                                                               \mathtt{let}\,v = \mathtt{RFAA}^{\mathrm{TSO}}(x,y,v')\,\mathtt{in}\,\mathtt{SetAdd}(d^{\mathtt{n}(y)},v)
```

Fig. 16: Implementation  $I_{W}$  of RDMA<sub>RMW</sub> into RDMA<sub>RMW</sub> TSO

*Proof.* By definition, we are given  $\mathcal{G} = \langle E, \mathsf{po}, \mathsf{stmp}, \mathsf{so} \rangle$  RDMA $_{\mathrm{RMW}}^{\mathrm{TSO}}$ -consistent (Definition 16) such that  $\langle \widetilde{v}, \langle E, \mathsf{po} \rangle \rangle \in \llbracket \lVert \widetilde{\mathsf{p}} \rVert_{I_{\mathsf{w}}} \rrbracket$ . Among others, it means  $\langle E, \mathsf{po} \rangle$  respects nodes and there exists well-formed  $\mathsf{v}_{\mathsf{R}}, \mathsf{v}_{\mathsf{w}}, \mathsf{rf}, \mathsf{mo}, \mathsf{nfo}, \mathsf{pf}, \mathsf{nd}$  rao such that ib is irreflexive, stmp is valid, so = iso  $\cup$  rf<sub>e</sub>  $\cup$  [ $\cup_n$  aNLW<sub>n</sub>]; pf  $\cup$  nfo  $\cup$  rb  $\cup$  mo  $\cup$  rao  $\cup$  ([aNRW]; iso<sup>-1</sup>; rao)  $\cup$  ([Inst]; ib), and hb  $\triangleq$  (ppo  $\cup$  so)<sup>+</sup> is irreflexive.

From Lemma 1, since  $\tilde{\mathbf{p}}$  uses only RDMA<sub>RMW</sub>, there is E', po', f such that  $\langle \tilde{v}, \langle E', \mathsf{po'} \rangle \rangle \in [\tilde{\mathbf{p}}]$  and  $\mathsf{abs}_{I_{\mathsf{W},\mathsf{RDMA}_{\mathsf{RMW}}}^{\mathsf{WAIT}}}(\langle E, \mathsf{po} \rangle, \langle E', \mathsf{po'} \rangle)$ . Note that this clearly implies  $\langle E, \mathsf{po} \rangle$  also respects nodes, as the implementation  $I_{\mathsf{W}}$  keeps the same locations. Our objective is to find stmp', so', and hb' such that  $\mathcal{G}' = \langle E', \mathsf{po'}, \mathsf{stmp'}, \mathsf{so'}, \mathsf{hb'} \rangle$  is  $\{\mathsf{RDMA}_{\mathsf{RMW}}^{\mathsf{WAIT}}\}$ -consistent (Definitions 2 and 8). To choose a valid function stmp', most values are forced. For remote compare-and-swap, we make the same choice as stmp. I.e. for each RCAS we assert it succeeds iff the corresponding RCAS<sup>TSO</sup> in its implementation succeeds. We will also pick  $\mathsf{hb'} \triangleq (\mathsf{ppo'} \cup \mathsf{so'})^+$  since there is no external constraints. Thus, we only need to carefully pick so' and show it works.

While our objective is not exactly local soundness (Definition 15), we still use a concretisation function  $g:\langle E',\mathsf{po'},\mathsf{stmp'}\rangle.\mathsf{SEvent}\to\mathcal{G}.\mathsf{SEvent}$  to then define so'

- For  $\mathbf{e}'=(t, \_, (\mathtt{Write}, (x, v), ()))$ , from the definition of the implementation  $I_{\mathtt{W}}$  and the abstraction f, there is some event  $\mathbf{e}=(t, \_, (\mathtt{Write}^{\mathrm{TSO}}, (x, v), ())) \in f^{-1}(\mathbf{e}')$ . We define  $g(\mathbf{e}', \mathtt{aCW}) = (\mathtt{e}, \mathtt{aCW})$ . For events calling Read, CAS, Mfence, and Rfence, we proceed similarly and let g map each subevent to their counterpart in the implementation.
- For e' = (t, -, (Get, (x, y, d), ())), there is some event  $e = (t, -, (Get^{TSO}, (x, y), (v))) \in f^{-1}(e')$ . We define  $g(e', aNRR_{n(y)}) = (e, aNRR_{n(y)})$  and  $g(e', aNLW_{n(y)}) = (e, aNLW_{n(y)})$ . We proceed similarly for Put, RCAS, and RFAA events.

• Finally for  $e' = (t, \_, (Wait, (d), ()))$ , there is in  $f^{-1}(e')$  some last event (in po order) of the form  $e = (t, \_, (SetIsEmpty, (d^N), true))$  confirming the set  $d^N$  tracking operations towards the last node N is empty. We define g(e', aWT) = (e, aMF).

We can see that  $g(\langle e', a' \rangle) = \langle e, a \rangle$  implies that f(e) = e' and that a is more restrictive than a'.

Each subevent in  $\mathcal{G}'.\mathcal{R}$  (resp.  $\mathcal{G}'.\mathcal{W}$ ) is mapped through g to a subevent in  $\mathcal{G}.\mathcal{R}$  (resp.  $\mathcal{G}.\mathcal{W}$ ) using the same stamp and location. Thus it is straightforward to define  $v_R'$ ,  $v_W'$ , rf', mo', nfo', and rao' by relying on their counterparts in  $\mathcal{G}$ . E.g.  $v_R'(s') \triangleq v_R(g(s'))$  and  $rf' \triangleq \{(s_1', s_2') \mid (g(s_1'), g(s_2')) \in rf\}$ . The well-formedness of  $v_R$ ,  $v_W$ , rf, mo, nfo, and rao trivially implies that of  $v_R'$ ,  $v_W'$ , rf', mo', nfo', and rao'. From this, we can define all the expected derived relations, including pfg', pfp', and  $ib' \triangleq (ippo' \cup iso' \cup rf' \cup pfg' \cup pfp' \cup nfo' \cup rb_i')^+$ . We then define  $so' \triangleq iso' \cup rf'_e \cup pfg' \cup nfo' \cup rb' \cup mo' \cup rao' \cup ([aNRW]; iso'^{-1}; rao') \cup ([Inst]; ib')$ , and as previously mentioned  $hb' \triangleq (ppo' \cup so')^+$ .

To show {RDMA<sub>RMW</sub>}-consistency, we are left to prove that  $\mathsf{ib}'$  and  $\mathsf{hb}'$  are irreflexive. For this, it is enough to show that  $g(\mathsf{ib}') \subseteq \mathsf{ib}$  and  $g(\mathsf{hb}') \subseteq \mathsf{hb} \triangleq (\mathsf{ppo} \cup \mathsf{so})^+$  since we know both  $\mathsf{ib}$  and  $\mathsf{hb}$  to be irreflexive.

For all subevent s', g(s') has a more restrictive stamp than s' (in most cases it is the same stamp, but for Wait the stamp aMF is more restrictive than aWT); this implies that  $g(\mathsf{ppo'}) \subseteq \mathsf{ppo}$ . Then, by definition, it is trivial to check that  $g(\mathsf{rf'}) \subseteq \mathsf{rf}$ ,  $g(\mathsf{mo'}) \subseteq \mathsf{mo}$ ,  $g(\mathsf{nfo'}) \subseteq \mathsf{nfo}$ ,  $g(\mathsf{ippo'}) \subseteq \mathsf{ippo}$ ,  $g(\mathsf{rf'}_e) \subseteq \mathsf{rf}_e$ ,  $g(\mathsf{iso'}) \subseteq \mathsf{iso}$ ,  $g(\mathsf{rb'}) \subseteq \mathsf{rb}$ , and  $g(\mathsf{rb'}) \subseteq \mathsf{rb}_i$ .

To finish the proof, we need the following crucial pieces:  $g(\mathsf{pfp'}) \subseteq \mathsf{ib}$ ,  $g(\mathsf{pfg'}) \subseteq \mathsf{ib}$ , and  $g(\mathsf{pfg'}) \subseteq \mathsf{hb}$ . In fact, it is enough to show that  $g(\mathsf{pfp'}) \subseteq \mathsf{ib}^?; \mathsf{pf}; \mathsf{ppo}^+$  and  $g(\mathsf{pfg'}) \subseteq \mathsf{pf}; \mathsf{ppo}^+$ . This is because  $\mathsf{pf}; \mathsf{ppo}^+ \subseteq \mathsf{ib}$ ,  $[\mathsf{aNLW}]; \mathsf{pf}; \mathsf{ppo}^+ \subseteq \mathsf{hb}$ , and the domain of  $g(\mathsf{pfg'})$  is included in  $\cup_n \mathcal{G}.\mathsf{aNLW}_n$  by definition.

Let us start with  $\operatorname{pfg}'$ , assuming  $((e_1',\operatorname{aNLW}_n),(e_2',\operatorname{aWT})) \in \operatorname{pfg}'$ . By definition they are of the form  $e_1' = (t, ., (., (., d), ()))$  and  $e_2' = (t, ., (\operatorname{Wait}, (d), ()))$ , for some t, d, and  $\operatorname{m}(e_1') \in \{\operatorname{Get},\operatorname{RCAS},\operatorname{RFAA}\}$ , with  $(e_1',e_2') \in \operatorname{po}'$  and n the remote node of this operation. By definition of the implementation and the abstraction,  $f^{-1}(e_1')$  contains two events  $e_1 = (t, ., (m, (...), (v)))$  with a similar method  $m \in \{\operatorname{Get}^{\operatorname{TSO}},\operatorname{RCAS}^{\operatorname{TSO}},\operatorname{RFAA}^{\operatorname{TSO}}\}$  and  $e_a = (t, ., (\operatorname{SetAdd}, (d^n, v), ()))$ , with  $e_1 \stackrel{\operatorname{po}}{\longrightarrow} e_a$ . Meanwhile  $f^{-1}(e_2')$  contains a last event  $e_2 = (t, ., (\operatorname{SetIsEmpty}, (d^N), \operatorname{true}))$  and an earlier event  $e_3 = (t, ., (\operatorname{SetIsEmpty}, (d^n), \operatorname{true}))$ , with  $e_3 \stackrel{\operatorname{po}^*}{\longrightarrow} e_2$ , confirming operations towards n are done (if n = N then  $e_2 = e_3$ ).

Since  $f(\mathbf{e}_a) = \mathbf{e}_1' \xrightarrow{\mathrm{po'}} \mathbf{e}_2' = f(\mathbf{e}_3)$  and f is an abstraction, we have  $\mathbf{e}_a \xrightarrow{\mathrm{po}} \mathbf{e}_3$ , i.e. the value v is added to  $d^n$  before the moment  $d^n$  is confirmed empty. By consistency (Definition 16), there is an in-between event  $\mathbf{e}_4 = (t, \neg, (\mathtt{SetRemove}, (d^n, v), ()))$  that removes this value, with  $\mathbf{e}_a \xrightarrow{\mathrm{po}} \mathbf{e}_4 \xrightarrow{\mathrm{po}} \mathbf{e}_3$ . From the definition of the implementation, such an event  $\mathbf{e}_4$  is immediately preceded (with maybe other  $\mathtt{SetRemove}$  in-between) by an event  $\mathbf{e}_p = (t, \neg, (\mathtt{Poll}, (n), (v)))$ . Now we argue that we necessarily have  $((\mathbf{e}_1, \mathtt{aNLW}_n), (\mathbf{e}_p, \mathtt{aWT})) \in \mathsf{pf}$ . From the well-formedness of  $\mathsf{pf}$ , we know that  $(\mathbf{e}_p, \mathtt{aWT})$  has a preimage  $(\mathsf{pf}$  is total and functional on its

range) and that this preimage outputs the value v. By consistency (Definition 16),  $e_1$ , with  $m(e_1) \in \{Get^{TSO}, RCAS^{TSO}, RFAA^{TSO}\}$ , is the only RDMA operation with output v. Thus  $(e_1, aNLW_n)$  is the preimage of  $(e_p, aWT)$  by pf.

Finally we have  $g(e'_1, aNLW_n) = (e_1, aNLW_n) \xrightarrow{pf} (e_p, aWT) \xrightarrow{ppo} (e_4, aMF) \xrightarrow{ppo^*} (e_3, aMF) \xrightarrow{ppo^*} (e_2, aWF) = g(e'_2, aWT)$ , which shows  $g(pfg') \subseteq pf; ppo^+$ .

For pfp' we have two cases. First, for a Put operation  $e'_1$ , having  $((e'_1, aNRW_n), (e'_2, aWT)) \in pfp'$  similarly implies  $g(e'_1, aNRW_n) \xrightarrow{pf; ppo^+} g(e'_2, aWT)$  for the same reasons. Second, for a successful remote RMW operation  $e'_1$ , having  $((e'_1, aNRW_n), (e'_2, aWT)) \in pfp'$ , with  $g(e'_1, aNRW_n) = (e_1, aNRW_n)$  instead implies  $(e_1, aNLW_n) \xrightarrow{pf; ppo^+} g(e'_2, aWT)$  for the same reasons. This is because the semantics of polls synchronises with the local write part of the operation, not with the remote write part. However, we do have  $g(e'_1, aNRW_n) \xrightarrow{ib} (e_1, aNLW_n)$  from the last component of ib, and thus  $g(pfg') \subseteq ib^2; pf; ppo^+$ .

Thus ib' and hb' are irreflexive, and  $\mathcal{G}'$  is  $\{\text{RDMA}_{RMW}^{WAIT}\}$ -consistent.

# $\mathbf{D}$ The RDMA $_{\mathrm{RMW}}^{\mathrm{TSO}}$ Memory Model

In this section, we present an operational (§D.1) and declarative model (§D.2) for RDMA<sup>TSO</sup> in the format of RDMA<sup>TSO</sup> as defined in [3], as well as an extension of their equivalence proof (§D.3 onwards).

The declarative format used in §D.2 (based on [3]) is slightly different from the one of C.1 (based on [4]), but they represent the same semantics.

### D.1 Operational Semantics

**Nodes and Threads.** We write  $\mathsf{Node} = \{1..N\}$  for the set of node identifiers, and  $\mathsf{Tid}$  for the set of thread identifiers. We write n (resp. t) to range over nodes (resp. threads), and given some node n we write  $\overline{n}$  to range over the set of all other nodes  $\mathsf{Node} \setminus \{n\}$ . Each thread runs on a particular node, so we write n(t) for the node the thread belongs to.

Note that the semantics of [3] assumes that the remote node  $\overline{n}$  of an operation is different from the local node n (i.e. they ignore loopback). As we extend their operational model, we keep their notations. However, as shown in [4], loopbacks are possible and follow exactly the same semantics.

**Memory.** Although all nodes can directly access all memory locations, whether an operation is towards local or remote memory is pivotal to our semantics, so we are always careful to note the node to which a memory location belongs. We write  $\mathsf{Loc}_n$  for the set of locations local to node n, and  $\mathsf{Loc} = \biguplus_n \mathsf{Loc}_n$  for the set of all locations. We use  $\mathsf{Loc}_{\overline{n}} = \mathsf{Loc} \setminus \mathsf{Loc}_n$  and write  $x^n, y^n, z^n$  for values in  $\mathsf{Loc}_n$ , respectively  $x^{\overline{n}}, y^{\overline{n}}, z^{\overline{n}}$  for  $\mathsf{Loc}_{\overline{n}}$ . When the node in question is sufficiently clear, we elide the superscript and instead simply write x or  $\overline{x}$  for local or remote locations respectively.

Values and Expressions. The language of expressions is standard and elided. We write  $v \in Val$  for values, with  $\mathbb{N} \subseteq Val$ , and  $e \in Exp$  for expressions. We write elocs(e) for the set of memory locations referenced in e, e[v/x] for the expression obtained by substituting all references to location x in e with value v, and [e] for the evaluation of e given it is closed, that is,  $elocs(e) = \emptyset$ . We use  $e^n$  for expressions where  $elocs(e^n) \subseteq Loc_n$ 

Commands and Programs. Commands are described by the  $C^n$  grammar below. CPU operations (CComm) are assignment, assumption of the value of a location, memory fence, compare-and-swap, and poll, which awaits the earliest completion notification of a remote operations towards  $\overline{n}$ .

RDMA operations (RComm) are either a 'get' of the form  $x:=\overline{y}$  which reads a remote location  $\overline{y}$  and writes its value to local location x, a 'put' ( $\overline{y}:=x$ ) which does the reverse, 'remote-CAS' (resp. 'remote-FAA') which executes a *remote* compare-and-swap (resp. fetch-and-add), and 'remote fence' which ensures all prior RDMA operations towards  $\overline{n}$  complete before any later RDMA operations towards  $\overline{n}$  execute. We note rRMW to cover both kind of remote read-modify-write operations, i.e. RCAS and RFAA.

Primitive operations (PComm) are CPU or RDMA operations, and commands (Comm) are the no-op, primitive operations, sequential composition (executes the first command, then the second), non-deterministic choice (executes one command or the other), and non-deterministic loop (executes the command some finite, possibly zero number of times).

A program P consists of a map from threads to commands, such that each  $t \in Tid$  is mapped to a command on n(t).

```
\begin{array}{ll} \mathsf{Comm} \ni \mathsf{C}^n ::= \mathtt{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*} & \mathsf{PComm} \ni \mathsf{c}^n ::= \mathsf{cc}^n \mid \mathsf{rc}^n \\ \mathsf{CComm} \ni \mathsf{cc}^n ::= x := e^n \mid \mathsf{assume}(x = v) \mid \mathsf{assume}(x \neq v) \mid \mathsf{mfence} \mid x := \mathsf{CAS}(y, e_1, e_2) \mid \mathsf{poll}(\overline{n}) \\ \mathsf{RComm} \ni \mathsf{rc}^n ::= x := \overline{y} \mid \overline{y} := x \mid x := \mathsf{RCAS}(\overline{y}, e_1, e_2) \mid x := \mathsf{RFAA}(\overline{y}, e) \mid \mathsf{rfence}(\overline{n}) \end{array}
```

**Store Buffers.** To permit the weak behaviours of TSO (i.e. write-read reordering), we assign each thread a *store buffer* B(t), which is a FIFO queue containing pending writes to memory by that thread. When a thread performs a CPU read, it reads the most recent entry for that location in its store buffer, if there is one, instead of the value in memory. The write at the head of the queue may be flushed to memory at any time, and mfence and CAS wait until the store buffer is empty before executing.

Queue Pairs. We follow the *simplified* operational model described in [3], and therefore consider a queue-pair structure comprising three FIFO queues:  $\mathbf{pipe}$ , which contains pending or in-progress RDMA operations;  $\mathbf{wb_R}$ , the remote write buffer, which contains pending writes to the memory of the remote node; and  $\mathbf{wb_L}$ , the local write buffer, which contains pending writes to the memory of the local node. The structure is shown in Fig. 17. Notice that under this simplified model, the transition between local and remote node in  $\mathbf{pipe}$  is continuous – we do not explicitly model the transition between local (yellow) and remote (pink) sides.

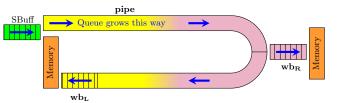


Fig. 17: Simple queue-pair structure.

**Remote Atomics.** To model the behaviour of RDMA atomic operations, we assign each node a *remote atomic lock* A(n), which is a boolean indicating whether an RDMA atomic is currently in progress *towards* that node.

Transitions of the Operational Semantics. We describe the rules governing the transitions between states, which comprise a program P, global memory M, store buffers B, queue pairs QP and remote atomic locks A.

Transitions take the form shown on the right, which should be read as: if  $\phi$  is true, then it is allowed for the system to transition from the state described by P...QP to the state described by P'...QP'.

$$\frac{\phi}{\mathsf{P},\mathsf{M},\mathsf{B},\mathsf{A},\mathsf{QP}\Rightarrow\mathsf{P}',\mathsf{M}',\mathsf{B}',\mathsf{A}',\mathsf{QP}'}$$

In practice, however, writing each transition rule in such a way would be verbose and hard to understand, as most transitions do not affect every part of the state. We can separate program transitions concerning P from hardware transitions concerning M, B, A, QP. In order to synchronise the two where necessary, we assign labels to certain transitions and require that a labelled program transition only occur if it is matched by a hardware transition with the same label (or vice-versa). Labels are of the form t:l where t is the thread executing at that step and  $l \in \mathsf{Lab}$  is the label of the operation. Silent transition, which affect only the program (resp. only hardware) are written with the empty label,  $\epsilon$ , and may be taken independently.

Fig. 18 shows the top-level rules of the operational semantics which govern this separation. We can henceforth consider the program and hardware transitions separately.

**Program Transitions.** Fig. 19 shows the program and command transitions (middle), labels (above) and expression rewriting rules (below). The transitions for non-remote commands are familiar from TSO. Notice that the transitions for get and put simply transition to skip with the relevant label; we know that this means there will be some relevant transition in the hardware. The transition to skip allows the program to continue executing, which we expect, as remote operations are handled asynchronously by the NIC.

This is similarly the case for the rules for remote-CAS and remote-FAA. The expressions involved are required to be closed, similarly to the rules for local write and CAS; the expressions must be evaluated before the transition.

$$\frac{P \xrightarrow{t:\varepsilon} P'}{P, M, B, A, QP \Rightarrow P', M, B, A, QP} \xrightarrow{M, B, A, QP \xrightarrow{t:\varepsilon} M', B', A', QP'} \frac{M, B, A, QP \xrightarrow{t:\varepsilon} M', B', A', QP'}{P, M, B, A, QP \xrightarrow{t:l} M', B', A', QP'} \frac{P \xrightarrow{t:l} P' \quad M, B, A, QP \xrightarrow{t:l} M', B', A', QP'}{P, M, B, A, QP \Rightarrow P', M', B', A', QP'}$$

Fig. 18: RDMA<sup>TSO</sup> operational semantics with the program and hardware transitions given in Fig. 19 and Fig. 20

It only makes sense to use a value, not an expression, in the label, since the corresponding hardware transition will only be concerned with values.

Hardware Domains. The upper section of Fig. 20 shows the hardware domains – that is, the states we are interested in *other* than the program. We have already described memory, store buffers, remote atomic locks and queue pairs, but note that the structures B, A, and QP are specifically maps from threads, nodes, and both, respectively, to the particular structures.

A remote atomic lock is a boolean,  $\bot$  (available) or  $\top$  (unavailable). A store buffer is a sequence of CPU writes and RDMA operations. A queue pair is a tuple of three sequences **pipe**,  $\mathbf{wb_R}$ , and  $\mathbf{wb_L}$ , where **pipe** may contain any of the operations described below except for a confirmation notification,  $\mathbf{wb_R}$  may contain NIC remote writes and NIC remote atomic writes, and  $\mathbf{wb_L}$  may contain NIC local writes and confirmation notifications.

- $y^{\overline{n}} := x^n$  denotes a put operation where the value of local memory location x is yet to be read (NIC local read);
- $y^{\overline{n}} := v$  denotes a NIC remote write of value v to remote location y, which occurs as the latter part of a put;
- ackp denotes the acknowledgement message returned by a put;
- $x^n := y^{\overline{n}}$  denotes a get operation where the value of the remote location y is yet to be read (pending NIC remote read)
- $x^n := v$  denotes a NIC local write of value v to local location x, which occurs as the latter part of a get or rRMW;
- RCAS $(z^n, x^{\overline{n}}, v, v')$  denotes a remote CAS towards remote location x, with expected value v, update value v', and returning to local location z;
- RFAA $(z^n, x^{\overline{n}}, v)$  similarly denotes a remote FAA towards x and returning to z, with increment value v;
- $y^{\overline{n}} :=_A v$  denotes a NIC remote write specifically in the case of an rRMW it is necessary for this to be disambiguated from the NIC remote write of a put, as we will see later;
- rfence( $\overline{n}$ ) denotes a remote fence towards node  $\overline{n}$ ;
- cn denotes a confirmation of a successful NIC remote write.

Hardware Transitions. All remote commands enter the queue-pair pipe via the thread's store buffer. When the program takes a transition step labelled with

Command transitions:

Fig. 19: The RDMA<sup>TSO</sup> program and command transitions

a remote CAS or FAA, the hardware takes a transition with a matching label, which adds that operation to the store buffer. The seventh transition rule allows remote commands at the head of the store buffer to enter the pipe of a queue pair, determined by their target node.

So far, we have seen that when an rRMW appears in the program, we can expect there to be a hardware transition which adds it to the store buffer, and later another hardware transition which removes it from the head of the store buffer and adds it to the suitable queue pair.

The final hardware transition introduces the queue-pair transitions, indicated by  $\rightarrow_{sqp}^4$ . When a particular queue pair takes a transition step, involving memory and the global remote atomic lock, the hardware takes a suitable corresponding transition. The queue-pair transitions merely involve a particular subset of the hardware states, so the relationship is straightforward. This separation is purely made for clarity and simplification of the queue-pair transition rules.

Queue-Pair Transitions. From Fig. 17, recall that remote operations enter the main **pipe** of the queue pair, then are suitably processed until they exit the pipe, possibly adding a write to  $\mathbf{wb_L}$  or  $\mathbf{wb_R}$  (or both). Note that the pipe grows to the left, so throughout,  $\alpha$  contains operations which are later in the program, while  $\beta$  contains earlier operations which have not yet completed.

The rules for remote fence, put, and get share a simple structure, where the premise for a transition either requires the operation to be at the head of the pipe  $\operatorname{sqp.pipe} = \alpha \cdot (\operatorname{operation})$ , or allows it to be in the middle of the pipe  $\operatorname{sqp.pipe} = \alpha \cdot (\operatorname{operation}) \cdot \beta$ , with some stipulation as to the operations allowed in  $\beta$ . In the prior case, the operation never executes before another, earlier operation; in the latter, it can execute before any operation in  $\beta$  which was issued before it. There may also be some requirement that buffers  $\operatorname{wb_L}$  or  $\operatorname{wb_R}$  contain no writes, due to PCIe guarantees:  $\operatorname{wb_L} \in \{\operatorname{cn}\}^*$  ( $\operatorname{wb_L}$  contains only confirmation notifications) or  $\operatorname{wb_R} = \epsilon$  (there are no operations in  $\operatorname{wb_R}$ ). Consider, for example, the first step of a put, which is a NIC local read described by rule 2. The value of location x is read from memory, so long as  $\operatorname{wb_L}$  has no pending writes and there are no other NIC local reads earlier in the pipe.

We can then describe the rules for rfence, put, and get at a high level:

Remote fence (rule 1) An refere may be removed from the pipe once it reaches the head (there are no earlier operations remaining to be processed). In combination with the fact that no other transition rule allows a step to be taken when there is an refere later in the pipe, this enforces the behaviour that all remote operations prior to an refere complete before it, and all later ones after it.

**Put** (rules 2-5) Rule 2: a NIC local read is performed, replacing the location x with its value in memory. Rule 3: the NIC remote write is sent to  $\mathbf{wb_R}$ , and an acknowledgement created in the pipe. Rule 4: the remote write is

<sup>&</sup>lt;sup>4</sup> SQP stands for simplified queue pair. We only considered the simplified three-buffer queue pair, so this disambiguation is technically unnecessary, but we maintain the notation for consistency with [3]

$$\begin{split} &\mathsf{M} \in \mathsf{Mem} \triangleq \mathsf{Loc} \to \mathsf{Val} \\ &\mathsf{A} \in \mathsf{RAMap} \triangleq \lambda n. \{\bot, \bot\} \\ &\mathsf{OP} \in \mathsf{SQPMap} \triangleq \lambda t. \{\lambda n(t), \mathsf{SQPair}_n^n) \\ &\mathsf{b} \in \mathsf{SBuff}_n \triangleq \{x^n := v, y^n := x^n, x^n := y^n, \mathsf{RCAS}(z^n, x^n, v), v'), \mathsf{RFA}(z^n, x^n, v), \mathsf{rfence}(n)\}^* \\ &\mathsf{sap} \in \mathsf{SQPair}_n^m \triangleq \mathsf{pipe}_n^m \times \mathsf{WBR}_n^m \times \mathsf{WBI}_n^m \\ &\mathsf{wb_L} \in \mathsf{WBL}_n^m \triangleq \{\mathsf{cn}, x^n := v\}^* \\ &\mathsf{wb_L} \in \mathsf{WBR}_n^m \triangleq \{y^m := v, y^n := a, v\}^* \\ &\mathsf{pipe} \in \mathsf{Pipe}_n^m \triangleq \left\{y^m := x^n, y^m := v, y^m := a, v, \mathsf{ack}_p, x^n := y^n, x^n := v\}^* \\ &\mathsf{pipe} \in \mathsf{Pipe}_n^m \triangleq \left\{y^m := v, y^n := a, v, \mathsf{ack}_p, x^n := y, x^n := v\}^* \\ &\mathsf{pipe} \in \mathsf{Pipe}_n^m \triangleq \left\{y^m := v, y^n := a, v, \mathsf{ack}_p, x^n := y, x^n := v\}^* \\ &\mathsf{pipe} \in \mathsf{Pipe}_n^m \triangleq \left\{y^m := v, y^n := a, v, \mathsf{ack}_p, x^n := y, x^n := v\}^* \\ &\mathsf{pipe} \in \mathsf{Pipe}_n^m \triangleq \left\{y^m := v, y^n := a, v, \mathsf{ack}_p, x^n := v, y^n := a, v\}^* \\ &\mathsf{pipe} \in \mathsf{Pipe}_n^m \triangleq \left\{y^m := v, y^n := a, v, \mathsf{ack}_p, x^n := y, x^n := v\}^* \\ &\mathsf{pipe} \in \mathsf{Pipe}_n^m \triangleq \left\{y^m := v, y^n := a, v, \mathsf{ack}_p, x^n := y, x^n := v\}^* \\ &\mathsf{pipe} \in \mathsf{Pipe}_n^m \triangleq \left\{y^m := v, y^n := a, v, \mathsf{ack}_p, x^n := y, x^n := v\}^* \\ &\mathsf{pipe} \in \mathsf{Pipe}_n^m \triangleq \left\{y^m := v, y^n := a, v, \mathsf{ack}_p, x^n := a, v, \mathsf{ack}_p, x^n$$

Fig. 20: RDMA<sup>TSO</sup> simplified hardware domains (above) and hardware transitions (below)

```
\mathbf{sqp.pipe} = \alpha \cdot (\mathtt{rfence}(\overline{n}))
                                                                                                             \overline{\mathsf{M},\mathsf{A},\mathbf{sqp}\to_{\mathbf{sqp}}\mathsf{M},\mathsf{A},\mathbf{sqp}[\mathbf{pipe}\mapsto\alpha]}
                               \begin{aligned} \mathbf{sqp.pipe} &= \alpha \cdot (\overline{y} := \underline{x}) \cdot \beta \quad \mathbf{wb_L} \in \big\{ \mathtt{cn} \big\}^* \\ \beta \in \big\{ y' := v', y' :=_A v', \overline{y'} := v', x' := \overline{y'}, \mathtt{RCAS}(z, \overline{x}, v, v'), \mathtt{RFAA}(z, \overline{x}, v), \mathtt{ack_p} \big\}^* \end{aligned}
                                                                              M, A, \mathbf{sqp} \rightarrow_{\mathsf{sqp}} M, A, \mathbf{sqp}[\mathbf{pipe} \mapsto \alpha \cdot (\overline{y} := \mathsf{M}(x)) \cdot \beta]
                                                               \begin{aligned} \mathbf{sqp.pipe} &= \alpha \cdot (\overline{y} := v) \cdot \beta \\ \beta \in \left\{ x' := \overline{y'}, x' := v', \mathsf{ack_p} \right\}^* \\ \mathbf{sqp'} &= \mathbf{sqp[pipe} \mapsto \alpha \cdot \mathsf{ack_p} \cdot \beta] [\mathbf{wb_R} \mapsto (\overline{y} := v) \cdot \mathbf{sqp.wb_R}] \end{aligned}
                                                                                                                                M, A, \mathbf{sqp} \rightarrow_{\mathbf{sqp}} M, A, \mathbf{sqp}'
                                                                                                                                                                                                                                         sqp.pipe = \alpha \cdot ack_p
                                                                                                                                                                                         \operatorname{\mathbf{sqp}}' = \operatorname{\mathbf{sqp}}[\operatorname{\mathbf{pipe}} \mapsto \alpha][\operatorname{\mathbf{wb_L}} \mapsto \operatorname{\mathtt{cn}} \cdot \operatorname{\mathbf{sqp.wb_L}}]
                                   \mathbf{sqp.wb_R} = \alpha \cdot (\overline{y} := v)
M, A, \operatorname{sqp} \to_{\operatorname{sqp}} M[\overline{y} \mapsto v], A, \operatorname{sqp}[\operatorname{\mathbf{wb}}_{\mathbf{R}} \mapsto \alpha]
                                                                                                                                                                                                                              M, A, \mathbf{sqp} \rightarrow_{sqp} M, A, \mathbf{sqp}'
                                                                                                                                                                            \beta \in \left\{ x' := \overline{y'}, x' := v', \mathtt{ack}_{\mathtt{p}} \right\}^*
                                                                    \mathbf{sqp.pipe} = \alpha \cdot (x := \overline{y}) \cdot \beta
                                                                             \operatorname{sqp.wb}_{\mathbf{R}} = \varepsilon \quad \operatorname{sqp'} = \operatorname{sqp}[\operatorname{\mathbf{pipe}} \mapsto \alpha \cdot (x := M(\overline{y})) \cdot \beta]
                                                                                                                                \overline{\mathsf{M},\mathsf{A},\mathbf{sqp} \to_{\mathsf{sqp}} \mathsf{M},\mathsf{A},\mathbf{s}} \overline{\mathbf{qp}}'
                                                                                                                                   \mathbf{sqp.pipe} = \alpha \cdot (x := v)
                                                                          \mathbf{sqp}' = \mathbf{sqp}[\mathbf{pipe} \mapsto \alpha][\mathbf{wb_L} \mapsto \mathbf{cn} \cdot (x := v) \cdot \mathbf{sqp.wb_L}]
                                                                                                                                M, A, \mathbf{sqp} \rightarrow_{\mathbf{sqp}} M, A, \mathbf{sqp}'
                                                                                                      \frac{\mathbf{sqp.wb_L} = \alpha \cdot (x := v) \cdot \beta \quad \beta \in \big\{ \mathsf{cn} \big\}^*}{\mathbf{sqp'} = \mathbf{sqp}[\mathbf{wb_L} \mapsto \alpha \cdot \beta]} \\ \frac{\mathbf{m, A, sqp} \rightarrow_{\mathsf{sqp}} \mathsf{M}[x \mapsto v], \mathsf{A, sqp'}}{\mathsf{M}[x \mapsto v], \mathsf{A, sqp'}}
                                                                                    \mathbf{sqp.pipe} = \alpha \cdot \mathtt{RCAS}(z, \overline{x}, v, v') \cdot \beta \quad \mathbf{sqp.wb_R} = \epsilon
                                                                        \mathsf{A}(n(\overline{x})) = \bot \quad \mathsf{M}(\overline{x}) \neq v \quad \beta \in \left\{x' := \overline{y'}, \overline{x'} := v', \mathsf{ack}_{\mathsf{p}}\right\}^*
                                                                                                      \operatorname{\mathbf{sqp}}' = \operatorname{\mathbf{sqp}}[\operatorname{\mathbf{pipe}} \mapsto \alpha \cdot (z := \operatorname{\mathsf{M}}(\overline{x})) \cdot \beta]
                                                                                                                                M, A, \mathbf{sqp} \rightarrow_{\mathsf{sqp}} M, A, \mathbf{sqp}'
                                                                                  \begin{aligned} \mathbf{sqp.pipe} &= \alpha \cdot \mathtt{RCAS}(z, \overline{x}, v, v') \cdot \beta \quad \mathbf{sqp.wb_R} = \epsilon \\ \mathsf{M}(\overline{x}) &= v \quad \beta \in \left\{ x' := \overline{y'}, x' := v', \mathtt{ack_p} \right\}^* \\ \mathsf{A}(n(\overline{x})) &= \bot \quad \mathsf{A}' = \mathsf{A}[n(\overline{x}) \mapsto \top] \end{aligned}
                                                              M, A, \operatorname{sqp} \to_{\operatorname{sqp}} M, A', \operatorname{sqp}[\operatorname{pipe} \mapsto \alpha \cdot (z := v) \cdot (\overline{x} :=_A v') \cdot \beta]
                                                                                         \begin{array}{l} \mathbf{sqp.pipe} = \alpha \cdot \mathtt{RFAA}(z,\overline{x},v) \cdot \beta \quad \mathbf{sqp.wb_R} = \epsilon \\ \mathsf{M}(\overline{x}) + v = v' \quad \beta \in \left\{x' := \overline{y'}, x' := v', \mathtt{ack_p}\right\}^* \\ \mathsf{A}(n(\overline{x})) = \bot \quad \mathsf{A}' = \mathsf{A}[n(\overline{x}) \mapsto \top] \end{array}
                                                              \mathsf{M}, \mathsf{A}, \mathbf{sqp} \to_{\mathsf{sqp}} \mathsf{M}, \mathsf{A}', \mathbf{sqp}[\mathbf{pipe} \mapsto \alpha \cdot (z := v) \cdot (\overline{x} :=_A v') \cdot \beta]
                                          \mathbf{sqp.pipe} = \alpha \cdot (\overline{x} :=_A v) \cdot \beta
                                         \mathbf{wb_R}' = (\overline{x} :=_A v) \cdot \mathbf{sqp.wb_R}\beta \in \left\{ x' := \overline{y'}, x' := v', \mathbf{ack_p} \right\}
                                                                                                                                                                                                                                \mathbf{sqp.wb_R} = \alpha \cdot (\overline{x} :=_A v)\mathbf{A}' = \mathbf{A}[n(\overline{x}) \mapsto \bot]
                                                                                                                                                                                                                                    \mathbf{sqp}' = \mathbf{sqp}[\mathbf{wb_R} \mapsto \alpha]
                   \mathbf{sqp}' = \mathbf{sqp}[\mathbf{pipe} \mapsto \alpha \cdot \beta][\mathbf{wb_R} \mapsto \mathbf{wb_R}']
                                                                                                                                                                                                                 \mathsf{M}, \overline{\mathsf{A}}, \overline{\mathbf{sqp}} \to_{\operatorname{sqp}} \mathsf{M}[\overline{x} \mapsto v], \mathsf{A}', \overline{\mathbf{sqp}}'
                                                 \mathsf{M},\mathsf{A},\mathbf{sqp}\to_{\mathsf{sqp}}\mathsf{M},\mathsf{A},\mathbf{sqp}'
```

Fig. 21: Queue-pair transitions of the simplified RDMA<sup>TSO</sup> operational semantics

committed to memory once it reaches the head of the queue. Rule 5: the acknowledgement in the pipe is converted to a confirmation notification in  $\mathbf{wb_L}$ , so that it can be polled.

Get (rules 6-8) Rule 6: a NIC remote read replaces the location  $\overline{y}$  with its value in memory. Rule 7: the NIC local write is sent to  $\mathbf{wb_L}$ , with a confirmation notification for the purpose of polling. Rule 8: the local write is committed to memory once there are no pending earlier writes in the queue.

Now, consider the rules for rRMWs. These rules are more complicated due to the need to check and update the remote atomic lock for the target node, which we see as  $A(n(\overline{x})) = \bot$  (the remote atomic lock for the target node is available), and  $A' = A[n(\overline{x}) \mapsto \top]$  (update the remote atomic lock for the target node to indicate it is busy). We also have distinct rules for success and failure of RCAS, depending on whether the remote memory location holds the expected value ( $M(\overline{x}) = v$  or  $M(\overline{x}) \neq v$ ).

The rules can then be interpreted as follows:

- (Rule 9) A failed RCAS read the remote memory location does not hold the expected value. This read can only occur when the remote atomic lock is available, otherwise it would violate the atomicity guarantee. The value of  $\overline{x}$  is read, and a NIC local write is added to the pipe to return that value to z. This is then handled by the same rules as for a get. The remote atomic lock is not obtained, since the remote location will not be written to.
- (Rule 10) A successful RCAS read the remote location contains the expected value. Once again, this requires that the remote atomic lock be available, and it is also obtained to ensure atomicity until the remote location is written to. A NIC local write is added to the pipe (similarly to 9), and a NIC remote atomic write to update the remote location is also added.
- (Rule 11) The remote read of an RFAA this is unconditionally successful. It is very similar to a successful RCAS, but the value for the NIC remote atomic write is calculated by adding v to the value of  $\overline{x}$  in memory.
- (Rule 12) A NIC remote atomic write in the pipe is processed into  $\mathbf{wb_R}$  similarly to a regular NIC remote write.
- (Rule 13) A NIC remote atomic write is committed to memory, and the remote atomic lock is released.

#### D.2 Declarative Semantics

A declarative semantics, in contrast to an operational one, describes only the events that occur in a system, not the state of the system itself. An execution is represented by a graph, with various relations over events. For example, given an event r, we say that it "reads-from" event w if r reads the value written to memory by w. We write  $(w,r) \in rf$  in this case. We then constrain these relations suitably to only allow execution graphs which make sense in the context of a program: considering rf again, we would naturally only allow  $(w,r) \in rf$  if the values of the read and write match.

Then, we know that an execution of a program is allowed if the graph of the execution is consistent. Contrast this with the operational semantics: there, our guarantee comes from the individual transition rules; here, it is due to the overall structure of the graph.

**Events and Executions.** An *execution* is a graph comprising a set of events and several relations over events; events are represented as graph nodes, and the relations are edges. An event has a unique identifier  $\iota$ , is created by a thread  $t \in \mathsf{Tid}$ , and has an event label  $l \in \mathsf{ELab}$  which describes the event.

**Definition 17 (Labels and events).** Each event label is associated to a node n. The set of event labels of node n is denoted by  $l \in \mathsf{ELab}_n$ , where l is a tuple with one of the following forms:

```
• (CPU) local read: l = 1R(x^n, v_r)
```

- (CPU) local write:  $l = 1W(x^n, v_w)$
- (CPU)  $CAS: l = CAS(x^n, v_r, v_w)$
- (CPU) memory fence: l = F
- (CPU) poll:  $l = P(\overline{n})$

- NIC local read:  $l = nlR(x^n, v_r, \overline{n})$
- $\bullet \ \mathit{NIC} \ \mathit{remote} \ \mathit{write:} \ l = \mathtt{nrW}(y^{\overline{n}}, v_w)$
- NIC remote read:  $l = nrR(y^{\overline{n}}, v_r)$
- NIC local write:  $l = nlW(x^n, v_w, \overline{n})$
- $NIC\ fence: l = nF(\overline{n})$
- NIC atomic remote read:  $l = \text{narR}(y^{\overline{n}}, v_r)$
- NIC atomic remote write:  $l = \text{narW}(y^{\overline{n}}, v_w)$

The set of event labels are defined  $\mathsf{ELab} \triangleq \bigcup_n \mathsf{ELab}_n$ . An event,  $\mathsf{e} \in \mathsf{Event}$ , is a triple  $(\iota, t, l)$ , where  $\iota \in \mathbb{N}$ ,  $t \in \mathsf{Tid}$  and  $l \in \mathsf{ELab}_{n(t)}$ .

We distinguish between events associated with the CPU (left) or NIC (right), with the prefix  ${\tt n}$  used for all NIC event labels. Note that a put, get, or rRMW is modelled by multiple events: a put  $\overline{x}:=y$  comprises a NIC local read event of type  ${\tt nlR}$  (on y) followed by a NIC remote write event  ${\tt nrW}$  (on  $\overline{x}$ ); conversely a get  $x:=\overline{y}$  comprises events of type  ${\tt nrR}$  (on  $\overline{y}$ ) and  ${\tt nlW}$  (on x). A successful rRMW (either successful RCAS or RFAA) is modelled by three events of type  ${\tt narR}$ ,  ${\tt narW}$  and  ${\tt nlW}$ , while a failed rRMW (RCAS only) is modelled by only  ${\tt narR}$  and  ${\tt nlW}$ .

For a given label l, we write  $\mathsf{type}(l), \mathsf{loc}(l), v_r(l), v_w(l), n(l)$  and  $\overline{n}(l)$  for the type, location, value read or written, and local or remote node, where applicable. For example, consider  $l = \mathsf{nlR}(x^n, v_r, \overline{n})$ :

```
• type(nlR(x^n, v_r, \overline{n})) = nlR
```

- $loc(nlR(x^n, v_r, \overline{n})) = x$
- $v_{\mathbf{r}}(\mathtt{nlR}(x^n, v_{\mathbf{r}}, \overline{n})) = v_{\mathbf{r}}$
- $v_{\rm w}({\tt nlR}(x^n,v_{\rm r},\overline{n}))$  is undefined
- $\bullet \ n(\mathtt{nlR}(x^n,v_{\mathbf{r}},\overline{n})) = n$
- $\overline{n}(\mathtt{nlR}(x^n, v_r, \overline{n})) = \overline{n}$

We write  $\iota(e)$ , t(e), l(e) for the relevant constituents of an event tuple  $e = (\iota, t, l)$ . We lift the functions on event labels to functions on events, for example  $\mathsf{type}(e) \triangleq \mathsf{type}(l(e))$ .

**Difference with** MOWGLI. The labels of this declarative semantics (à la [3]) roughly corresponds to the stamps of MOWGLI (à la [4]) in the main paper. E.g.

a NIC local read has a label nlR here and corresponds to the stamp (family) anlR in Fig. 9. However, there is one major discrepancy. The declarative semantics of this section distinguishes between NIC remote writes performed by put operations (label nrW) and performed by rRMW operations (label narW), while they both correspond to the single stamp  $aNRW_n$ .

The main reason is that this semantics, by decomposing operations into multiple events, creates a po ordering between the remote write and local write parts of a (successful) remote RMW. As such, we cannot enforce a ppo ordering between the two parts (narW and nlW) as they might not finish in order, but we can enforce a ppo ordering between the remote write of a put and later local writes (nrW and nlW), making the semantics more straightforward. With MOWGLI, each operation generates a single event, and there is no po ordering between subevents of the same event. Thus we can add a dependency between aNRW<sub>n</sub> and aNLW<sub>n</sub> (cell G10 in Fig. 9), and it will not create an internal dependency within the same rRMW operation.

A secondary reason is that the two labels (nrW and narW) correspond to different behaviours of the operational semantics. Making the distinction renders the equivalence proof more tractable.

Issue and Observation Points. Some types of events do not occur instantaneously: for example, a local write event 1W first enters the store before, before later being committed to memory. We therefore distinguish between the point at which an event is *issued* by the CPU or NIC, and the point at which it is *observed*, when its effect becomes visible in memory. An event is *instantaneous* if it either has no visible effect on memory, or if it affects memory immediately, as is the case for a local CAS operation. For instantaneous events, the issue and observation points coincide.

**Notation.** Once again, we follow and extend the notation of [3]. For a set A and relations  $r, r_1, r_2$ , we write:

```
r^+ for the transitive closure of r;

r^{-1} for the inverse of r;

r|_{A} \triangleq r \cap (A \times A) for the restriction of r to set A;

[A] \triangleq \{(a,a) \mid a \in A\} for the identity relation

r_1; r_2 \triangleq \{(a,b) \mid \exists c.(a,c) \in r_1 \land (c,b) \in r_2\} for relational composition;

r|_{\text{imm}} \triangleq r \setminus (r;r) for the immediate edges in r, when r is a strict partial order.
```

For a set of events E, location x and label type X, we also define:

```
\begin{split} E_x &\triangleq \{\mathbf{e} \in E \mid \mathbf{loc}(\mathbf{e}) = x\}, \text{ the events towards } x; \\ E.\mathbf{X} &\triangleq \{\mathbf{e} \in E \mid \mathbf{type}(\mathbf{e}) = \mathbf{X}\}, \text{ the events of type X}; \\ E.\mathcal{R} &\triangleq E.\mathbf{1R} \cup E.\mathbf{CAS} \cup E.\mathbf{n1R} \cup E.\mathbf{nrR} \cup E.\mathbf{narR}, \text{ the set of reads}; \\ E.\mathcal{W} &\triangleq E.\mathbf{1W} \cup E.\mathbf{CAS} \cup E.\mathbf{n1W} \cup E.\mathbf{nrW} \cup E.\mathbf{narW}, \text{ the set of writes}; \\ E.\mathbf{Inst} &\triangleq E \setminus (E.\mathbf{1W} \cup E.\mathbf{n1W} \cup E.\mathbf{nrW} \cup E.\mathbf{narW}), \text{ the set of instantaneous events}. \end{split}
```

Finally, we define the following relations:

```
\begin{array}{ll} \textbf{Same-location:} & \mathsf{sloc} \triangleq \left\{ (\mathsf{e},\mathsf{e}') \in \mathsf{Event}^2 \mid \mathsf{loc}(\mathsf{e}) = \mathsf{loc}(\mathsf{e}') \right\} \\ \textbf{Same-thread:} & \mathsf{sthd} \triangleq \left\{ (\mathsf{e},\mathsf{e}') \in \mathsf{Event}^2 \mid t(\mathsf{e}) = t(\mathsf{e}') \right\} \\ \textbf{Same-queue-pair:} & \mathsf{sqp} \triangleq \left\{ (\mathsf{e},\mathsf{e}') \in \mathsf{Event}^2 \mid t(\mathsf{e}) = t(\mathsf{e}') \land \overline{n}(\mathsf{e}) = \overline{n}(\mathsf{e}') \right\} \end{array}
```

Note that these relations are all symmetric, and  $sqp \subseteq sthd$ . Given events E, we write E.sloc for  $sloc|_E$ , likewise for E.sthd and E.sqp.

**Definition 18 (Pre-executions).** A pre-execution is a tuple  $G = \langle E, po, rf, mo, pf, nfo, rao \rangle$ , where:

- $E \subseteq \text{Event}$  is the set of events and includes a set of initialisation events,  $E^0 \subseteq E$ , comprising a single write with label  $1\mathbb{W}(x,0)$  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  $E^0 \times (E \setminus E^0) \subseteq \text{po}$ .
- rf  $\subseteq E.W \times E.R$  is the 'reads-from' relation on events of the same location with matching values; i.e.  $(a,b) \in \mathsf{rf} \Rightarrow (a,b) \in \mathsf{sloc} \wedge v_w(a) = v_r(b)$ . Moreover, rf is total and functional on its range: every read in E.R is related to exactly one write in E.W.
- $\operatorname{mo} \triangleq \bigcup_{x \in \mathsf{Loc}} \operatorname{mo}_x$  is the 'modification-order', where each  $\operatorname{mo}_x$  is a strict total order on  $E.\mathcal{W}_x$  with  $E_x^0 \times (E.\mathcal{W}_x \setminus E_x^0) \subseteq \operatorname{mo}_x$  describing the order in which writes on x reach the memory.
- pf ⊆ (E.nlW ∪ E.nrW) × E.P is the 'polls-from' relation, relating earlier (in program-order) 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).
- 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 $\cup E.$ narR),  $b \in (E.$ nrW $\cup E.$ narW), then  $(a,b) \in$  nfo $\cup$  nfo $^{-1}$ .
- rao  $\triangleq \bigcup_{n \in \mathsf{Node}} \mathsf{rao}_n$  is the 'remote-atomic-order', where each rao<sub>n</sub> is a strict total order on  $\{e \mid e \in E.\mathtt{narR} \land \overline{n}(e) = n\}$  describing the order in which remote atomics towards n are executed.

The definitions of po, rf and mo are familiar from TSO, while pf and nfo are introduced in [3]. As mentioned previously, nfo represents the PCIe guarantee that a NIC local read flushes pending NIC remote writes on the same queue pair, and likewise for NIC local reads/writes. We introduce rao, which totally orders NIC remote atomic reads towards a given node and help enforce the rRMW atomicity guarantee.

**Derived Relations.** Given a pre-execution  $\langle E, po, rf, mo, pf, nfo, rao \rangle$ , we define the following *derived* relations:

•  $\mathsf{rb} \triangleq (\mathsf{rf}^{-1}; \mathsf{mo}) \setminus [E]$  is the *reads-before* relation, relating each read r to writes that are  $\mathsf{mo}$ -after the write from which r reads.

- rf<sub>i</sub>  $\triangleq$  [1W]; (rf  $\cap$  sthd); [1R] is the rf-buffer relation, which includes rf edges only for CPU operations on the same thread, which thus share a store buffer; therefore when  $w \xrightarrow{\text{rf}_i} r$ , it may be that the write w is not yet visible (committed to memory) when it is read by r, since CPU reads check the store buffer.
- $\mathsf{rf}_{\mathsf{e}} \triangleq \mathsf{rf} \setminus \mathsf{rf}_{\mathsf{i}}$  is the  $\mathsf{rf}_{\mathsf{i}}$ -complement: if  $w \xrightarrow{\mathsf{rf}_{\mathsf{e}}} r$ , then r only occurs after w is observable.
- $rb_i \triangleq [1R]; (rb \cap sthd); [1W]$  is the rb-buffer relation, analogously.
- ar  $\triangleq$  [narW]; (po $|_{\mathrm{imm}}^{-1}$ ) is the *atomic-write-to-read* relation, connecting the remote write of a successful rRMW to their corresponding read.

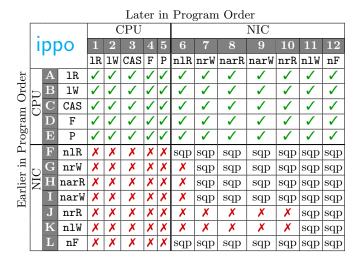
Note that these derived relations contain no additional information. We introduce them for ease and brevity of notation.

**Preserved Program Order.** We identify which events in po are *issued* in order, and which are *observed* in order. The observation point of an event is no earlier than its issue point, so two events in po are only observed in order if they are issued in order. Furthermore, when the po-earlier event is instantaneous, the events are observed in order if and only if they are issued in order.

We therefore define two relations: ippo, the issue-preserved-program-order relation, and oppo, the observation-preserved-program-order relation, where oppo  $\subseteq$  ippo  $\subseteq$  po. The tables in Fig. 22 show these relations. Each row indicates the po-earlier event, while each column indicates that which is po-later. A cell labelled  $\checkmark$  indicates the event pair is in ippo (resp. oppo) and must be issued (resp. observed) in program order, while  $\checkmark$  indicates they are not in ippo/oppo and may be issued/observed out of program order. The label sqp indicates that the events are in ippo/oppo if they are events on the same queue pair.

We can observe high-level reordering rules by looking at each quadrant of the two tables, which partition the event pairs by their categorisation as CPU or NIC events. The top left quadrant contains pairs of CPU events. Observe that CPU events are always issued in program order, and only an earlier CPU write may be observed out of order, as all other CPU events are instantaneous. The bottom left quadrants shows that an earlier NIC event may always be issued or observed after a later CPU event, matching our intuition that NIC events execute concurrently, as if in a separate thread; conversely the top right shows that earlier CPU events always complete before later NIC events. In the bottom right quadrant, we can see that a pair of NIC events are only ordered if they are on the same queue pair.

The relations ippo and oppo differ in only six cells. A CPU write may be buffered and hence not observed by a later CPU read or poll (B1 and B5). Other CPU writes and CAS or fence operations go via the store buffer, so earlier writes will be observed first. Similarly, a remote fence may be observed before an earlier NIC remote (atomic) write (resp. local), if that write is buffered in  $\mathbf{wb_R}$  (resp.  $\mathbf{wb_L}$ ) (G12 and I12, resp. K12). Finally, a po-later nlW may be observed before a po-earlier narW (I11). This occurs specifically in the case where both are created by the same rRMW, because the writes are sent to  $\mathbf{wb_L}$  and  $\mathbf{wb_R}$  respectively and may be committed in either order.



Later in Program Order															
				CPU				NIC							
	oppo			1	2	3	4	5	6	7	8	9	10	11	12
Earlier in Program Order				1R	lW	CAS	F	Р	nlR	nrW	narR	narW	nrR	nlW	nF
	CPU	A	1R	<b>\</b>	1	1	1	✓	<b>✓</b>	1	<b>√</b>	<b>✓</b>	1	<b>✓</b>	1
		В	1W	X	1	<b>✓</b>	1	X	<b>✓</b>	<b>✓</b>	<b>✓</b>	<b>✓</b>	1	<b>✓</b>	1
		$\mathbf{C}$	CAS	<b>✓</b>	1	<b>✓</b>	1	✓	<b>✓</b>	<b>✓</b>	<b>✓</b>	<b>✓</b>	1	<b>✓</b>	<b>✓</b>
		D	F	<b>\</b>	1	<b>✓</b>	1	✓	✓	<b>✓</b>	<b>✓</b>	<b>✓</b>	1	<b>✓</b>	<b>✓</b>
		E	P	>	✓	<b>\</b>	✓	✓	<b>✓</b>	<b>\</b>	<b>✓</b>	>	1	>	<b>✓</b>
	NIC	F	nlR	X	X	Х	Х	X	$\operatorname{sqp}$	$\operatorname{sqp}$	sqp	$\operatorname{sqp}$	$\operatorname{sqp}$	$\operatorname{sqp}$	$\operatorname{sqp}$
		G	nrW	X	Х	Х	Х	X	Х	$\operatorname{sqp}$	sqp	sqp	sqp	$\operatorname{sqp}$	X
		Η	narR	X	Х	Х	Х	Х	Х	$\operatorname{sqp}$	sqp	sqp	$\operatorname{sqp}$	$\operatorname{sqp}$	$\operatorname{sqp}$
		Ι	narW	X	X	X	X	X	X	$\operatorname{sqp}$	sqp	$\operatorname{sqp}$	$\operatorname{sqp}$	X	X
		J	nrR	X	X	X	X	X	X	X	X	X	X	$\operatorname{sqp}$	$\operatorname{sqp}$
		K	nlW	X	X	X	X	X	X	X	X	X	X	$\operatorname{sqp}$	X
		L	nF	X	X	Х	Х	X	$\operatorname{sqp}$	$\operatorname{sqp}$	sqp	$\operatorname{sqp}$	$\operatorname{sqp}$	$\operatorname{sqp}$	$\operatorname{sqp}$

Fig. 22: The RDMA<sup>TSO</sup> ordering constraints on ippo (above) and oppo (below), 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.

**Definition 19 (Executions).** A pre-execution  $G = \langle E, po, rf, mo, pf, nfo, rao \rangle$  is well-formed if the following hold for all  $w, r, w_1, w_2, p_2$ :

- 1) Poll events poll-from the oldest non-polled remote operation on the same queue pair:
  - $\begin{array}{l} \textit{if } w_1 \in G. \texttt{nlW} \cup G. \texttt{nrW} \textit{ and } w_1 \xrightarrow{\texttt{po} \cap \texttt{sqp}} w_2 \xrightarrow{\texttt{pf}} p_2, \textit{ then there exists } p_1 \textit{ such that } w_1 \xrightarrow{\texttt{pf}} p_1 \xrightarrow{\texttt{po}} p_2. \end{array}$
- 2) 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  $w \in G.nrW$  then  $(r, w) \in po|_{imm}$  for some  $r \in G.nlR$ . The case  $w \in G.nlW$  is handled by (6) below.
- 3) Read and write events of a put (resp. get) have matching values: if  $(r,w) \in G.po|_{imm}$ ,  $type(r) \in \{nlR,nrR\}$  and  $type(w) \in \{nlW,nrW\}$ , then  $v_r(r) = v_w(w)$ .
- 4) Each rRMW operation corresponds to either an atomic remote read followed by a local write, or an atomic remote read, followed by an atomic remote write, followed by a local write: 1) if  $r \in G$ .narR then  $(r, w_1) \in \mathsf{po}|_{imm}$  for some  $w_1 \in G$ .narW  $\cup G$ .nlW, and if  $w_1 \in G$ .narW then  $(w_1, w_2) \in \mathsf{po}|_{imm}$  for some  $w_2 \in G$ .nlW, and 2) if  $w_1 \in G$ .narW then  $(r, w_1) \in \mathsf{po}|_{imm}$  for some  $r \in \mathsf{narR}$ , and  $(w_1, w_2) \in \mathsf{po}|_{imm}$  for some  $w_2 \in \mathsf{nlW}$ . The case for  $w_2 \in \mathsf{nlW}$  is handled by (6) below.
- 5) Remote atomic read and local write events of an rRMW have matching values: if  $(r,w) \in G.po|_{imm}$ , type(r) = narR and type(w) = nlW, then  $v_r(r) = v_w(w)$ ; and if  $(r,w_1), (w_1,w_2) \in G.po|_{imm}$ , type(r) = narR,  $type(w_1) = narW$  and  $type(w_2) = nlW$ , then  $v_r(r) = v_w(w_2)$ .
- 6. (2) and (4) auxiliary in the case of  $w \in nlW$ . If  $w \in G.nlW$  then either:
  - 1)  $(r, w) \in po|_{imm} \text{ for some } r \in G.nrR \text{ or }$
  - 2)  $(r, w) \in po|_{imm}$  for some  $r \in G$ .narR or
  - 3)  $(r, w'), (w', w) \in po|_{imm}$  for some  $r \in G$ .narR and  $w' \in G$ .narW.

An execution is a pre-execution (Def. 18) that is well-formed.

Given an execution G, we write G.E, G.mo, G.ippo and so forth to project the components and derived relations of G. When the execution is question is clear, we simply write E, mo or similar.

**Definition 20** (RDMA<sup>TSO</sup>-consistency). An execution  $\langle E, po, rf, mo, pf, nfo, rao \rangle$  is RDMA<sup>TSO</sup>-consistent iff 1) ib is irreflexive; and 2) ob is irreflexive, where:

```
\begin{split} \text{ib} &\triangleq \big( \text{ippo} \cup \text{rf} \cup \text{pf} \cup \text{nfo} \cup \text{rb}_i \cup \big( \text{ob}; [\texttt{Inst}] \big) \big)^+ & \text{(`issued-before')} \\ \text{ob} &\triangleq \big( \text{oppo} \cup \text{rf}_e \cup \big( [\texttt{nlW}]; \text{pf} \big) \cup \text{nfo} \cup \text{rb} \cup \text{mo} \cup \text{rao} \cup \big( \text{ar}; \text{rao} \big) \cup \big( [\texttt{Inst}]; \text{ib} \big) \big)^+ \\ & \text{(`observed-before')} \end{split}
```

These relations extend ippo and oppo respectively to describe the issue and observation orders across threads and nodes. They are required to be irreflexive, i.e. an event cannot be issued or observed before itself.

The remaining components of ib are (a) rf: if  $w \xrightarrow{rf} r$  then w was at least issued (if not observed) before r – recall that if the read and write are both CPU events on the same thread, w may not be observable; (b) pf: similarly  $w \xrightarrow{pf} p$  only if w was issued before p; (c) nfo: NIC events arrive in  $\mathbf{wb_L/wb_R}$  in the order they are issued; (d)  $\mathbf{rb_i}$ : if  $r \xrightarrow{\mathbf{rb_i}} w$ , then r must be issued before w, otherwise r would read from w or an mo-later w'; (e) ob; [Inst]: in general, an event is observed no earlier than it is issued, and for an instantaneous event, the two points coincide. Thus  $e \xrightarrow{ob} e'$  implies  $e \xrightarrow{ib} e'$  when e' is instantaneous. As noted in [3], this last component is optional and does not modify the semantics.

On the other hand, for ob we have (a)  $rf_e$ : if  $w \xrightarrow{rf_e} r$  then w was committed to memory before r, since r cannot read from the store buffer of another thread; (b) [nlW]; pf: NIC local writes cannot be polled until they are committed to memory; (c) nfo: NIC events are observed in the same order they arrive in  $\mathbf{wb_L/wb_R}$ ; (d) rb: if  $r \xrightarrow{rb} w$ , then w was not observed before r, otherwise it would have been committed to memory before r; (e) mo: if  $w \xrightarrow{mo} w'$ , then w was observed in memory before w'; (f) rao: remote atomic reads are (issued and) completed in the defined order; (g) ar; rao: if  $w \xrightarrow{ar} r \xrightarrow{rao} r'$ , then we have that r and w are the read and write of the same rRMW operation, thus w must be observed before the rao-later r' to ensure atomicity. (h) [Inst]; ib: by a similar logic to above, we know that the ib-earlier instantaneous event is also observed earlier, since its issue and observation points coincide.

**Semantics of a Program.** Given a program P, we can generate an event graph (E, po), by a standard process, which we describe below. We then choose any rf, mo, pf, nfo, rao such that the execution is consistent. The semantics of P are the set of consistent executions of P.

Thread to Event Graph. Given a thread identifier  $t \in \text{Tid}$  and a sequence of labels  $l_1, \ldots, l_n \in \text{ELab}$ , we define the *event graphs of* t as  $(\{e_1, \ldots, e_n\}, po) \in G^t(l_1, \ldots, l_n)$  where: (a)  $l(e_i) = l_i$  for all  $1 \le i \le n$ ; (b)  $\iota(e_i) \ne \iota(e_j)$  for all  $1 \le i \le n$ ; (c)  $t(e_i) = t$  for all  $1 \le i \le n$ ; (d)  $t(e_i) = t$  for all  $t(e_i) = t$  for all

**Initial Event Graph.** Given a set of locations Loc, we define  $G_{init} = (E_0, \emptyset)$ , such that for each  $x \in \text{Loc}$  there is exactly one  $e \in E_0$  with  $l(e) = 1 \mathbb{W}(x, 0)$ , and every event in  $E_0$  has a unique identifier. We call  $E_0$  the set of *initialisation* events.

**Sequential Composition.** For two event graphs  $G_1$  and  $G_2$ , we define their sequential composition  $G_1$ ;  $G_2 = (E, po)$  where

$$E \triangleq G_1.E \uplus G_2.E$$
 
$$\mathsf{po} \triangleq G_1.\mathsf{po} \cup G_2.\mathsf{po} \cup (G_1.E \times G_2.E)$$

Note that all events in  $G_2$  are ordered po-after every event in  $G_1$ . Sequential composition is defined only where the set of events of each graph are disjoint, i.e.  $G_1.E \cap G_2.E = \emptyset$ .

**Parallel Composition.** We define parallel composition by  $G_1 \parallel G_2 = (E, po)$  where

$$E \triangleq G_1.E \uplus G_2.E$$
  
po  $\triangleq G_1.po \cup G_2.po$ 

Note that the events of each graph are not po-ordered with respect to one another. We also require that the event sets be disjoint. As this operation is commutative and associative, it is straightforward to lift it to sets of graphs, which we denote by  $\parallel \mathcal{G}$ , where  $\mathcal{G}$  is a set of event graphs.

**Program to Event Graph.** A program P generates G if  $G = G_{init}$ ; ( $\|_{t \in \mathsf{Tid}} G_t$ ) and there is a set of sequences  $s_t \in S$  such that  $\mathsf{P}(t) \mapsto s_t$  and  $G_t \in G^t(s_t)$  for all  $t \in \mathsf{Tid}$ .

The operation  $C \rightarrow s$  relates a sequential program C to a sequence of labels s it generates. The definition is standard and show in Fig. 23. Note that RDMA operations generate multiple events, and for local and remote CAS operations, we distinguish between success and failure cases.

**Theorem 7.** The operational and declarative semantics of RDMA $_{\rm RMW}^{\rm TSO}$  are equivalent.

*Proof.* See Appendix D.3 onwards, extending the proof of [3].

## D.3 Annotated Labels and Inference Rules

On top of the 12 labels presented in Appendix D.2, we create six new labels:  $\mathtt{Put}(\overline{y},x)$ ,  $\mathtt{Get}(x,\overline{y})$ ,  $\mathtt{RCAS}(z,\overline{x},v,u)$ ,  $\mathtt{RFAA}(z,\overline{x},u)$ ,  $\mathtt{nlEX}(\overline{n})$ , and  $\mathtt{nrEX}(\overline{n})$ . These labels can also be used to create events (when bundled with an event identifier and a thread identifier).

We note  $E^{\mathrm{ext}}$  the extended set of all events, including the six new labels. Recall that  $\mathcal{R} = \mathtt{lR} \cup \mathtt{CAS} \cup \mathtt{nlR} \cup \mathtt{nrR} \cup \mathtt{narR} \subseteq E^{\mathrm{ext}}$  and  $\mathcal{W} = \mathtt{lW} \cup \mathtt{CAS} \cup \mathtt{nlW} \cup \mathtt{nrW} \cup \mathtt{narW} \subseteq E^{\mathrm{ext}}$ . We also note  $\mathtt{nEX} = \mathtt{nlEX} \cup \mathtt{nrEX}$  and  $\mathtt{rRMW} = \mathtt{RCAS} \cup \mathtt{RFAA}$ .

For annotated labels, we reuse most names from labels, but they are different entities. For instance we note  $r \in \mathtt{1R}$  for an event with label  $\mathtt{1R}$ , while  $\lambda = \mathtt{1R}\langle \ldots \rangle$  is an annotated label.

We use  $\mathsf{type}(\lambda)$  to denote the type of the annotated label (1R, 1W, CAS, F, Push, NIC, nlR, nrR, nlW, nrW, CN, P, nF, B,  $\mathcal{E}$ ). We use  $r(\lambda), w(\lambda), u(\lambda), a(\lambda), f(\lambda), p(\lambda), e(\lambda), \ldots$  to access the elements of a  $\lambda \in \mathsf{ALabel}$  where applicable. Also, we note  $t(\lambda)$  for the thread of the first argument of  $\lambda$ .

The annotated program transitions (Fig. 25) use an additional annotated label CASF $\langle r, w \rangle$  with  $r \in \mathtt{IR}$  and  $w \in \mathcal{W}$  to represent a failed CAS operation. This case is then translated into two labels (a memory fence and a local read) when creating a path in §D.4. Also, note that the annotated domains (e.g. the store buffers and the queue pairs) contain events, not annotated labels.

$$\begin{array}{c} C \leadsto C' & C' \rightarrowtail s \\ \hline C \rightarrowtail s & \frac{C_1 \rightarrowtail s_1 & C_2 \rightarrowtail s_2}{C_1; C_2 \rightarrowtail s_1, s_2} & \frac{\mathsf{elocs}(e) = \emptyset}{x := e \rightarrowtail \mathsf{IW}(x, \llbracket e \rrbracket)} \\ & \frac{\mathsf{elocs}(e_{\mathrm{old}}) = \mathsf{elocs}(e_{\mathrm{new}}) = \emptyset}{z := \llbracket e_{\mathrm{old}} \rrbracket \rightarrowtail s} \\ \hline z := \mathsf{CAS}(x, e_{\mathrm{old}}, e_{\mathrm{new}}) \rightarrowtail \mathsf{CAS}(x, \llbracket e_{\mathrm{old}} \rrbracket, \llbracket e_{\mathrm{new}} \rrbracket), s} \\ & \frac{\mathsf{elocs}(e_{\mathrm{old}}) = \mathsf{elocs}(e_{\mathrm{new}}) = \emptyset}{z := \mathsf{CAS}(x, e_{\mathrm{old}}, e_{\mathrm{new}}) \rightarrowtail \mathsf{F}, \mathsf{1R}(x, v), s} & \frac{v' \neq v}{\mathsf{assume}(x \neq v) \rightarrowtail \mathsf{1R}(x, v')} \\ \hline \hline x := y^{\overline{n}} \rightarrowtail \mathsf{nrR}(y^{\overline{n}}, v), \mathsf{n1W}(x, v, \overline{n}) & y^{\overline{n}} := x \rightarrowtail \mathsf{n1R}(x, v, \overline{n}), \mathsf{nrW}(y^{\overline{n}}, v)} \\ \hline x := v \mapsto \mathsf{nF}(\overline{n}) & \overline{\mathsf{pol1}(\overline{n})} \rightarrowtail \mathsf{P}(\overline{n}) & \overline{\mathsf{skip}} \rightarrowtail \epsilon \\ \hline elocs(e_{old}) = \mathsf{elocs}(e_{new}) = \emptyset & v \neq \llbracket e_{old} \rrbracket \\ \hline z := \mathsf{RCAS}(x^{\overline{n}}, e_{old}, e_{new}) \rightarrowtail \mathsf{narR}(x^{\overline{n}}, v), \mathsf{n1W}(z, v, \overline{n}) \\ \hline elocs(e_{old}) = \mathsf{elocs}(e_{new}) = \emptyset \\ \hline z := \mathsf{RFAA}(x^{\overline{n}}, e_{old}, e_{new}) \rightarrowtail \mathsf{narR}(x^{\overline{n}}, [e_{old}]], \mathsf{narW}(x^{\overline{n}}, \llbracket e_{new} \rrbracket), \mathsf{n1W}(z, \llbracket e_{old} \rrbracket, \overline{n}) \\ \hline elocs(e) = \emptyset & v' = v + \llbracket e \rrbracket \\ \hline z := \mathsf{RFAA}(x^{\overline{n}}, e) \rightarrowtail \mathsf{narR}(x^{\overline{n}}, v), \mathsf{narW}(x^{\overline{n}}, v'), \mathsf{n1W}(z, v, \overline{n}) \\ \hline \end{array}$$

Fig. 23: Label Sequences Construction

# initialisation. Given a program P, let

```
\begin{array}{lll} \mathsf{M}_0 \in \mathsf{AMem} & \text{s.t. } \forall x \in \mathsf{Loc. } \mathsf{M}_0(x) = init_x \text{ with } l(init_x) \triangleq \mathsf{1W}(x,0) \\ \mathsf{b}_0 \in \mathsf{ASBuff} & \mathsf{b}_0 \triangleq \varepsilon \\ \mathsf{B}_0 \in \mathsf{ASBMap} & \mathsf{B}_0 \triangleq \lambda t. \mathsf{b}_0 \\ \mathsf{A}_0 \in \mathsf{RAMap} & \mathsf{A}_0 \triangleq \lambda t. \bot \\ \mathsf{qp}_0 \in \mathsf{AQPair} & \mathsf{qp}_0 \triangleq \langle \varepsilon, \varepsilon, \varepsilon \rangle \\ \mathsf{QP}_0 \in \mathsf{AQPMap} & \mathsf{QP}_0 \triangleq \lambda t. \lambda n. \mathsf{qp}_0 \end{array}
```

```
\lambda \in \mathsf{ALabel}
\lambda \triangleq | \operatorname{lR}\langle r, w \rangle
                                                          where r \in \mathtt{lR}, w \in \mathcal{W}, \mathtt{eq_{loc\&v}}(r, w)
         | lW\langle w \rangle
                                                          where w \in \mathtt{lW}
         | \operatorname{CAS}\langle u, w \rangle
                                                          where u \in \mathtt{CAS}, w \in \mathcal{W}, \mathtt{eq_{loc\&v}}(u, w)
         | F\langle f \rangle
                                                          where f \in \mathbf{F}
                                                          where a \in (\mathtt{Put} \cup \mathtt{Get} \cup \mathtt{RCAS} \cup \mathtt{RFAA} \cup \mathtt{nF})
          |\operatorname{Push}\langle a\rangle
         |\operatorname{NIC}\langle a\rangle
                                                          where a \in (\mathtt{Put} \cup \mathtt{Get} \cup \mathtt{RCAS} \cup \mathtt{RFAA} \cup \mathtt{nF})
         | \operatorname{nlR}\langle r, w, a, w' \rangle
                                                          where r \in nlR, w \in \mathcal{W}, a \in Put, w' \in nrW, eq_{loc&v}(r, w),
                                                                       loc_r(a) = loc(r), loc_w(a) = loc(w'), v_r(r) = v_w(w')
         |\operatorname{nrR}\langle r, w, a, w'\rangle
                                                          where r \in nrR, w \in \mathcal{W}, a \in Get, w' \in nlW, eq_{loc&v}(r, w),
                                                                       \mathsf{loc}_r(a) = \mathsf{loc}(r), \mathsf{loc}_w(a) = \mathsf{loc}(w'), v_r(r) = v_w(w')
                                                         where r \in \mathtt{narR}, w \in \mathcal{W}, a \in \mathtt{RRMW}, w' \in \mathtt{nlW}, w'' \in \mathtt{narW},
         |\operatorname{narR}\langle r, w, a, w', w''\rangle
                                                                       eq_{loc&v}(r, w), loc_r(a) = loc(r) = loc(w''),
                                                                       loc_w(a) = loc(w'), v_r(r) = v_w(w'),
                                                                      a \in \mathtt{RCAS} \implies v_{\mathrm{r}}(r) = v_{\mathrm{e}}(a) \wedge v_{\mathrm{w}}(w'') = v_{\mathrm{u}}(a)
                                                                      a \in RFAA \implies v_{\mathbf{w}}(w'') = v_{\mathbf{r}}(r) + v(a)
                                                          where r \in \mathtt{narR}, w \in \mathcal{W}, a \in \mathtt{rRMW}, w' \in \mathtt{nlW}, \mathtt{eq}_{\mathtt{loc\&v}}(r, w),
         | \operatorname{naF} \langle r, w, a, w' \rangle
                                                                       loc_r(a) = loc(r), loc_w(a) = loc(w'),
                                                                       v_{\mathrm{r}}(r) = v_{\mathrm{w}}(w^l), v_{\mathrm{r}}(r) \neq v_{\mathrm{e}}(a)
         | \ \mathtt{nlW} \langle w, e \rangle
                                                          where w \in \mathtt{nlW}, e \in \mathtt{nlEX}, \mathtt{sameqp}(w, e)
          |\operatorname{nrW}\langle w,e\rangle
                                                          where w \in nrW, e \in nrEX, sameqp(w, e)
         |\operatorname{narW}\langle w \rangle
                                                          where w \in \mathtt{narW}
         | \mathsf{CN}\langle e \rangle
                                                          where e \in \mathtt{nrEX}
         | P\langle p, e \rangle
                                                          where p \in P, e \in nEX, sameqp(p, e)
         \mid nF\langle f \rangle
                                                          where f \in nF
         \mid \mathsf{B}\langle w \rangle
                                                          where w \in \mathcal{W}
         \mid \mathcal{E}\langle t \rangle
                                                          where t \in \mathsf{Tid}
                                   \mathsf{eq}_{\mathsf{loc}\&\mathsf{v}}(r,w) \ \triangleq \ \mathsf{loc}(r) = \mathsf{loc}(w) \land v_{\mathsf{r}}(r) = v_{\mathsf{w}}(w)
                                  sameqp(e, e') \triangleq t(e) = t(e') \land \overline{n}(e) = \overline{n}(e')
```

Fig. 24: Annotated Labels

$$\begin{array}{c} \mathbf{Program \ transitions: Prog} \xrightarrow{\mathsf{ALabel} \uplus \{\mathsf{CASF}\}} \mathsf{Prog} \\ \mathbf{Command \ transitions: Comm} \xrightarrow{\mathsf{ALabel} \uplus \{\mathsf{CASF}\}} \mathsf{Comm} \end{array}$$

$$\begin{array}{c} \frac{\mathsf{C}_1 \overset{\lambda}{\to} \mathsf{C}_1'}{\mathsf{C}_1; \mathsf{C}_2 \overset{\lambda}{\to} \mathsf{C}_1'; \mathsf{C}_2} & \frac{i \in \{1,2\}}{\mathsf{ckip}} \mathsf{C} & \frac{i \in \{1,2\}}{\mathsf{C}_1 + \mathsf{C}_2} \overset{\mathcal{E}(t)}{\to} \mathsf{C}_i & \frac{\mathcal{E}(t)}{\mathsf{C}^* \to \mathsf{skip}} \\ \\ \frac{\mathsf{C} \overset{\mathcal{E}(t)}{\to} \mathsf{C}_i'}{\mathsf{C}} & \frac{\mathsf{C} \overset{\mathcal{E}(t)}{\to} \mathsf{C}_i'} & \frac{\mathsf{elocs}(e) = \emptyset & w = (\iota, t, \mathsf{1W}(x, \llbracket e \rrbracket))}{x := e} \overset{\mathsf{IW}(w)}{\to} \mathsf{skip} \\ \\ \frac{\mathsf{elocs}(e_{\mathrm{old}}) = \mathsf{elocs}(e_{\mathrm{new}}) = \emptyset & v \neq \llbracket e_{\mathrm{old}} \rrbracket & r = (\iota, t, \mathsf{1R}(x, v))}{x := e} & \frac{\mathsf{elocs}(e_{\mathrm{old}}) = \mathsf{elocs}(e_{\mathrm{new}}) = \emptyset & v \neq \llbracket e_{\mathrm{old}} \rrbracket & r = (\iota, t, \mathsf{1R}(x, v))}{z := \mathsf{CAS}(x, e_{\mathrm{old}}, e_{\mathrm{new}}) & \frac{\mathsf{CASF}(r, w)}{\mathsf{CASF}(r, w)} \neq z := v} \\ \\ \frac{\mathsf{elocs}(e_{\mathrm{old}}) = \mathsf{elocs}(e_{\mathrm{new}}) = \emptyset & u = (\iota, t, \mathsf{CAS}(x, \llbracket e_{\mathrm{old}} \rrbracket, \llbracket e_{\mathrm{new}} \rrbracket))}{z := \mathsf{CAS}(x, e_{\mathrm{old}}, e_{\mathrm{new}}) & \frac{\mathsf{CAS}(u, w)}{z} \neq z := \llbracket e_{\mathrm{old}} \rrbracket} \\ \frac{f = (\iota, t, \mathsf{F})}{\mathsf{mfence}} & \frac{a = (\iota, t, \mathsf{Get}(x, \overline{y}))}{x := \overline{y}} & \frac{a = (\iota, t, \mathsf{Push}(a)}{z} \Rightarrow \mathsf{skip} \\ \\ \frac{a = (\iota, t, \mathsf{nF}(\overline{n}))}{\mathsf{rfence}(\overline{n})} & \frac{v = \llbracket e_{\mathrm{old}} \rrbracket}{v} & u = \llbracket e_{\mathrm{new}} \rrbracket & a = (\iota, t, \mathsf{RCAS}(z, \overline{x}, v, u)) \\ & u = \llbracket e_{\mathrm{new}} \rrbracket & a = (\iota, t, \mathsf{RCAS}(z, \overline{x}, v, u)) \\ \\ \frac{v = \llbracket e_{\mathrm{old}} \rrbracket}{v} & u = \llbracket e_{\mathrm{new}} \end{split} & \frac{p = (\iota, t, \mathsf{Polh}(a)}{v} \Rightarrow \mathsf{skip} \\ \\ \frac{\mathsf{elocs}(e) = \emptyset}{z} & u = \llbracket e \rrbracket & a = (\iota, t, \mathsf{RFAA}(z, \overline{x}, u)) & p = (\iota, t, \mathsf{Polh}(a)}{v} \Rightarrow \mathsf{skip} \\ \\ \frac{\mathsf{elocs}(e) = \mathbb{Q}_{\mathsf{N}} & u = \mathbb{Q}_$$

Fig. 25: RDMA<sup>TSO</sup> program and command transitions for the annotated semantics

$$\begin{split} \mathsf{M} \in \mathsf{AMem} &\triangleq \{ m \in \mathsf{Loc} \to \mathcal{W} \mid \forall x \in \mathsf{Loc.loc}(m[x]) = x \} & \mathsf{B} \in \mathsf{ASBMap} \triangleq \mathsf{Tid} \to \mathsf{ASBuff} \\ \mathsf{A} \in \mathsf{RAMap} &\triangleq \lambda n. \, \{\bot, \top\} & \mathsf{QP} \in \mathsf{AQPMap} \triangleq \mathsf{Tid} \to (\mathsf{Node} \to \mathsf{AQPair}) \\ \mathsf{b} \in \mathsf{ASBuff} &\triangleq (\mathsf{1W} \cup \mathsf{Get} \cup \mathsf{Put} \cup \mathsf{nF} \cup \mathsf{RCAS} \cup \mathsf{RFAA})^* & \mathsf{sqp} \in \mathsf{AQPair} \triangleq \mathsf{APipe} \times \mathsf{AWBR} \times \mathsf{AWBL} \\ & \mathsf{pipe} \in \mathsf{APipe} \triangleq (\mathsf{Get} \cup \mathsf{Put} \cup \mathsf{nF} \cup \mathsf{nrW} \cup \mathsf{narW} \cup \mathsf{nrEX} \cup \mathsf{nlW} \cup \mathsf{RCAS} \cup \mathsf{RFAA})^* \\ & \mathsf{wb}_{\mathbf{R}} \in \mathsf{AWBR} \triangleq (\mathsf{nrW}, \mathsf{narW})^* & \mathsf{wb}_{\mathbf{L}} \in \mathsf{AWBL} \triangleq (\mathsf{nlW} \cup \mathsf{nlEX} \cup \mathsf{nrEX})^* \end{split}$$

$$\frac{\mathsf{B}' = \mathsf{B}[t(w) \mapsto w \cdot \mathsf{B}(t(w))]}{\mathsf{M}, \mathsf{B}, \mathsf{A}, \mathsf{QP}} \xrightarrow{\mathsf{IW}(w)} \mathsf{M}, \mathsf{B}', \mathsf{A}, \mathsf{QP} \xrightarrow{\mathsf{IR}(r,w)} \mathsf{M}, \mathsf{B}', \mathsf{A}, \mathsf{QP} \xrightarrow{\mathsf{IR}(r,w)} \mathsf{M}, \mathsf{B}, \mathsf{A}, \mathsf{QP} \xrightarrow{\mathsf{IR}(r,w)} \mathsf{M}, \mathsf{$$

$$\text{with} \quad (\mathsf{M} \blacktriangleleft \alpha)(x) = \begin{cases} \mathsf{M}[x] & \alpha = \varepsilon \\ w & \alpha = w \cdot \beta \wedge w \in \mathcal{W} \wedge \mathsf{loc}(w) = x \\ (\mathsf{M} \blacktriangleleft \beta)(x) & \alpha = e \cdot \beta \wedge (e \not\in \mathcal{W} \vee \mathsf{loc}(e) \neq x) \end{cases}$$

 $\overline{\text{Fig. 26: RDMA}^{\text{TSO}}}$  hardware domains and hardware transitions for the annotated semantics

```
\mathbf{pipe} = \alpha \cdot f \qquad f = (\iota, t, \mathtt{nF}(n))
                                                                         \overline{\mathsf{M},\mathsf{A},\langle\mathbf{pipe},\mathbf{wb_R},\mathbf{wb_L}\rangle\xrightarrow{\mathsf{nF}\langle f\rangle}_{\mathsf{sqp}}\mathsf{M},\mathsf{A},\langle\alpha,\mathbf{wb_R},\mathbf{wb_L}\rangle}
              \mathbf{pipe} = \alpha \cdot a \cdot \beta \qquad a = (\iota_a, t, \mathtt{Put}(\overline{y}, x)) \qquad \mathsf{M}(x) = w \qquad r = (\iota_r, t, \mathtt{nlR}(x, v_{\mathrm{w}}(w), n(\overline{y})))
w' = (\iota_{w'}, t, \operatorname{nrW}(\overline{y}, v_{\operatorname{w}}(w))) \qquad \beta \in (\operatorname{nrW} \cup \operatorname{narW} \cup \operatorname{Get} \cup \operatorname{nlW} \cup \operatorname{RCAS} \cup \operatorname{RFAA} \cup \operatorname{nrEX})^* \qquad \operatorname{wb_L} \in \operatorname{nEX}^*
                                               \mathsf{M}, \mathsf{A}, \langle \mathbf{pipe}, \mathbf{wb_R}, \mathbf{wb_L} \rangle \xrightarrow{\mathtt{nlR}\langle r, w, a, w' \rangle}_{\mathsf{sqp}} \mathsf{M}, \mathsf{A}, \langle \alpha \cdot w' \cdot \beta, \mathbf{wb_R}, \mathbf{wb_L} \rangle
                                                                                                                                                           e = (\iota_e, t, \mathtt{nrEX}(n(\overline{y}))) \qquad eta \in (\mathtt{Get} \cup \mathtt{nlW} \cup \mathtt{nrEX})^*
                                                                    w = (\iota_w, t, \mathtt{nrW}(\overline{y}, v))
                                                     M, A, \langle \mathbf{pipe}, \mathbf{wb_R}, \mathbf{wb_L} \rangle \xrightarrow{\operatorname{nrW}\langle w, e \rangle} \operatorname{sqp} M, A, \langle \alpha \cdot e \cdot \beta, w \cdot \mathbf{wb_R}, \mathbf{wb_L} \rangle
                                                                                                                           \mathbf{wb_R} = \alpha \cdot w \qquad w \in \mathtt{nrW}
                                                   \mathsf{M}, \mathsf{A}, \langle \mathbf{pipe}, \mathbf{wb}_{\mathbf{L}} \rangle \xrightarrow{\mathsf{B}\langle w \rangle}_{\mathsf{sqp}} \mathsf{M}[\mathsf{loc}(w) \mapsto w], \mathsf{A}, \langle \mathbf{pipe}, \alpha, \mathbf{wb}_{\mathbf{L}} \rangle
                                                                                                                            \mathbf{pipe} = \alpha \cdot e \qquad e \in \mathtt{nrEX}
                                                                    \mathsf{M}, \mathsf{A}, \langle \mathbf{pipe}, \mathbf{wb_R}, \mathbf{wb_L} \rangle \xrightarrow[]{\mathsf{CN}\langle e \rangle}_{\mathsf{sqp}} \mathsf{M}, \mathsf{A}, \langle \alpha, \mathbf{wb_R}, e \cdot \mathbf{wb_L} \rangle
                       \begin{array}{cccc} \mathbf{pipe} = \alpha \cdot a \cdot \beta & a = (\iota_a, t, \mathtt{Get}(x, \overline{y})) & \mathsf{M}(\overline{y}) = w & r = (\iota_r, t, \mathtt{nrR}(\overline{y}, v_{\mathbf{w}}(w), v_{\mathbf{w}}(w), n(\overline{y}))) & \beta \in (\mathtt{Get} \cup \mathtt{nlW} \cup \mathtt{nrEX})^* & \mathbf{wb_R} = \varepsilon \end{array}
                                                                                                                                                                                                                            r = (\iota_r, t, \text{nrR}(\overline{y}, v_w(w)))
                                               \mathsf{M}, \mathsf{A}, \langle \mathbf{pipe}, \mathbf{wb_R}, \mathbf{wb_L} \rangle \xrightarrow{\mathbf{nrR}\langle r, w, a, w' \rangle}_{\mathsf{sqp}} \mathsf{M}, \mathsf{A}, \langle \alpha \cdot w' \cdot \beta, \mathbf{wb_R}, \mathbf{wb_L} \rangle
                                                           \mathbf{pipe} = \alpha \cdot w \qquad w = (\iota_w, t, \mathtt{nlW}(x, v, n)) \qquad e = (\iota_e, t, \mathtt{nlEX}(n))
                                                           M, A, \langle \mathbf{pipe}, \mathbf{wb_R}, \mathbf{wb_L} \rangle \xrightarrow{\mathtt{nlW}\langle w, e \rangle}_{\mathtt{sqp}} M, A, \langle \alpha, \mathbf{wb_R}, e \cdot w \cdot \mathbf{wb_L} \rangle
                                                                                                \mathbf{wb_L} = \alpha \cdot w \cdot \beta w \in \mathtt{nlW} \beta \in \mathtt{nEX}^*
                                             \mathsf{M}, \mathsf{A}, \overline{\langle \mathbf{pipe}, \mathbf{wb_R}, \mathbf{wb_L}\rangle \ \frac{\mathsf{B}\langle w\rangle}{\mathsf{sqp}} \ \mathsf{M}[\mathsf{loc}(w) \mapsto w], \mathsf{A}, \langle \mathbf{pipe}, \mathbf{wb_R}, \alpha \cdot \beta\rangle}
           \begin{aligned} \mathbf{pipe} &= \alpha \cdot a \cdot \beta \quad \mathbf{wb_R} = \varepsilon \quad \mathsf{M}(\overline{x}) = w \\ v_{\mathbf{w}}(w) &\neq v \qquad \mathsf{A}(n(\overline{x})) = \bot \qquad a = (\iota_a, t, \mathtt{RCAS}(z, \overline{x}, v, u)) \\ \underline{r = (\iota_r, t, \mathtt{narR}(\overline{x}, v_{\mathbf{w}}(w)))} \qquad w' = (\iota_{w'}, t, \mathtt{nlW}(z, v_{\mathbf{w}}(w), n(\overline{x}))) \qquad \beta \in (\mathtt{Get} \cup \mathtt{nlW} \cup \mathtt{nrEX})^* \end{aligned}
                                                M, A, \langle \mathbf{pipe}, \mathbf{wb_R}, \mathbf{wb_L} \rangle \xrightarrow{\mathrm{naF}\langle r, w, a, w' \rangle}_{\mathrm{sqp}} M, A, \langle \alpha \cdot w' \cdot \beta, \mathbf{wb_R}, \mathbf{wb_L} \rangle
            \begin{aligned} & \mathbf{pipe} = \alpha \cdot a \cdot \beta \quad \mathbf{wb_R} = \varepsilon \quad \mathbf{M}(\overline{x}) = w \\ v_{\mathbf{w}}(w) = v \quad & \mathbf{A}(n(\overline{x})) = \bot \quad a = (\iota_a, t, \mathtt{RCAS}(z, \overline{x}, v, u)) \quad r = (\iota_r, t, \mathtt{narR}(\overline{x}, v_{\mathbf{w}}(w))) \\ w'' = (\iota_{w''}, t, \mathtt{narW}(\overline{x}, u)) \quad & w' = (\iota_{w'}, t, \mathtt{nlW}(z, v_{\mathbf{w}}(w), n(\overline{x}))) \quad \beta \in (\mathtt{Get} \cup \mathtt{nlW} \cup \mathtt{nreX})^* \end{aligned}
            \overbrace{\mathsf{M},\mathsf{A},\langle\mathbf{pipe},\mathbf{wb_R},\mathbf{wb_L}\rangle\xrightarrow{\mathsf{narR}\langle r,w,a,w',w''\rangle}_{\mathsf{sqp}}\mathsf{M},\mathsf{A}[n(\overline{x})\mapsto\top],\langle\alpha\cdot w'\cdot w''\cdot\beta,\mathbf{wb_R},\mathbf{wb_L}\rangle}
                  \begin{array}{c} \mathbf{pipe} = \alpha \cdot a \cdot \beta \quad \mathbf{wb_R} = \varepsilon \quad \mathsf{M}(\overline{x}) = w \\ v_\mathbf{w}(w) + v = u \quad \mathsf{A}(n(\overline{x})) = \bot \quad a = (\iota_a, t, \mathtt{RFAA}(z, \overline{x}, v)) \qquad r = (\iota_r, t, \mathtt{narR}(\overline{x}, v_\mathbf{w}(w))) \\ w'' = (\iota_{w''}, t, \mathtt{narW}(\overline{x}, u))) \quad w' = (\iota_{w'}, t, \mathtt{nlW}(z, v_\mathbf{w}(w))) \qquad \beta \in (\mathtt{Get} \cup \mathtt{nlW} \cup \mathtt{nrEX})^* \end{array}
            \overline{\mathsf{M},\mathsf{A},\langle\mathsf{pipe},\mathsf{wb}_{\mathbf{R}},\mathsf{wb}_{\mathbf{L}}\rangle} \xrightarrow{\mathsf{narR}\langle r,w,a,w',w''\rangle}_{\mathsf{sop}} \mathsf{M},\mathsf{A}[n(\overline{x})\mapsto \top],\langle\alpha\cdot w'\cdot w''\cdot\beta,\mathsf{wb}_{\mathbf{R}},\mathsf{wb}_{\mathbf{L}}\rangle
                                               \mathbf{pipe} = \alpha \cdot w \cdot \beta \qquad \beta \in \left( \mathtt{Get} \cup \mathtt{nlW} \cup \mathtt{nrEX} \right)^* \qquad w = \left( \iota_w, t, \mathtt{narW}(\overline{x}, v) \right)
                                                            M, A, \langle \mathbf{pipe}, \mathbf{wb_R}, \mathbf{wb_L} \rangle \xrightarrow{\underline{\mathsf{narW}} \langle w \rangle}_{\mathsf{sqp}} M, A, \langle \alpha \cdot \beta, w \cdot \mathbf{wb_R}, \mathbf{wb_L} \rangle
                                                                                                    \mathbf{wb_R} = \alpha \cdot w w = (\iota_w, t, \mathtt{narW}(\overline{x}, v))
                                M, A, \langle \mathbf{pipe}, \mathbf{wb_R}, \mathbf{wb_L} \rangle \xrightarrow{B\langle w \rangle}_{sqp} M[loc(w) \mapsto w], A[n(\overline{x}) \mapsto \bot], \langle \mathbf{pipe}, \alpha, \mathbf{wb_L} \rangle
```

Fig. 27: Annotated 3 Buffers NIC Semantics

# D.4 Paths, Gluing, and Other Definitions

We define a path as:  $\pi \in \mathsf{Path} \triangleq (\mathsf{ALabel} \setminus \mathcal{E}\langle t \rangle)^*$ 

We define Annotated Operational Semantics Gluing with the following rules.

$$\frac{\mathsf{P} \xrightarrow{\mathcal{E}\langle t \rangle} \mathsf{P}'}{\mathsf{P}, \mathsf{M}, \mathsf{B}, \mathsf{A}, \mathsf{QP}, \pi \Rightarrow \mathsf{P}', \mathsf{M}, \mathsf{B}, \mathsf{A}, \mathsf{QP}, \pi}$$

$$\frac{\mathsf{P} \xrightarrow{\lambda} \mathsf{P}' \qquad \mathsf{M}, \mathsf{B}, \mathsf{A}, \mathsf{QP} \xrightarrow{\lambda} \mathsf{M}', \mathsf{B}', \mathsf{A}, \mathsf{QP}'}{\lambda \in (\mathsf{1R} \cup \mathsf{1W} \cup \mathsf{CAS} \cup \mathsf{F} \cup \mathsf{Push} \cup \mathsf{P}) \qquad \mathsf{fresh}(\lambda, \pi)} \xrightarrow{\mathsf{P}, \mathsf{M}, \mathsf{B}, \mathsf{A}, \mathsf{QP}, \pi \Rightarrow \mathsf{P}', \mathsf{M}', \mathsf{B}', \mathsf{A}, \mathsf{QP}', \lambda \cdot \pi}$$

$$\frac{\mathsf{M}, \mathsf{B}, \mathsf{A}, \mathsf{QP} \xrightarrow{\lambda} \mathsf{M}', \mathsf{B}', \mathsf{A}', \mathsf{QP}'}{\lambda \in (\mathsf{NIC} \cup \mathsf{n1R} \cup \mathsf{nrR} \cup \mathsf{n1W} \cup \mathsf{nrW} \cup \mathsf{naF} \cup \mathsf{narR} \cup \mathsf{narW} \cup \mathsf{CN} \cup \mathsf{nF} \cup \mathsf{B})} \xrightarrow{\mathsf{fresh}(\lambda, \pi)} \xrightarrow{\mathsf{P}, \mathsf{M}, \mathsf{B}, \mathsf{A}, \mathsf{QP}, \pi \Rightarrow \mathsf{P}, \mathsf{M}', \mathsf{B}', \mathsf{A}', \mathsf{QP}', \lambda \cdot \pi}$$

$$\frac{\mathsf{P} \xrightarrow{\mathsf{CASF}\langle r, w \rangle}{\mathsf{P}, \mathsf{M}, \mathsf{B}, \mathsf{A}, \mathsf{QP}} \xrightarrow{\lambda_1} \mathsf{M}, \mathsf{B}, \mathsf{A}, \mathsf{QP} \xrightarrow{\lambda_2} \mathsf{M}, \mathsf{B}, \mathsf{A}, \mathsf{QP} \qquad \mathsf{fresh}(\lambda_1, \pi) \qquad \mathsf{fresh}(\lambda_2, \pi)}{\mathsf{P}, \mathsf{M}, \mathsf{B}, \mathsf{A}, \mathsf{QP}, \pi \Rightarrow \mathsf{P}', \mathsf{M}, \mathsf{B}, \mathsf{A}, \mathsf{QP}, \lambda_2 \cdot \lambda_1 \cdot \pi} \xrightarrow{\mathsf{Fresh}(\lambda_2, \pi)} \xrightarrow{\mathsf{P}, \mathsf{M}, \mathsf{B}, \mathsf{A}, \mathsf{QP}, \pi \Rightarrow \mathsf{P}', \mathsf{M}, \mathsf{B}, \mathsf{A}, \mathsf{QP}, \lambda_2 \cdot \lambda_1 \cdot \pi}$$

Two annotated labels are non-conflicting  $(\lambda_1 \bowtie \lambda_2)$  if they are of a different type or if their relevant arguments are disjoints. An annotated label is fresh if it does not conflict with any previous annotated label.

 $\mathbf{Relevant}: \mathsf{ALabel} \to 2^{E^{\mathrm{ext}}}$ 

$$\lambda_1 \bowtie \lambda_2 \triangleq \mathsf{type}(\lambda_1) \neq \mathsf{type}(\lambda_2) \vee \mathbf{Relevant}(\lambda_1) \cap \mathbf{Relevant}(\lambda_2) = \emptyset$$

$$\begin{aligned} &\mathsf{fresh}(\lambda,\pi) \; \triangleq \; \forall \lambda' \in \pi, \; \lambda \bowtie \lambda' \\ &\mathsf{nodup}(\pi) \; \triangleq \; \forall \pi_2,\lambda,\pi_1. \; \pi = \pi_2 \cdot \lambda \cdot \pi_1 \; \Longrightarrow \; \mathsf{fresh}(\lambda,\pi_1) \end{aligned}$$

 $\mathbf{Relevant}(\lambda)$  are the arguments that are important to consider to avoid duplicating events. The excluded events are the write operations we lookup when reading. For instance:

- Having both  $1R\langle r_1, w \rangle$  and  $1R\langle r_2, w \rangle$  during an execution is fine, since w can be looked up any number of time.
- Having both  $nlR\langle r_1, w_1, a, e_1 \rangle$  and  $nlR\langle r_2, w_2, a, e_2 \rangle$  during an execution is problematic, since it means the put operation a is being run twice.

## Completeness.

$$\operatorname{complete}(\pi) \triangleq \forall a, w', e, r, w, f, w''. \\ \operatorname{lW}\langle w \rangle \in \pi \implies \operatorname{B}\langle w \rangle \in \pi \\ \wedge \operatorname{Push}\langle a \rangle \in \pi \implies \operatorname{NIC}\langle a \rangle \in \pi \\ \wedge \operatorname{NIC}\langle f \rangle \in \pi \wedge f \in \operatorname{nF} \implies \operatorname{nF}\langle f \rangle \in \pi \\ \wedge \operatorname{NIC}\langle a \rangle \in \pi \wedge a \in \operatorname{Put} \implies \exists r, w, w'. \operatorname{n1R}\langle r, w, a, w' \rangle \in \pi \\ \wedge \operatorname{NIC}\langle a \rangle \in \pi \wedge a \in \operatorname{Get} \implies \exists r, w, w'. \operatorname{narR}\langle r, w, a, w' \rangle \in \pi \\ \wedge \operatorname{NIC}\langle a \rangle \in \pi \wedge a \in \operatorname{RFAA} \implies \exists r, w, w'. \operatorname{narR}\langle r, w, a, w', w'' \rangle \in \pi \\ \wedge \operatorname{NIC}\langle a \rangle \in \pi \wedge a \in \operatorname{RCAS} \implies \begin{pmatrix} \exists r, w, w'. \operatorname{narR}\langle r, w, a, w' \rangle \in \pi \\ \vee \exists r, w, w'. w''. \operatorname{narR}\langle r, w, a, w' \rangle \in \pi \end{pmatrix} \\ \wedge \operatorname{n1R}\langle r, w, a, w' \rangle \in \pi \implies \exists e. \operatorname{nrW}\langle w', e \rangle \in \pi \\ \wedge \operatorname{narR}\langle r, w, a, w', w'' \rangle \in \pi \implies \exists e. \operatorname{nlW}\langle w', e \rangle \in \pi \\ \wedge \operatorname{narR}\langle r, w, a, w', w'' \rangle \in \pi \implies \exists e. \operatorname{nlW}\langle w', e \rangle \in \pi \\ \wedge \operatorname{narR}\langle r, w, a, w', w'' \rangle \in \pi \implies \exists e. \operatorname{nlW}\langle w', e \rangle \in \pi \\ \wedge \operatorname{narR}\langle r, w, a, w' \rangle \in \pi \implies \exists e. \operatorname{nlW}\langle w', e \rangle \in \pi \\ \wedge \operatorname{narR}\langle r, w, a, w' \rangle \in \pi \implies \exists e. \operatorname{nlW}\langle w', e \rangle \in \pi \\ \wedge \operatorname{narR}\langle w, e \rangle \in \pi \implies \operatorname{B}\langle w \rangle \in \pi \wedge \operatorname{CN}\langle e \rangle \in \pi \\ \wedge \operatorname{narW}\langle w, e \rangle \in \pi \implies \operatorname{B}\langle w \rangle \in \pi \wedge \operatorname{CN}\langle e \rangle \in \pi \\ \wedge \operatorname{narW}\langle w \rangle \in \pi \implies \operatorname{B}\langle w \rangle \in \pi \wedge \operatorname{CN}\langle e \rangle \in \pi$$

Informal: every pending operation is done and (most) buffers are empty. Note that some nEX (i.e., completion notifications) might still be in  $wb_L$ .

For a path  $\pi$  without duplicate (e.g. if  $\mathsf{nodup}(\pi)$  holds), we define the total ordering of its annotated labels as follows. Note that the early part of the path is on the right.

$$\lambda_1 \prec_{\pi} \lambda_2 \triangleq \exists \pi_1, \pi_2, \pi_3 \text{ s.t. } \pi = \pi_3 \cdot \lambda_2 \cdot \pi_2 \cdot \lambda_1 \cdot \pi_1$$

# Backward Completeness. (with ordering)

 $\mathsf{backComp}(\pi) \triangleq \forall a, w', e, r, w, f, p, w''.$ 

Poll Order.

$$\mathsf{pollOrder}(\pi) \, \triangleq \, \forall e_1, e_2. \, \begin{pmatrix} \mathsf{sameqp}(e_1, e_2) \\ \land \, \lambda_1 \in \{ \mathsf{nlW} \langle \_, e_1 \rangle, \mathsf{CN} \langle e_1 \rangle \} \\ \land \, \lambda_2 \in \{ \mathsf{nlW} \langle \_, e_2 \rangle, \mathsf{CN} \langle e_2 \rangle \} \\ \land \, \lambda_1 \prec_\pi \, \lambda_2 \\ \land \, \mathsf{P} \langle \_, e_2 \rangle \in \pi \end{pmatrix} \implies \mathsf{P} \langle \_, e_1 \rangle \prec_\pi \mathsf{P} \langle \_, e_2 \rangle$$

#### Flush Order.

```
bufFlushOrd(\pi) \triangleq
          \forall w_1, w_2 \in \mathtt{IW}. \ \left( \begin{array}{l} t(w_1) = t(w_2) \implies \\ (\mathsf{B}\langle w_2 \rangle \in \pi \wedge \mathtt{IW}\langle w_1 \rangle \prec_\pi \mathtt{IW}\langle w_2 \rangle) \iff \mathsf{B}\langle w_1 \rangle \prec_\pi \mathsf{B}\langle w_2 \rangle \end{array} \right)
 \land \forall a_1, a_2 \in (\texttt{Get} \cup \texttt{Put} \cup \texttt{nF} \cup \texttt{RCAS} \cup \texttt{RFAA}).
             \int t(a_1) = t(a_2) \implies
            \left( (\operatorname{NIC}\langle a_2 \rangle \in \pi \wedge \operatorname{Push}\langle a_1 \rangle \prec_{\pi} \operatorname{Push}\langle a_2 \rangle) \iff \operatorname{NIC}\langle a_1 \rangle \prec_{\pi} \operatorname{NIC}\langle a_2 \rangle \right)
  \land \forall a_1 \in (\texttt{Get} \cup \texttt{Put} \cup \texttt{nF} \cup \texttt{RCAS} \cup \texttt{RFAA}), w_2 \in \texttt{1W}.
              \begin{pmatrix} t(a_1) = t(w_2) \implies \\ (\mathsf{B}\langle w_2\rangle \in \pi \land \mathsf{Push}\langle a_1\rangle \prec_{\pi} \mathsf{IW}\langle w_2\rangle) \iff \mathsf{NIC}\langle a_1\rangle \prec_{\pi} \mathsf{B}\langle w_2\rangle \\ \land (\mathsf{NIC}\langle a_1\rangle \in \pi \land \mathsf{IW}\langle w_2\rangle \prec_{\pi} \mathsf{Push}\langle a_1\rangle) \iff \mathsf{B}\langle w_2\rangle \prec_{\pi} \mathsf{NIC}\langle a_1\rangle \end{pmatrix} 
\land \forall w \in \mathsf{IW}, f \in \mathsf{F}. \; \mathsf{IW}\langle w \rangle \prec_{\pi} \mathsf{F}\langle f \rangle \land t(w) = t(f) \implies \mathsf{B}\langle w \rangle \prec_{\pi} \mathsf{F}\langle f \rangle
 \land \ \forall w \in \mathtt{IW}, u \in \mathtt{CAS}. \ \mathtt{IW} \langle w \rangle \prec_{\pi} \mathtt{CAS} \langle u, \_ \rangle \land t(w) = t(u) \implies \mathtt{B} \langle w \rangle \prec_{\pi} \mathtt{CAS} \langle u, \_ \rangle
 \wedge \ \forall w \in \mathtt{nlW}, r \in \mathtt{nlR}.
           (\mathtt{nlW}\langle w, \_ \rangle \prec_{\pi} \mathtt{nlR}\langle r, \_, \_, \_ \rangle \land \mathtt{sameqp}(w, r)) \implies \mathsf{B}\langle w \rangle \prec_{\pi} \mathtt{nlR}\langle r, \_, \_, \_ \rangle
 \land \forall w \in \mathtt{nrW}, r \in \mathtt{nrR}.
           (\mathtt{nrW}\langle w, \_ \rangle \prec_\pi \mathtt{nrR}\langle r, \_, \_, \_ \rangle \wedge \mathsf{sameqp}(w, r)) \implies \mathsf{B}\langle w \rangle \prec_\pi \mathtt{nrR}\langle r, \_, \_, \_ \rangle
  \land \ \forall w \in \mathtt{nrW}, r \in \mathtt{narR}.
           (\operatorname{nrW}\langle w, \_ \rangle \prec_{\pi} \operatorname{naF}\langle r, \_, \_, \_ \rangle \wedge \operatorname{sameqp}(w, r)) \implies \mathsf{B}\langle w \rangle \prec_{\pi} \operatorname{naF}\langle r, \_, \_, \_ \rangle
  \wedge \ \forall w \in \mathtt{nrW}, r \in \mathtt{narR}.
           (\mathtt{nrW}\langle w, \_ \rangle \prec_\pi \mathtt{narR}\langle r, \_, \_, \_, \_ \rangle \wedge \mathsf{sameqp}(w,r)) \implies \mathsf{B}\langle w \rangle \prec_\pi \mathtt{narR}\langle r, \_, \_, \_, \_ \rangle
  \wedge \ \forall w \in \mathtt{narW}, r \in \mathtt{nrR}.
           (\mathtt{narW}\langle w \rangle \prec_{\pi} \mathtt{nrR}\langle r, \_, \_, \_ \rangle \wedge \mathtt{sameqp}(w,r)) \implies \mathsf{B}\langle w \rangle \prec_{\pi} \mathtt{nrR}\langle r, \_, \_, \_ \rangle
   \land \ \forall w \in \mathtt{narW}, r \in \mathtt{narR}.
           (\mathtt{narW}\langle w \rangle \prec_{\pi} \mathtt{naF}\langle r, \_, \_, \_ \rangle \wedge \mathsf{sameqp}(w,r)) \implies \mathsf{B}\langle w \rangle \prec_{\pi} \mathtt{naF}\langle r, \_, \_, \_ \rangle
 \land \forall w \in \mathtt{narW}, r \in \mathtt{narR}.
           (\text{narW}\langle w \rangle \prec_{\pi} \text{narR}\langle r, \_, \_, \_, \_ \rangle \land \text{samegp}(w, r)) \implies \mathsf{B}\langle w \rangle \prec_{\pi} \text{narR}\langle r, \_, \_, \_, \_ \rangle
```

## NIC Order.

```
\mathsf{nicActOrder}(\pi) \triangleq \forall a_1, a_2. \ \mathsf{NIC}(a_1) \prec_{\pi} \mathsf{NIC}(a_2) \land \mathsf{samegp}(a_1, a_2) \Longrightarrow
                (a_1 \in \mathtt{nF} \land a_2 \in \mathtt{Get} \land \mathtt{nrR}\langle -, -, a_2, - \rangle \in \pi \implies \mathtt{nF}\langle a_1 \rangle \prec_{\pi} \mathtt{nrR}\langle -, -, a_2, - \rangle)
  \land (a_1 \in \mathtt{nF} \land a_2 \in \mathtt{Put} \land \mathtt{nlR}\langle \_, \_, a_2, \_ \rangle \in \pi \implies \mathtt{nF}\langle a_1 \rangle \prec_{\pi} \mathtt{nlR}\langle \_, \_, a_2, \_ \rangle)
  \land (a_1 \in nF \land a_2 \in RCAS \land naF\langle_{-,-,}a_2,_{-}\rangle \in \pi \implies nF\langle a_1 \rangle \prec_{\pi} naF\langle_{-,-,}a_2,_{-,}\rangle)
  \land (a_1 \in \mathsf{nF} \land a_2 \in \mathsf{rRMW} \land \mathsf{narR}\langle \neg, \neg, a_2, \neg, \neg \rangle \in \pi \implies \mathsf{nF}\langle a_1 \rangle \prec_{\pi} \mathsf{narR}\langle \neg, \neg, a_2, \neg, \neg \rangle)
   \land (a_1 \in \mathtt{nF} \land a_2 \in \mathtt{nF} \land \mathtt{nF} \langle a_2 \rangle \in \pi \implies \mathtt{nF} \langle a_1 \rangle \prec_{\pi} \mathtt{nF} \langle a_2 \rangle)
  \land \ (a_1 \in \mathsf{Get} \land a_2 \in \mathsf{nF} \land \mathsf{nF} \langle a_2 \rangle \in \pi \implies \mathsf{nrR} \langle \neg, \neg, a_1, w_1 \rangle \prec_{\pi} \mathsf{nlW} \langle w_1, \neg \rangle \prec_{\pi} \mathsf{nF} \langle a_2 \rangle)
                \left( \begin{array}{l} a_1 \in \operatorname{Put} \wedge a_2 \in \operatorname{nF} \wedge \operatorname{nF} \langle a_2 \rangle \in \pi \\ \Longrightarrow \operatorname{nlR} \langle \_, \_, a_1, w_1 \rangle \prec_{\pi} \operatorname{nrW} \langle w_1, e_1 \rangle \prec_{\pi} \operatorname{CN} \langle e_1 \rangle \prec_{\pi} \operatorname{nF} \langle a_2 \rangle \end{array} \right) 
                     \begin{pmatrix} a_1 \in \mathtt{RCAS} \wedge a_2 \in \mathtt{nF} \wedge \mathtt{nF} \langle a_2 \rangle \in \pi \\ \Longrightarrow \mathtt{narR} \langle {}_{-}, {}_{-}, a_1, w_1, w_2 \rangle \prec_{\pi} \mathtt{narW} \langle w_2 \rangle \prec_{\pi} \mathtt{nlW} \langle w_1, {}_{-} \rangle \prec_{\pi} \mathtt{nF} \langle a_2 \rangle \\ \lor \mathtt{naF} \langle {}_{-}, {}_{-}, a_1, w_1 \rangle \prec_{\pi} \mathtt{nlW} \langle w_1, {}_{-} \rangle \prec_{\pi} \mathtt{nF} \langle a_2 \rangle \end{pmatrix} 
                 \begin{pmatrix} a_1 \in \mathtt{RFAA} \wedge a_2 \in \mathtt{nF} \wedge \mathtt{nF} \langle a_2 \rangle \in \pi \\ \Longrightarrow \mathtt{narR} \langle ., ., a_1, w_1, w_2 \rangle \prec_{\pi} \mathtt{narW} \langle w_2 \rangle \prec_{\pi} \mathtt{nlW} \langle w_1, . \rangle \prec_{\pi} \mathtt{nF} \langle a_2 \rangle \end{pmatrix} 
                  \begin{pmatrix} a_1 \in \mathsf{Get} \wedge a_2 \in \mathsf{Get} \wedge \mathsf{nrR}\langle \_, \_, a_2, w_2 \rangle \prec_\pi \mathsf{nlW}\langle w_2, \_ \rangle \\ \Longrightarrow \mathsf{nrR}\langle \_, \_, a_1, w_1 \rangle \prec_\pi \mathsf{nlW}\langle w_1, \_ \rangle \prec_\pi \mathsf{nlW}\langle w_2, \_ \rangle \\ \end{pmatrix} 
  \land \  \left( \begin{matrix} a_1 \in \mathsf{Get} \land a_2 \in \mathsf{Put} \land \mathsf{nlR}\langle \_, \_, a_2, w_2 \rangle \prec_\pi \mathsf{nrW}\langle w_2, e_2 \rangle \prec_\pi \mathsf{CN}\langle e_2 \rangle \\ \Longrightarrow \mathsf{nrR}\langle \_, \_, a_1, w_1 \rangle \prec_\pi \mathsf{nlW}\langle w_1, \_ \rangle \prec_\pi \mathsf{CN}\langle e_2 \rangle \end{matrix} \right) 
                 \left( \begin{array}{l} a_1 \in \operatorname{Get} \wedge a_2 \in \operatorname{RCAS} \wedge \left( \begin{array}{l} \operatorname{narR} \langle \_, \_, a_2, w_2, \_ \rangle \prec_{\pi} \operatorname{nlW} \langle w_2, \_ \rangle \\ \vee \operatorname{naF} \langle \_, \_, a_2, w_2 \rangle \prec_{\pi} \operatorname{nlW} \langle w_2, \_ \rangle \end{array} \right) \right)   \Longrightarrow \operatorname{nrR} \langle \_, \_, a_1, w_1 \rangle \prec_{\pi} \operatorname{nlW} \langle w_1, \_ \rangle \prec_{\pi} \operatorname{nlW} \langle w_2, \_ \rangle 
  \land \ \left( \begin{array}{l} a_1 \in \mathtt{Get} \land a_2 \in \mathtt{RFAA} \land \mathtt{narR}\langle \_, \_, a_2, w_2, \_\rangle \prec_\pi \mathtt{nlW}\langle w_2, \_\rangle \\ \Longrightarrow \mathtt{nrR}\langle \_, \_, a_1, w_1\rangle \prec_\pi \mathtt{nlW}\langle w_1, \_\rangle \prec_\pi \mathtt{nlW}\langle w_2, \_\rangle \end{array} \right.
                \begin{pmatrix} a_1 \in \mathtt{Put} \land a_2 \in \mathtt{Get} \land \mathtt{nrR} \langle \_, \_, a_2, \_ \rangle \in \pi \\ \Longrightarrow \mathtt{nlR} \langle \_, \_, a_1, w_1 \rangle \prec_{\pi} \mathtt{nrW} \langle w_1, \_ \rangle \prec_{\pi} \mathtt{nrR} \langle \_, \_, a_2, \_ \rangle 
                 \begin{pmatrix} a_1 \in \operatorname{Put} \wedge a_2 \in \operatorname{Get} \wedge \operatorname{nrR}\langle \_, \_, a_2, w_2 \rangle \prec_{\pi} \operatorname{nlW}\langle w_2, \_\rangle \\ \Longrightarrow \operatorname{nlR}\langle \_, \_, a_1, w_1 \rangle \prec_{\pi} \operatorname{nrW}\langle w_1, e_1 \rangle \prec_{\pi} \operatorname{CN}\langle e_1 \rangle \prec_{\pi} \operatorname{nlW}\langle w_2, \_\rangle \end{pmatrix} 
                \left( \begin{array}{l} a_1 \in \mathtt{Put} \wedge a_2 \in \mathtt{RCAS} \wedge \mathtt{naF}\langle \_, \_, a_2, w_2 \rangle \in \pi \\ \Longrightarrow \ \mathtt{n1R}\langle \_, \_, a_1, w_1 \rangle \prec_\pi \mathtt{nrW}\langle w_1, e_1 \rangle \prec_\pi \mathtt{naF}\langle \_, \_, a_2, w_2 \rangle \right) 
                \land (a_1 \in \mathtt{Put} \land a_2 \in \mathtt{Put} \land \mathtt{nlR}\langle \_, \_, a_2, \_\rangle \in \pi \implies \mathtt{nlR}\langle \_, \_, a_1, \_\rangle \prec_{\pi} \mathtt{nlR}\langle \_, \_, a_2, \_\rangle)
                 \begin{pmatrix} a_1 \in \mathtt{Put} \land a_2 \in \mathtt{Put} \land \mathtt{nlR} \langle \_, \_, a_2, w_2 \rangle \prec_{\pi} \mathtt{nrW} \langle w_2, \_ \rangle \\ \Longrightarrow \mathtt{nlR} \langle \_, \_, a_1, w_1 \rangle \prec_{\pi} \mathtt{nrW} \langle w_1, \_ \rangle \prec_{\pi} \mathtt{nrW} \langle w_2, \_ \rangle \\ \end{pmatrix} 
 \wedge \ \left( \begin{array}{l} a_1 \in \mathtt{Put} \wedge a_2 \in \mathtt{Put} \wedge \mathtt{nlR} \langle \_, \_, a_2, w_2 \rangle \prec_{\pi} \mathtt{nrW} \langle w_2, e_2 \rangle \prec_{\pi} \mathtt{CN} \langle e_2 \rangle \\ \Longrightarrow \ \mathtt{nlR} \langle \_, \_, a_1, w_1 \rangle \prec_{\pi} \mathtt{nrW} \langle w_1, e_1 \rangle \prec_{\pi} \mathtt{CN} \langle e_1 \rangle \prec_{\pi} \mathtt{CN} \langle e_2 \rangle \end{array} \right)
```

## NIC Atomicity.

$$\begin{split} \operatorname{nicAtomicity}(\pi) \; & \triangleq \; \forall a_1, a_2, r, w. \\ \left( \begin{array}{c} \lambda_1 = \operatorname{narR}\langle r_1, .., a_1, .., w \rangle \\ \wedge \; \lambda_2 \in \{\operatorname{naF}\langle .., .., a_2, .. \rangle, \operatorname{narR}\langle .., .., a_2, .., .. \rangle\} \\ \wedge \; a_1, a_2 \in \operatorname{rRMW} \; \wedge \; \overline{n}(a_1) = \overline{n}(a_2) \; \wedge \; \lambda_1 \prec_{\pi} \lambda_2 \end{array} \right) \; \Longrightarrow \; \operatorname{B}\langle w \rangle \prec_{\pi} \lambda_2 \end{split}$$

## Read Order.

$$\begin{array}{lll} \mathsf{wfrd}(\pi) & \triangleq & \forall \pi_2, r, w, \pi_1. \ \pi = \pi_2 \cdot \mathsf{1R} \langle r, w \rangle \cdot \pi_1 \implies \mathsf{wfrdCPU}(r, w, \pi_1) \\ & \wedge \ \forall \pi_2, u, w, \pi_1. \ \pi = \pi_2 \cdot \mathsf{CAS} \langle u, w \rangle \cdot \pi_1 \implies \mathsf{wfrdCPU}(u, w, \pi_1) \\ & \wedge \ \forall \pi_2, r, w, \pi_1. \ \pi = \pi_2 \cdot \mathsf{n1R} \langle r, w, ., . \rangle \cdot \pi_1 \implies \mathsf{wfrdNIC}(r, w, \pi_1) \\ & \wedge \ \forall \pi_2, r, w, \pi_1. \ \pi = \pi_2 \cdot \mathsf{nrR} \langle r, w, ., . \rangle \cdot \pi_1 \implies \mathsf{wfrdNIC}(r, w, \pi_1) \\ & \wedge \ \forall \pi_2, r, w, \pi_1. \ \pi = \pi_2 \cdot \mathsf{naF} \langle r, w, ., ., . \rangle \cdot \pi_1 \implies \mathsf{wfrdNIC}(r, w, \pi_1) \\ & \wedge \ \forall \pi_2, r, w, \pi_1. \ \pi = \pi_2 \cdot \mathsf{narR} \langle r, w, ., ., . \rangle \cdot \pi_1 \implies \mathsf{wfrdNIC}(r, w, \pi_1) \end{array}$$

$$\begin{split} \mathsf{wfrdCPU}(r,w,\pi) \; &\triangleq \; \begin{pmatrix} \exists \pi_2,\lambda,\pi_1. \; \pi = \pi_2 \cdot \lambda \cdot \pi_1 \\ \land \; \lambda \in \{\mathsf{B}\langle w \rangle, \mathsf{CAS}\langle w, \lrcorner \rangle \} \\ \land \; \{\mathsf{B}\langle w' \rangle, \mathsf{CAS}\langle w', \lrcorner \rangle \in \pi_2 \mid \mathsf{loc}(w') = \mathsf{loc}(r) \} = \emptyset \\ \land \; \left\{ w' \mid \mathsf{lW}\langle w' \rangle \in \pi \land \mathsf{B}\langle w' \rangle \notin \pi \land \\ \land \; \mathsf{loc}(w') = \mathsf{loc}(r) \land t(w') = t(r) \end{cases} \right\} = \emptyset \\ \lor \; \left\{ \exists \pi_2,\lambda,\pi_1. \; \pi = \pi_2 \cdot \lambda \cdot \pi_1 \\ \land \; \lambda = \mathsf{lW}\langle w \rangle \land t(w) = t(r) \land \mathsf{B}\langle w \rangle \notin \pi_2 \\ \land \; \{\mathsf{lW}\langle w' \rangle \in \pi_2 \mid \mathsf{loc}(w') = \mathsf{loc}(r) \land t(w') = t(r) \} = \emptyset \right) \\ \lor \; \left\{ \begin{matrix} w = init_{\mathsf{loc}(w)} \land \\ \mathsf{B}\langle w' \rangle, \mathsf{CAS}\langle w', \lrcorner \rangle \in \pi, \; \middle| \; \mathsf{loc}(w') = \mathsf{loc}(r) \land t(w'') = t(r) \\ \mathsf{lW}\langle w'' \rangle \in \pi \end{matrix} \right\} = \emptyset \\ \end{split} \\ \mathsf{wfrdNIC}(r,w,\pi) \; \triangleq \; \left\{ \begin{matrix} \exists \pi_2,\lambda,\pi_1. \; \pi = \pi_2 \cdot \lambda \cdot \pi_1 \\ \land \; \lambda \in \{\mathsf{B}\langle w \rangle,\mathsf{CAS}\langle w, \lrcorner \rangle \} \\ \land \; \{\mathsf{B}\langle w' \rangle,\mathsf{CAS}\langle w', \lrcorner \rangle \in \pi_2 \mid \mathsf{loc}(w') = \mathsf{loc}(r) \} = \emptyset \\ \end{cases} \\ \lor \; \left\{ \begin{matrix} \exists w : \mathsf{loc}(w) \land \mathsf{loc}(w') = \mathsf{loc}(w') = \mathsf{loc}(r) \end{Bmatrix} \right\} \\ \lor \; \left\{ \begin{matrix} \exists w : \mathsf{loc}(w) \land \mathsf{loc}(w', \neg) \in \pi_2 \mid \mathsf{loc}(w') = \mathsf{loc}(r) \end{Bmatrix} \right\} \\ \lor \; \left\{ \begin{matrix} \exists w : \mathsf{loc}(w) \land \mathsf{loc}(w', \neg) \in \pi_2 \mid \mathsf{loc}(w') = \mathsf{loc}(r) \end{Bmatrix} \right\} \\ \lor \; \left\{ \begin{matrix} \exists w : \mathsf{loc}(w) \land \mathsf{loc}(w', \neg) \in \pi_2 \mid \mathsf{loc}(w') = \mathsf{loc}(r) \end{Bmatrix} \right\} \\ \lor \; \left\{ \begin{matrix} \exists w : \mathsf{loc}(w) \land \mathsf{loc}(w', \neg) \in \pi_2 \mid \mathsf{loc}(w') = \mathsf{loc}(r) \end{Bmatrix} \right\} \\ \lor \; \left\{ \begin{matrix} \exists w : \mathsf{loc}(w) \land \mathsf{loc}(w', \neg) \in \pi_2 \mid \mathsf{loc}(w') = \mathsf{loc}(r) \end{Bmatrix} \right\} \\ \lor \; \left\{ \begin{matrix} \exists w : \mathsf{loc}(w) \land \mathsf{loc}(w', \neg) \in \pi_2 \mid \mathsf{loc}(w') = \mathsf{loc}(r) \end{Bmatrix} \right\} \\ \lor \; \left\{ \begin{matrix} \exists w : \mathsf{loc}(w) \land \mathsf{loc}(w', \neg) \in \pi_2 \mid \mathsf{loc}(w') = \mathsf{loc}(r) \end{Bmatrix} \right\} \\ \lor \; \left\{ \begin{matrix} \exists w : \mathsf{loc}(w) \land \mathsf{loc}(w', \neg) \in \pi_2 \mid \mathsf{loc}(w$$

# Well-formed path.

$$\begin{split} \mathsf{wfp}(\pi) &\triangleq & \mathsf{nodup}(\pi) \\ & \land & \mathsf{backComp}(\pi) \\ & \land & \mathsf{bufFlushOrd}(\pi) \\ & \land & \mathsf{pollOrder}(\pi) \\ & \land & \mathsf{nicActOrder}(\pi) \\ & \land & \mathsf{nicAtomicity}(\pi) \\ & \land & \mathsf{wfrd}(\pi) \end{split}$$

#### Definition 21.

$$\begin{split} \mathsf{wf}(\mathsf{M},\mathsf{B},\mathsf{A},\mathsf{QP},\pi) &\triangleq & \mathsf{wfp}(\pi) \\ &\wedge \ \forall x \in \mathsf{Loc.} \ \mathsf{M}(x) = \mathsf{read}(\pi,x) \\ &\wedge \ \forall t \in \mathsf{Tid.B}(t) = \mathsf{mksbuff}(\varepsilon,t,\pi) \\ &\wedge \ \forall n \in \mathsf{Node.A}(n) = \mathsf{chkatm}(n,\pi) \\ &\wedge \ \forall t \in \mathsf{Tid.} \forall \overline{n} \in (\mathsf{Node} \setminus \{n(t)\}). \begin{pmatrix} \mathsf{QP}(t)(\overline{n}).\mathbf{pipe} = \mathsf{mkpipe}(\varepsilon,t,\overline{n},\pi) \\ \mathsf{QP}(t)(\overline{n}).\mathbf{wb_R} = \mathsf{mkwbR}(\varepsilon,t,\overline{n},\pi) \\ \mathsf{QP}(t)(\overline{n}).\mathbf{wb_L} = \mathsf{mkwbL}(\varepsilon,t,\overline{n},\pi) \end{pmatrix} \end{split}$$

Where, the functions read, mksbuff, chkatm, mkpipe, mkwbR, and mkwbL are defined below.

$$\begin{split} \operatorname{read}(\lambda \cdot \pi, x) &\triangleq \begin{cases} w & \lambda \in \{\mathsf{B}\langle w \rangle, \mathsf{CAS}\langle w, \lrcorner \rangle\} \wedge \mathsf{loc}(w) = x \\ \operatorname{read}(\varepsilon, x) & \text{otherwise} \end{cases} \\ \operatorname{read}(\varepsilon, x) &\triangleq init_x \end{split}$$

$$\mathsf{mksbuff}(\mathsf{b},t,\varepsilon) \triangleq \mathsf{b}$$
 
$$\mathsf{mksbuff}(\mathsf{b},t,\pi \cdot \lambda) \triangleq \begin{cases} \mathsf{mksbuff}(w \cdot \mathsf{b},t,\pi) & \lambda = \mathsf{1W}\langle w \rangle \wedge t(w) = t \wedge \mathsf{B}\langle w \rangle \notin \pi \\ \mathsf{mksbuff}(a \cdot \mathsf{b},t,\pi) & \lambda = \mathsf{Push}\langle a \rangle \wedge \mathsf{NIC}\langle a \rangle \notin \pi \wedge t(a) = t \\ \mathsf{mksbuff}(\mathsf{b},t,\pi) & \mathsf{otherwise} \end{cases}$$

$$\mathsf{chkatm}(n,\pi) \triangleq \begin{cases} \bot & \forall w. \begin{pmatrix} \mathsf{narR}\langle \_, \_, a, \_, w \rangle \in \pi \\ \land \, \overline{n}(a) = n \end{pmatrix} \implies \mathsf{B}\langle w \rangle \in \pi \\ \top & \mathsf{otherwise} \end{cases}$$

$$\mathsf{mkpipe}(\mathbf{pipe}, t, \overline{n}, \varepsilon) \triangleq \mathbf{pipe}$$

$$\mathsf{mkpipe}(\mathbf{pipe},t,n,\varepsilon) = \mathbf{pipe}$$

$$\begin{pmatrix} \mathsf{mkpipe}(a \cdot \mathbf{pipe},t,\overline{n},\pi) & \mathsf{if} & \begin{pmatrix} t(\lambda) = t \wedge \overline{n}(\lambda) = \overline{n} \wedge \lambda = \mathsf{NIC}\langle a \rangle \\ \wedge \mathsf{nIR}\langle_{-}, -, a, -\rangle \not\in \pi \wedge \mathsf{nrR}\langle_{-}, -, a, -\rangle \not\in \pi \\ \wedge \mathsf{nr}\langle a \rangle \not\in \pi \wedge \mathsf{nrr}\langle_{-}, -, a, -\rangle \not\in \pi \end{pmatrix}$$

$$\mathsf{mkpipe}(w \cdot \mathbf{pipe},t,\overline{n},\pi) & \mathsf{if} & \begin{pmatrix} t(\lambda) = t \wedge \overline{n}(\lambda) = \overline{n} \wedge \lambda = \mathsf{NIC}\langle a \rangle \\ \wedge \mathsf{nIR}\langle_{-}, -, a, -, -\rangle \not\in \pi \end{pmatrix}$$

$$\mathsf{mkpipe}(e \cdot \mathbf{pipe},t,\overline{n},\pi) & \mathsf{if} & \begin{pmatrix} t(\lambda) = t \wedge \overline{n}(\lambda) = \overline{n} \wedge \lambda = \mathsf{NIC}\langle a \rangle \\ \wedge \mathsf{nIR}\langle_{-}, -, a, w \rangle \in \pi \wedge \mathsf{nrW}\langle w, -\rangle \not\in \pi \end{pmatrix}$$

$$\mathsf{mkpipe}(w \cdot \mathbf{pipe},t,\overline{n},\pi) & \mathsf{if} & \begin{pmatrix} t(\lambda) = t \wedge \overline{n}(\lambda) = \overline{n} \wedge \lambda = \mathsf{NIC}\langle a \rangle \\ \wedge \mathsf{nrR}\langle_{-}, -, a, w \rangle \in \pi \wedge \mathsf{nrW}\langle w, -\rangle \not\in \pi \end{pmatrix}$$

$$\mathsf{mkpipe}(w \cdot \mathbf{pipe},t,\overline{n},\pi) & \mathsf{if} & \begin{pmatrix} t(\lambda) = t \wedge \overline{n}(\lambda) = \overline{n} \wedge \lambda = \mathsf{NIC}\langle a \rangle \\ \wedge \mathsf{nrR}\langle_{-}, -, a, w \rangle \in \pi \wedge \mathsf{nIW}\langle w, -\rangle \not\in \pi \end{pmatrix}$$

$$\mathsf{mkpipe}(w \cdot w' \cdot \mathbf{pipe},t,\overline{n},\pi) & \mathsf{if} & \begin{pmatrix} t(\lambda) = t \wedge \overline{n}(\lambda) = \overline{n} \wedge \lambda = \mathsf{NIC}\langle a \rangle \\ \wedge \mathsf{nar}\langle_{-}, -, a, w \rangle \in \pi \wedge \mathsf{nIW}\langle w, -\rangle \not\in \pi \end{pmatrix}$$

$$\mathsf{mkpipe}(\mathbf{pipe},t,\overline{n},\pi) & \mathsf{if} & \begin{pmatrix} t(\lambda) = t \wedge \overline{n}(\lambda) = \overline{n} \wedge \lambda = \mathsf{NIC}\langle a \rangle \\ \wedge \mathsf{nar}\langle_{-}, -, a, w \rangle \in \pi \wedge \mathsf{narW}\langle w' \rangle \not\in \pi \end{pmatrix}$$

$$\mathsf{mkpipe}(\mathbf{pipe},t,\overline{n},\pi) & \mathsf{otherwise}$$

 $\mathsf{mkwbR}(\mathbf{wb}_{\mathbf{R}}, t, \overline{n}, \varepsilon) \triangleq \mathbf{wb}_{\mathbf{R}}$ 

$$\mathsf{mkwbR}(\mathbf{wb_R}, t, \overline{n}, \pi \cdot \lambda) \triangleq \begin{cases} \mathsf{mkwbR}(w \cdot \mathbf{wb_R}, t, \overline{n}, \pi) & \text{if } \begin{pmatrix} t(\lambda) = t \wedge \overline{n}(\lambda) = \overline{n} \wedge \mathsf{B}\langle w \rangle \notin \pi \\ \wedge \lambda \in \{\mathsf{nrW}\langle w, \lrcorner \rangle, \mathsf{narW}\langle w \rangle\} \end{cases} \\ \mathsf{mkwbR}(\mathbf{wb_R}, t, \overline{n}, \pi) & \text{otherwise} \end{cases}$$

 $\mathsf{mkwbL}(\mathbf{wb_L}, t, \overline{n}, \varepsilon) \triangleq \mathbf{wb_L}$ 

$$\mathsf{mkwbL}(\mathbf{w}\mathbf{b_L},t,\overline{n},\pi) \triangleq \begin{cases} \mathsf{mkwbL}(e \cdot w \cdot \mathbf{w}\mathbf{b_L},t,\overline{n},\pi) & \text{if } \begin{pmatrix} t(\lambda) = t \wedge \ \overline{n}(\lambda) = \overline{n} \wedge \ \lambda = \mathsf{nlW}\langle w,e \rangle \\ \wedge \ \mathsf{B}\langle w \rangle \notin \pi \wedge \mathsf{P}\langle_{-},e \rangle \notin \pi \end{cases} \\ \mathsf{mkwbL}(e \cdot \mathbf{w}\mathbf{b_L},t,\overline{n},\pi) & \text{if } \begin{pmatrix} t(\lambda) = t \wedge \ \overline{n}(\lambda) = \overline{n} \wedge \ \lambda = \mathsf{nlW}\langle w,e \rangle \\ \wedge \ \mathsf{B}\langle w \rangle \in \pi \wedge \mathsf{P}\langle_{-},e \rangle \notin \pi \end{cases} \\ \mathsf{mkwbL}(e \cdot \mathbf{w}\mathbf{b_L},t,\overline{n},\pi) & \text{if } \begin{pmatrix} t(\lambda) = t \wedge \ \overline{n}(\lambda) = \overline{n} \wedge \ \lambda = \mathsf{nlW}\langle w,e \rangle \\ \wedge \ \mathsf{B}\langle w \rangle \in \pi \wedge \mathsf{P}\langle_{-},e \rangle \notin \pi \end{cases} \\ \mathsf{mkwbL}(\mathbf{w}\mathbf{b_L},t,\overline{n},\pi) & \text{otherwise} \end{cases}$$

**Theorem 8.** For all  $P, P', M, M', B, B', A, A', QP, QP', \pi, \pi'$ :

- wf(M<sub>0</sub>, B<sub>0</sub>, A<sub>0</sub>, QP<sub>0</sub>,  $\varepsilon$ );
- if P, M, B, A, QP,  $\pi \Rightarrow$  P', M', B', A', QP',  $\pi'$  and wf(M, B, A, QP,  $\pi$ ) then wf(M', B', A', QP',  $\pi'$ );
- if  $P, M_0, B_0, A_0, QP_0, \varepsilon \Rightarrow^* (\lambda t.skip), M, B_0, A_0, QP, \pi \text{ such that for all } t, \overline{n} \text{ we have } QP(t)(\overline{n}) = \langle \varepsilon, \varepsilon, nEX^* \rangle, \text{ then } wf(M, B_0, A_0, QP, \pi) \text{ and } complete(\pi).$

The proof of the first part follows trivially from the definitions of  $M_0$ ,  $B_0$ ,  $A_0$ , and  $QP_0$ . The second part is proved by induction on the structure of  $\Rightarrow$ . The last part follows from the previous two parts and induction on the length of  $\Rightarrow^*$ , as well as how the definition of wf on empty store buffers and queue pairs (regardless of nex in wb<sub>L</sub>) implies complete( $\pi$ ).

#### D.5 From Annotated Semantics to Declarative Semantics

We define

$$\mathsf{getEG}(\pi) \triangleq \begin{cases} (\mathsf{Event}, \mathsf{po}, \mathsf{rf}, \mathsf{pf}, \mathsf{mo}, \mathsf{nfo}, \mathsf{rao}) & \text{if } \mathsf{wfp}(\pi) \land \mathsf{complete}(\pi) \\ \mathsf{undefined} & \mathsf{otherwise} \end{cases}$$

with

Event 
$$\triangleq$$
 Event<sub>0</sub>  $\cup$  {getA( $\lambda$ ) |  $\lambda \in \pi$ }

Recall that  $\mathsf{Event}_0$  is the set of initialisation events  $\{init_x \mid x \in \mathsf{Loc}\}$ , where  $l(init_x) = \mathsf{1W}(x,0)$ 

$$getA: ALabel \rightarrow Event$$

We define  $\mathsf{getl}\lambda(\cdot,\pi)$  and  $\mathsf{getO}\lambda(\cdot,\pi)$  to perform the reverse operation of  $\mathsf{getA}$ . In the case of write events,  $\mathsf{getl}\lambda(\cdot,\pi)$  retrieves the first label sending the write to the buffer, while  $\mathsf{getO}\lambda(\cdot,\pi)$  retrieves the second label committing the write to memory.

$$\operatorname{getI}\lambda(-,\pi), \operatorname{getO}\lambda(-,\pi) : \{\operatorname{getA}(\lambda) \mid \lambda \in \pi\} \to \operatorname{\mathsf{ALabel}}$$

For all  $\lambda \in \pi$ :

- if  $type(\lambda) \in \{lR, CAS, F, P, nlR, nrR, narR, naF, nF\},$ then  $getl\lambda(getA(\lambda), \pi) \triangleq getO\lambda(getA(\lambda), \pi) \triangleq \lambda;$
- if  $type(\lambda) \in \{lW, nlW, nrW, narW\}$ , then  $getl\lambda(getA(\lambda), \pi) \triangleq \lambda$  while  $getO\lambda(getA(\lambda), \pi) \triangleq B(\lambda)$ .
- if  $\lambda = \mathsf{B}\langle w \rangle$ , then from  $\mathsf{backComp}(\pi)$  there is  $\lambda' \prec_{\pi} \lambda$  such that  $\mathsf{type}(\lambda) \in \{\mathsf{1W}, \mathsf{n1W}, \mathsf{nrW}, \mathsf{narW}\}$  and  $\mathsf{getA}(\lambda') = \mathsf{getA}(\lambda) = w$ . From the previous case, we have  $\mathsf{getI}\lambda(w,\pi) \triangleq \lambda'$  and  $\mathsf{getO}\lambda(w,\pi) \triangleq \lambda$ .

From this we define two relations IB and OB on Event total on all meaningful events by copying the ordering in  $\pi$ .

$$\mathsf{IB} \triangleq \{(e_1, e_2) \mid \mathsf{getI}\lambda(e_1, \pi) \prec_{\pi} \mathsf{getI}\lambda(e_2, \pi)\} \cup (\mathsf{Event}_0 \times (\mathsf{Event} \setminus \mathsf{Event}_0))$$

$$\mathsf{OB} \triangleq \{(e_1, e_2) \mid \mathsf{getO}\lambda(e_1, \pi) \prec_{\pi} \mathsf{getO}\lambda(e_2, \pi)\} \cup (\mathsf{Event}_0 \times (\mathsf{Event} \setminus \mathsf{Event}_0))$$

From  $\mathsf{wfp}(\pi)$ , IB and OB are transitive and irreflexive. Note: we could make IB and OB total by adding an arbitrary total order on  $\mathsf{Event}_0$ .

$$\mathsf{rf} \triangleq \left\{ (w,r) \;\middle|\; \begin{array}{l} \mathsf{1R}\langle r,w\rangle \in \pi \vee \mathsf{n1R}\langle r,w,\_,\_\rangle \in \pi \vee \mathsf{nrR}\langle r,w,\_,\_\rangle \in \pi \vee \mathsf{CAS}\langle r,w\rangle \in \pi \\ \vee \; \mathsf{narR}\langle r,w,\_,\_\rangle \in \pi \vee \mathsf{naF}\langle r,w,\_,\_\rangle \in \pi \end{array} \right\}$$

$$\mathsf{pf} \triangleq \left\{ (w,p) \;\middle|\; \begin{array}{l} \mathsf{nlW}\langle w,e \rangle \prec_{\pi} \mathsf{P}\langle p,e \rangle \\ \vee \; \mathsf{nrW}\langle w,e \rangle \prec_{\pi} \mathsf{P}\langle p,e \rangle \end{array} \right\}$$

$$\lambda \in \{ \mathsf{1R} \langle e, \_ \rangle, \mathsf{1W} \langle e \rangle, \mathsf{CAS} \langle e, \_ \rangle, \mathsf{Push} \langle e \rangle, \mathsf{P} \langle e, \_ \rangle, \mathsf{F} \langle e \rangle \}$$
 
$$\lambda \subset \mathsf{Push} \langle a \rangle \wedge \begin{pmatrix} \lambda \subset_{\pi} \mathsf{n1R} \langle e, \_, a, \_ \rangle \\ \vee \lambda \subset_{\pi} \mathsf{n1R} \langle -, \_, a, e \rangle \\ \vee \lambda \subset_{\pi} \mathsf{nrR} \langle -, \_, a, e \rangle \\ \vee \lambda \subset_{\pi} \mathsf{nrR} \langle -, \_, a, e \rangle \\ \vee \lambda \subset_{\pi} \mathsf{naF} \langle -, \_, a, e \rangle \\ \vee \lambda \subset_{\pi} \mathsf{narR} \langle -, \_, a, e \rangle \\ \vee \lambda \subset_{\pi} \mathsf{narR} \langle -, \_, a, e, \_ \rangle \\ \vee \lambda \subset_{\pi} \mathsf{narR} \langle -, \_, a, e, \_ \rangle \\ \vee \lambda \subset_{\pi} \mathsf{narR} \langle -, \_, a, e, \_ \rangle \end{pmatrix}$$

$$\mathsf{po} \triangleq \begin{pmatrix} \mathsf{Event}_0 \times (\mathsf{Event} \setminus \mathsf{Event}_0) \\ \cup \left\{ (e_1, e_2) \middle| \begin{array}{c} \lambda_1 \prec_\pi \lambda_2 \wedge t(\lambda_1) = t(\lambda_2) \\ \wedge \lambda_1 \text{ generates } e_1 \text{ in } \pi \\ \wedge \lambda_2 \text{ generates } e_2 \text{ in } \pi \end{array} \right\} \\ \cup \left\{ (r, w) \middle| \begin{array}{c} \mathsf{nlR} \langle r, \ldots, w \rangle \in \pi \\ \vee \mathsf{narR} \langle r, \ldots, w \rangle \in \pi \\ \vee \mathsf{narR} \langle r, \ldots, w \rangle \in \pi \\ \vee \mathsf{narR} \langle r, \ldots, w \rangle \in \pi \end{array} \right\} \\ \cup \left\{ (w_1, w_2) \middle| \mathsf{narR} \langle -, \ldots, w_2, w_1 \rangle \right\} \in \pi \end{pmatrix}$$

$$\mathsf{nfo} \triangleq \begin{pmatrix} \{(r,w) \mid \mathsf{sameqp}(r,w) \land \mathsf{nlR}\langle r, \_, \_, \_\rangle \prec_{\pi} \mathsf{nlW}\langle w, \_\rangle \prec_{\pi} \mathsf{B}\langle w \rangle \} \\ \cup \ \{(r,w) \mid \exists \lambda_r, \lambda_w.\mathsf{sameqp}(r,w) \land \lambda_r \prec_{\pi} \lambda_w \prec_{\pi} \mathsf{B}\langle w \rangle \} \\ \cup \ \{(w,r) \mid \mathsf{sameqp}(w,r) \land \mathsf{nlW}\langle w, \_\rangle \prec_{\pi} \mathsf{B}\langle w \rangle \prec_{\pi} \mathsf{nlR}\langle r, \_, \_, \_\rangle \} \\ \cup \ \{(w,r) \mid \exists \lambda_r, \lambda_w.\mathsf{sameqp}(w,r) \land \lambda_w \prec_{\pi} \mathsf{B}\langle w \rangle \prec_{\pi} \lambda_r \} \\ \mathsf{where} \ \lambda_r \in \{\mathsf{nrR}\langle r', \ldots \rangle, \mathsf{naF}\langle r', \ldots \rangle, \mathsf{narR}\langle r', \ldots \rangle \} \\ \lambda_w \in \{\mathsf{nrW}\langle w, \_\rangle, \mathsf{narW}\langle w \rangle \} \\ \end{pmatrix}$$

$$\operatorname{rao} \triangleq \left( \left\{ (r_1, r_2) \;\middle|\; \overline{n}(a_1) = \overline{n}(a_2) \land \begin{pmatrix} \lambda_1 \prec_{\pi} \lambda_2 \\ \land \; \lambda_1 \in \{\operatorname{naF}\langle r_1, a_1, \ldots \rangle, \operatorname{narR}\langle r_1, a_1, \ldots \rangle\} \\ \land \; \lambda_2 \in \{\operatorname{naF}\langle r_2, a_2, \ldots \rangle, \operatorname{narR}\langle r_2, a_2, \ldots \rangle\} \end{pmatrix} \right) \right)$$

From an execution graph  $E = \text{getEG}(\pi)$ , we use the definitions of the paper to define oppo, ippo, rf<sub>i</sub>, rf<sub>e</sub>, rb<sub>i</sub>, ar<sub>i</sub>, ob, and ib.

$$\textbf{Lemma 2.} \ w \in \texttt{nlW} \implies \exists r. \left( \begin{array}{c} r \in \texttt{nrR} \wedge (r,w) \in \texttt{po}|_{imm} \\ \vee \ r \in \texttt{narR} \wedge (r,w) \in \texttt{po}|_{imm} \cup (\texttt{po}|_{imm})^2 \end{array} \right)$$

*Proof.* By definition of po, we can only have such  $w \in nlW$  if there is some  $\lambda = Push\langle a \rangle$  which generates w in  $\pi$ . Then we can consider the cases of a such that  $Push\langle a \rangle$  generates some  $w \in nlW$ . Either:

- $a \in Put$ , then there is some  $r \in nrR$  with  $(r, w) \in po|_{imm}$
- $a \in \text{RCAS} \cup \text{RFAA}$ , then there is some  $r \in \text{narR}$  with either  $(r, w) \in \text{po}|_{\text{imm}}$  (in the case of a failed RCAS) or  $(r, w) \in (\text{po}|_{\text{imm}})^2$  (in the case of a successful RCAS or RFAA)

**Theorem 9.** getEG( $\pi$ ) is well-formed.

*Proof.* We need to check the conditions of a pre-execution (Def. 18) and of well-formedness (Def. 19). For the pre-execution conditions:

- Checking Event<sup>0</sup> × (Event \ Event<sup>0</sup>) ⊆ po: by definition.
- Checking po is a union of strict partial orders each on one thread: If  $t(e_1) \neq t(e_2)$ , then  $(e_1, e_2) \notin po$  and  $(e_2, e_1) \notin po$  by definition. If  $t(e_1) = t(e_2)$ , then either  $(e_1, e_2) \in po$  or  $(e_2, e_1) \in po$ . This comes from the second case of the definition of po: if there is  $\lambda_1$  and  $\lambda_2$  such that  $\lambda_i$  generates  $e_i$  in  $\pi$ , then either  $\lambda_1 \prec_{\pi} \lambda_2$  or  $\lambda_2 \prec_{\pi} \lambda_1$ .
- Checking that rf is functional on its range: If  $r \in \mathcal{R} \subseteq \{ \text{getA}(\lambda) \mid \lambda \in \pi \}$ , then we have either  $\text{lR}\langle r, \_\rangle$ ,  $\text{nlR}\langle r, \_, \_, \_\rangle$ ,  $\text{nrR}\langle r, \_, \_, \_\rangle$ ,  $\text{naF}\langle r, \_, \_, \_\rangle$ , or  $\text{narR}\langle r, \_, \_, \_\rangle$  in  $\pi$ , and r have at least one antecedent.

If  $(w,r) \in \text{rf}$ , let us assume  $r \in \text{nlR}$ , then by definition  $\text{nlR}\langle r, w, \_, \_ \rangle \in \pi$ . Since  $\text{nodup}(\pi)$ , for all  $w' \neq w$ , we have  $\text{nlR}\langle r, w', \_, \_ \rangle \notin \pi$ , and syntactically we cannot write  $\text{lR}\langle r, \_ \rangle$  or  $\text{nrR}\langle r, \_, \_, \_ \rangle$ , so  $(w',r) \notin \text{rf}$ . Similarly,  $r \in \text{lR}$ ,  $r \in \text{nrR}$ ,  $r \in \text{nrR}$  only have one antecedent.

- Checking that rf relates events on the same location with matching values: By syntactic definition of the annotated labels 1R, n1R, nrR, naF and narR, e.g.,  $1R\langle r,w\rangle \implies eq_{loc\&v}(r,w)$ .
- Checking that mo is a union of strict total orders for writes on each variables: By definition of mo, given that we have  $\mathsf{complete}(\pi)$ , e.g., if  $\mathsf{1W}\langle w\rangle \in \pi$  then  $\mathsf{B}\langle w\rangle \in \pi$ .
- Checking that  $\operatorname{pf} \subseteq \operatorname{po} \cap \operatorname{sqp}$ : If  $(w,p) \in \operatorname{pf}$  with  $w \in \operatorname{nlW}$  (resp.  $\operatorname{nrW}$ ), then we have  $\operatorname{nlW}\langle w,e \rangle \prec_{\pi} \operatorname{P}\langle p,e \rangle$ . There is  $\lambda$  such that  $\lambda$  generates w in  $\pi$ , and we have  $\lambda \prec_{\pi} \operatorname{nlW}\langle w,e \rangle \prec_{\pi} \operatorname{P}\langle p,e \rangle$ . Also, t(p)=t(w) and  $\overline{n}(p)=\overline{n}(e)=\overline{n}(w)$ , so we have  $(w,p) \in \operatorname{po}$  and  $(w,p) \in \operatorname{sqp}$ .
- Checking that pf is functional on its domain: If  $(w,p) \in \mathsf{pf}$  with  $w \in \mathsf{nlW}$  (resp.  $\mathsf{nrW}$ ), then we have  $\mathsf{nlW}\langle w,e \rangle \prec_{\pi} \mathsf{P}\langle p,e \rangle$ . From  $\mathsf{nodup}(\pi)$ , for all  $p' \neq p$  we have  $\mathsf{P}\langle p,e \rangle \notin \pi$ , so w has at most one image.
- Checking that pf is total and functional on its range: If  $p \in \text{Event}$ , then there is  $e \in \text{nlex}$  (resp. nrex) such that  $P\langle p, e \rangle \in \pi$ . From backComp $(\pi)$  there is  $w \in \text{nlw}$  (resp. nrw) such that  $\text{nlw}\langle w, e \rangle \prec_{\pi} P\langle p, e \rangle$ , and so  $(w, p) \in \text{pf}$ . From  $\text{nodup}(\pi)$ , e cannot be used in another nlw (resp. nrw) annotated label, and p has exactly one antecedent.
- Checking that for all  $(a,b) \in \operatorname{sqp}$ ,  $a \in \operatorname{nrR} \cup \operatorname{naF} \cup \operatorname{narR}$ ,  $b \in \operatorname{nrW} \cup \operatorname{narW}$ , (resp.  $\operatorname{nlR/nlW}$ ) then  $(a,b) \in \operatorname{nfo} \cup \operatorname{nfo}^{-1}$ :

  By definition of nfo, given that  $\operatorname{bufFlushOrd}(\pi)$  forbids such interleavings as  $\operatorname{nrW}\langle w, \_ \rangle \prec_{\pi} \operatorname{nrR}\langle r, \_, \_, \_ \rangle \prec_{\pi} \operatorname{B}\langle w \rangle$  (resp.  $\operatorname{nlW}$  and  $\operatorname{nlR}$ ) when  $\operatorname{sameqp}(r,w)$ .
- Checking that rao is a union of strict total orders for remote atomic reads: By definition of rao.

For the well-formedness conditions:

- (1) Let us assume  $(w_1, w_2) \in po \cap sqp$  and  $(w_2, p_2) \in pf$ . The three events are on the same thread and queue pair.
  - If  $w_1 \in \mathtt{nlW}$ , then by  $\mathtt{complete}(\pi)$  there is a chain  $\mathtt{Push}\langle a_1 \rangle \prec_\pi \mathtt{NIC}\langle a_1 \rangle \prec_\pi \mathtt{nR} \prec_\pi \mathtt{nlW}\langle w_1, e_1 \rangle$  for some  $\mathtt{nR} \in \{\mathtt{nrR}\langle_-, -, a_1, w_1 \rangle, \mathtt{naF}\langle_-, -, a_1, w_1 \rangle, \mathtt{narR}\langle_-, -, a_1, w_1, \rangle\};$  if  $w_1 \in \mathtt{nrW}$ , there is instead a chain  $\mathtt{Push}\langle a_1 \rangle \prec_\pi \mathtt{NIC}\langle a_1 \rangle \prec_\pi \mathtt{nlR}\langle_-, -, a_1, w_1 \rangle \prec_\pi \mathtt{nrW}\langle w_1, e_1 \rangle \prec_\pi \mathtt{CN}\langle e_1 \rangle$ . Similarly there is a chain for  $w_2$ . By  $(w_1, w_2) \in \mathtt{po}$  we have  $\mathtt{Push}\langle a_1 \rangle \prec_\pi \mathtt{Push}\langle a_2 \rangle$ , and by  $\mathtt{bufFlushOrd}(\pi)$  we have  $\mathtt{NIC}\langle a_1 \rangle \prec_\pi \mathtt{NIC}\langle a_2 \rangle$ .
  - Let us call  $\lambda_1$  the last annotated label on the chain for  $w_1$ , i.e., either  $\mathtt{nlW}\langle w_1, e_1 \rangle$  or  $\mathsf{CN}\langle e_1 \rangle$ . Similarly,  $\lambda_2$  is the last annotated label on the chain for  $w_2$ . There are four cases to consider, but in all four  $\mathsf{nicActOrder}(\pi)$  implies  $\lambda_1 \prec_{\pi} \lambda_2$ . Then, from  $\mathsf{pollOrder}(\pi)$ , there is  $p_1$  such that  $\mathsf{P}\langle p_1, e_1 \rangle \prec_{\pi} \mathsf{P}\langle p_2, e_2 \rangle$ . By definitions, we have both  $(w_1, p_1) \in \mathsf{pf}$  and  $(p_1, p_2) \in \mathsf{po}$ .
- (2) If  $r \in \mathtt{nlR}$ , then there is  $w \in \mathtt{nrW}$  (taken from  $\mathtt{nlR}\langle r, \_, \_, w \rangle$ ) such that  $(r, w) \in \mathtt{po}|_{\mathrm{imm}}$ . This is by the last case of definition of po, since there is  $\lambda_a$  such that we have both  $\lambda_a$  generates r in  $\pi$  and  $\lambda_a$  generates w in  $\pi$ . Similarly for  $\mathtt{nrR/nlW}$  and  $\mathtt{nrW/nlR}$ .

- (3) If  $(r, w) \in \text{po}|_{\text{imm}}$ ,  $\text{type}(r) \in \{\text{nlR}, \text{nrR}\}$ , and  $\text{type}(w) \in \{\text{nlW}, \text{nrW}\}$ , then  $(r, w) \in \text{po}$  comes from the third case of the definition of po, and we have either  $\text{nlR}\langle r, \_, \_, w \rangle$  or  $\text{nrR}\langle r, \_, \_, w \rangle$  in  $\pi$ . In both cases, we have  $v_r(r) = v_w(w)$  by syntactic definition of the annotated labels.
- (4) (a) If  $r \in \mathtt{narR}$ , then either: There is  $\mathtt{naF}\langle r, ., ., w \rangle \in \pi$  such that  $w \in \mathtt{nlW}$  and  $(r, w) \in \mathtt{po}|_{\mathrm{imm}}$ . This follows from the second case definition of po. There is  $\mathtt{narR}\langle r, ., ., w_2, w_1 \rangle \in \pi$  such that  $w_1 \in \mathtt{narW}$ ,  $w_2 \in \mathtt{nlW}$ , and  $(r, w_1), (w_1, w_2) \in \mathtt{po}|_{\mathrm{imm}}$ . This follows from the second and third cases of the definition of po since there is  $\lambda_a$  which generates  $r, w_1$  and  $w_2$  in  $\pi$ . (b) If  $w \in \mathtt{narW}$  then  $(r, w), (w, w') \in \mathtt{po}|_{\mathrm{imm}}$  with  $r \in \mathtt{narR}$  and  $w' \in \mathtt{nlW}$  comes from the second case definition of po.
- (5) If  $(r,w) \in G.\mathsf{po}|_{\mathrm{imm}}$ ,  $\mathsf{type}(r) = \mathsf{narR}$  and  $\mathsf{type}(w) = \mathsf{nlW}$ , then (r,w) comes from the second case definition of po and we have  $\mathsf{naF}\langle r, ., ., w \rangle \in \pi$ . Then  $v_{\mathrm{r}}(r) = v_{\mathrm{w}}(w)$  by the syntax of annotated labels. If  $(r,w_1), (w_1,w_2) \in G.po|_{\mathrm{imm}}$ ,  $\mathsf{type}(r) = \mathsf{narR}$ ,  $\mathsf{type}(w_1) = \mathsf{narW}$  and  $\mathsf{type}(w_2) = \mathsf{nlW}$ , then  $(r,w_1)$  comes from the second case definition of po and  $(w_1,w_2)$  from the third case, so we have  $\mathsf{narR}\langle r, ., ., w_2, w_1 \rangle \in \pi$ . Then  $v_{\mathrm{r}}(r) = v_{\mathrm{w}}(w_2)$  by the syntax of annotated labels.
- (6) Comes from Lem. 2.

Lemma 3. OB;  $[Inst] \subseteq IB$  and [Inst];  $IB \subseteq OB$ .

*Proof.* If  $(e_1, e_2) \in \mathsf{OB}$ ; [Inst], then  $\mathsf{getO}\lambda(e_1, \pi) \prec_{\pi} \mathsf{getO}\lambda(e_2, \pi) = \mathsf{getI}\lambda(e_2, \pi)$ .

- If  $e_1 \in \text{Inst}$ , then  $\text{getO}\lambda(e_1, \pi) = \text{getI}\lambda(e_1, \pi)$ , so we have  $\text{getI}\lambda(e_1, \pi) \prec_{\pi} \text{getI}\lambda(e_2, \pi)$  and  $(e_1, e_2) \in \text{IB}$ .
- If  $e_1 \in \{1 \text{W}, \text{nlW}, \text{nrW}, \text{narW}\}$ , there is  $\lambda$  such that  $\mathsf{type}(\lambda) \in \{1 \text{W}, \text{nlW}, \text{nrW}, \text{narW}\}$ ,  $\mathsf{getA}(\lambda) = e_1$ , and  $\mathsf{getI}\lambda(e_1, \pi) = \lambda \prec_{\pi} \mathsf{B}\langle e_1 \rangle = \mathsf{getO}\lambda(e_1, \pi)$ . By transitivity we again have  $\mathsf{getI}\lambda(e_1, \pi) \prec_{\pi} \mathsf{getI}\lambda(e_2, \pi)$  and  $(e_1, e_2) \in \mathsf{IB}$ .

With a similar reasoning, we can see that [Inst];  $IB \subseteq OB$ .

**Theorem 10.** getEG( $\pi$ ) is consistent.

*Proof.* From Definition 20, we need to check that both ib and ob are irreflexive. Since IB and OB are irreflexive, it is enough to show that  $ib \subseteq IB$  and  $ob \subseteq OB$ . The explicit definition using limits is the following (where  $rf_e \triangleq (rf \setminus rf_i)$  includes  $(rf \cap sqp)$  since we assume the PCIe guarantees hold):

```
\begin{split} & \mathsf{ib}^0 \triangleq (\mathsf{ippo} \ \cup \ \mathsf{rf} \ \cup \ \mathsf{pf} \ \cup \ \mathsf{rb_i} \ \cup \ \mathsf{nfo})^+ \\ & \mathsf{ob}^0 \triangleq (\mathsf{oppo} \ \cup \ \mathsf{rf_e} \ \cup \ [\mathtt{nlW}]; \mathsf{pf} \ \cup \ \mathsf{rb} \ \cup \ \mathsf{nfo} \ \cup \ \mathsf{mo} \ \cup \ \mathsf{rao} \ \cup \ \mathsf{ar}; \mathsf{rao})^+ \\ & \mathsf{ib}^{n+1} \triangleq (\mathsf{ib}^n \cup \mathsf{ob}^n; [\mathtt{Inst}])^+ \\ & \mathsf{ob}^{n+1} \triangleq (\mathsf{ob}^n \cup [\mathtt{Inst}]; \mathsf{ib}^n)^+ \\ & \mathsf{ib} \triangleq \lim_{n \to \infty} \mathsf{ib}^n \\ & \mathsf{ob} \triangleq \lim_{n \to \infty} \mathsf{ob}^n \end{split}
```

It is then enough to show that  $\mathsf{ib}^0 \subseteq \mathsf{IB}$  and  $\mathsf{ob}^0 \subseteq \mathsf{OB}$ . Using Lemma 3 above, we can check the induction case:

```
\mathsf{ib}^{n+1} = (\mathsf{ib}^n \cup \mathsf{ob}^n; [\mathtt{Inst}])^+ \subseteq (\mathsf{ib}^n \cup \mathsf{OB}; [\mathtt{Inst}])^+ \subseteq (\mathsf{IB} \cup \mathsf{IB})^+ = \mathsf{IB}
\mathsf{ob}^{n+1} = (\mathsf{ob}^n \cup [\mathtt{Inst}]; \mathsf{ib}^n)^+ \subseteq (\mathsf{ob}^n \cup [\mathtt{Inst}]; \mathsf{IB})^+ \subseteq (\mathsf{OB} \cup \mathsf{OB})^+ = \mathsf{OB}
```

Since IB and OB are transitive, we need to check the components of  $ib^0$  and  $ob^0$ . There are twelve cases to verify.

- Checking ippo  $\subseteq$  IB. Let  $E^{\text{cpu}} = \{1\text{R}, 1\text{W}, \text{CAS}, \text{F}, \text{P}\}$  and  $E^{\text{nic}} = \{\text{n1R}, \text{nrR}, \text{narR}, \text{naF}, \text{n1W}, \text{nrW}, \text{narW}, \text{nF}\}$ .  $[E^{\text{cpu}}]; \text{po} \subseteq \text{IB}$  by definition of po and IB:  $E^{\text{cpu}}$  are the events for which the same annotated label is used to define po and IB, i.e.,  $\forall e \in E^{\text{cpu}}, \text{getl} \lambda(e, \pi)$  generates e in  $\pi$ . To check that  $[E^{\text{nic}}]; \text{ippo}; [E^{\text{nic}}] \subseteq \text{IB}$ , there are 36 cases to consider. They are all trivially satisfied by  $\text{nicActOrder}(\pi)$  and  $\text{backComp}(\pi)$ .
- Checking oppo ⊆ OB.
   From above we have [Inst]; oppo ⊆ [Inst]; ippo ⊆ [Inst]; IB ⊆ OB.
   [1W]; po; [Event \ (1R ∪ P)] ⊆ OB by using bufFlushOrd(π).
   For the remaining cases:
- (G7) [nrW];  $(po \cap sqp)$ ;  $[nrW] \subseteq OB$  comes from  $nicActOrder(\pi)$  (i.e.,  $nrW\langle ... \rangle \prec_{\pi} nrW\langle ... \rangle)$  and  $bufFlushOrd(\pi)$  (i.e.,  $B\langle ... \rangle \prec_{\pi} B\langle ... \rangle$ ).
- (G8) [nrW];  $(po \cap sqp)$ ;  $[narR] \subseteq OB$  comes from  $nicActOrder(\pi)$  (i.e.,  $nrW\langle ... \rangle \prec_{\pi} narR\langle ... \rangle$ ) and  $bufFlushOrd(\pi)$  (i.e.,  $B\langle ... \rangle \prec_{\pi} narR\langle ... \rangle$ ).
- (G9) If  $e_1 \in \operatorname{nrW}$ ,  $e_3 \in \operatorname{narW}$ , and  $(e_1, e_3) \in (\operatorname{po} \cap \operatorname{sqp})$ , then from Def. 19 there is  $e_2 \in \operatorname{narR}$  such that  $(e_2, e_3) \in \operatorname{po}|_{\operatorname{imm}}$  and thus  $(e_1, e_2) \in (\operatorname{po} \cap \operatorname{sqp})$ . From case G8 above, we have  $(e_1, e_2) \in \operatorname{OB}$ . From backComp $(\pi)$ , we have  $(e_2, e_3) \in [\operatorname{Inst}]$ ;  $\operatorname{IB} \subseteq \operatorname{OB}$ . Thus  $[\operatorname{nrW}]$ ;  $(\operatorname{po} \cap \operatorname{sqp})$ ;  $[\operatorname{narW}] \subseteq \operatorname{OB}$ .
- (G10) [nrW];  $(po \cap sqp)$ ;  $[nrR] \subseteq OB$  comes from  $nicActOrder(\pi)$  (i.e.,  $nrW\langle ... \rangle \prec_{\pi} nrR\langle ... \rangle$ ) and  $bufFlushOrd(\pi)$  (i.e.,  $nrW\langle ... \rangle \prec_{\pi} B\langle ... \rangle \prec_{\pi} nrR\langle ... \rangle$ ).
- (G11) If  $e_1 \in \mathtt{nrW}$ ,  $e_3 \in \mathtt{nlW}$ , and  $(e_1, e_3) \in (\mathtt{po} \cap \mathtt{sqp})$ , then from Def. 19 there is  $e_2 \in (\mathtt{narR} \cup \mathtt{nrR})$  such that  $(e_2, e_3) \in \mathtt{po}|_{\mathrm{imm}}^{\{1,2\}}$  and thus  $(e_1, e_2) \in (\mathtt{po} \cap \mathtt{sqp})$ . Then  $(e_1, e_2) \subseteq \mathsf{OB}$  comes from cases G9 and G10 respectively. From  $\mathtt{backComp}(\pi)$ , we have  $(e_2, e_3) \in [\mathtt{Inst}]$ ;  $\mathtt{IB} \subseteq \mathsf{OB}$ . Thus  $[\mathtt{nrW}]$ ;  $(\mathtt{po} \cap \mathtt{sqp})$ ;  $[\mathtt{nlW}] \subseteq \mathsf{OB}$ .
  - (I7)  $[\mathtt{narW}]; (\mathtt{po} \cap \mathtt{sqp}); [\mathtt{nrW}] \subseteq \mathsf{OB} \text{ comes from nicActOrder}(\pi) \text{ (i.e., narW}\langle \ldots \rangle \prec_{\pi} \mathtt{nrW}\langle \ldots \rangle) \text{ and bufFlushOrd}(\pi) \text{ (i.e., } \mathsf{B}\langle \ldots \rangle \prec_{\pi} \mathsf{B}\langle \ldots \rangle).$
  - (I8) [narW];  $(po \cap sqp)$ ;  $[narR] \subseteq OB$  follows from Def. 19,  $nicActOrder(\pi)$  and  $bufFlushOrd(\pi)$  by similar reasoning to I7.
  - (I9) [narW];  $(po \cap sqp)$ ;  $[narW] \subset OB$  follows similarly to I7.
- (I10) [narW];  $(po \cap sqp)$ ;  $[nrR] \subseteq OB$  follows similarly to I7.
- (K11) [nlW];  $(po \cap sqp)$ ;  $[nlW] \subseteq OB$  comes from  $nicActOrder(\pi)$  (i.e.,  $nlW\langle ... \rangle \prec_{\pi} nlW\langle ... \rangle$ ) and  $bufFlushOrd(\pi)$  (i.e.,  $B\langle ... \rangle \prec_{\pi} B\langle ... \rangle$ ).
  - Checking  $\mathsf{rf}_{\mathsf{e}} \subseteq \mathsf{OB}$ . If  $(w,r) \in \mathsf{rf}_{\mathsf{e}}$ , there is  $\pi_1$  and  $\pi_2$  such that  $\pi = \pi_2 \cdot \mathsf{getO}\lambda(r,\pi) \cdot \pi_1$ , and we use  $\mathsf{wfrd}(\pi)$ .

- If  $r \in 1\mathbb{R}$ , we have  $\mathsf{wfrdCPU}(r, w, \pi_1)$ . The definition allow for three different cases. In the first case,  $\lambda \in \{\mathsf{B}\langle w \rangle, \mathsf{CAS}\langle w, \_ \rangle\}$  is in  $\pi_1$ ; we have  $\lambda = \mathsf{getO}\lambda(w,\pi) \prec_{\pi} \mathsf{getO}\lambda(r,\pi)$  and so  $(w,r) \in \mathsf{OB}$ . In the second case, we have  $\lambda = 1\mathbb{W}\langle w \rangle$  and t(w) = t(r); so  $(w,r) \in [1\mathbb{W}]$ ;  $(\mathsf{rf} \cap \mathsf{sthd})$ ;  $[1\mathbb{R}] = \mathsf{rf}_i$ , which contradicts  $(w,r) \in \mathsf{rf}_{\mathsf{e}} = \mathsf{rf} \backslash \mathsf{rf}_i$ . In the third case,  $w = init_x$  for some location x, so  $(w,r) \in \mathsf{Event}_0 \times (\mathsf{Event} \backslash \mathsf{Event}_0) \subseteq \mathsf{OB}$ .
- If  $r \in CAS$ , similarly to above, except the second case of wfrdCPU $(r, w, \pi_1)$  is not possible because of bufFlushOrd $(\pi)$ : B $\langle w \rangle \notin \pi_1$  while CAS acts as a memory fence.
- If  $r \in nlR$ , we have  $wfrdNIC(r, w, \pi_1)$ , with two possibilities. In the first case,  $\lambda \in \{B\langle w \rangle, CAS\langle w, \_ \rangle\}$  is in  $\pi_1$ ; we have  $\lambda = getO\lambda(w, \pi) \prec_{\pi} getO\lambda(r, \pi)$  and so  $(w, r) \in OB$ . In the second case,  $w = init_x$  for some location x, so  $(w, r) \in Event_0 \times (Event \setminus Event_0) \subset OB$ .
- If  $r \in nrR$  or narR, similarly to above.
- Checking rf  $\subseteq$  IB. From above we have rf<sub>e</sub> = rf<sub>e</sub>; [Inst]  $\subseteq$  OB; [Inst]  $\subseteq$  IB. If  $(w,r) \in$  rf<sub>i</sub>  $\subseteq$  [1W]; rf; [1R], then there is  $1R\langle r,w\rangle \in \pi$ . There is  $\pi_1$  and  $\pi_2$  such that  $\pi = \pi_2 \cdot 1R\langle r,w\rangle \cdot \pi_1$ . So by  $wfrd(\pi)$  we have  $wfrdCPU(r,w,\pi_1)$  which implies  $1W\langle w\rangle \prec_{\pi} 1R\langle r,w\rangle$  and  $(w,r) \in IB$ .
- Checking [nlW]; pf  $\subseteq$  OB. If  $(w,p) \in$  pf with  $w \in$  nlW, then there exists e such that nlW $\langle w,e \rangle \prec_{\pi} P\langle p,e \rangle$ . From backComp $(\pi)$ , we have nlW $\langle w,e \rangle \prec_{\pi} B\langle w \rangle \prec_{\pi} P\langle p,e \rangle$  and so  $(w,p) \in$  OB.
- Checking  $\mathsf{pf} \subseteq \mathsf{IB}$ . If  $(w,p) \in \mathsf{pf}$ , then there exists e such that either  $\mathsf{nlW}\langle w,e \rangle \prec_\pi \mathsf{P}\langle p,e \rangle$  or  $\mathsf{nrW}\langle w,e \rangle \prec_\pi \mathsf{P}\langle p,e \rangle$ . In both cases we immediately have  $(w,p) \in \mathsf{IB}$ .
- Checking  $\mathsf{rb_i} \subseteq \mathsf{IB}$ . If  $(r, w') \in \mathsf{rb_i}$  then  $r \in \mathsf{IR}$ ,  $w' \in \mathsf{IW}$ , t(r) = t(w'), and there exists w such that  $(w, r) \in \mathsf{rf}$  and  $(w, w') \in \mathsf{mo}$ . There is  $\pi_4$  and  $\pi_3$  such that  $\pi = \pi_4 \cdot \mathsf{IR} \langle r, w \rangle \cdot \pi_3$ . So by  $\mathsf{wfrd}(\pi)$  we have  $\mathsf{wfrdCPU}(r, w, \pi_3)$ , and there is three cases to consider.
  - In the first case,  $\pi_3 = \pi_2 \cdot \lambda_w \cdot \pi_1$ , with  $\lambda_w \in \{ \mathsf{B}\langle w \rangle, \mathsf{CAS}\langle w, \_ \rangle \}$ , and  $\mathsf{B}\langle w' \rangle \notin \pi_2$ . Since  $(w,w') \in \mathsf{mo}$  we have  $\mathsf{B}\langle w' \rangle \notin \pi_1$ , an so  $\mathsf{B}\langle w' \rangle \notin \pi_3$ . The last condition of the first case then gives us  $\mathsf{IW}\langle w' \rangle \notin \pi_3$ , which implies  $(r,w') \in \mathsf{IB}$ .
  - In the second case,  $\pi_3 = \pi_2 \cdot \lambda_w \cdot \pi_1$ , with  $\lambda_w = 1 \mathbb{W}\langle w \rangle$ , thread $(w) = \mathsf{thread}(r)$ , and  $\mathbb{B}\langle w \rangle \notin \pi_3$ . Then w and w' are on the same thread, and by  $\mathsf{bufFlushOrd}(\pi)$  and  $(w,w') \in \mathsf{mo}$  we have  $1 \mathbb{W}\langle w \rangle \prec_{\pi} 1 \mathbb{W}\langle w' \rangle$  and  $1 \mathbb{W}\langle w' \rangle \notin \pi_1$ . The last condition of the second case gives us  $1 \mathbb{W}\langle w' \rangle \notin \pi_2$ , so  $1 \mathbb{W}\langle w' \rangle \notin \pi_3$  and  $(r,w') \in \mathsf{IB}$ .
  - In the last case,  $w = init_x$  for some location x, and we immediately get  $\mathbb{IW}\langle w' \rangle \notin \pi_3$ , which implies  $(r, w') \in \mathbb{IB}$ .
- Checking  $\operatorname{rb} \subseteq \operatorname{OB}$ . If  $(r,w') \in \operatorname{rb}$ , then there exists w such that  $(w,r) \in \operatorname{rf}$  and  $(w,w') \in \operatorname{mo}$ . By definition of rf, there is  $\pi_4$  and  $\pi_3$  such that  $\pi = \pi_4 \cdot \lambda_r \cdot \pi_3$ , with  $\lambda_r \in \{\operatorname{1R}\langle r,w\rangle, \operatorname{CAS}\langle r,w\rangle, \operatorname{n1R}\langle r,w, \ldots, \ldots\rangle, \operatorname{narR}\langle r,w, \ldots, \ldots\rangle\}$ . So by  $\operatorname{wfrd}(\pi)$  we have either  $\operatorname{wfrdNIC}(r,w,\pi_3)$  or  $\operatorname{wfrdCPU}(r,w,\pi_3)$ , and there are five cases to consider.

- In the first case of wfrdNIC $(r, w, \pi_3)$ ,  $\pi_3 = \pi_2 \cdot \text{getO}\lambda(w, \pi) \cdot \pi_1$ , and  $\text{getO}\lambda(w', \pi) \notin \pi_2$ . Since  $(w, w') \in \text{mo}$  we have  $\text{getO}\lambda(w', \pi) \notin \pi_1$ , and thus  $\text{getO}\lambda(w', \pi) \notin \pi_3$ . So  $\text{getO}\lambda(w', \pi) \in \pi_4$  and  $(r, w') \in \text{OB}$ .
- In the last case  $\mathsf{wfrdNIC}(r, w, \pi_3)$ ,  $w = init_x$  for some location x, and we immediately have  $\mathsf{getO}\lambda(w', \pi) \notin \pi_3$ , which implies  $(r, w') \in \mathsf{OB}$ .
- For the first case of  $\mathsf{wfrdCPU}(r, w, \pi_3)$ , same reasoning as for the first case of  $\mathsf{wfrdNIC}$ .
- For the second case of  $\mathsf{wfrdCPU}(r, w, \pi_3), \ \pi_3 = \pi_2 \cdot \mathsf{getl}\lambda(w, \pi) \cdot \pi_1$ , with  $\mathsf{thread}(w) = \mathsf{thread}(r)$ , and  $\mathsf{getO}\lambda(w, \pi) \notin \pi_3$ . So  $\mathsf{getO}\lambda(w, \pi) \in \pi_4$ , and since  $(w, w') \in \mathsf{mo}$  we have  $\mathsf{getO}\lambda(w', \pi) \in \pi_4$  as well, and  $(r, w') \in \mathsf{OB}$ .
- For the last case of  $\mathsf{wfrdCPU}(r, w, \pi_3)$ , same reasoning as for the last case of  $\mathsf{wfrdNIC}$ .
- Checking nfo ⊆ IB.
  By definition of nfo.
  Checking nfo ⊆ OB.
- By definition of nfo.

   Checking mo ⊆ OB.
- By definition of mo, as what matters are the  $init_x$ ,  $B\langle w \rangle$ , and  $CAS\langle w, \_ \rangle$  events.
- Checking rao ⊆ OB.
   By definition of rao.
- Checking ar; rao  $\subseteq$  OB. If  $(w, r_1) \in \operatorname{ar}$  then  $\operatorname{narR}\langle r_1, a_1, ..., w \rangle \in \pi$  for some  $a_1 \in \operatorname{rRMW}$ , and if  $(r_1, r_2) \in \operatorname{rao}$  then  $\operatorname{narR}\langle r_1, ..., a_1, ..., w \rangle \prec_{\pi} \lambda_r$  for some  $\lambda_r \in \{\operatorname{naF}\langle r_2, ..., a_2, ... \rangle, \operatorname{narR}\langle r_2, ..., a_2, ... \rangle\}$ , with  $\overline{n}(a_1) = \overline{n}(a_2)$ . Then using nicAtomicity( $\pi$ ) we have that  $\operatorname{B}\langle w \rangle \prec_{\pi} \lambda_r$ .

### D.6 From Declarative Semantics to Annotated Semantics

From a program P and a well-formed consistent execution graph  $G = (\mathsf{Event}, \mathsf{po}, \mathsf{rf}, \mathsf{pf}, \mathsf{mo}, \mathsf{nfo}, \mathsf{rao}),$  where  $(\mathsf{Event}, \mathsf{po})$  is generated by P, we want to reconstruct an annotated semantics execution.

**Theorem 11.** ib and ob can be extended into total relations IB and OB on Event such that:

- IB and OB are irreflexive and transitive;
- OB;  $[Inst] \subseteq IB$  and [Inst];  $IB \subseteq OB$ .

*Proof.* We show that if ib is not already total we can extend it (and maybe ob) into a strictly bigger relation satisfying the constraints of the theorem. Let us assume that there is  $(a,b) \in \mathsf{Event}^2$  such that  $(a,b) \notin \mathsf{ib}$  and  $(b,a) \notin \mathsf{ib}$ . We arbitrarily decide to include (a,b) in our relation and we define  $\mathsf{ib}' = (\mathsf{ib} \cup \{(a,b)\})^+$  and  $\mathsf{ob}' = (\mathsf{ob} \cup [\mathsf{Inst}]; \mathsf{ib}')^+$ .

Clearly  $\mathsf{ib'}$  and  $\mathsf{ob'}$  are transitive,  $\mathsf{ib'}$  is irreflexive, and  $[\mathtt{Inst}]$ ;  $\mathsf{ib'} \subseteq \mathsf{ob'}$ . We need to prove the following two facts:  $\mathsf{ob'}$  is still irreflexive; and  $\mathsf{ob'}$ ;  $[\mathtt{Inst}] \subseteq \mathsf{ib'}$ .

First, let us check that  $(ob\cup[Inst];ib')^+$  is irreflexive. Since ob and ([Inst];ib') are both transitive and irreflexive, a cycle would only be possible by alternating between the two components, so it is enough to show that  $(ob;([Inst];ib'))^+$  is irreflexive. But  $(ob;([Inst];ib'))^+ = ((ob;[Inst]);ib')^+ \subseteq (ib;ib')^+ \subseteq ib'$  is irreflexive. Thus ob' is irreflexive.

Then, we need to check that ob'; [Inst]  $\subseteq$  ib'. Using some rewriting, ob' =  $(ob \cup [Inst]; ib')^+ = ob \cup (ob^*; ([Inst]; ib'))^+; ob^*$ . We know ob; [Inst]  $\subseteq$  ib', which also implies ob\*; [Inst]  $\subseteq$  ib'\*. So ob'; [Inst] = ob; [Inst]  $\cup$  ((ob\*; [Inst]); ib')+; (ob\*; [Inst])  $\subseteq$  ib'  $\cup$  (ib'\*; ib')+; ib'\*  $\subseteq$  ib'.

Once  ${\sf ib}$  is a total relation on  ${\sf Event}$ , we can similarly freely extend  ${\sf ob}$  into a total relation.

We use Theorem 11 above to extend ib and ob into total relations IB and OB. Since (Event, po) is derived from P, by Appendix D.2 we have that for all  $t \in \mathsf{Tid}$  there are  $s_t$  and  $G_t$  such that  $G_t \in G^t(s_t)$ ,  $\mathsf{P}(t) \mapsto s_t$  and (Event, po) =  $G_{init}$ ; ( $\|_{t \in \mathsf{Tid}} G_t$ ). We consider each premise of the form  $C \mapsto s$ , where C is a primitive command, to generate new events and annotated labels.

- If  $s = r \in 1\mathbb{R}$ , from well-formedness conditions, there is w such that  $(w, r) \in \mathsf{rf}$  and  $\mathsf{eq}_{\mathsf{loc}\&v}(r, w)$ . We create an annotated label  $\mathsf{lR}\langle r, w \rangle$ .
- If s = u, s' where  $u \in CAS$ , from well-formedness conditions, there is w such that  $(w, u) \in \mathsf{rf}$  and  $\mathsf{eq}_{\mathsf{loc}\&v}(u, w)$ . We create an annotated label  $\mathsf{CAS}\langle u, w \rangle$ , then process s'.
- If s = f, r, s' where  $f \in F$ ,  $r \in 1R$ , and  $w \in 1W$ , from well-formedness conditions, there is w' such that  $(w', r) \in \mathsf{rf}$  and  $\mathsf{eq}_{\mathsf{loc}\&v}(r, w')$ . We create annotated labels  $\mathsf{F}\langle f \rangle$ ,  $\mathsf{1R}\langle r, w' \rangle$ ,  $\mathsf{1W}\langle w \rangle$  and  $\mathsf{B}\langle w \rangle$ , then process s'.
- If  $s = w \in 1W$ , we create annotated labels  $1W\langle w \rangle$  and  $B\langle w \rangle$ .
- If  $s = f \in F$ , we create annotated labels F(f).

- If s = r, w where  $r \in nlR$  and  $w \in nrW$ , we create two events  $a \in Put$  and  $e \in nrEX$ , and the annotated labels  $Push\langle a \rangle$ ,  $NlC\langle a \rangle$ ,  $nlR\langle r, w', a, w \rangle$  (where  $(w',r) \in rf$ ),  $nrW\langle w,e \rangle$ ,  $B\langle w \rangle$ , and  $CN\langle e \rangle$ . If there is p such that  $(w,p) \in pf$ , we also create an annotated label  $P\langle p,e \rangle$ . To simplify later definition, we also extend po such that the event a is placed just before r, and e just after w. I.e., let  $po' = po \cup \{(e',a) \mid (e',r) \in po\} \cup \{(a,e') \mid (r,e') \in po^*\}$  and redefine  $po = po' \cup \{(e',e) \mid (e',w) \in po'^*\} \cup \{(e,e') \mid (w,e') \in po'\}$ . Note: from well-formedness conditions, every nlR and every nrW are part of such a pair.
- If s = r, w where  $r \in \text{nrR}$  and  $w \in \text{nlW}$ , we similarly create  $a \in \text{Get}$ ,  $e \in \text{nlEX}$ ,  $\text{Push}\langle a \rangle$ ,  $\text{NIC}\langle a \rangle$ ,  $\text{nrR}\langle \dots \rangle$ ,  $\text{nlW}\langle \dots \rangle$ ,  $\text{B}\langle \dots \rangle$ , and potentially  $\text{P}\langle \dots \rangle$ .
- If s = r, w where  $r \in \text{narR}$  and  $w \in \text{nlW}$ , we have C of the form  $z := \text{RCAS}(\overline{x}, e, e')$ , so we use the values  $\llbracket e \rrbracket$  and  $\llbracket e' \rrbracket$  to create  $a \in \text{RCAS}$ ,  $\text{Push}\langle a \rangle$ ,  $\text{NIC}\langle a \rangle$ ,  $\text{naF}\langle \ldots \rangle$ ,  $\text{nlW}\langle \ldots \rangle$ ,  $\text{B}\langle \ldots \rangle$ , and potentially  $\text{P}\langle \ldots \rangle$ .
- If  $s = r, w_1, w_2$  where  $r \in \text{narR}$ ,  $w_1 \in \text{narW}$ ,  $w_2 \in \text{nlW}$ , we have C either of the form  $z := \text{RFAA}(\overline{x}, e)$  or  $z := \text{RCAS}(\overline{x}, e_1, e_2)$ , so we create  $a \in \text{RFAA}$  or  $a \in \text{RCAS}$  accordingly, and  $\text{Push}\langle a \rangle$ ,  $\text{NIC}\langle a \rangle$ ,  $\text{narR}\langle \ldots \rangle$ ,  $\text{narW}\langle w_1 \rangle$ ,  $\text{nlW}\langle w_2, \ldots \rangle$ ,  $\text{B}\langle w_1 \rangle$ ,  $\text{B}\langle w_2 \rangle$  and potentially  $\text{P}\langle \ldots \rangle$ .
- If  $s = f \in nF$ , we create the annotated labels Push(f), NIC(f), and nF(f).
- We ignore  $s = p \in P$ , as this is already handled by our earlier cases.

Then, we use IB and OB to reconstruct a partial path from these annotated labels. We define a path  $\pi_0$  such that:

- $\pi_0 \in (\mathsf{ALabel} \setminus (\mathsf{Push} \cup \mathsf{NIC} \cup \mathsf{CN}))^*$
- $\operatorname{getI}\lambda(e_1,\pi_0) \prec_{\pi_0} \operatorname{getI}\lambda(e_2,\pi_0) \iff (e_1,e_2) \in \operatorname{IB}$
- $\bullet \ \operatorname{getO}{\lambda}(e_1,\pi_0) \prec_{\pi_0} \operatorname{getO}{\lambda}(e_2,\pi_0) \iff (e_1,e_2) \in \operatorname{OB}$
- $\forall w \in \{1\text{W}, \text{nlW}, \text{nrW}, \text{narW}\}, \text{getI}\lambda(w, \pi_0) \prec_{\pi_0} \text{getO}\lambda(w, \pi_0)$

This is possible from the properties of IB and OB. For pairs of annotated labels not ordered by IB or OB, we decide to order  $1 \mathbb{W}\langle w \rangle / n 1 \mathbb{W}\langle w, \_ \rangle / n r \mathbb{W}\langle w, \_ \rangle / n r \mathbb{W}\langle w \rangle$  first and  $\mathbb{B}\langle w \rangle$  last. Note that the annotated labels  $Push\langle ... \rangle$ ,  $NIC\langle ... \rangle$ , and  $CN\langle ... \rangle$  not covered by IB/OB are not yet integrated in  $\pi_0$ .

Then we extend  $\pi_0$  to add annotated labels not considered by the declarative semantics. We use the following extension function that introduces a new annotated label as early as possible after a set of dependencies.

$$\mathsf{extend}(\pi,\lambda,S) \triangleq \begin{cases} \pi_2 \cdot \lambda \cdot \lambda' \cdot \pi_1 & \text{if } \pi = \pi_2 \cdot \lambda' \cdot \pi_1 \wedge \lambda' \in S \wedge \pi_2 \cap S = \emptyset \\ \pi \cdot \lambda & \text{if } \pi \cap S = \emptyset \end{cases}$$

We define a new function to recover the first annotated label corresponding to an event:

$$E^{\text{ext}} \triangleq \text{Event} \cup (\text{Get} \cup \text{Put} \cup \text{RCAS} \cup \text{RFAA} \cup \text{nlEX} \cup \text{nrEX})$$

$$getCPU : E^{ext} \rightarrow ALabel$$

$$\mathtt{getCPU}(e) \triangleq \begin{cases} \mathtt{getI}\lambda(e,\pi_0) & \text{if } e \in E^{\mathrm{cpu}} = \{\mathtt{lR},\mathtt{lW},\mathtt{CAS},\mathtt{F},\mathtt{P}\} \\ \mathtt{Push}\langle e \rangle & \text{if } e \in \{\mathtt{Put},\mathtt{Get},\mathtt{RCAS},\mathtt{RFAA},\mathtt{nF}\} \\ \mathtt{undefined} & \mathrm{otherwise} \end{cases}$$

And a similar function for events emptying a CPU buffer:

$$\mathsf{getTSO}: E^{\mathsf{ext}} \rightharpoonup \mathsf{ALabel}$$
 
$$\mathsf{getTSO}(e) \triangleq \begin{cases} \mathsf{B}\langle e \rangle & \text{if } e \in \mathsf{IW} \\ \mathsf{NIC}\langle e \rangle & \text{if } e \in \{\mathsf{Put}, \mathsf{Get}, \mathsf{RCAS}, \mathsf{RFAA}, \mathsf{nF}\} \\ \mathsf{undefined} & \mathsf{otherwise} \end{cases}$$

Let us consider  $(a_1, \ldots, a_n) = \mathsf{Event} \cap \{\mathsf{Put}, \mathsf{Get}, \mathsf{RCAS}, \mathsf{RFAA}, \mathsf{nF}\}$  in po order, i.e., if i < j then  $(a_j, a_i) \notin \mathsf{po}$ . We extend  $\pi_0$  successively until we get  $\pi_n$ :

- We introduce Push as early as possible: Let  $\pi' = \mathsf{extend}(\pi_{i-1}, \mathsf{Push}\langle a_i \rangle, \{\mathsf{getCPU}(e) \mid (e, a_i) \in \mathsf{po}\})$
- We introduce NIC as early as possible: Let  $\pi'' = \mathsf{extend}(\pi', \mathsf{NIC}(a_i), \{\mathsf{Push}(a_i)\} \cup \{\mathsf{getTSO}(e) \mid (e, a_i) \in \mathsf{po}\})$
- If  $a_i \in \operatorname{Put}$ , there is  $e_i \in \operatorname{nrEX}$  such that  $\operatorname{nlR}\langle -, -, a_i, w \rangle \prec_{\pi_0} \operatorname{nrW}\langle w, e_i \rangle$ . We also introduce  $\operatorname{CN}$ : Let  $S = \{\operatorname{nrW}\langle w, e_i \rangle\} \cup \{\operatorname{nlW}\langle -, e \rangle \mid (e, e_i) \in \operatorname{po} \cap \operatorname{sqp}\} \cup \{\operatorname{CN}\langle e \rangle \mid (e, e_i) \in \operatorname{po} \cap \operatorname{sqp}\}$ , we pose  $\pi_i = \operatorname{extend}(\pi'', \operatorname{CN}\langle e_i \rangle, S)$ . Otherwise, i.e.  $a_i \notin \operatorname{Put}$ , we simply have  $\pi_i = \pi''$

Finally,  $\pi = \pi_n$  is our path for an annotated semantics reduction. We clearly have  $\mathsf{complete}(\pi)$  by definition. Our goal is then to prove that  $\mathsf{wfp}(\pi)$  holds. It is composed of seven properties. Note that we already have the existence of the relevant annotated labels, and we need to show that the ordering constraints are respected.

#### nodup

 $\mathsf{nodup}(\pi)$  directly comes from the definition of annotated labels. There is no conflict in event usage.

## backComp

Here are a couple lemmas showing that the new annotated labels are not placed too late and do not disturb the expected ordering.

**Lemma 4.** For all  $a \in \{\text{Put}, \text{Get}, \text{RCAS}, \text{RFAA}, \text{nF}\}\ and\ b \in \text{Event},\ if\ (a,b) \in \text{po}^*,\ then\ \text{Push}\langle a \rangle \prec_{\pi} \text{getl}\lambda(b,\pi_0).$ 

Proof. We take an arbitrary b, and proceed for a in po order, i.e., we can assume it holds for  $e \in \{\text{Put}, \text{Get}, \text{RCAS}, \text{RFAA}, \text{nF}\}$  such that  $(e, a) \in \text{po}$ . By definition,  $\text{Push}\langle a \rangle$  comes from an extension  $\pi'' = \text{extend}(\pi', \text{Push}\langle a \rangle, \{\text{getCPU}(e) \mid (e, a) \in \text{po}\})$  and has been placed either first—and the result is trivial—or just after some getCPU(e) with  $(e, a) \in \text{po}$ . If  $e \in \{\text{Put}, \text{Get}, \text{RCAS}, \text{RFAA}, \text{nF}\}$ , we have  $\text{Push}\langle e \rangle \prec_{\pi''}$   $\text{Push}\langle a \rangle \prec_{\pi''}$   $\text{getl}\lambda(b, \pi_0)$  by induction hypothesis. If  $e \in E^{\text{cpu}} = \{\text{1R}, \text{1W}, \text{CAS}, \text{F}, \text{P}\}$ , we have  $\text{getl}\lambda(e, \pi_0) \prec_{\pi''} \text{Push}\langle a \rangle \prec_{\pi''}$   $\text{getl}\lambda(b, \pi_0)$  since  $(e, b) \in \text{ippo} \subseteq \text{IB}$ .

**Lemma 5.**  $\forall a \in \{\text{Put}, \text{Get}, \text{RCAS}, \text{RFAA}, \text{nF}\}, \forall b \in \{\text{nF}, \text{nrR}, \text{nlR}, \text{narR}, \text{1W}\}, if (a, b) \in \text{po*}, then \text{NIC}\langle a \rangle \prec_{\pi} \text{getO}\lambda(b, \pi_0).$ 

*Proof.* We take an arbitrary  $b \in \{nF, nrR, nlR, narR\}$ , and proceed for a in poorder, i.e., we can assume it holds for  $e \in \{Put, Get, RCAS, RFAA, nF\}$  such that  $(e, a) \in po$ . By definition,  $NIC\langle a \rangle$  comes from an extension  $\pi'' = extend(\pi', NIC\langle a \rangle, S)$ , with  $S = \{Push\langle a \rangle\} \cup \{getTSO(e) \mid (e, a) \in po\}$ , and has been placed just after some  $\lambda \in S$ .

- If  $\lambda = \text{Push}\langle a \rangle$ , then we have  $\lambda \prec_{\pi''} \text{NIC}\langle a \rangle \prec_{\pi''} \text{getO}\lambda(b, \pi_0)$  using Lemma 4 above, since  $\text{getI}\lambda(b, \pi_0) = \text{getO}\lambda(b, \pi_0)$  or  $\text{getI}\lambda(b, \pi_0) \prec_{\pi''} \text{getO}\lambda(b, \pi_0)$ .
- If  $\lambda = \text{getTSO}\langle e \rangle = \text{NIC}\langle e \rangle$  for some  $e \in \{\text{Put}, \text{Get}, \text{RCAS}, \text{RFAA}, \text{nF}\}$ , then we have  $\lambda \prec_{\pi''} \text{NIC}\langle a \rangle \prec_{\pi''} \text{getO}\lambda(b, \pi_0)$  by induction hypothesis.
- If  $\lambda = \mathsf{getTSO}\langle e \rangle = \mathsf{B}\langle e \rangle$  for some  $e \in \mathsf{IW}$ , then we have  $\mathsf{B}\langle e \rangle \prec_{\pi''} \mathsf{NIC}\langle a \rangle \prec_{\pi''} \mathsf{getO}\lambda(b,\pi_0)$  since  $(e,b) \in \mathsf{oppo} \subseteq \mathsf{OB}$ .

**Lemma 6.** For all w, e, p, if  $nrW(w, e) \in \pi$  and  $P(p, e) \in \pi$ , then  $CN(e) \prec_{\pi} P(p, e)$ .

*Proof.* Once again, we proceed for e in po order, i.e., we can assume the result holds for  $e' \in \mathtt{nrEX}$  such that  $(e',e) \in \mathtt{po}$ .  $\mathsf{CN}\langle e \rangle$  is inserted in some operation  $\pi'' = \mathsf{extend}(\pi', \mathsf{CN}\langle e \rangle, S)$ , with  $S = \{\mathtt{nrW}\langle w, e \rangle\} \cup \{\mathtt{nlW}\langle -, e' \rangle \mid (e', e) \in \mathtt{po} \cap \mathtt{sqp}\} \cup \{\mathsf{CN}\langle e' \rangle \mid (e', e) \in \mathtt{po} \cap \mathtt{sqp}\}$ . It is then placed just after some label  $\lambda \in S$ .

- If  $\lambda = \text{nrW}(w, e)$ , we have  $\lambda \prec_{\pi''} \text{CN}(e) \prec_{\pi''} \text{P}(p, e)$  because  $(w, p) \in \text{pf} \subseteq \text{IB}$ .
- If  $\lambda = \mathsf{CN}\langle e' \rangle$  with  $(e', e) \in \mathsf{po} \cap \mathsf{sqp}$ , then there is some w' such that  $(w', w) \in \mathsf{po} \cap \mathsf{sqp}$  and  $\mathsf{nrW}\langle w', e' \rangle \in \pi'$ . From well-formedness condition number 1 (see Definition 19), there is some p' such that  $(w', p') \in \mathsf{pf}$  and  $(p', p) \in \mathsf{po}$ . By induction hypothesis, we have  $\mathsf{CN}\langle e' \rangle \prec_{\pi'} \mathsf{P}\langle p', e' \rangle$ , and from  $(p', p) \in \mathsf{IB}$  we have  $\mathsf{P}\langle p', e' \rangle \prec_{\pi'} \mathsf{P}\langle p, e \rangle$ . In the end, we have the result  $\mathsf{CN}\langle e' \rangle \prec_{\pi''} \mathsf{CN}\langle e \rangle \prec_{\pi''} \mathsf{P}\langle p, e \rangle$ .
- If  $\lambda = \mathtt{nlW}\langle w', e' \rangle$  with  $(e', e) \in \mathtt{po} \cap \mathtt{sqp}$ , then we also have  $(w', w) \in \mathtt{po} \cap \mathtt{sqp}$ , so from well-formedness condition number 1 (see Definition 19), there is some p' such that  $(w', p') \in \mathtt{pf}$  and  $(p', p) \in \mathtt{po}$ . We have  $\mathtt{nlW}\langle w', e' \rangle \prec_{\pi''} \mathsf{CN}\langle e \rangle \prec_{\pi''} \mathsf{P}\langle p', e' \rangle \prec_{\pi''} \mathsf{P}\langle p, e \rangle$ .

We can then check that we have  $\mathsf{backComp}(\pi)$ :

- $1 \text{W} \langle w \rangle \prec_{\pi} \text{B} \langle w \rangle$  comes from the third property when defining  $\pi_0$ ; similarly for nlW, nrW and narW.
- $\operatorname{Push}\langle a \rangle \prec_{\pi} \operatorname{NIC}\langle a \rangle$  comes from the extension process.
- NIC $\langle f \rangle \prec_{\pi} nF \langle f \rangle$  comes from Lemma 5; similarly for NIC $\langle a \rangle \prec_{\pi} nlR/nrR/naF/narR\langle \ldots \rangle$ .
- $nlR(r, w, a, w') \prec_{\pi} nrW(w', e)$  comes from  $(r, w') \in ippo \subseteq IB$ ; similarly for nrR/nlW, naF/nlW, narR/nlW and narR/narW.
- $nrW(w,e) \prec_{\pi} CN(e)$  comes from the extension process
- $\mathtt{nlW}\langle w, e \rangle \prec_{\pi} \mathsf{B}\langle w \rangle \prec_{\pi} \mathsf{P}\langle p, e \rangle$  comes from  $(w, p) \in [\mathtt{nlW}]$ ;  $\mathsf{pf} \subseteq \mathsf{OB}$ .
- $\mathsf{CN}\langle e \rangle \prec_{\pi} \mathsf{P}\langle p, e \rangle$  comes from Lemma 6.

Thus we have  $\mathsf{backComp}(\pi)$ .

#### bufFlushOrd

- $1 \mathbb{W}\langle w_1 \rangle \prec_{\pi} 1 \mathbb{W}\langle w_2 \rangle \iff \mathbb{B}\langle w_1 \rangle \prec_{\pi} \mathbb{B}\langle w_2 \rangle$  when  $t(w_1) = t(w_2)$  comes the fact that  $[1 \mathbb{W}]$ ; po;  $[1 \mathbb{W}] \subseteq (\mathbb{IB} \cup \mathbb{OB})$ , so both sides are true if and only if  $(w_1, w_2) \in \mathbb{P}$  similarly for  $\mathbb{N}$  and  $\mathbb{N}$  narW on the same queue pair.
- When  $t(a_1) = t(a_2)$ ,  $\operatorname{Push}\langle a_1 \rangle \prec_{\pi} \operatorname{Push}\langle a_2 \rangle \iff \operatorname{NIC}\langle a_1 \rangle \prec_{\pi} \operatorname{NIC}\langle a_2 \rangle \iff (a_1, a_2) \in \operatorname{po}$  from the definition of the extension process (to define  $\pi_n$ ).
- For  $a \in \{\text{Put}, \text{Get}, \text{nF}, \text{RCAS}, \text{RFAA}\}, w \in \text{IW}, \text{ such that } t(a) = t(w)$ :
  - If  $(w, a) \in po$ , then  $1 \mathbb{W} \langle w \rangle \prec_{\pi} Push \langle a \rangle$  and  $B \langle w \rangle \prec_{\pi} NIC \langle a \rangle$  from the definition of the extension process.
  - If  $(a, w) \in po$ , then  $Push\langle a \rangle \prec_{\pi} 1 \mathbb{W}\langle w \rangle$  and  $NIC\langle a \rangle \prec_{\pi} B\langle w \rangle$  from Lemmas 4 and 5.
- When t(w) = t(f),  $1 \mathbb{W} \langle w \rangle \prec_{\pi} \mathbb{F} \langle f \rangle$  implies  $(w, f) \in \text{po (since [F]; po; [1W] } \subseteq \text{ippo } \subseteq \text{IB})$ , which implies  $\mathbb{B} \langle w \rangle \prec_{\pi} \mathbb{F} \langle f \rangle$  (since [1W]; po; [F]  $\subseteq \text{oppo } \subseteq \text{OB}$ ); similarly for CAS.
- If  $w \in nlW$ ,  $r \in nlR$ , and sameqp(w,r), then from the definition of pre-executions (see condition 6 of Definition 18), either  $(w,r) \in nfo$  or  $(r,w) \in nfo$ . If  $nlW\langle w, \_ \rangle \prec_{\pi} nlR\langle r, \_, \_, \_ \rangle$ , then  $(r,w) \notin nfo$  (since  $nfo \subseteq lB$ ) and  $(w,r) \in nfo$ . Thus,  $B\langle w \rangle \prec_{\pi} nlR\langle r, \_, \_, \_ \rangle$  (since  $nfo \subseteq OB$ ); similarly for  $w \in \{nrW, narW\}$  and  $r \in \{nrR, narR\}$ .

Thus we have  $bufFlushOrd(\pi)$ .

# pollOrder

Lemma 7. For all  $e_1$ ,  $e_2 \in \{nlEX, nrEX\}$ , such that  $sameqp(e_1, e_2)$ , let  $\lambda_1 \in \{nlW\langle_-, e_1\rangle, CN\langle e_1\rangle\}$ ,  $\lambda_2 \in \{nlW\langle_-, e_2\rangle, CN\langle e_2\rangle\}$ , then  $(e_1, e_2) \in po \iff \lambda_1 \prec_\pi \lambda_2$ .

*Proof.* By symmetry, we only need to show  $(e_1,e_2) \in po \implies \lambda_1 \prec_{\pi} \lambda_2$ . Once again, we proceed for  $e_1$  in po order, i.e., we can assume the result holds for  $e' \in nEX$  such that  $(e',e_1) \in po$ .

- If  $\lambda_1 = \mathtt{nlW}\langle w_1, e_1 \rangle$  and  $\lambda_2 = \mathtt{nlW}\langle w_2, e_2 \rangle$ , then  $(e_1, e_2) \in \mathtt{po}$  implies  $(w_1, w_2) \in \mathtt{(po} \cap \mathtt{sqp})$ , so  $(w_1, w_2) \in \mathtt{ippo} \subseteq \mathtt{IB}$  and  $\lambda_1 \prec_{\pi} \lambda_2$ .
- If  $\lambda_1 = \mathtt{nlW}\langle w_1, e_1 \rangle$  and  $\lambda_2 = \mathsf{CN}\langle e_2 \rangle$ , then by definition of the extension process we have  $\lambda_1 \prec_{\pi} \lambda_2$ .
- If  $\lambda_1 = \mathsf{CN}\langle e_1 \rangle$  and  $\lambda_2 = \mathsf{nlW}\langle w_2, e_2 \rangle$ , then  $\lambda_1$  is inserted in some operation  $\pi'' = \mathsf{extend}(\pi', \mathsf{CN}\langle e_1 \rangle, S)$ , with  $S = \{\mathsf{nrW}\langle -, e_1 \rangle\} \cup \{\mathsf{nlW}\langle -, e' \rangle \mid (e', e_1) \in \mathsf{po} \cap \mathsf{sqp}\} \cup \{\mathsf{CN}\langle e' \rangle \mid (e', e_1) \in \mathsf{po} \cap \mathsf{sqp}\}$ . It is then placed just after some label  $\lambda \in S$ .
  - If  $\lambda = \text{nrW}(w_1, e_1)$ , we have  $\lambda \prec_{\pi''} \lambda_1 \prec_{\pi''} \lambda_2$  because  $(w_1, w_2) \in \text{ippo} \subseteq \mathsf{IB}$ .
  - If  $\lambda = \mathsf{CN}\langle e' \rangle$  or  $\lambda = \mathsf{nlW}\langle -, e' \rangle$ , with  $(e', e_1) \in \mathsf{po} \cap \mathsf{sqp}$ , then by induction hypothesis  $\lambda \prec_{\pi''} \lambda_1 \prec_{\pi''} \lambda_2$ .
- If  $\lambda_1 = \mathsf{CN}\langle e_1 \rangle$  and  $\lambda_2 = \mathsf{CN}\langle e_2 \rangle$ , then by definition of the extension process we have  $\lambda_1 \prec_{\pi} \lambda_2$ .

Let us assume we have  $e_1, e_2, p_2, \lambda_1, \lambda_2$  such that  $\mathsf{sameqp}(e_1, e_2), \lambda_1 \in \{\mathsf{nlW}\langle -, e_1 \rangle, \mathsf{CN}\langle e_1 \rangle\}, \lambda_2 \in \{\mathsf{nlW}\langle -, e_2 \rangle, \mathsf{CN}\langle e_2 \rangle\}, \lambda_1 \prec_{\pi} \lambda_2, \text{ and } \mathsf{P}\langle p_2, e_2 \rangle \in \pi.$ 

From the creation of the events  $e_1$  and  $e_2$ , there is some  $w_1, w_2 \in \{\text{nlW}, \text{nrW}\}$  such that  $(w_i, e_i) \in \text{po}|_{\text{imm}}$ . From Lemma 7, we have  $(e_1, e_2) \in \text{po}$  and thus

 $(w_1, w_2) \in (po \cap sqp)$ . By definition, we also have  $(w_2, p_2) \in pf$ . From well-formedness condition number 1 (see Definition 19), there is some  $p_1$  such that  $(w_1, p_1) \in pf$  and  $(p_1, p_2) \in po$ . Thus we have  $P\langle p_1, e_1 \rangle \prec_{\pi} P\langle p_2, e_2 \rangle$  as required to prove  $pollOrder(\pi)$ .

#### nicActOrder

Let  $a_1$  and  $a_2$  such that  $NIC\langle a_1\rangle \prec_{\pi} NIC\langle a_2\rangle$  and  $sameqp(a_1, a_2)$ . From the definition of the extension process, we have  $(a_1, a_2) \in po$ .

- If  $a_1 \in nF$  or  $a_2 \in nF$ , then most of the required results hold by definition of ippo. The only exception is  $\mathsf{CN}\langle e \rangle \prec_{\pi} \mathsf{nF}\langle a_2 \rangle$  which holds (by induction on e in po order) because all the dependencies of  $\mathsf{CN}\langle e \rangle$  are before  $\mathsf{nF}\langle a_2 \rangle$  by ippo.
- If  $(a_1 \in \text{Get} \land a_2 \in \text{Get})$ , the result holds by ippo.
- If  $(a_1 \in \text{Get} \land a_2 \in \text{Put})$ , the result holds by Lemma 7.
- If  $(a_1 \in \text{Get} \land a_2 \in \text{RCAS} \cup \text{RFAA})$ , the results hold by ippo.
- If  $(a_1 \in \text{Put} \land a_2 \in \text{Get})$ , the first result holds by ippo, the second by Lemma 7.
- If  $(a_1 \in \text{Put} \land a_2 \in \text{Put})$ , the first two results hold by ippo, the last one by Lemma 7.
- If  $(a_1 \in Put \land a_2 \in RCAS \cup RFAA)$ , the results hold by ippo.
- If  $(a_1 \in \mathtt{RCAS} \cup \mathtt{RFAA} \land a_2 \in \mathtt{Get})$ , the results hold by ippo.
- If  $(a_1 \in \text{RCAS} \cup \text{RFAA} \land a_2 \in \text{Put})$ , the first result holds by ippo, the latter by Lemma 7.
- If  $(a_1, a_2 \in RCAS \cup RFAA)$ , the first result holds by ippo, the latter by Lemma 7.

Thus we have  $nicActOrder(\pi)$ .

## nicAtomicity

For every  $a_1, a_2 \in \text{rRMW}$  where  $\overline{n}(a_1) = \overline{n}(a_2)$ , if  $\text{narR}\langle r_1, a_1, ..., w \rangle \prec_\pi \lambda_r$  where  $\lambda_r \in \{\text{naF}\langle r_2, ..., a_2, ... \rangle, \text{narR}\langle r_2, ..., a_2, ... \rangle\}$ , then from the extension process we have  $(r_1, w) \in \text{po}|_{\text{imm}}$ , and  $w \in \text{narW}$ , so  $(w, r_1) \in \text{ar}$ . Then we need to show that  $(r_1, r_2) \in \text{rao}$ . Suppose, for contradiction, that  $(r_1, r_2) \notin \text{rao}$ . By definition of rao, for each node n, rao<sub>n</sub> is a total order on  $\{e \in \text{narR} \mid \overline{n}(e) = n\}$ . Thus we have either  $(r_1, r_2) \in \text{rao}$  or  $(r_2, r_1) \in \text{rao}$ , and by assumption the prior is not the case so  $(r_2, r_1) \in \text{rao} \subseteq \text{OB}$ . However, since  $\text{narR}\langle r_1, \ldots \rangle \prec_\pi \lambda_r$ , we have  $(r_1, r_2) \in \text{OB}$ , which is a contradiction, as OB is irreflexive. Therefore reject our original assumption. Thus  $(r_1, r_2) \in \text{rao}$ , then we have  $(w, r_2) \in \text{ar}$ ; rao  $\subseteq \text{OB}$ , so  $\mathbb{B}\langle w \rangle \prec_\pi \lambda_r$ . Thus we have nicAtomicity $(\pi)$ .

# wfrd

Let us assume we have  $\pi = \pi_4 \cdot \lambda_r \cdot \pi_3$ , with  $\lambda_r \in \{1R\langle r, w \rangle, CAS\langle r, w \rangle, n1R\langle r, w, \_, \_ \rangle, nrR\langle r, w, \_, \_ \rangle, naF\langle r, w, \_, \_ \rangle, narR\langle r, w, \_, \_, \_ \rangle \}$ . In all cases we have  $(w, r) \in \mathsf{rf}$ . Another important fact is that  $\forall w', (w, w') \in \mathsf{mo} \implies (r, w') \in \mathsf{rb}$ .

- If  $\lambda_r = 1R\langle r, w \rangle$ , we need to show wfrdCPU $(r, w, \pi_3)$ .
  - If  $w = init_{\mathsf{loc}(w)}$ , then we need to check that  $\{\mathsf{B}\langle w' \rangle, \mathsf{CAS}\langle w', \_ \rangle \in \pi_3 \mid \mathsf{loc}(w') = \mathsf{loc}(r)\} = \emptyset$  and  $\{\mathsf{lW}\langle w'' \rangle \in \pi_3 \mid \mathsf{loc}(w'') = \mathsf{loc}(r) \land t(w'') = t(r)\} = \emptyset$ . For the first, such a w' would imply  $(r, w') \in \mathsf{rb} \subseteq \mathsf{OB}$ , which contradicts the ordering with  $\lambda_r$ . For the second, such an w'' would imply  $(r, w'') \in \mathsf{rb}_{\mathsf{i}} \subseteq \mathsf{IB}$ , and  $\lambda_r \prec_\pi \mathsf{lW}\langle w'' \rangle$  which similarly contradicts the ordering with  $\lambda_r$ .

- If  $w \in 1\mathbb{W}$ , t(w) = t(r), and  $\mathsf{B}\langle w \rangle \notin \pi_3$ . From  $(w,r) \in \mathsf{rf}_{\mathsf{i}} \subseteq \mathsf{IB}$ , we have  $\lambda_w = 1\mathbb{W}\langle w \rangle \prec_\pi \lambda_r$ , i.e.,  $\pi_3 = \pi_2 \cdot \lambda_w \cdot \pi_1$ . We need to show that  $\{1\mathbb{W}\langle w' \rangle \in \pi_2 \mid \mathsf{loc}(w') = \mathsf{loc}(r) \wedge t(w') = t(r)\} = \emptyset$ . Such a w' would imply  $(w,w') \in \mathsf{po}$  (from  $[1\mathbb{W}]$ ;  $\mathsf{po}$ ;  $[1\mathbb{W}] \subseteq \mathsf{ippo} \subseteq \mathsf{IB}$ , and the execution graph forcing either  $(w,w') \in \mathsf{po}$  or  $(w',w) \in \mathsf{po}$ ,  $(w,w') \in \mathsf{mo}$  (from  $[1\mathbb{W}]$ ;  $\mathsf{po}$ ;  $[1\mathbb{W}] \subseteq \mathsf{oppo} \subseteq \mathsf{OB}$ , and well-formedness conditions forcing either  $(w,w') \in \mathsf{mo}$  or  $(w',w) \in \mathsf{mo}$ ), and  $(r,w') \in \mathsf{rb}_{\mathsf{i}} \subseteq \mathsf{IB}$  would contradicts the ordering with  $\lambda_r$ .
- Else we have  $\lambda_w \in \pi_3$ , with  $\lambda_w \in \{\mathsf{B}\langle w \rangle, \mathsf{CAS}\langle w, \_ \rangle\}$ . If  $w \in \mathsf{IW}$  and t(w) = t(r), this is the remaining subcase, else it comes from  $(w,r) \in \mathsf{rf}_\mathsf{e} \subseteq \mathsf{OB}$ . Thus we have  $\pi_3 = \pi_2 \cdot \lambda_w \cdot \pi_1$ , and we need to check two properties. First, we check that  $\{\mathsf{B}\langle w' \rangle, \mathsf{CAS}\langle w', \_ \rangle \in \pi_2 \mid \mathsf{loc}(w') = \mathsf{loc}(r)\} = \emptyset$ . It holds because such a w' would again imply  $(r, w') \in \mathsf{rb} \subseteq \mathsf{OB}$  and contradict the ordering with  $\lambda_r$ . Second, we check that  $\{w' \mid \mathsf{IW}\langle w' \rangle \in \pi_3 \wedge \mathsf{B}\langle w' \rangle \notin \pi_3 \wedge \mathsf{B}\langle w' \rangle \in \mathsf{IW}(w') = \mathsf{IV}(w') = \mathsf{I$
- If  $\lambda_r = \mathsf{CAS}\langle r, w \rangle$ , we similarly check that  $\mathsf{wfrdCPU}(r, w, \pi_3)$  holds. The difference is that cases that previously contradicted  $(\mathsf{rb_i} \subseteq \mathsf{IB})$  now contradict  $\mathsf{bufFlushOrd}(\pi)$  that forces the buffer of t(r) to be empty when performing  $\lambda_r$ .
- If  $\lambda_r = \text{nlR}\langle r, w, \_, \_ \rangle$ , we need to show wfrdNIC $(r, w, \pi_3)$ .
  - If  $w = init_{\mathtt{loc}(w)}$ , then we need to check that  $\{\mathsf{B}\langle w'\rangle, \mathtt{CAS}\langle w', \_\rangle \in \pi_3 \mid \mathtt{loc}(w') = \mathtt{loc}(r)\} = \emptyset$ . Such a w' would imply  $(r, w') \in \mathsf{rb} \subseteq \mathsf{OB}$ , which contradicts the ordering with  $\lambda_r$ .
  - Else we have  $\lambda_w \in \pi_3$ , with  $\lambda_w \in \{B\langle w \rangle, CAS\langle w, \_\rangle\}$ . This comes from  $(w,r) \in \mathsf{rf}_{\mathsf{e}} \subseteq \mathsf{OB}$ . Thus we have  $\pi_3 = \pi_2 \cdot \lambda_w \cdot \pi_1$ , and we need to check that
    - $\{B\langle w'\rangle, CAS\langle w', \_\rangle \in \pi_2 \mid loc(w') = loc(r)\} = \emptyset$ . It holds because such a w' would again imply  $(r, w') \in rb \subseteq OB$  and contradict the ordering with  $\lambda_r$ .
- If  $\lambda_r = \operatorname{nrR}\langle r, w, \_, \_ \rangle$ ,  $\operatorname{naF}\langle r, w, \_, \_ \rangle$  or  $\operatorname{narR}\langle r, w, \_, \_, \_ \rangle$ , we similarly check that  $\operatorname{wfrdNIC}(r, w, \pi_3)$  for the same reasons.

Thus we have  $\mathsf{wfrd}(\pi)$ .

**Theorem 12.** Let G be a well-formed consistent execution graph generated from a program P. Let  $\pi$  be the path obtained from G by the process defined above. Then there is M', QP' (such that forall  $t, \overline{n}$  we have  $QP'(t)(\overline{n}) = \langle \varepsilon, \varepsilon, nEX^* \rangle$ ), and an equivalent path  $\pi'$  (producing the same outcome as  $\pi$ ) such that  $P, M_0, B_0, A_0, QP_0, \varepsilon \Rightarrow^* (\lambda t.skip), M', B_0, A_0, QP', <math>\pi'$ .

*Proof.* From above, we have  $\mathsf{wf}(\pi)$ . This shows that the program configuration can perform the events described by the annotated labels of  $\pi$ . The remaining part of the proof is simply to check that the command rewritings used when deriving the execution graph from P (see Fig. 23) can be used as  $\mathcal{E}$  transitions in the annotated semantics for P, which follows from the definitions.

## D.7 Operational Semantics and Annotated Semantics

We define forgetful functions from annotated configurations to operational configurations. For memories, we replace the write event by the value written. For labels within annotated configurations, we drop some arguments to recover the data structure of the operational semantics.

$$\label{eq:matter_model} \begin{split} [\![\cdot]\!]_{\mathrm{M}} : \mathsf{AMem} &\to \mathsf{Mem} \\ [\![\mathsf{M}]\!]_{\mathrm{M}} &\triangleq \lambda x. v_{\mathrm{w}}(\mathsf{M}(x)) \end{split}$$

$$\begin{split} \llbracket \cdot \rrbracket_{\mathrm{op}} : E^{\mathrm{ext}} &\rightharpoonup \left\{ \begin{aligned} y^{\overline{n}} &:= x^{n}, y^{\overline{n}} := v, \mathrm{ack_{p}}, x^{n} := y^{\overline{n}}, x^{n} := v, \\ x &:= \mathrm{RCAS}(y^{n}, v, v'), x := \mathrm{RFAA}(y^{n}, v), \mathrm{cn, rfence}(\overline{n}) \end{aligned} \right\} \\ & & & & & & & & & & & & & & & \\ \lVert \mathrm{IW}(x, v_{\mathrm{w}}) \rVert_{\mathrm{op}} & \triangleq x := v_{\mathrm{w}} & & & & & & & & & \\ \lVert \mathrm{nIW}(x, v_{\mathrm{w}}, \overline{n}) \rVert_{\mathrm{op}} & \triangleq \overline{y} := v_{\mathrm{r}} & & & & & & & & \\ \lVert \mathrm{nIW}(x, v_{\mathrm{w}}, \overline{n}) \rVert_{\mathrm{op}} & \triangleq x := v_{\mathrm{w}} & & & & & & & & \\ \lVert \mathrm{nIW}(x, v_{\mathrm{w}}, \overline{n}) \rVert_{\mathrm{op}} & \triangleq x := v_{\mathrm{w}} & & & & & & & \\ \lVert \mathrm{nIW}(x, v_{\mathrm{w}}, \overline{n}) \rVert_{\mathrm{op}} & \triangleq x := v_{\mathrm{w}} & & & & & & & \\ \lVert \mathrm{nIW}(x, v_{\mathrm{w}}, \overline{n}) \rVert_{\mathrm{op}} & \triangleq x := v_{\mathrm{w}} & & & & & & \\ \lVert \mathrm{nIW}(x, v_{\mathrm{w}}, \overline{n}) \rVert_{\mathrm{op}} & \triangleq x := v_{\mathrm{w}} & & & & & & \\ \lVert \mathrm{nIW}(x, v_{\mathrm{w}}, \overline{n}) \rVert_{\mathrm{op}} & \triangleq x := v_{\mathrm{w}} & & & & & & \\ \lVert \mathrm{nIW}(x, v_{\mathrm{w}}, \overline{n}) \rVert_{\mathrm{op}} & \triangleq x := v_{\mathrm{w}} & & & & & \\ \lVert \mathrm{nIW}(x, v_{\mathrm{w}}, \overline{n}) \rVert_{\mathrm{op}} & \triangleq x := v_{\mathrm{w}} & & & & & \\ \lVert \mathrm{nIW}(x, v_{\mathrm{w}}, \overline{n}) \rVert_{\mathrm{op}} & \triangleq x := v_{\mathrm{w}} & & & & & \\ \lVert \mathrm{nIW}(x, v_{\mathrm{w}}, \overline{n}) \rVert_{\mathrm{op}} & \triangleq x := v_{\mathrm{w}} & & & & & \\ \lVert \mathrm{nIW}(x, v_{\mathrm{w}}, \overline{n}) \rVert_{\mathrm{op}} & \triangleq x := v_{\mathrm{w}} & & & & & \\ \lVert \mathrm{nIW}(x, v_{\mathrm{w}}, \overline{n}) \rVert_{\mathrm{op}} & \triangleq x := v_{\mathrm{w}} & & & & & \\ \lVert \mathrm{nIW}(x, v_{\mathrm{w}}, \overline{n}) \rVert_{\mathrm{op}} & \triangleq x := v_{\mathrm{w}} & & & & & \\ \lVert \mathrm{nIW}(x, v_{\mathrm{w}}, \overline{n}) \rVert_{\mathrm{op}} & \triangleq x := v_{\mathrm{w}} & & & & \\ \lVert \mathrm{nIW}(x, v_{\mathrm{w}}, \overline{n}) \rVert_{\mathrm{op}} & \triangleq x := v_{\mathrm{w}} & & & & \\ \lVert \mathrm{nIW}(x, v_{\mathrm{w}}, \overline{n}) \rVert_{\mathrm{op}} & \triangleq x := v_{\mathrm{w}} & & & & \\ \lVert \mathrm{nIW}(x, v_{\mathrm{w}}, \overline{n}) \rVert_{\mathrm{op}} & \triangleq x := v_{\mathrm{w}} & & & & \\ \lVert \mathrm{nIW}(x, v_{\mathrm{w}}, \overline{n}) \rVert_{\mathrm{op}} & \triangleq x := v_{\mathrm{w}} & & & \\ \lVert \mathrm{nIW}(x, v_{\mathrm{w}}, \overline{n}) \rVert_{\mathrm{op}} & \triangleq x := v_{\mathrm{w}} & & & \\ \lVert \mathrm{nIW}(x, v_{\mathrm{w}}, \overline{n}) \rVert_{\mathrm{op}} & \triangleq x := v_{\mathrm{w}} & & & \\ \lVert \mathrm{nIW}(x, v_{\mathrm{w}}, \overline{n}) \rVert_{\mathrm{op}} & \triangleq x := v_{\mathrm{w}} & & & \\ \lVert \mathrm{nIW}(x, v_{\mathrm{w}}, \overline{n}) \rVert_{\mathrm{op}} & \triangleq x := v_{\mathrm{w}} & & & \\ \lVert \mathrm{nIW}(x, v_{\mathrm{w}}, \overline{n}) \rVert_{\mathrm{op}} & \triangleq x := v_{\mathrm{w}} & & & \\ \lVert \mathrm{nIW}(x, v_{\mathrm{w}}, \overline{n}) \rVert_{\mathrm{op}} & \triangleq x := v_{\mathrm{w}} & & & \\ \lVert \mathrm{nIW}(x, v_{\mathrm{w}}, \overline{n}) \rVert_{\mathrm{op}} & \triangleq x := v_{\mathrm{$$

$$\begin{split} \llbracket \cdot \rrbracket_{\mathrm{opl}} : E^{\mathrm{ext}} & \rightharpoonup \left\{ \begin{aligned} y^{\overline{n}} &:= x^n, y^{\overline{n}} := v, x^n := y^{\overline{n}}, x^n := v, \\ x &:= \mathtt{RCAS}(y^n, v, v'), x := \mathtt{RFAA}(y^n, v), \mathtt{cn}, \mathtt{rfence}(\overline{n}) \end{aligned} \right\} \\ & \qquad \qquad \begin{bmatrix} l \rrbracket_{\mathrm{opl}} &= \begin{cases} \mathtt{cn} & \text{if } l = \mathtt{nrEX}(\overline{n}) \\ \llbracket l \rrbracket_{\mathrm{op}} & \text{otherwise} \end{aligned} \end{split}$$

The labels that cannot appear in a well-formed annotated configuration are not mapped. For put operations, the operational semantics uses both  $(ack_p)$  and (cn) while the annotated semantics uses the label nrEX, so the mapping is different for labels in  $wb_L$ .

 $[\![\cdot]\!]_{\text{op}}$  and  $[\![\cdot]\!]_{\text{opl}}$  are extended to lists in an obvious way.

We then extend this to configurations as expected. We overload notations to simplify the formulas.

For  $\mathsf{qp} = \langle \mathbf{pipe}, \mathbf{wb_R}, \mathbf{wb_L} \rangle \in \mathsf{AQPair}$ , we define  $\llbracket \mathsf{qp} \rrbracket \triangleq \langle \llbracket \mathbf{pipe} \rrbracket_{\mathrm{op}}, \llbracket \mathbf{wb_R} \rrbracket_{\mathrm{op}}, \llbracket \mathbf{wb_L} \rrbracket_{\mathrm{opl}} \rangle$ . For  $\mathsf{QP} \in \mathsf{AQPMap}$ , we define  $\llbracket \mathsf{QP} \rrbracket \triangleq \lambda t. \lambda \overline{n}. \llbracket \mathsf{QP}(t)(\overline{n}) \rrbracket$ . For  $\mathsf{B} \in \mathsf{ASBMap}$ , we define  $\llbracket \mathsf{B} \rrbracket \triangleq \lambda t. \llbracket \mathsf{B}(t) \rrbracket_{\mathrm{op}}$ .

**Theorem 13.** For all  $P, P' \in Prog$ ,  $M, M' \in AMem$ ,  $B, B' \in ASBMap$ ,  $A, A' \in RAMap$ ,  $QP, QP' \in AQPMap$ ,  $\pi, \pi' \in Path$ ,  $if P, M, B, A, QP, \pi \Rightarrow P', M', B', A', QP', \pi'$  and  $wf(M, B, A, QP, \pi)$ , then  $P, [\![M]\!]_M, [\![B]\!], A, [\![QP]\!] \Rightarrow P', [\![M']\!]_M, [\![B']\!], A', [\![QP']\!].$ 

*Proof.* By straightforward induction on  $\Rightarrow$ .

**Theorem 14.** For all M  $\in$  AMem, M"  $\in$  Mem, B  $\in$  ASBMap, B"  $\in$  SBMap, A, A'  $\in$  RAMap, QP  $\in$  AQPMap, QP"  $\in$  QPMap, and  $\pi \in$  Path, if P,  $[\![M]\!]_M$ ,  $[\![B]\!]$ , A,  $[\![QP]\!] \Rightarrow$  P', M", B", A', QP" and wf(M, B, A, QP,  $\pi$ ), then there exists M'  $\in$  AMem, B'  $\in$  ASBMap, QP'  $\in$  AQPMap, and  $\pi' \in$  Path such that  $[\![M']\!]_M = M$ ",  $[\![B']\!] = B$ ",  $[\![QP']\!] = QP''$ , and P, M, B, A, QP,  $\pi \Rightarrow P'$ , M', B', A', QP',  $\pi'$ .

*Proof.* By straightforward induction on  $\Rightarrow$ . In some cases, the reduction enforces a specific annotated label  $\lambda$  and we have  $\pi' = \lambda \cdot \pi$ ; we then need wf(M, B, A, QP,  $\pi$ ) to check that  $\lambda$  is fresh enough for  $\pi$ .

# Theorem 15 (Operational and Annotated Semantics Equivalence). For all program P.

- [M<sub>0</sub>]<sub>M</sub>, [B<sub>0</sub>], A<sub>0</sub>, and [QP<sub>0</sub>] are the initialisation for the operational semantics:
- $If P, M_0, B_0, A_0, QP_0, \varepsilon \Rightarrow^* P', M', B', A', QP', \pi' then P, \llbracket M_0 \rrbracket_M, \llbracket B_0 \rrbracket, A_0, \llbracket QP_0 \rrbracket \Rightarrow^* P', \llbracket M' \rrbracket_M, \llbracket B' \rrbracket, A', \llbracket QP' \rrbracket$
- If P,  $[M_0]_M$ ,  $[B_0]$ ,  $A_0$ ,  $[QP_0]$   $\Rightarrow^* P'$ , M'', B'', A', QP'' then there exists  $M' \in AMem$ ,  $B' \in ASBMap$ ,  $QP' \in AQPMap$ , and  $\pi' \in Path$  such that  $[M']_M = M''$ , [B'] = B'', [QP'] = QP'', and P,  $M_0$ ,  $B_0$ ,  $A_0$ ,  $QP_0$ ,  $\varepsilon \Rightarrow^* P'$ , M', B', A', QP',  $\pi'$ .

*Proof.* The first point comes from unfolding the definitions. The other two are proved by straightforward induction on  $\Rightarrow^*$  and using Theorems 13 and 14. The condition  $wf(M, B, A, QP, \pi)$  is obtained by applying Theorem 8 when needed.