

Imperial College London

MENG INDIVIDUAL PROJECT

IMPERIAL COLLEGE LONDON

DEPARTMENT OF COMPUTING

Semantics of RDMA Remote Read-Modify-Write Operations

Author:
Max Stuppel

Supervisors:
Dr. Azalea Raad
Dr. Guillaume Ambal

Second Marker:
Prof. Alastair F. Donaldson

June 13, 2025

Abstract

Remote direct memory access (RDMA) allows each computer in a network to directly access the memory of every other machine, without involving the operating system of either side. This improves data transfer speeds, but reduces synchronisation, making it harder to reason about program behaviours. Formal memory models allow us to precisely describe which program behaviours are allowed, and to verify that a particular program only observes behaviours that the programmer expects.

The RDMA^{TSO} memory model describes the semantics of RDMA, but it lacks support for remote read-modify-write operations. Read-modify-write (RMW) operations are the foundation for all synchronisation in single-machine concurrent programming, and RDMA remote RMWs allow the same techniques to be lifted to the network level, facilitating high-speed, fine-grained concurrent programming.

In this report, we extend RDMA^{TSO} with support for remote RMWs. We provide an intuitive account of remote RMW behaviours, which has been reviewed and confirmed by an expert from NVIDIA, a leading vendor of RDMA systems. We formalise these behaviours through two models, one operational and the other declarative, and prove that the two models are equivalent. This work expands the range of programs for which RDMA^{TSO} can describe a formal semantics, providing a robust mechanism to understand and reason about RDMA programs involving remote RMWs.

Acknowledgements

I owe great thanks to my supervisors, Azalea Raad and Guillaume Ambal, who have kindly and patiently supported me throughout this project. This work has been not only deeply interesting in itself, but also personally enlightening, in that it has played a significant role in my decision to pursue further research. I am very grateful to have had the opportunity to work on this topic.

I would also like to thank my second marker, Alastair Donaldson, for his diligent review of my interim report, which supported the final introduction and background sections seen here. Thanks also to Haggai Eran, whose expertise in RDMA has been invaluable in deepening and clarifying my understanding of this technology, and to Dan Iorga, for helping me better understand the Alloy modelling language.

Finally, I want to thank all of my friends and family who have supported me throughout this project and my degree. I am particularly grateful to my parents, who instilled in me the curiosity and ambition that have driven me at Imperial, and to my partner Laura, whose continual encouragement and support has helped me immeasurably.

Contents

1	Introduction	4
2	Background	6
2.1	Concurrency	6
2.2	Synchronisation	6
2.3	Memory Models	7
2.4	RDMA	9
2.5	Remote RMWs	10
3	Overview	12
3.1	RDMA Concurrency	12
3.2	Remote RMW Behaviours	13
4	Operational Semantics	17
4.1	States of the Operational Semantics	17
4.1.1	Memory	17
4.1.2	Program State	17
4.1.3	Store Buffers	18
4.1.4	RDMA Operations	18
4.2	Transitions of the Operational Semantics	19
4.2.1	Program Transitions	20
4.3	Hardware Domains	20
4.3.1	Hardware Transitions	21
4.3.2	Queue-Pair Transitions	21
5	Declarative Semantics	27
5.1	Events and Executions	27
5.2	Semantics of a Program	33
6	Equivalence	36
6.1	Structure of the Proof	36
6.2	Annotated Semantics	36
6.3	Annotated Semantics to Declarative Semantics	38
6.4	Declarative Semantics to Annotated Semantics	38
6.5	Operational and Annotated Semantics	39
7	Discussion	40
7.1	Conclusion	40
7.2	Limitations	40
7.3	Future Work	40
A	Declarations	44
A.1	Use of Generative AI	44
A.2	Ethical Considerations	44
A.3	Sustainability	44
A.4	Availability of Materials	44

B	Annotated Semantics	45
B.1	Annotated Labels and Inference Rules	45
B.2	Paths, Gluing, and Other Definitions	50
B.3	From Annotated Semantics to Declarative Semantics	59
B.4	From Declarative Semantics to Annotated Semantics	67
B.5	Operational Semantics and Annotated Semantics	75
C	Encoding the Declarative Model in Alloy	77
C.1	Prototype Encoding	77
C.2	Example Litmus Test	82

1 Introduction

Concurrency refers to the ability of a system to execute multiple tasks (threads) at once, either by executing them *simultaneously* (in parallel) on multiple processor cores, or by time-sharing a single core and frequently context switching between tasks. All modern CPUs feature a multi-core design, so to maximise performance a program must distribute its work across multiple threads, which can then run concurrently: if a program can fully utilise all cores on a four-core multi-processor, then it would run in roughly one quarter of the time compared to if it executed the same tasks sequentially on a single core. Since the processor, not the programmer, decides how to distribute threads across cores and when to interrupt a thread to allow another to execute, the programmer must assume that the instructions of different threads may be executed in an arbitrary order. We call the rules which govern the order in which program instructions may be executed the *memory model*, or concurrent semantics, of that processor.

Remote direct memory access (RDMA) technologies enable a machine to directly read and write the memory of another machine within the network, bypassing the operating system of each and thereby reducing the number of CPU cycles consumed in transferring data [1]. This is particularly useful in massively parallel computer clusters (e.g. cloud computing, big data and scientific computation), where this enables high-throughput, low-latency networking [2]. While older implementations of RDMA, such as *InfiniBand* [3], required specialised hardware, newer implementations of RDMA, such as *RDMA over Converged Ethernet* [4], operate on top of ordinary networking infrastructure. This has brought costs in line with traditional networking [5], leading to widespread adoption of RDMA as of 2018 [6].

Remote read (get) and write (put) operations are dispatched to the network interface card (NIC), which executes transfers in parallel as the program continues on the CPU. Due to the independent operation of CPU and NIC for local and remote operations respectively, instructions may not be executed in the order the programmer intuitively expects, leading to unintended behaviours; this is comparable to the operation of shared memory in a multi-processor. However, there is the additional caveat that the concurrent behaviour is multi-tiered: a program consists of several threads, with those threads distributed across a number of machines (nodes) in the network (inter-node concurrency), and each node executing several threads (intra-node concurrency). To write *correct* programs for RDMA networks, we need to know the order in which local (CPU) and remote (NIC) operations are executed, how they can be reordered, and when the effects of those operations become visible, both to local and remote threads.

To this end, similarly to how the concurrent semantics of Intel-x86 [7] or ARM [8] have been formalised, there has been some progress in formally describing the semantics of RDMA. Namely, two models exist in the literature: coreRMA [9] and RDMA^{TSO} [10]. The coreRMA model is subject to several significant limitations, most notably that it models intra-node concurrency as following sequential consistency: an unrealistic and simplified memory model which no modern processor adheres to by default. RDMA^{TSO} is instead based on the memory model of x86 processors, a reasonable assumption as Intel and AMD processors that use this architecture are widespread. However, RDMA^{TSO} does not model the behaviour of every remote operation. While it does describe get and put operations, as well as remote *fences* useful for synchronisation, it omits remote *read-modify-write* (RMW, also known as atomic) operations.

Read-modify-write operations are an essential tool for concurrent programming [11],

and RDMA specifies remote RMWs to allow similar techniques to be applied at the network level. These operations read, modify, and write a remote memory location, while guaranteeing that no other remote RMW acts on that memory location between the read and write. This enables the implementation of shared higher-level synchronisation tools, such as locks and barriers. We propose extensions to the formal models described by RDMA^{TSO} , to also specify the behaviour of such remote atomic operations.

Contributions and Outline

In §2, we review the existing literature on memory models and RDMA. In §3, we overview the RDMA^{TSO} memory model and present our intuitive account of the weak behaviours of remote RMWs, which we illustrate through a number of examples. In §4, we develop the operational semantics of RDMA^{TSO} with new rules for remote RMWs, and in §5, we extend the declarative model. In §6, we discuss how the two models are proven equivalent. Finally in §7, we conclude with a discussion of the limitations of this project and future work on this topic.

2 Background

2.1 Concurrency

Since the end of the “free ride” for processor speed enjoyed during the era of Moore’s law, the microprocessor landscape has become dominated by multi-core processors (MCPs) [12], which provide faster (effective) speeds than uniprocessors due to the distribution of work across multiple cores, which execute in parallel. In order to maximise the utilisation of these MCPs, and to minimise program execution time, programs must be written concurrently so that their work can be efficiently distributed across processor cores.

However, concurrent programming is widely considered to be challenging [13]. This is largely due to the non-deterministic nature of concurrent program execution: instructions from each program thread may interleave arbitrarily, which can lead to “Heisenbugs” [14] – errors which are only present in a subset of the possible interleavings, and hence appear seemingly at random while evading reliable reproduction. A typical cause of these types of bugs is a *data race*, where two program threads access the same memory location without synchronisation, with at least one of the accesses being a write. Depending on which thread happens to execute first, different program behaviours may be observed.

For example, consider the litmus test in Fig. 2.1a, where two threads attempt to increment the counter x concurrently: only if the two threads execute sequentially – that is, one thread executes to completion before the other commences – will the expected result of $x = 2$ be achieved; any other interleaving will cause one thread to overwrite the other, resulting in $x = 1$. This is equivalent to each thread executing $x := x + 1$, since that instruction would also be separated into a read and a write step; we have explicitly used two operations in our example for clarity.

2.2 Synchronisation

So, how can we avoid data races in our programs? We need a way to access data that guarantees *mutual exclusion*: the guarantee that so long as one thread is performing an operation on a shared memory location, no other thread may access it until the operation is complete. Modern processors provide *Read-Modify-Write* (RMW) instructions that facilitate this. An example of an RMW is *Compare-and-Swap* (CAS), which accepts three parameters and returns a value: $a := \text{CAS}(x, e, u)$. Given a memory location x containing a value v , and values e (expected) and u (update), it checks if $v = e$, and if so assigns u to location x ; otherwise, it does not alter memory. It returns v – so in our example $a = v$ after the command returns – which allows the caller to check whether or not the swap was performed. Crucially, the hardware guarantees that this happens *atomically*, without any other thread accessing location x during the operation. In fact, CAS is sufficiently powerful to implement all other RMWs [11], and can be used to implement higher-level synchronisation tools such as locks and barriers [15, Chs. 7.2 and 17.2].

Another example of an RMW is *Fetch-and-Add* (FAA): $a := \text{FAA}(x, i)$, which given a memory location x containing value v , atomically assigns $v + i$ to location x , returning v . Unlike CAS, we might not know what value is stored at x before writing to it, but for use cases such as counters it allows a simpler implementation than CAS. However, if we needed to enforce a strict guarantee that the counter value stays within some bound, for example, FAA would not be suitable since the update occurs unconditionally.

$x = 0$	
$a := x$	$b := x$
$x := a + 1$	$x := b + 1$

(a) $x = 1, x = 2$ ✓

$x = 0$	
$a := \text{FAA}(x, 1)$	$b := \text{FAA}(x, 1)$

(b) $x = 1$ ✗ $x = 2$ ✓

Figure 2.1: Left: unsafe incrementation of the counter results in a data race — whichever thread reads first must write before the other reads or writes, or else its effect will be lost. Right: safe implementation using FAA guarantees the effects of both threads will be visible.

Notation. Columns (separated by $|$) denote distinct threads, variables x, y denote shared memory locations, a, b are thread-local variables (assumed to be initially 0), and the first line indicates the initial values of memory locations. Captions indicate allowed outcomes for variables with ✓, and disallowed outcomes with ✗.

Fig. 2.1b shows a safe implementation for concurrently implementing a counter using FAA. Since FAA guarantees that no other operation may occur on x between the read and write, we can be sure that both incrementations will be visible, resulting in a final state of $x = 2$.

2.3 Memory Models

Given the ability to synchronise between threads and ensure mutual exclusion where necessary, it seems that it would not be too challenging for programmers to reason about potential interleavings of instructions from each program thread, then use the tools available to write correct programs. However, the real world is significantly more complicated. The interleaving model, called *Sequential Consistency* (SC), is not realistic: no modern computer processor provides the guarantee that instructions within a thread are executed in order [16]. Instead, they only provide *weak* memory models.

To be precise, under sequential consistency, it is guaranteed that instructions within a thread execute in order, while instructions from different threads may interleave in an arbitrary execution order. A weak memory model, then, is one which, in addition to permitting any execution ordering allowed by SC, also permits some other orders – which is to say that the guarantee that instructions within a thread execute in order is *relaxed* in some way. It may also be, as we will see in the case of RDMA, that we cannot simply define one order over the instructions, as the order in which they are issued and that in which they become observable may be different.

There are numerous weak memory models, however we will focus on *Total Store Ordering* (TSO). TSO is a description of program execution on multiprocessors with store buffers, originally formulated to describe SPARC V8/9 processors [17]. It permits certain reorderings of instructions that are not allowed in SC, namely:

Write-read reordering (on different locations): a later read may be reordered *before* an earlier write, provided it is on a different memory location. This leads to weak behaviour known as *store-buffering*.

Write-to-read forwarding (on the same location): the value of an earlier write can be *copied* into a later read.

$x = y = 0$	$x = y = 0$	$x = y = 0$	$x = y = 0$
$x := 1$ $y := 1$ $a := y$ $b := x$	$x := 1$ $y := 1$ mfence mfence $a := y$ $b := x$	$a := y$ $b := x$ $x := 1$ $y := 1$	$x := 1$ $a := y$ $y := 1$ $b := x$
(a) $a = b = 0$ ✓	(b) $a = b = 0$ ✗	(c) $a = b = 1$ ✗	(d) $(a, b) = (1, 0)$ ✗

Figure 2.2: TSO litmus tests. Reproduced from [10] with permission.

Write-assignment reordering: a later assignment (which does not access memory) may be reordered *before* an earlier write.

Store-buffering is so-called because it can be modelled by assuming a buffer between a thread’s private registers and shared memory, which behaves as a FIFO queue for write (store) instructions. Write instructions do not immediately affect memory, but rather enter the store buffer and execute at some later point, once they have reached the head of the queue. This facilitates write-read reordering, as a later read may not see an earlier write because it is still in the buffer; similarly, an assignment may not observe a buffered write. Write-to-read forwarding can be interpreted as a read viewing the latest value from the thread’s store buffer, rather than reading from memory.

The TSO model has also been adapted to describe x86 processors, by Owens et al. [7]. In this paper, the authors define the concurrent semantics of the x86 architecture, which they call x86-TSO due to its similarity to the SPARC memory model. Notably, x86-TSO supports a *memory fence* instruction, which flushes a thread’s store buffer. Henceforth, we use TSO to refer to the x86 memory model, rather than the original SPARC definition.

Inevitably, it is more challenging for programmers to reason about weak memory models than SC; indeed, TSO is a relatively simple model, yet there may be far more possible execution orders to consider due to the reordering within threads. In order to make these weak models easier to understand, it is useful to provide a *formal semantics* that programmers may use to build their intuition, or which may be used to write tools to automatically check programs for bugs.

Operational semantics are well established in the literature for defining the step-by-step execution of a program by providing a set of *transition rules*. These rules specify how the program’s state can be updated at any given execution step, depending on whether particular logical properties hold at that point. For any valid program execution, one may write down a logical trace of the execution using the transition rules, called a derivation tree. Hence, an operational semantics gives a low-level description of program execution, which may be used for formal reasoning, and can help programmers understand program behaviour.

Declarative semantics, on the other hand, describe program executions as a graph: we use the program instructions (events) as nodes in the graph, and define several relations over the events, representing dependencies. For example, we might define the *reads-from* (**rf**) relation: if a read r reads the value written by an earlier write w , then we add an **rf** dependency from w to r . We might also impose certain restrictions on the relations in our graph, for example that our graph should not contain any cycles. Such execution graphs are particularly useful for understanding the properties and correctness of a program, and are much more concise than operational semantics; however they may be less useful for gaining an intuitive understanding of how programs execute.

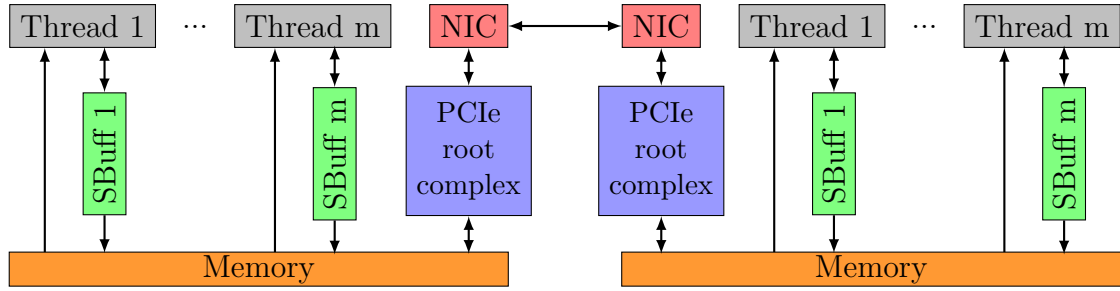


Figure 2.3: RDMA network overview [10, page 10]. Communication between NICs occurs in parallel to ongoing execution of local threads.

2.4 RDMA

Conventionally, communication over the network between separate machines requires the involvement of the operating system, which is responsible for the network and transport layers of the network stack. In common usage, the $\sim 1\text{ms}$ latency introduced due to encoding/decoding the TCP/IP stack is a negligible factor in the overall latency, however in modern computer clusters the propagation delay is far lower. Thus, it is desirable to use a thinner network stack which bypasses the OS, so that remote reads and writes use fewer CPU cycles. This has led to widespread adoption of *Remote Direct Memory Access* (RDMA), which allows computers in the network to directly transfer data without the need for caches, switches, or any involvement by the CPU – the operation may be handled solely by the Network Interface Card (NIC), which transfers messages directly between main memory and the wire. Transfers are executed in parallel with system operations, reducing latency and increasing throughput [6].

In particular, as these network operations are executed in parallel to other operations on the CPU, they do not observe sequential consistency; hence it is necessary to define a memory model for RDMA so that programmers may understand the weak behaviours that can be observed, and use that knowledge to write correct programs and libraries for RDMA networks.

The first attempt to provide a formal semantics of RDMA programs was coreRMA [9]. However, this work is subject to a number of significant shortcomings, as described by Ambal et al. [10]. They identify four limitations of the coreRMA model: (1) the coreRMA authors assume that the concurrency on each node follows *sequential consistency* – this assumption is unrealistic as all modern CPU architectures (in particular, x86 and ARM, which are ubiquitous) offer a weaker model by default; (2) they provide only a *declarative* model, while operational models are required for certain applications, such as invariant synthesis [18], and may provide better intuition than declarative models, as they are explicitly intended to be easy to understand [19]; (3) the model is not validated against existing implementations, and they could not observe any of the weak behaviours permitted by coreRMA; (4) coreRMA is neither stronger nor weaker than the RDMA specification, admitting certain behaviours not allowed by the specification, while prohibiting behaviours which the specification permits.

Due to these limitations, Ambal et al. define RDMA^{TSO} , a memory model for RDMA which improves on coreRMA in a number of ways: it assumes an x86 architecture, hence modelling individual nodes as observing TSO; they provide both operational and declarative models; and they have validated that all weak behaviours permitted by the model

$x = 0$			$x = 0$
$a := \text{CAS}(x, 0, 7) \mid b := \text{FAA}(x, 2)$	$a := \text{nCAS}(x^3, 0, 7)$	$b := \text{nFAA}(x^3, 2)$	
(a) $x = 2, x = 9$ ✓ $x = 7$ ✗	(b) $x = 2, x = 9$ ✓ $x = 7$ ✗		
		$x = 0$	
	$a := \text{nCAS}(x^2, 0, 7)$	$b := \text{FAA}(x, 2)$	
(c) $x = 2, x = 7, x = 9$ ✓			

Figure 2.4: Remote RMWs are non-atomic with respect to local operations: this results in behaviours not visible when only using local RMWs. Remote RMWs are still atomic with respect to each other. We note remote CAS (resp. FAA) by **nCAS** (resp. **nFAA**), where **n** stands for NIC. Note that for RDMA litmus tests we use a double vertical line to indicate threads on separate nodes, and x^n to denote the variable x on foreign node n .

are either observed on actual hardware, or permitted by the specification but not implemented in current systems. Thus, due to the multiple models being proven equivalent, extensive hardware testing, and confirmation from hardware experts, RDMA^{TSO} may be observed to be a very accurate model of concurrency possible on RDMA systems running on the x86 architecture.

2.5 Remote RMWs

Despite the advantages over coreRMA, the RDMA^{TSO} model is arguably incomplete, as it does not describe all operations provided by RDMA. It provides the semantics of all local operations (as in TSO), and the semantics of remote read (get) and remote write (put) operations, as well as the semantics of remote fences instructions used for synchronisation. It does not provide, however, the semantics of remote read-modify-write (rRMW) instructions defined in the RDMA specification. Remote RMWs have the potential to be a very useful tool for synchronisation in distributed algorithms over RDMA networks, however they are subject to a peculiar set of properties.

Under TSO, (local) RMWs are not subject to any reordering: this makes them easy to reason about, since we know precisely what value the memory location in question held immediately before and immediately after the execution of the RMW, and, since the RMW executes as an atomic unit, we know that there was no intermediary value held or operation executed on that location. RDMA remote RMWs are much more challenging to understand and use by comparison.

Firstly, rRMWs have complex reordering rules with respect to other operations in the same thread, as we will detail later. Secondly, and more importantly, rRMWs only guarantee atomicity with respect to *other rRMWs* operations, *not* local operations on the target node, nor other types of remote operations towards that node. Although an option is available on some hardware to execute rRMWs with a higher level of atomicity, this is not globally supported, nor enabled by default. Thus, RDMA programmers will need to operate under the default assumption that other operations (in particular, local or remote writes) may occur between the read and write of a remote RMW.

Consider the examples in Fig. 2.4. Fig. 2.4a shows how the strong atomicity guarantees

of local RMWs require the effects of both operations to be visible: either the CAS succeeds first, setting x to 7, after which the FAA adds two to it to achieve $x = 9$; or the FAA occurs first, setting $x = 2$, after which the CAS fails. No other outcome is possible, as the execution of the two RMWs may not overlap – we certainly cannot observe $x = 7$, as this would mean we have somehow lost or overwritten the FAA. When using multiple remote RMWs as in Fig. 2.4b, we observe similar behaviour due to the guarantee that an rRMW is atomic with respect to other rRMWs operations. However in Fig. 2.4c, we see that such an outcome is possible when using a mixture of local and remote RMWs: the $x = 7$ result is achieved when the (local) FAA is executed between the read and write of the remote CAS. Such an execution is permitted due to the weak atomicity guarantees of remote RMWs.

For these reasons, a formal definition of the semantics of RDMA remote RMWs will be particularly helpful to RDMA programmers, and to that end we will present suitable extensions to RDMA^{TSO} such that programmers may confidently reason about all aspects of RDMA programs.

3 Overview

3.1 RDMA Concurrency

An interesting consequence of RDMA operations being handled by the NIC separately to the CPU is that it creates two distinct tiers of concurrency. In general, the CPU does not wait for NIC operations to complete – this would greatly slow down program execution, as despite network communication being relatively fast (compared to TCP/IP), it is still much slower than local operations on the CPU. Therefore, we can consider CPU (intra-node) concurrency (governed by TSO) separately to network-level (inter-node) concurrency. The reordering rules between CPU and NIC operations are very simple: an earlier NIC operations can *always* be reordered *after* a later CPU operation.

The CPU reordering rules have already been discussed in §2.3. Reordering rules for RDMA operations on a single thread (except rMWs) are described in [10] and shown in Fig. 3.1. To summarise:

- CPU operations always complete before later RDMA operations (Fig. 3.1a)
- RDMA operations may be reordered after later CPU operations, as if executing concurrently in a separate thread (Fig. 3.1b)
- Poll awaits the earliest (unpolled) RDMA operation towards a given node (Figs. 3.1c and 3.1d)
- RDMA operations towards distinct nodes may be reordered (Figs. 3.1e and 3.1f)
- An earlier put always completes before a later get towards the same node (Fig. 3.1g)
- An earlier get may be reordered after a later put towards the same node (Fig. 3.1h)
- RDMA operations always complete before a later remote fence towards the same node

We are more interested in RDMA programs running on multiple nodes. Fig. 3.2 shows a number of litmus tests for inter-node concurrency. Store buffering, as seen in TSO, is similarly possible under RDMA (Fig. 3.2a). Additional weak behaviours are also possible, including load buffering (Fig. 3.2b); load buffering is analogous to store buffering, but concerning unobserved reads rather than writes (in our case, the read is a get operation).

Furthermore, whilst polling can prevent most weak behaviours, including load buffering (Fig. 3.2c), store buffering cannot be prevented due to a difference in how a poll awaits the completion of a put as opposed to a get. In the case of a get, the poll only completes once the value read on the remote node has been written to local memory. A put, however, may be considered complete while the write of the value to remote memory is still pending: this is because the confirmation of a put is returned when the operation is buffered, not when it is complete. In practice, the write is almost always de-buffered before the confirmation has travelled across the network, but the behaviour of Fig. 3.2d is allowed by the architecture specification.

However, it is still possible to prevent store buffering. To our knowledge, RDMA is implemented on hardware which meets the PCIe standard [20]. In certain instances, PCIe provides stronger guarantees than the RDMA specification, which we will always assume

<table><tr><td>$x=0$</td><td>$z=0$</td></tr><tr><td>$x := 1$ $z^2 := x$</td><td></td></tr></table>	$x=0$	$z=0$	$x := 1$ $z^2 := x$		<table><tr><td>$x=0$</td><td>$z=0$</td></tr><tr><td>$z^2 := x$ $x := 1$</td><td></td></tr></table>	$x=0$	$z=0$	$z^2 := x$ $x := 1$		<table><tr><td>$x=0$</td><td>$z=0$</td></tr><tr><td>$z^2 := x$ poll(2) $x := 1$</td><td></td></tr></table>	$x=0$	$z=0$	$z^2 := x$ poll (2) $x := 1$		<table><tr><td>$x=0$</td><td>$z=0$</td></tr><tr><td>$z^2 := x$ $z^2 := x$ poll(2) $x := 1$</td><td></td></tr></table>	$x=0$	$z=0$	$z^2 := x$ $z^2 := x$ poll (2) $x := 1$		<table><tr><td>$x=0$</td><td>$z=1$</td><td>$y=2$</td></tr><tr><td>$x := z^2$ $x := y^3$</td><td></td><td></td></tr></table>	$x=0$	$z=1$	$y=2$	$x := z^2$ $x := y^3$		
$x=0$	$z=0$																									
$x := 1$ $z^2 := x$																										
$x=0$	$z=0$																									
$z^2 := x$ $x := 1$																										
$x=0$	$z=0$																									
$z^2 := x$ poll (2) $x := 1$																										
$x=0$	$z=0$																									
$z^2 := x$ $z^2 := x$ poll (2) $x := 1$																										
$x=0$	$z=1$	$y=2$																								
$x := z^2$ $x := y^3$																										
(a) $z=0$ ✗ $z=1$ ✓	(b) $z=0$ ✓ $z=1$ ✓	(c) $z=0$ ✓ $z=1$ ✗	(d) $z=0$ ✓ $z=1$ ✓	(e) $x=1$ ✓ $x=2$ ✓																						
<table><tr><td>$x=0$</td><td>$z=1$</td><td>$y=2$</td></tr><tr><td>$y^3 := x$ $x := z^2$</td><td></td><td></td></tr></table>	$x=0$	$z=1$	$y=2$	$y^3 := x$ $x := z^2$			<table><tr><td>$x=1$</td><td>$y=z=0$</td></tr><tr><td>$z^2 := x$ $x := y^2$</td><td></td></tr></table>	$x=1$	$y=z=0$	$z^2 := x$ $x := y^2$		<table><tr><td>$x=1$</td><td>$y=z=0$</td></tr><tr><td>$x := y^2$ $z^2 := x$</td><td></td></tr></table>	$x=1$	$y=z=0$	$x := y^2$ $z^2 := x$		<table><tr><td>$x=1$</td><td>$y=z=0$</td></tr><tr><td>$x := y^2$ rfence (2) $z^2 := x$</td><td></td></tr></table>	$x=1$	$y=z=0$	$x := y^2$ rfence (2) $z^2 := x$						
$x=0$	$z=1$	$y=2$																								
$y^3 := x$ $x := z^2$																										
$x=1$	$y=z=0$																									
$z^2 := x$ $x := y^2$																										
$x=1$	$y=z=0$																									
$x := y^2$ $z^2 := x$																										
$x=1$	$y=z=0$																									
$x := y^2$ rfence (2) $z^2 := x$																										
(f) $y=0$ ✓ $y=1$ ✓	(g) $z=0$ ✗ $z=1$ ✓	(h) $z=0$ ✓ $z=1$ ✓	(i) $z=0$ ✓ $z=1$ ✗																							

Figure 3.1: RDMA litmus tests showing reordering rules for a single node. Column (separated by ||) denote distinct nodes, and the top line of each column gives the initial values for shared memory locations on that node. Reproduced from [10] with permission.

$y=0$	$x=0$	$x=0$	$y=0$	$x=0$	$y=0$	$y=0$	$x=0$	$y=w=0$	$x=z=0$
$x^2 := 1$	$y^1 := 1$	$a := y^2$	$b := x^1$	$a := y^2$	$b := x^1$	$x^2 := 1$	$y^1 := 1$	$x^2 := 1$	$y^1 := 1$
$a := y$	$b := x$	$x := 1$	$y := 1$	poll (2)	poll (1)	poll (2)	poll (1)	$c := z^2$	$d := w^1$
				$x := 1$	$y := 1$	$a := y$	$b := x$	poll (2)	poll (1)
								poll (2)	poll (1)
								$a := y$	$b := x$
(a) $a=b=0$ ✓	(b) $a=b=1$ ✓	(c) $a=b=1$ ✗	(d) $a=b=0$ ✓	(e) $a=b=0$ ✗					

Figure 3.2: RDMA litmus tests for concurrency across multiple nodes. Reproduced from [10] with permission.

are available due to its ubiquity. The standard guarantees that a read cannot complete before a write which was issued earlier. Therefore, a NIC remote read (get) must flush all pending NIC remote writes (puts) to memory¹, so by polling a get which follows a put, we guarantee that the NIC remote write of the put has been committed to memory (Fig. 3.2e).

3.2 Remote RMW Behaviours

We now present an intuitive account of rRMW behaviours.

Atomicity. Remote RMWs are atomic only with respect to each other – not with respect to any other memory operations. This is in contrast to standard RMWs, which are atomic with respect to all other operations. We will need to consider that either CPU operations on the remote side, or other remote operations (puts and gets) from any node can result in non-atomic behaviour for the rRMW.

Reordering. Due to the two-tiered concurrency of RDMA, all remote operations may be reordered after a later local operation – this equally applies to rRMWs. From [10],

¹The converse is also true: the NIC local read of a put flushes pending NIC local writes due to gets. However this is not particularly useful, since poll is well behaved with respect to get.

we know that an earlier get can be reordered after a later put. We can say that get/put reordering is allowed. Other reorderings involving any combination of get and put are not allowed: get/get¹, put/get and put/put. Remote RMWs behave much like puts, so that get/rRMW is allowed, and all other reorderings are disallowed.

The litmus tests in Fig. 3.3 illustrate these behaviours. Figs. 3.3a to 3.3d show that if either a put or get reads or writes the location targeted by a remote CAS, non-atomic behaviour may be observed. In contrast, Fig. 3.3e shows that multiple rRMWs are atomic with respect to each other. Figs. 3.3f to 3.3i show that rRMWs never reorder with respect to puts or other rRMWs, and an earlier rRMW completes before a later get. Fig. 3.3j shows that an earlier get may be reordered after a later rRMW.

Discussion: Remote RMW Behaviours

Early in the project, it was believed that rRMWs had stronger atomicity guarantees, and weaker reordering behaviours.

Atomicity. We initially believed that rRMWs were atomic with respect to all other remote operations (i.e. puts and gets, in addition to other rRMWs). This confusion came about due to overloaded use of the term “atomic access” in the InfiniBand specification [21]. Section 9.4.5.1 reads:

Atomicity of the read/modify/write on the responder’s node by the ATOMIC Operation shall be assured in the presence of concurrent atomic accesses [by the NIC].

The correct interpretation of this quote is that rRMWs guarantee atomicity only with respect to rRMWs issued by other threads towards the same node – that is to say, we take “atomic access” to mean an RDMA remote read-modify-write. However, the term may also be used in the context of a simple read or write being atomic, meaning that while a memory location is being written to, it cannot be read – for example, see (non-)atomic accesses in C/C++ [22]. Without this guarantee, it is possible to read a memory location as containing a value which was never written to it, even after the memory location has been initialised. This “atomicity” is a far more ubiquitous guarantee, and puts and gets would both be considered “atomic accesses” in this sense. Fortunately, some other parts of the document help to clarify this: section 10.6.7.2 mentions that

AtomicOps are not required by the architecture to be atomic with respect to RDMA writes

which reinforces that they are only atomic with respect to other rRMWs².

Reordering. We also believed that an earlier rRMW could be reordered *after* a later rRMW. This was due to Table 76, which shows which pairs of operations require a fence to guarantee in-order semantics. It shows that when the first and second operations are both RDMA Atomics, a fence is required. This seems to imply that without a fence, a pair of rRMWs may execute out of order. In fact, such behaviour is not possible – section 9.4.5.2 reads:

¹Technically, a pair of gets may be partially reordered. They may read remote memory out of order, but must write to local memory in order.

²The term ‘AtomicOp’ in the quote refers to an RDMA remote RMW.

$x = 0$		$y = 1$
$x := 2$ $y^3 := x$	$a := \text{nCAS}(x^1, 0, 3)$	

(a) $(x, y) = (3, 2)$ ✓

		$x = 0$
$x^3 := 1$	$a := \text{nCAS}(x^3, 0, 2)$	

(b) $x = 2$ ✓

		$x = 0$
$a := x^3$	$b := \text{nCAS}(x^3, 0, 2)$	$x := 1$

(c) $(a, x) = (1, 2)$ ✓

$x = 0$		$y = 1$
$x := y^3$	$a := \text{nCAS}(x^1, 0, 2)$	

(d) $x = 2$ ✓

		$x = 0$
$a := \text{nFAA}(x^3, 1)$	$b := \text{nCAS}(x^3, 0, 2)$	

(e) $x = 1$ ✓ $x = 3$ ✓ $x = 2$ ✗

	$x = 0$
$a := \text{nCAS}(x^2, 1, 2)$ $b := \text{nCAS}(x^2, 0, 1)$	

(f) $x = 2$ ✗

	$x = 0$
$a := \text{nCAS}(x^2, 1, 2)$ $x^2 := 1$	

(g) $x = 2$ ✗

	$x = 0$
$x^2 := 2$ $a := \text{nCAS}(x^2, 0, 1)$	$b := x$

(h) $b = 1$ ✗

	$x = 0$
$a := \text{nCAS}(x^2, 0, 1)$ $b := x^2$	

(i) $b = 0$ ✗

	$x = 0$
$a := x^2$ $b := \text{nCAS}(x^2, 0, 1)$	

(j) $a = 1$ ✓

Figure 3.3: Litmus tests illustrating atomicity and reordering of rRMWs with respect to other remote operations.

the ATOMIC Operation only accesses memory after prior (non-RDMA READ) requests access memory and before subsequent requests access memory

which clearly enforces that rRMWs are in-order with respect to each other.

It may be that the apparently erroneous table entry is due to completion of an rRMW being defined not with respect to the rRMW accessing memory, but rather the response packet it returns. This is described in section 9.4.5.2:

For the requestor [the thread issuing the rRMW], an ATOMIC Operation is considered complete when the response packet returns.

So if we wish for a later rRMW to wait for the response of an earlier rRMW, we indeed require a fence, however the order of memory accesses is maintained regardless.

4 Operational Semantics

At its core, an operational semantics consists of a logical model of the state of a machine or system, and a set of rules describing how that system may transition from one state to the next. The transition rules can be composed into *derivations*, which describe the execution of a program. If we can build a derivation tree for a given execution, then we know the execution is possible for our system, because it is composed from allowed transitions.

In our case, the system is a network of machines (nodes) communicating using RDMA, and the states we are interested in are (a) the state of the program, i.e. which commands have been executed, and which remain to be executed, for each program thread; (b) the state of memory on each node; (c) the state of the TSO store buffer of each thread; and (d) the state of any network communications. Under RDMA, each thread has access to a *queue pair* for communicating with each other node in the network, so we are interested in the state of the queue pairs of each thread.

4.1 States of the Operational Semantics

Throughout, we follow and extend the notation of [10].

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

4.1.1 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 Loc_n for the set of locations local to node n , and $\text{Loc} = \bigsqcup_n \text{Loc}_n$ for the set of all locations. We use $\text{Loc}_{\bar{n}} = \text{Loc} \setminus \text{Loc}_n$ and write x^n, y^n, z^n for values in Loc_n , respectively $x^{\bar{n}}, y^{\bar{n}}, z^{\bar{n}}$ for $\text{Loc}_{\bar{n}}$. When the node in question is sufficiently clear, we elide the superscript and instead simply write x or \bar{x} for local or remote locations respectively.

4.1.2 Program State

Values and Expressions. The language of expressions is standard and elided. We write $v \in \text{Val}$ for values, with $\mathbb{N} \subseteq \text{Val}$, and $e \in \text{Exp}$ for expressions. We write $\text{elocs}(e)$ 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 $\llbracket e \rrbracket$ for the evaluation of e given it is *closed*, that is, $\text{elocs}(e) = \emptyset$. We use e^n for expressions where $\text{elocs}(e^n) \subseteq \text{Loc}_n$.

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-set, and poll, which awaits the earliest completion notification of a remote operations towards \bar{n} .

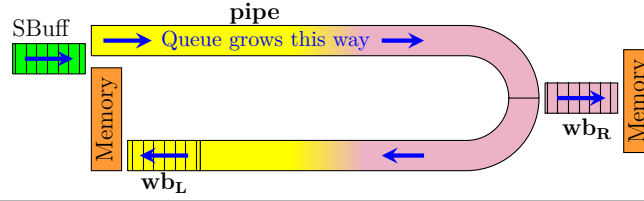


Figure 4.1: Simple queue-pair structure.

RDMA operations (RComm) are either a ‘get’ of the form $x := \bar{y}$ which reads a remote location \bar{y} and writes its value to local location x , a ‘put’ ($\bar{y} := x$) which does the reverse, ‘remote-CAS’¹ (resp. ‘remote-FAA’) which executes a *remote atomic* compare-and-set (resp. fetch-and-add), and ‘remote fence’ which ensures all prior RDMA operations towards \bar{n} complete before any later RDMA operations towards \bar{n} execute.

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 $t \in \text{Tid}$ is mapped to a command on $n(t)$.

$$\begin{aligned} \text{Comm} \ni C^n &::= \text{skip} \mid c^n \mid C_1^n, C_2^n \mid C_1^n + C_2^n \mid C^{n*} & \text{PComm} \ni c^n &::= \text{cc}^n \mid \text{rc}^n \\ \text{CComm} \ni \text{cc}^n &::= x := e^n \mid \text{assume}(x = v) \mid \text{assume}(x \neq v) \mid \text{mfence} \mid x := \text{CAS}(y, e_1, e_2) \mid \text{poll}(\bar{n}) \\ \text{RComm} \ni \text{rc}^n &::= x := \bar{y} \mid \bar{y} := x \mid x := \text{nCAS}(\bar{y}, e_1, e_2) \mid x := \text{nFAA}(\bar{y}, e) \mid \text{rfence}(\bar{n}) \end{aligned}$$

4.1.3 Store Buffers

To implement the weak behaviour of TSO described in §2.3, 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.

4.1.4 RDMA Operations

Queue Pairs. We follow the *simplified* operational model described in [10], and therefore consider a queue-pair structure comprising three FIFO queues: **pipe**, which contains pending or in-progress RDMA operations; **wb_R**, the remote write buffer, which contains pending writes to the memory of the remote node; and **wb_L**, the local write buffer, which contains pending writes to the memory of the local node. The structure is shown in Fig. 4.1. Notice that under this simplified model, the transition between local and remote node in **pipe** is continuous – we do not explicitly model the transition between local (yellow) and remote (pink) sides.

¹Note, we use the prefix **n** for ‘NIC’, rather than **r** for ‘remote’, for consistency with the notation used later in §5.

$$\begin{array}{c}
\frac{P \xrightarrow{t:\epsilon} P'}{P, M, B, A, QP \Rightarrow P', M, B, A, QP} \quad \frac{M, B, A, QP \xrightarrow{t:\epsilon} M', B', A', QP'}{P, M, B, A, QP \Rightarrow P, M', B', A', QP'} \\
\\
\frac{P \xrightarrow{t:l} P' \quad M, B, A, QP \xrightarrow{t:l} M', B', A', QP'}{P, M, B, A, QP \Rightarrow P', M', B', A', QP'}
\end{array}$$

Figure 4.2: RDMA^{TSO} operational semantics with the program and hardware transitions given in Fig. 4.3 and Fig. 4.4

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.

Discussion: States of the Operational Semantics

Extending the language of commands to accommodate the new remote atomic operations is straightforward – we just need to include the new RDMA commands. We also need to introduce a remote atomic lock for each node, which will be used in the relevant transition rules to guarantee atomicity. Without this structure, it would be necessary to check all queue-pairs towards a particular \bar{n} for in-progress rRMWs, which would be semantically equivalent but much more verbose. See §4.3.2 for more discussion of the structure A as it relates to the transition rules specifically.

4.2 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 ϕ is true, then it is allowed for the system to transition from the state described by $P \dots QP$ to the state described by $P' \dots QP'$.

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 $l \in \mathbf{Lab}$ is the label of the operation. *Silent transition*, which affect only the program (resp. only hardware) are written with the empty label, ϵ , and may be taken independently.

Fig. 4.2 shows the top-level rules of the operational semantics which govern this separation. We can henceforth consider the program and hardware transitions separately.

4.2.1 Program Transitions

Fig. 4.3 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, highlighted in yellow. 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. 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.

4.3 Hardware Domains

The upper section of Fig. 4.4 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, \perp (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**, **wb_R** and **wb_L**, where **pipe** may contain any of the operations described below except for a confirmation notification, **wb_R** may contain NIC remote writes and NIC remote atomic writes, and **wb_L** may contain NIC local writes and confirmation notifications.

- $y^{\bar{n}} := x^n$ denotes a put operation where the value of local memory location x is yet to be read (NIC local read);
- $y^{\bar{n}} := v$ denotes a NIC remote write of value v to remote location y , which occurs as the latter part of a put;
- **ack_p** denotes the acknowledgement message returned by a put;
- $x^n := y^{\bar{n}}$ 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;
- **nCAS**($z^n, x^{\bar{n}}, v, v'$) denotes a remote (NIC) CAS towards remote location x , with expected value v , update value v' , and returning to local location z ;
- **nFAA**($z^n, x^{\bar{n}}, v$) similarly denotes a NIC FAA towards x and returning to z , with increment value v ;
- $y^{\bar{n}} :=_A v$ 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 in §4.3.2;

- **rfence** \bar{n} denotes a remote fence towards node \bar{n} ;
- **cn** denotes a confirmation of a successful NIC remote write.

4.3.1 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 a remote CAS or FAA, the hardware takes a transition with a matching label, which adds that operation to the store buffer – once again these are highlighted in yellow, in the lower part of Fig. 4.4. 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} ². 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.

4.3.2 Queue-Pair Transitions

The queue-pair transitions model most of the interesting behaviours discussed in §3.2. The new transition rules for rRMWs are highlighted in yellow. From Fig. 4.1, 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 **wb_L** or **wb_R** (or both). Note that the pipe grows to the left, so throughout, α contains operations which are later in the program, while β 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 **sqp.pipe** = $\alpha \cdot (\text{operation})$, or allows it to be in the middle of the pipe **sqp.pipe** = $\alpha \cdot (\text{operation}) \cdot \beta$, with some stipulation as to the operations allowed in β . In the prior case, the operation never executes before another, earlier operation; in the latter, it can execute before any operation in β which was issued before it. There may also be some requirement that buffers **wb_L** or **wb_R** contain no writes, due to PCIe guarantees: **wb_L** $\in \{\text{cn}\}^*$ (**wb_L** contains only confirmation notifications) or **wb_R** = ϵ (there are no operations in **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 **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 rfence may be removed from the pipe once it reaches the head (there are no earlier operations remaining to be processed). In combination

²SQP 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 [10]

Program transitions: $\text{Prog} \xrightarrow{\text{Tid}:\text{Lab}\Psi\{\varepsilon\}} \text{Prog}$	Command transitions: $\text{Comm} \xrightarrow{\text{Lab}\Psi\{\varepsilon\}} \text{Comm}$
$\text{Lab} \triangleq \bigcup_n \text{Lab}_n \quad l \in \text{Lab}_n \triangleq \left\{ \begin{array}{l} \text{lw}(x^n, v), \text{lr}(x^n, v), \text{CASS}(x^n, v_1, v_2), \text{CASF}(x^n, v), \\ \text{F}, \text{P}(\bar{n}), \text{Get}(x^n, \bar{y}^n), \text{Put}(\bar{y}^n, x^n), \text{rF}(\bar{n}), \\ \text{nCAS}(y, x^{\bar{n}}, v_1, v_2), \text{nFAA}(y, x^{\bar{n}}, v) \end{array} \mid \begin{array}{l} x, y \in \text{Loc}, \\ v, v_1, v_2 \in \text{Val} \end{array} \right\}$	
<hr/> $\frac{C_1 \xrightarrow{l} C'_1}{C_1; C_2 \xrightarrow{l} C'_1; C_2} \quad \frac{}{\text{skip}; C \xrightarrow{\varepsilon} C} \quad \frac{i \in \{1, 2\}}{C_1 + C_2 \xrightarrow{\varepsilon} C_i} \quad \frac{}{C^* \xrightarrow{\varepsilon} \text{skip}} \quad \frac{}{C^* \xrightarrow{\varepsilon} C; C^*}$ $\frac{C \rightsquigarrow C'}{C \xrightarrow{\varepsilon} C'} \quad \frac{\text{elocs}(e) = \emptyset}{x := e \xrightarrow{\text{lw}(x, \llbracket e \rrbracket)} \text{skip}} \quad \frac{\text{elocs}(e_{\text{old}}) = \text{elocs}(e_{\text{new}}) = \emptyset \quad v \neq \llbracket e_{\text{old}} \rrbracket}{z := \text{CAS}(x, e_{\text{old}}, e_{\text{new}}) \xrightarrow{\text{CASF}(x, v)} z := v}$ $\frac{\text{elocs}(e_{\text{old}}) = \text{elocs}(e_{\text{new}}) = \emptyset}{z := \text{CAS}(x, e_{\text{old}}, e_{\text{new}}) \xrightarrow{\text{CASS}(x, \llbracket e_{\text{old}} \rrbracket, \llbracket e_{\text{new}} \rrbracket)} z := \llbracket e_{\text{old}} \rrbracket} \quad \frac{}{\text{mfence} \xrightarrow{\text{F}} \text{skip}}$ $\frac{}{x := \bar{y} \xrightarrow{\text{Get}(x, \bar{y})} \text{skip}} \quad \frac{}{\bar{y} := x \xrightarrow{\text{Put}(\bar{y}, x)} \text{skip}}$ <div style="background-color: #ffffcc; padding: 10px; margin: 10px 0;"> $\frac{\text{elocs}(e_{\text{old}}) = \text{elocs}(e_{\text{new}}) = \emptyset \quad v = \llbracket e_{\text{old}} \rrbracket \quad v' = \llbracket e_{\text{new}} \rrbracket}{z := \text{nCAS}(\bar{x}, e_{\text{old}}, e_{\text{new}}) \xrightarrow{\text{nCAS}(z, \bar{x}, v, v')} \text{skip}}$ </div> <div style="background-color: #ffffcc; padding: 10px; margin: 10px 0;"> $\frac{\text{elocs}(e) = \emptyset \quad v = \llbracket e \rrbracket}{z := \text{nFAA}(\bar{x}, e) \xrightarrow{\text{nFAA}(z, \bar{x}, v)} \text{skip}}$ </div> $\frac{}{\text{rfence } \bar{n} \xrightarrow{\text{rF}(\bar{n})} \text{skip}} \quad \frac{}{\text{poll}(\bar{n}) \xrightarrow{\text{P}(\bar{n})} \text{skip}}$ $\frac{}{\text{assume}(x = v) \xrightarrow{\text{lr}(x, v)} \text{skip}} \quad \frac{v \neq v'}{\text{assume}(x \neq v') \xrightarrow{\text{lr}(x, v)} \text{skip}} \quad \frac{\text{P}(t) \xrightarrow{l} C}{\text{P} \xrightarrow{t:l} \text{P}[t \mapsto C]}$ <hr/> $x := e \rightsquigarrow \text{assume}(y = v); x := e[v/y] \quad \text{for } y \in \text{elocs}(e), v \in \text{Val}$ $z := \text{CAS}(x, e_{\text{old}}, e_{\text{new}}) \rightsquigarrow \text{assume}(y = v); z := \text{CAS}(x, e_{\text{old}}[v/y], e_{\text{new}}) \quad \text{for } y \in \text{elocs}(e_{\text{old}}), v \in \text{Val}$ $z := \text{CAS}(x, e_{\text{old}}, e_{\text{new}}) \rightsquigarrow \text{assume}(y = v); z := \text{CAS}(x, e_{\text{old}}, e_{\text{new}}[v/y]) \quad \text{for } y \in \text{elocs}(e_{\text{new}}), v \in \text{Val}$ $z := \text{nCAS}(\bar{x}, e_{\text{old}}, e_{\text{new}}) \rightsquigarrow \text{assume}(y = v); z := \text{nCAS}(\bar{x}, e_{\text{old}}[v/y], e_{\text{new}}) \quad \text{for } y \in \text{elocs}(e_{\text{old}}), v \in \text{Val}$ $z := \text{nCAS}(\bar{x}, e_{\text{old}}, e_{\text{new}}) \rightsquigarrow \text{assume}(y = v); z := \text{nCAS}(\bar{x}, e_{\text{old}}, e_{\text{new}}[v/y]) \quad \text{for } y \in \text{elocs}(e_{\text{new}}), v \in \text{Val}$ $z := \text{nFAA}(\bar{x}, e) \rightsquigarrow \text{assume}(y = v); z := \text{nFAA}(\bar{x}, e[v/y]) \quad \text{for } y \in \text{elocs}(e), v \in \text{Val}$ <hr/>	

Figure 4.3: The RDMA^{TSO} program and command transitions

$$\begin{aligned}
& \mathbf{M} \in \mathbf{Mem} \triangleq \mathbf{Loc} \rightarrow \mathbf{Val} & \mathbf{B} \in \mathbf{SBMap} \triangleq \lambda t \in \mathbf{Tid}. \mathbf{SBuff}_{n(t)} \\
& \mathbf{A} \in \mathbf{RAMap} \triangleq \lambda n. \{ \perp, \top \} & \mathbf{QP} \in \mathbf{SQPMap} \triangleq \lambda t. (\lambda \bar{n}(t). \mathbf{SQPair}_{\bar{n}}^{\bar{n}}) \\
& \mathbf{b} \in \mathbf{SBuff}_n \triangleq \{ x^n := v, y^{\bar{n}} := x^n, x^n := y^{\bar{n}}, \mathbf{nCAS}(z^n, x^{\bar{n}}, v, v'), \mathbf{nFAA}(z^n, x^{\bar{n}}, v), \mathbf{rfence} \bar{n} \}^* \\
& \mathbf{sqp} \in \mathbf{SQPair}_{\bar{n}}^{\bar{n}} \triangleq \mathbf{Pipe}_{\bar{n}}^{\bar{n}} \times \mathbf{WBR}_{\bar{n}}^{\bar{n}} \times \mathbf{WBL}_{\bar{n}}^{\bar{n}} \\
& \mathbf{wb}_L \in \mathbf{WBL}_{\bar{n}}^{\bar{n}} \triangleq \{ \mathbf{cn}, x^n := v \}^* & \mathbf{wb}_R \in \mathbf{WBR}_{\bar{n}}^{\bar{n}} \triangleq \{ y^{\bar{n}} := v, y^{\bar{n}} :=_A v \}^* \\
& \mathbf{pipe} \in \mathbf{Pipe}_{\bar{n}}^{\bar{n}} \triangleq \left\{ \begin{array}{l} y^{\bar{n}} := x^n, y^{\bar{n}} := v, y^{\bar{n}} :=_A v, \mathbf{ack}_p, x^n := y^{\bar{n}}, x^n := v, \\ \mathbf{nCAS}(z^n, x^{\bar{n}}, v, v'), \mathbf{nFAA}(z^n, x^{\bar{n}}, v), \mathbf{rfence} \bar{n} \end{array} \right\}^*
\end{aligned}$$

$$\begin{aligned}
& \frac{\mathbf{B}' = \mathbf{B}[t \mapsto (x := v) \cdot \mathbf{B}(t)]}{\mathbf{M}, \mathbf{B}, \mathbf{A}, \mathbf{QP} \xrightarrow{t:\mathbf{W}(x,v)} \mathbf{M}, \mathbf{B}', \mathbf{A}, \mathbf{QP}} & \frac{(\mathbf{M} \triangleleft \mathbf{B}(t))(x) = v}{\mathbf{M}, \mathbf{B}, \mathbf{A}, \mathbf{QP} \xrightarrow{t:\mathbf{LR}(x,v)} \mathbf{M}, \mathbf{B}, \mathbf{A}, \mathbf{QP}} & \frac{\mathbf{B}(t) = \varepsilon \quad \mathbf{M}(x) = v_1}{\mathbf{M}, \mathbf{B}, \mathbf{A}, \mathbf{QP} \xrightarrow{t:\mathbf{CASS}(x,v_1,v_2)} \mathbf{M}[x \mapsto v_2], \mathbf{B}, \mathbf{A}, \mathbf{QP}} \\
& \frac{\mathbf{B}(t) = \varepsilon \quad \mathbf{M}(x) = v}{\mathbf{M}, \mathbf{B}, \mathbf{A}, \mathbf{QP} \xrightarrow{t:\mathbf{CASF}(x,v)} \mathbf{M}, \mathbf{B}, \mathbf{A}, \mathbf{QP}} & \frac{\mathbf{B}(t) = \varepsilon}{\mathbf{M}, \mathbf{B}, \mathbf{A}, \mathbf{QP} \xrightarrow{t:\mathbf{F}} \mathbf{M}, \mathbf{B}, \mathbf{A}, \mathbf{QP}} & \frac{\mathbf{B}(t) = \mathbf{b} \cdot (x := v)}{\mathbf{M}, \mathbf{B}, \mathbf{A}, \mathbf{QP} \xrightarrow{t:\varepsilon} \mathbf{M}[x \mapsto v], \mathbf{B}[t \mapsto \mathbf{b}], \mathbf{A}, \mathbf{QP}} \\
& \frac{\mathbf{B}(t) = \mathbf{b} \cdot \mathbf{rc}^n \quad \mathbf{rc}^n \in \{ x := y^{\bar{n}}, y^{\bar{n}} := x, \mathbf{nCAS}(z, \bar{x}, v, v'), \mathbf{nFAA}(z, \bar{x}, v), \mathbf{rfence} \bar{n} \}}{\mathbf{QP}(t)(\bar{n}) = \mathbf{sqp} \quad \mathbf{sqp}' = \mathbf{sqp}[\mathbf{pipe} \mapsto \mathbf{rc}^n \cdot \mathbf{sqp}.\mathbf{pipe}]} \\
& \frac{}{\mathbf{M}, \mathbf{B}, \mathbf{A}, \mathbf{QP} \xrightarrow{t:\varepsilon} \mathbf{M}, \mathbf{B}[t \mapsto \mathbf{b}], \mathbf{A}, \mathbf{QP}[t \mapsto \mathbf{QP}(t)[\bar{n} \mapsto \mathbf{sqp}']]} \\
& \frac{\mathbf{B}' = \mathbf{B}[t \mapsto (x := \bar{y}) \cdot \mathbf{B}(t)]}{\mathbf{M}, \mathbf{B}, \mathbf{A}, \mathbf{QP} \xrightarrow{t:\mathbf{Get}(x,\bar{y})} \mathbf{M}, \mathbf{B}', \mathbf{A}, \mathbf{QP}} & \frac{\mathbf{B}' = \mathbf{B}[t \mapsto (\bar{y} := x) \cdot \mathbf{B}(t)]}{\mathbf{M}, \mathbf{B}, \mathbf{A}, \mathbf{QP} \xrightarrow{t:\mathbf{Put}(\bar{y},x)} \mathbf{M}, \mathbf{B}', \mathbf{A}, \mathbf{QP}} \\
& \frac{\mathbf{B}' = \mathbf{B}[t \mapsto \mathbf{nCAS}(z, \bar{x}, v, v') \cdot \mathbf{B}(t)]}{\mathbf{M}, \mathbf{B}, \mathbf{A}, \mathbf{QP} \xrightarrow{t:\mathbf{nCAS}(z,\bar{x},v,v')} \mathbf{M}, \mathbf{B}', \mathbf{A}, \mathbf{QP}} & \frac{\mathbf{B}' = \mathbf{B}[t \mapsto \mathbf{nFAA}(z, \bar{x}, v) \cdot \mathbf{B}(t)]}{\mathbf{M}, \mathbf{B}, \mathbf{A}, \mathbf{QP} \xrightarrow{t:\mathbf{nFAA}(z,\bar{x},v)} \mathbf{M}, \mathbf{B}', \mathbf{A}, \mathbf{QP}} \\
& \frac{\mathbf{B}' = \mathbf{B}[t \mapsto (\mathbf{rfence} \bar{n}) \cdot \mathbf{B}(t)]}{\mathbf{M}, \mathbf{B}, \mathbf{A}, \mathbf{QP} \xrightarrow{t:\mathbf{RF}(\bar{n})} \mathbf{M}, \mathbf{B}', \mathbf{A}, \mathbf{QP}} \\
& \frac{\mathbf{QP}(t)(\bar{n}) = \mathbf{sqp} \quad \mathbf{sqp}.\mathbf{wb}_L = \alpha \cdot \mathbf{cn} \quad \mathbf{sqp}' = \mathbf{sqp}[\mathbf{wb}_L \mapsto \alpha]}{\mathbf{M}, \mathbf{B}, \mathbf{A}, \mathbf{QP} \xrightarrow{t:\mathbf{P}(\bar{n})} \mathbf{M}, \mathbf{B}, \mathbf{A}, \mathbf{QP}[t \mapsto \mathbf{QP}(t)[\bar{n} \mapsto \mathbf{sqp}']]} & \frac{\mathbf{M}, \mathbf{A}, \mathbf{QP}(t)(\bar{n}) \rightarrow_{\mathbf{sqp}} \mathbf{M}', \mathbf{A}', \mathbf{sqp} \quad (\text{Fig. 4.5})}{\mathbf{M}, \mathbf{B}, \mathbf{A}, \mathbf{QP} \xrightarrow{t:\varepsilon} \mathbf{M}', \mathbf{B}, \mathbf{A}', \mathbf{QP}[t \mapsto \mathbf{QP}(t)[\bar{n} \mapsto \mathbf{sqp}']]}
\end{aligned}$$

$$\text{with } (\mathbf{M} \triangleleft \alpha)(x) \triangleq \begin{cases} v & \text{if } \alpha = \beta \cdot (x := v) \cdot - \wedge \forall v'. x := v' \notin \beta \\ \mathbf{M}(x) & \text{if } \forall v. x := v \notin \alpha \end{cases}$$

Figure 4.4: RDMA^{TSO} simplified hardware domains (above) and hardware transitions (below)

with the fact that no other transition rule allows a step to be taken when there is an rfence later in the pipe, this enforces the behaviour that all remote operations prior to an rfence 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 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 \bar{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, highlighted in yellow. 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 $A(n(\bar{x})) = \perp$ (the remote atomic lock for the target node is available), and $A' = A[n(\bar{x}) \mapsto \top]$ (update the remote atomic lock for the target node to indicate it is busy). We also have distinct rules for success and failure of nCAS, depending on whether the remote memory location holds the expected value ($M(\bar{x}) = v$ or $M(\bar{x}) \neq v$).

The rules can then be interpreted as follows:

(Rule 9) A failed nCAS 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 \bar{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 nCAS 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 nFAA – this is unconditionally successful. It is very similar to a successful nCAS, but the value for the NIC remote atomic write is calculated by adding v to the value of \bar{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.

$$\begin{array}{c}
\frac{\text{sqp.pipe} = \alpha \cdot (\text{rfence } \bar{n})}{M, A, \text{sqp} \rightarrow_{\text{sqp}} M, A, \text{sqp}[\text{pipe} \mapsto \alpha]} \\
\\
\frac{\text{sqp.pipe} = \alpha \cdot (\bar{y} := x) \cdot \beta \quad \text{wb}_L \in \{\text{cn}\}^* \quad \beta \in \{y' := v', y' :=_A v', \bar{y}' := v', x' := \bar{y}', \text{nCAS}(z, \bar{x}, v, v'), \text{nFAA}(z, \bar{x}, v), \text{ack}_p\}^*}{M, A, \text{sqp} \rightarrow_{\text{sqp}} M, A, \text{sqp}[\text{pipe} \mapsto \alpha \cdot (\bar{y} := M(x)) \cdot \beta]} \\
\\
\frac{\text{sqp.pipe} = \alpha \cdot (\bar{y} := v) \cdot \beta \quad \beta \in \{x' := \bar{y}', x' := v', \text{ack}_p\}^* \quad \text{sqp}' = \text{sqp}[\text{pipe} \mapsto \alpha \cdot \text{ack}_p \cdot \beta][\text{wb}_R \mapsto (\bar{y} := v) \cdot \text{sqp.wb}_R]}{M, A, \text{sqp} \rightarrow_{\text{sqp}} M, A, \text{sqp}'} \\
\\
\frac{\text{sqp.wb}_R = \alpha \cdot (\bar{y} := v) \quad M, A, \text{sqp} \rightarrow_{\text{sqp}} M[\bar{y} \mapsto v], A, \text{sqp}[\text{wb}_R \mapsto \alpha]}{\text{sqp}' = \text{sqp}[\text{pipe} \mapsto \alpha][\text{wb}_L \mapsto \text{cn} \cdot \text{sqp.wb}_L]} \quad \frac{\text{sqp.pipe} = \alpha \cdot \text{ack}_p \quad \text{sqp}' = \text{sqp}[\text{pipe} \mapsto \alpha][\text{wb}_L \mapsto \text{cn} \cdot \text{sqp.wb}_L]}{M, A, \text{sqp} \rightarrow_{\text{sqp}} M, A, \text{sqp}'} \\
\\
\frac{\text{sqp.pipe} = \alpha \cdot (x := \bar{y}) \cdot \beta \quad \beta \in \{x' := \bar{y}', x' := v', \text{ack}_p\}^* \quad \text{sqp.wb}_R = \epsilon \quad \text{sqp}' = \text{sqp}[\text{pipe} \mapsto \alpha \cdot (x := M(\bar{y})) \cdot \beta]}{M, A, \text{sqp} \rightarrow_{\text{sqp}} M, A, \text{sqp}'} \\
\\
\frac{\text{sqp.pipe} = \alpha \cdot (x := v) \quad \text{sqp}' = \text{sqp}[\text{pipe} \mapsto \alpha][\text{wb}_L \mapsto \text{cn} \cdot (x := v) \cdot \text{sqp.wb}_L]}{M, A, \text{sqp} \rightarrow_{\text{sqp}} M, A, \text{sqp}'} \quad \frac{\text{sqp.wb}_L = \alpha \cdot (x := v) \cdot \beta \quad \beta \in \{\text{cn}\}^* \quad \text{sqp}' = \text{sqp}[\text{wb}_L \mapsto \alpha \cdot \beta]}{M, A, \text{sqp} \rightarrow_{\text{sqp}} M[x \mapsto v], A, \text{sqp}'} \\
\\
\frac{\text{sqp.pipe} = \alpha \cdot \text{nCAS}(z, \bar{x}, v, v') \cdot \beta \quad \text{sqp.wb}_R = \epsilon \quad A(n(\bar{x})) = \perp \quad M(\bar{x}) \neq v \quad \beta \in \{x' := \bar{y}', x' := v', \text{ack}_p\}^* \quad \text{sqp}' = \text{sqp}[\text{pipe} \mapsto \alpha \cdot (z := M(\bar{x})) \cdot \beta]}{M, A, \text{sqp} \rightarrow_{\text{sqp}} M, A, \text{sqp}'} \\
\\
\frac{\text{sqp.pipe} = \alpha \cdot \text{nCAS}(z, \bar{x}, v, v') \cdot \beta \quad \text{sqp.wb}_R = \epsilon \quad M(\bar{x}) = v \quad \beta \in \{x' := \bar{y}', x' := v', \text{ack}_p\}^* \quad A(n(\bar{x})) = \perp \quad A' = A[n(\bar{x}) \mapsto \top]}{M, A, \text{sqp} \rightarrow_{\text{sqp}} M, A', \text{sqp}[\text{pipe} \mapsto \alpha \cdot (z := v) \cdot (\bar{x} :=_A v') \cdot \beta]} \\
\\
\frac{\text{sqp.pipe} = \alpha \cdot \text{nFAA}(z, \bar{x}, v) \cdot \beta \quad \text{sqp.wb}_R = \epsilon \quad M(\bar{x}) + v = v' \quad \beta \in \{x' := \bar{y}', x' := v', \text{ack}_p\}^* \quad A(n(\bar{x})) = \perp \quad A' = A[n(\bar{x}) \mapsto \top]}{M, A, \text{sqp} \rightarrow_{\text{sqp}} M, A', \text{sqp}[\text{pipe} \mapsto \alpha \cdot (z := v) \cdot (\bar{x} :=_A v') \cdot \beta]} \\
\\
\frac{\text{sqp.pipe} = \alpha \cdot (\bar{x} :=_A v) \cdot \beta \quad \text{wb}_R' = (\bar{x} :=_A v) \cdot \text{sqp.wb}_R \quad \beta \in \{x' := \bar{y}', x' := v', \text{ack}_p\}^* \quad \text{sqp}' = \text{sqp}[\text{pipe} \mapsto \alpha \cdot \beta][\text{wb}_R \mapsto \text{wb}_R']}{M, A, \text{sqp} \rightarrow_{\text{sqp}} M, A, \text{sqp}'} \quad \frac{\text{sqp.wb}_R = \alpha \cdot (\bar{x} :=_A v) \quad A' = A[n(\bar{x}) \mapsto \perp] \quad \text{sqp}' = \text{sqp}[\text{wb}_R \mapsto \alpha]}{M, A, \text{sqp} \rightarrow_{\text{sqp}} M[\bar{x} \mapsto v], A', \text{sqp}'}
\end{array}$$

Figure 4.5: Queue-pair transitions of the simplified RDMA^{TSO} operational semantics

$y=0$	$x=0$
$x^2 := 1$ $\text{poll}(2)$ $a := y$	$y^1 := 1$ $\text{poll}(1)$ $b := x$

(a) $(a, b) = (0, 0)$ ✓

$y=0$	$x=0$
$c := \text{nFAA}(x^2, 1)$ $\text{poll}(2)$ $a := y$	$d := \text{nFAA}(y^1, 1)$ $\text{poll}(1)$ $b := x$

(b) $(a, b) = (0, 0)$ ✓

Figure 4.6: Remote RMWs exhibit the same weak behaviours as puts, including store buffering even in the presence of polls.

Discussion: Queue-Pair Transitions

The atomicity requirement motivated the addition of **A**, however it previously took on a different form. Initially, rather than extending **pipe** and introducing **A**, only **A** was present. When an rRMW was in progress, the description of it would be held in **A**, rather than a boolean. This simplified the atomicity checks, however made reordering much more complicated, as other rules would need to check **A** as well as β for reordering purposes. Using a boolean instead, and maintaining the in-progress rRMW in the pipe, made the semantics much clearer.

It is also noteworthy that the local and remote write of an rRMW may be processed in either order. In particular, this means that local write may reach memory while the remote write is still in **wb_R**. The local write indicates the completion of the rRMW for the local thread, so this seems problematic: the rRMW is considered complete for the local side before it has finished on the remote side. In fact, this is because rRMWs allow all the same weak behaviours as puts, other than those excluded by the atomicity guarantee. As discussed in §3.1, some of these cannot be prevented by a poll as we would expect. See Fig. 4.6 which replicates the store buffering example from Fig. 3.2d, and shows the same behaviour with rRMWs.

5 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 $(w, r) \in \text{rf}$ 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 $(w, r) \in \text{rf}$ 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.

5.1 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 ι , is created by a thread $t \in \text{Tid}$, and has an event label $l \in \text{ELab}$ which describes the event.

Definition 1 (Labels and events). Each event label is associated to a node n . The set of *event labels of node n* is denoted by $l \in \text{ELab}_n$, where l is a tuple with one of the following forms:

- (CPU) local read: $l = \text{1R}(x^n, v_r)$
- (CPU) local write: $l = \text{1W}(x^n, v_w)$
- (CPU) CAS: $l = \text{CAS}(x^n, v_r, v_w)$
- (CPU) memory fence: $l = \text{F}$
- (CPU) poll: $l = \text{P}(\bar{n})$
- NIC local read: $l = \text{n1R}(x^n, v_r, \bar{n})$
- NIC remote write: $l = \text{nrW}(y^{\bar{n}}, v_w)$
- NIC remote read: $l = \text{nrR}(y^{\bar{n}}, v_r)$
- NIC local write: $l = \text{n1W}(x^n, v_w, \bar{n})$
- NIC fence: $l = \text{nF}(\bar{n})$
- NIC atomic remote read:
 $l = \text{narR}(y^{\bar{n}}, v_r)$
- NIC atomic remote write:
 $l = \text{narW}(y^{\bar{n}}, v_w)$

The set of event labels are defined $\text{ELab} \triangleq \bigcup_n \text{ELab}_n$.

An *event*, $e \in \text{Event}$, is a triple (ι, t, l) , where $\iota \in \mathbb{N}$, $t \in \text{Tid}$ and $l \in \text{ELab}_{n(t)}$.

We distinguish between events associated with the CPU (left) or NIC (right), with the prefix **n** used for all NIC event labels. Note that a put, get, or rRMW is modelled by multiple events: a put $\bar{x} := y$ comprises a NIC local read event of type **n1R** (on y) followed by a NIC remote write event **nrW** (on \bar{x}); conversely a get $x := \bar{y}$ comprises events of type **nrR** (on \bar{y}) and **n1W** (on x). A successful rRMW (either successful nCAS or nFAA) is modelled by three events of type **narR**, **narW** and **n1W**, while a failed rRMW (nCAS only) is modelled by only **narR** and **n1W**.

For a given label l , we write $\text{type}(l)$, $\text{loc}(l)$, $v_r(l)$, $v_w(l)$, $n(l)$ and $\bar{n}(l)$ for the type, location, value read or written, and local or remote node, where applicable. For example, consider $l = \text{n1R}(x^n, v_r, \bar{n})$:

- $\text{type}(\text{n1R}(x^n, v_r, \bar{n})) = \text{n1R}$
- $\text{loc}(\text{n1R}(x^n, v_r, \bar{n})) = x$
- $v_r(\text{n1R}(x^n, v_r, \bar{n})) = v_r$
- $v_w(\text{n1R}(x^n, v_r, \bar{n}))$ is undefined
- $n(\text{n1R}(x^n, v_r, \bar{n})) = n$
- $\bar{n}(\text{n1R}(x^n, v_r, \bar{n})) = \bar{n}$

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

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 [10]. 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
- $r_1; r_2 \triangleq \{(a, b) \mid \exists c. (a, c) \in r_1 \wedge (c, b) \in r_2\}$ 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:

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

Finally, we define the following relations:

- Same-location:** $\text{sloc} \triangleq \{(\mathbf{e}, \mathbf{e}') \in \text{Event}^2 \mid \text{loc}(\mathbf{e}) = \text{loc}(\mathbf{e}')\}$
- Same-thread:** $\text{sthd} \triangleq \{(\mathbf{e}, \mathbf{e}') \in \text{Event}^2 \mid t(\mathbf{e}) = t(\mathbf{e}')\}$
- Same-queue-pair:** $\text{sqp} \triangleq \{(\mathbf{e}, \mathbf{e}') \in \text{Event}^2 \mid t(\mathbf{e}) = t(\mathbf{e}') \wedge \bar{n}(\mathbf{e}) = \bar{n}(\mathbf{e}')\}$

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

		$x = y = 0$
$a := \text{nCAS}(x^3, 0, 7)$	$b := \text{nCAS}(y^3, 0, 42)$	$x := 1$ $y := 2$ $y := 3$ $x := 4$

(a) $(x, b) = (7, 2)$ ✗

Figure 5.1: Remote RMWs are guaranteed to be mutually exclusive whenever they are towards the same node, even if they target different locations.

Definition 2 (Pre-executions). A *pre-execution* is a tuple $G = \langle E, \text{po}, \text{rf}, \text{mo}, \text{pf}, \text{nfo}, \text{rao} \rangle$, 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 $1W(x, 0)$ for each $x \in \text{Loc}$.
- $\text{po} \subseteq E \times E$ is the ‘*program order*’ relation defined as a disjoint union of strict total orders, each ordering the events of one thread, with $E^0 \times (E \setminus E^0) \subseteq \text{po}$.
- $\text{rf} \subseteq E.\mathcal{W} \times E.\mathcal{R}$ is the ‘*reads-from*’ relation on events of the same location with matching values; i.e. $(a, b) \in \text{rf} \Rightarrow (a, b) \in \text{sloc} \wedge v_w(a) = v_r(b)$. Moreover, rf is total and functional on its range: every read in $E.\mathcal{R}$ is related to exactly one write in $E.\mathcal{W}$.
- $\text{mo} \triangleq \bigcup_{x \in \text{Loc}} \text{mo}_x$ is the ‘*modification-order*’, where each mo_x is a strict total order on $E.\mathcal{W}_x$ with $E_x^0 \times (E.\mathcal{W}_x \setminus E_x^0) \subseteq \text{mo}_x$ describing the order in which writes on x reach the memory.
- $\text{pf} \subseteq (E.\text{n1W} \cup E.\text{nrW}) \times E.\text{P}$ is the ‘*polls-from*’ relation, relating earlier (in program-order) NIC writes to later poll operations on the *same queue pair*; i.e. $\text{pf} \subseteq \text{po} \cap \text{sqp}$. 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 $E.\text{P}$ polls from exactly one NIC write).
- $\text{nfo} \subseteq E.\text{sqp}$ is the ‘*NIC flush order*’, such that for all $(a, b) \in E.\text{sqp}$, if $a \in E.\text{n1R}, b \in E.\text{n1W}$, then $(a, b) \in \text{nfo} \cup \text{nfo}^{-1}$, and if $a \in (E.\text{nrR} \cup E.\text{narR}), b \in (E.\text{nrW} \cup E.\text{narW})$, then $(a, b) \in \text{nfo} \cup \text{nfo}^{-1}$.
- $\text{rao} \triangleq \bigcup_{n \in \text{Node}} \text{rao}_n$ is the ‘*remote-atomic-order*’, where each rao_n is a strict total order on $\{e \mid e \in E.\text{narR} \wedge \bar{n}(e) = n\}$ 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 [10]. The reason nfo is important is subtle: recall from §3.1 that we model the requirement from PCIe that a NIC local read flushes pending NIC remote writes on the same queue pair, and likewise for NIC local reads/writes. We will therefore need to maintain an order over these events to check that the property holds.

We introduce rao , which totally orders NIC remote atomic reads towards a given node. This will help us to check the atomicity guarantee.

Discussion: Pre-executions

At first glance, it might appear that **rao** should be a total order *per-location*, rather than per-thread; indeed, it was so in an earlier draft of the semantics. However, this results in an atomicity guarantee which is too weak. Consider the example shown in Fig. 5.1. When we observe the outcome $x = 7$, this means that every local write on node 3 occurred between the read and write of the nCAS towards x : in particular, $y = 2$ is only ever true within this critical period, not outside it. Thus if $x = 7$ then we cannot have $b = 2$, as this would imply the two rRMWs were not mutually exclusive. If atomicity were only guaranteed per-location, then this would be allowed, but this is forbidden by the specification.

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

- $\text{rb} \triangleq (\text{rf}^{-1}; \text{mo}) \setminus [E]$ is the *reads-before* relation, relating each read r to writes that are **mo**-after the write from which r reads.
- $\text{rf}_b \triangleq [1W]; (\text{rf} \cap \text{sthd}); [1R]$ 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 $w \xrightarrow{\text{rf}_b} r$, 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.
- $\text{rf}_{\bar{b}} \triangleq \text{rf} \setminus \text{rf}_b$ is the rf_b -complement: if $w \xrightarrow{\text{rf}_{\bar{b}}} r$, then r only occurs after w is observable.
- $\text{rb}_b \triangleq [1R]; (\text{rb} \cap \text{sthd}); [1W]$ is the *rb-buffer* relation, analogously.
- $\text{ar} \triangleq [\text{narW}]; (\text{po}|_{\text{imm}}^{-1})$ is the *atomic-write-to-read* relation, connecting the remote write of an 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. In order to model the reordering rules of §3.1, we must 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 $\text{oppo} \subseteq \text{ippo} \subseteq \text{po}$. The tables in Fig. 5.2 show these relations, with the new cells due to rRMW events highlighted. Each row indicates the po -earlier event, while each column indicates that which is po -later. A cell labelled **✓** indicates the event pair is in **ippo** (resp. **oppo**) and must be issued (resp. observed) in program order, while **✗** 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

ippo			CPU					NIC						
			1	2	3	4	5	6	7	8	9	10	11	12
			1R	1W	CAS	F	P	nlR	nrW	narR	narW	nrR	nlW	nF
CPU	A	1R	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
	B	1W	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
	C	CAS	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
	D	F	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
	E	P	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
NIC	F	nlR	✗	✗	✗	✗	✗	sqp	sqp	sqp	sqp	sqp	sqp	sqp
	G	nrW	✗	✗	✗	✗	✗	✗	sqp	sqp	sqp	sqp	sqp	sqp
	H	narR	✗	✗	✗	✗	✗	✗	sqp	sqp	sqp	sqp	sqp	sqp
	I	narW	✗	✗	✗	✗	✗	✗	sqp	sqp	sqp	sqp	sqp	sqp
	J	nrR	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	sqp	sqp
	K	nlW	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	sqp	sqp
	L	nF	✗	✗	✗	✗	✗	sqp	sqp	sqp	sqp	sqp	sqp	sqp

oppo			CPU					NIC						
			1	2	3	4	5	6	7	8	9	10	11	12
			1R	1W	CAS	F	P	nlR	nrW	narR	narW	nrR	nlW	nF
CPU	A	1R	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
	B	1W	✗	✓	✓	✓	✗	✓	✓	✓	✓	✓	✓	✓
	C	CAS	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
	D	F	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
	E	P	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
NIC	F	nlR	✗	✗	✗	✗	✗	sqp	sqp	sqp	sqp	sqp	sqp	sqp
	G	nrW	✗	✗	✗	✗	✗	✗	sqp	sqp	sqp	sqp	sqp	✗
	H	narR	✗	✗	✗	✗	✗	✗	sqp	sqp	sqp	sqp	sqp	sqp
	I	narW	✗	✗	✗	✗	✗	✗	sqp	sqp	sqp	sqp	✗	✗
	J	nrR	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	sqp	sqp
	K	nlW	✗	✗	✗	✗	✗	✗	✗	✗	✗	✗	sqp	✗
	L	nF	✗	✗	✗	✗	✗	sqp	sqp	sqp	sqp	sqp	sqp	sqp

Figure 5.2: The RDMA^{TSO} ordering constraints on **ippo** (above) and **oppo** (below), where ✓ denotes that instructions are ordered (and cannot be reordered), ✗ denotes they are not ordered (and may be reordered), and sqp denotes they are ordered iff they are on the same queue pair. Highlighted: new cells due to rRMWs.

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 **wb_R** (resp. **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 **wb_L** and **wb_R** respectively and may be committed in either order.

Definition 3 (Executions). A pre-execution $G = \langle E, \text{po}, \text{rf}, \text{mo}, \text{pf}, \text{nfo}, \text{rao} \rangle$ 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-pollled remote operation on the same queue pair:
if $w_1 \in G.\text{nlW} \cup G.\text{nrW}$ and $w_1 \xrightarrow{\text{po} \cap \text{sqp}} w_2 \xrightarrow{\text{pf}} p_2$, then there exists p_1 such that $w_1 \xrightarrow{\text{pf}} p_1 \xrightarrow{\text{po}} p_2$.
2. Each put (resp. get) operation corresponds to two events: a read and a write with the read immediately preceding the write in po: (a) if $r \in G.\text{nlR}$ (resp. $r \in G.\text{nrR}$), then $(r, w) \in \text{po}|_{\text{imm}}$ for some $w \in G.\text{nrW}$ ($w \in G.\text{nlW}$); and (b) if $w \in G.\text{nrW}$ then $(r, w) \in \text{po}|_{\text{imm}}$ for some $r \in G.\text{nlR}$. The case $w \in G.\text{nlW}$ is handled by (6) below.
3. Read and write events of a put (resp. get) have matching values:
if $(r, w) \in G.\text{po}|_{\text{imm}}$, $\text{type}(r) \in \{\text{nlR}, \text{nrR}\}$ and $\text{type}(w) \in \{\text{nlW}, \text{nrW}\}$, 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: (a) if $r \in G.\text{narR}$ then $(r, w_1) \in \text{po}|_{\text{imm}}$ for some $w_1 \in G.\text{narW} \cup G.\text{nlW}$, and if $w_1 \in G.\text{narW}$ then $(w_1, w_2) \in \text{po}|_{\text{imm}}$ for some $w_2 \in G.\text{nlW}$, and (b) if $w_1 \in G.\text{narW}$ then $(r, w_1) \in \text{po}|_{\text{imm}}$ for some $r \in \text{narR}$, and $(w_1, w_2) \in \text{po}|_{\text{imm}}$ for some $w_2 \in \text{nlW}$. The case for $w_2 \in \text{nlW}$ is handled by (6) below.
5. Remote atomic read and local write events of an rRMW have matching values: if $(r, w) \in G.\text{po}|_{\text{imm}}$, $\text{type}(r) = \text{narR}$ and $\text{type}(w) = \text{nlW}$, then $v_r(r) = v_w(w)$; and if $(r, w_1), (w_1, w_2) \in G.\text{po}|_{\text{imm}}$, $\text{type}(r) = \text{narR}$, $\text{type}(w_1) = \text{narW}$ and $\text{type}(w_2) = \text{nlW}$, then $v_r(r) = v_w(w_2)$.
6. (2) and (4) auxiliary in the case of $w \in \text{nlW}$. If $w \in G.\text{nlW}$ then either:
 - (a) $(r, w) \in \text{po}|_{\text{imm}}$ for some $r \in G.\text{nrR}$ or
 - (b) $(r, w) \in \text{po}|_{\text{imm}}$ for some $r \in G.\text{narR}$ or
 - (c) $(r, w'), (w', w) \in \text{po}|_{\text{imm}}$ for some $r \in G.\text{narR}$ and $w' \in G.\text{narW}$.

An *execution* is a pre-execution (Def. 2) that is well-formed.

Given an execution G , we write $G.E$, $G.\text{mo}$, $G.\text{ippo}$ 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 4 (RDMA^{TSO} -consistency). An execution $\langle E, \text{po}, \text{rf}, \text{mo}, \text{pf}, \text{nfo}, \text{rao} \rangle$ is RDMA^{TSO} -consistent iff 1. ib is irreflexive; and 2. ob is irreflexive, where:

$$\begin{aligned} \text{ib} &\triangleq \left(\text{ippo} \cup \text{rf} \cup \text{pf} \cup \text{nfo} \cup \text{rb}_b \cup \text{rao} \cup (\text{ob}; [\text{Inst}]) \right)^+ && \text{('issued-before')} \\ \text{ob} &\triangleq \left(\text{oppo} \cup \text{rf}_b \cup ([\text{n1W}]; \text{pf}) \cup \text{nfo} \cup \text{rb} \cup \text{mo} \cup (\text{ar}; \text{rao}) \cup ([\text{Inst}]; \text{ib}) \right)^+ && \text{('observed-before')} \end{aligned}$$

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 $w \xrightarrow{\text{rf}} r$ 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 $w \xrightarrow{\text{pf}} p$ only if w was issued before p ; (c) nfo : NIC events arrive in wb_L/wb_R in the order they are issued; (d) rb_b : if $r \xrightarrow{\text{rb}_b} w$, then r must be issued before w , otherwise r would read from w or an mo -later w' ; (e) rao : remote atomic reads are issued in the defined order; (f) $\text{ob}; [\text{Inst}]$: in general, an event is observed no earlier than it is issued, and for an instantaneous event, the two points coincide. Thus $e \xrightarrow{\text{ob}} e'$ implies $e \xrightarrow{\text{ib}} e'$ when e' is instantaneous.

On the other hand, for ob we have (a) rf_b : if $w \xrightarrow{\text{rf}_b} r$ then w was committed to memory before r , since r cannot read from the store buffer of another thread; (b) $[\text{n1W}]; \text{pf}$: 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 wb_L/wb_R ; (d) rb : if $r \xrightarrow{\text{rb}} w$, then w was not observed before r , otherwise it would have been committed to memory before r ; (e) mo : if $w \xrightarrow{\text{mo}} w'$, then w was observed in memory before w' ; (f) $\text{ar}; \text{rao}$: if $w \xrightarrow{\text{ar}} r \xrightarrow{\text{rao}} r'$, 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. (g) $[\text{Inst}]; \text{ib}$: 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.

5.2 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 $\text{rf}, \text{mo}, \text{pf}, \text{nfo}, \text{rao}$ 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 $l_1, \dots, l_n \in \text{ELab}$, we define the *event graphs of t* as $(\{e_1, \dots, e_n\}, \text{po}) \in G^t(l_1, \dots, l_n)$ where: (a) $l(e_i) = l_i$ for all $1 \leq i \leq n$; (b) $\iota(e_i) \neq \iota(e_j)$ for all $1 \leq i < j \leq n$; (c) $t(e_i) = t$ for all $1 \leq i \leq n$; (d) $\text{po} = \{(e_i, e_j) \mid 1 \leq i < j \leq n\}$.

Initial Event Graph. Given a set of locations Loc , we define $G_{\text{init}} = (E_0, \emptyset)$, such that for each $x \in \text{Loc}$ there is exactly one $e \in E_0$ with $l(e) = 1W(x, 0)$, 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 $G_1; G_2 = (E, \text{po})$ where

$$\begin{aligned} E &\triangleq G_1.E \uplus G_2.E \\ \text{po} &\triangleq G_1.\text{po} \cup G_2.\text{po} \cup (G_1.E \times G_2.E) \end{aligned}$$

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 $G_1 \parallel G_2 = (E, \text{po})$ where

$$\begin{aligned} E &\triangleq G_1.E \uplus G_2.E \\ \text{po} &\triangleq G_1.\text{po} \cup G_2.\text{po} \end{aligned}$$

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 $G = G_{\text{init}}; (\parallel_{t \in \text{Tid}} G_t)$ and there is a set of sequences $s_t \in S$ such that $P(t) \mapsto s_t$ and $G_t \in G^t(s_t)$ for all $t \in \text{Tid}$.

The operation $C \mapsto s$ relates a sequential program C to a sequence of labels s it generates. The definition is standard and show in Fig. 5.3. Note that RDMA operations generate multiple events, and for local and remote CAS operations, we distinguish between success and failure cases.

Theorem 1. The operational and declarative semantics of RDMA^{TSO} are equivalent.

Proof. See §6

□

$$\begin{array}{c}
\frac{C \rightsquigarrow C' \quad C' \rightsquigarrow s}{C \rightsquigarrow s} \quad \frac{C_1 \rightsquigarrow s_1 \quad C_2 \rightsquigarrow s_2}{C_1; C_2 \rightsquigarrow s_1, s_2} \quad \frac{\text{elocs}(e) = \emptyset}{x := e \rightsquigarrow \text{lw}(x, \llbracket e \rrbracket)} \\
\\
\frac{\text{elocs}(e_{\text{old}}) = \text{elocs}(e_{\text{new}}) = \emptyset \quad z := \llbracket e_{\text{old}} \rrbracket \rightsquigarrow s}{z := \text{CAS}(x, e_{\text{old}}, e_{\text{new}}) \rightsquigarrow \text{CAS}(x, \llbracket e_{\text{old}} \rrbracket, \llbracket e_{\text{new}} \rrbracket), s} \\
\\
\frac{\text{elocs}(e_{\text{old}}) = \text{elocs}(e_{\text{new}}) = \emptyset \quad v \neq \llbracket e_{\text{old}} \rrbracket \quad z := v \rightsquigarrow s}{z := \text{CAS}(x, e_{\text{old}}, e_{\text{new}}) \rightsquigarrow \text{F}, \text{lr}(x, v), s} \quad \overline{\text{mfence} \rightsquigarrow \text{F}} \quad \overline{\text{assume}(x = v) \rightsquigarrow \text{lr}(x, v)} \\
\\
\overline{\text{assume}(x \neq v) \rightsquigarrow \text{lr}(x, v')} \quad \overline{x := y^{\bar{n}} \rightsquigarrow \text{nrR}(y^{\bar{n}}, v), \text{nlW}(x, v, \bar{n})} \\
\\
\overline{y^{\bar{n}} := x \rightsquigarrow \text{nlR}(x, v, \bar{n}), \text{nrW}(y^{\bar{n}}, v)} \quad \overline{\text{rfence}(\bar{n}) \rightsquigarrow \text{nF}(\bar{n})} \quad \overline{\text{poll}(\bar{n}) \rightsquigarrow \text{P}(\bar{n})} \\
\\
\overline{\text{skip} \rightsquigarrow \epsilon} \quad \frac{\text{elocs}(e_{\text{old}}) = \text{elocs}(e_{\text{new}}) = \emptyset \quad v \neq \llbracket e_{\text{old}} \rrbracket}{z := \text{nCAS}(x^{\bar{n}}, e_{\text{old}}, e_{\text{new}}) \rightsquigarrow \text{narR}(x^{\bar{n}}, v), \text{nlW}(z, v, \bar{n})} \\
\\
\frac{\text{elocs}(e_{\text{old}}) = \text{elocs}(e_{\text{new}}) = \emptyset}{z := \text{nCAS}(x^{\bar{n}}, e_{\text{old}}, e_{\text{new}}) \rightsquigarrow \text{narR}(x^{\bar{n}}, \llbracket e_{\text{old}} \rrbracket), \text{narW}(x^{\bar{n}}, \llbracket e_{\text{new}} \rrbracket), \text{nlW}(z, \llbracket e_{\text{old}} \rrbracket, \bar{n})} \\
\\
\frac{\text{elocs}(e) = \emptyset \quad v' = v + \llbracket e \rrbracket}{z := \text{nFAA}(x^{\bar{n}}, e) \rightsquigarrow \text{narR}(x^{\bar{n}}, v), \text{narW}(x^{\bar{n}}, v'), \text{nlW}(z, v, \bar{n})}
\end{array}$$

Figure 5.3: Label Sequences Construction

6 Equivalence

The proof of Theorem 1 can be found in Appendix B, which directly adapts the proof given in [10]. However, this proof is necessarily verbose, so this section provides an overview of the structure, highlights, and notable additions to the proof.

6.1 Structure of the Proof

The goal of the proof is to show that the operational and declarative models, given in §4 and §5 respectively, are equivalent.

By ‘equivalent’, we mean that given a program P , the set of executions of P given by either semantics each describe the same program behaviours, such that each execution derived by the operational semantics has a corresponding execution graph given by the declarative semantics, and vice-versa.

Therefore, we aim to show that (a) given an execution graph G generated by a program P , there is an equivalent derivation D of an execution of P ; and (b) given a derivation D of an execution of P , there is an equivalent execution graph G generated by P .

We note that a derivation D of an execution takes the form

$$P, M_0, B_0, A_0, QP_0 \Rightarrow^* \lambda t. \text{skip}, M, B_0, A_0, \lambda t. \lambda \bar{n}. \langle \epsilon, \epsilon, \text{cn}^* \rangle$$

where \Rightarrow^* is the reflexive transitive closure of \Rightarrow defined in Fig. 4.2.

Note that M_0 and similar denote initial states. A derivation must finish with the empty program (all threads contain only `skip`), all store buffers empty, all remote atomic locks released, and queue pairs containing no pending operations except for confirmation notifications, as it is allowed for these not to be polled.

However, it will be very challenging to convert between operational and declarative semantics directly, because the prior is concerned primarily with states while the latter only represents events. For the proof in the direction from operational to declarative (resp. declarative to operational), we would first need to derive the events that are implicit in the transitions between states (resp. the states implied by the events), then remove information about states (resp. events) that is no longer relevant.

In other words, for the proof in either direction, we first need to build a view that contains all the information of both semantics, then remove information to arrive at the appropriate form. There would be no reason not to use the same intermediary form for both directions. Therefore, we define an *annotated* semantics, as a direct extension of the operational semantics, where states are extended with event labels.

We then aim to prove that 1. the operational semantics is equivalent to the annotated semantics and 2. the declarative semantics is equivalent to the annotated semantics. Naturally, each point will require a proof in either direction.

Ensuring that the annotated semantics is defined as a direct extension of the operational one will make the proof of the two directions for the first point straightforward, while the latter point will be more complex.

6.2 Annotated Semantics

The annotated semantics is defined in §B.1. The states of the operational semantics are extended to contain events, rather than values. For example, a memory location

will contain a *particular* write event, rather than a mere value. This will help us to associate related events, in order to build the relations of the declarative semantics. The set of labels used is an extension of those defined for the declarative semantics, with the introduction of labels representing the operations put, get, nCAS and nFAA, and RDMA acknowledgement messages, as these are not explicitly modelled by labels in an execution graph, only by relations.

We then define annotated labels, which will be the labels for the transitions of the annotated semantics; each comprises a label type and tuple of events. The annotated labels explicitly connect transitions that concern the same event. Recall that for RDMA operations in the operational semantics, we expect to see a number of related transitions as it is processed into, through, and out of a queue-pair. In the annotated semantics, we directly connect those transitions by generating an event in the relevant program transition, then including that event in the label of each related hardware and queue-pair transition.

The transitions of the annotated semantics correspond *directly* to those of the operational semantics, but with events instead of values checked in most places. These transitions have been suitably extended to match the corresponding extensions to the operational semantics.

Paths. We extend the annotated semantics with a path π , which is a sequence of annotated labels. This serves as a history of transitions taken, as each time a labelled transition occurs, the label is added to π . This means that, given a particular state of the operational semantics, we are able to use π to determine every transition taken up to that point. We maintain a number of invariants to ensure that a path is well-formed, and to identify if it corresponds to a complete execution. A number of these invariants need to be updated due to the introduction of rRMWs, and we also introduce a new invariant to model the atomicity guarantee.

Path Conditions. A path is *complete* if there are no pending operations. We can check if an operation is complete by looking for all the relevant labels in the path. For example, if there was a put operation in the program, there should be labels in the path corresponding to the operation being pushed to the store-buffer, pushed to the queue-pair, a NIC local read, NIC remote write, acknowledgement, and committing the NIC remote write to memory.

A path is *backwards-complete* if labels do not appear out of thin air. For example, if there is a NIC remote write, there must be a corresponding prior NIC local read, since a NIC remote write can only appear as part of a put operation. Likewise there must be an earlier label pushing the operation to the queue pair, and so on.

A path is *well-formed* if it is backwards-complete and observes a number of ordering rules.

- Poll order: for any two writes on the same queue-pair, if the later write is polled, then the earlier write must be polled first;
- Flush order: any pair of writes that go via the same buffer (i.e. store buffer or queue-pair write buffer) must be committed to memory in the order they were buffered;
- NIC order: for a pair of operations which are pushed to the same queue-pair, describes the minimum progress achieved by the first operation, given the progress of the second operation;

- NIC atomicity: an earlier successful rRMW must have had its remote write committed to remote memory before a later rRMW reads;
- Read order: a read should observe the latest write that is visible to it, i.e. for NIC reads, the latest write in memory, and for CPU reads, the latest write in the store buffer, if there is one, or else the latest write in memory.

We then define well-formedness for **M**, **B**, **A**, and **QP** with respect to π , which requires that events in π are reflected in those states. For example, the write event at a location in **M** should match the latest write event to that location which has been flushed in π .

We then observe the result that well-formedness is preserved by the transitions of the annotated semantics.

6.3 Annotated Semantics to Declarative Semantics

For this part of the proof, we show that a well-formed and complete path π yields a consistent execution. We define $\text{getEG}(\pi)$ which finds E , po , rf , pf , mo , nfo and rao . This is straightforward, as by construction π provides an explicit total order for every event.

We then prove that $\text{getEG}(\pi)$ is a pre-execution and meets the conditions for well-formedness and consistency. Showing that it is a well-formed pre-execution is straightforward, while consistency is more interesting. The conceit of the proof is to show that each component of ib and ob is a subset of the total orders on issue and observation points given by π . The notable new cases are due to the extensions to ippo and oppo , and the new component of ob , ar ; rao . The proofs for the new cases of ippo and oppo rely heavily on the extended definition of NIC order, the cases of which broadly correspond to cells in the ippo/oppo tables. In the case of ar ; rao , the new invariant regarding NIC atomicity is used.

6.4 Declarative Semantics to Annotated Semantics

Now, we need to do the converse: given a consistent execution $G = \langle E, \text{po}, \text{rf}, \text{pf}, \text{mo}, \text{nfo}, \text{rao} \rangle$, we aim to construct a well-formed path π . Using π , it is then straightforward to reconstruct the rest of the annotated semantics derivation.

To build π , we extend ib and ob into total orders. We update this process to accommodate rRMWs, which requires some additional work not present in the prior version of the proof. The event graph does not contain information about failed nCAS operations, nor does it distinguish between a successful nCAS and an nFAA. Therefore, we need to look at the derivation from **P** to G , as described in §5.2. By checking the derivation, we can retrieve the information that was lost during the creation of the event graph.

We then need to show that each path invariant holds. The most interesting case is the proof of NIC atomicity. We argue that when a successful rRMW with read r_1 and write w_1 appears in the path, there is no other rRMW with read r_1 towards the same node which reads between r_1 and w_1 . Such a read would force $r_1 \xrightarrow{\text{rao}} r_2$, and thus $w_1 \xrightarrow{\text{ar}; \text{rao}} r_2$, which make ib and ob reflexive. This is not allowed by the consistency conditions.

6.5 Operational and Annotated Semantics

The annotated semantics was deliberately designed as a direct extension of the operational semantics in order to make this side of the proof much easier. The rules extending a derivation of the operational semantics to one of the annotated semantics are straightforward, as is the respective process for restricting an annotated semantics derivation to one of the operational semantics. We extend these definitions with rules in the case of rRMWs, which leads to the result we require.

7 Discussion

7.1 Conclusion

This project gives the first intuitive and formal descriptions of the behaviour of RDMA remote RMWs, outside of a hardware manual. We extend the RDMA^{TSO} memory model to provide operational and declarative semantics, enabling formal reasoning and verification of programs reliant on these operations. By formalising remote RMW behaviours, we have clarified ambiguities in the hardware specification, enabling programmers to understand them in simple terms. In the future, we hope this will lead to more reliable RDMA program and libraries, and potentially to the development of automated verification tools for RDMA.

We are confident that the models given in this report align with intended behaviour for rRMWs given by the InfiniBand technical specification [21]. The behaviours given by these models have been checked against the specification, and confirmed by an expert from NVIDIA, a leading vendor of RDMA hardware.

7.2 Limitations

Although we have sought to confirm the accuracy of our model, we have not conducted independent validation on RDMA hardware. Only an extensive testing campaign would reveal whether real networks implement the specification faithfully. As is the case with all memory models, there may be some behaviours which are permitted but not observed. It was noted in [10] that store buffering with puts was not observed when the earlier put was polled. NVIDIA have confirmed that their current hardware implements a stronger guarantee for polled NIC remote writes, which means the same behaviour is also not possible when rRMWs are used. Therefore, the examples shown in Fig. 4.6 are not possible under current NVIDIA hardware, but this may change in future implementations. In general, we cannot be sure the weak behaviours we have described are possible on current hardware, unless we observe it through empirical validation.

Our work is also limited to RDMA implementations where individual machines adhere to x86-TSO. While x86 is the most common architecture used for RDMA, there are other possibilities; recent research has shown the viability of ARM processors for high-performance computing workloads [23, 24, 25]. This could lead to wider use of RDMA on ARM, which does not observe TSO.

7.3 Future Work

Conformance Testing. The models presented here could be used to develop conformance tests, to confirm that implementation of RDMA faithfully adhere to the specification. One approach to this would be to follow the method shown in [26], by encoding the declarative model in Alloy, a tool for bounded model checking. This is then used to find program executions which are on the boundary of being allowed or disallowed. Appendix C contains a prototype encoding of the declarative semantics, and an example encoding of a litmus test. The next step would be to extend this with a framework for encoding litmus tests, as the process is too tedious and verbose to be done manually.

It is equally possible to encode the operational semantics using other bounded model checkers, such as CMBC [27].

Other Architectures. It would be interesting to investigate the semantics of RDMA in the presence of other architectures such as ARMv8. As the semantics of RDMA operations is independent from the semantics of local operations, it may be possible to develop a memory model for RDMA which is *parametric* in the semantics of the underlying architecture. This would make the semantics durable with respect to changes in the hardware with which RDMA is used.

Verifying Programs and Libraries. The main utility of these models is to verify RDMA libraries and programs. This may be performed manually using the event graph construction we have described, but developing a tool to automate this would enable faster and more extensive verification. One approach is to extend an existing verification tool such as Memalloy [28], which is extensible to arbitrary memory models, such as RDMA^{TSO} .

References

- [1] M. Su et al. RFP: A remote fetching paradigm for RDMA-accelerated systems. *CoRR*, abs/1512.07805, 2015. arXiv: 1512.07805. URL: <http://arxiv.org/abs/1512.07805>.
- [2] N. Boscia and H. S. Sidhu. Comparison of 40G RDMA and Traditional Ethernet Technologies. Technical report NAS-2014-01, NASA Advanced Supercomputing Division, NASA Ames Research Center, Moffett Field, CA, 2014. URL: https://www.nas.nasa.gov/assets/nas/pdf/papers/NAS_Technical_Report_NAS-2014-01.pdf.
- [3] Introduction to InfiniBand. White Paper 2003WP, Mellanox Technologies Inc., 2003. URL: https://network.nvidia.com/pdf/whitepapers/IB_Intro_WP_190.pdf.
- [4] Annex A16: RDMA over Converged Ethernet (RoCE). Technical report, InfiniBand Trade Association, Beaverton, OR, 2010. Available upon request.
- [5] R. Gerstenberger et al. Enabling highly scalable remote memory access programming with MPI-3 one sided. *Commun. ACM*, 61(10):106–113, 2018. DOI: 10.1145/3264413.
- [6] A. Shpiner et al. RoCE rocks without PFC: detailed evaluation. In M. Alizadeh and Y. Zhu, editors, *Proceedings of the Workshop on Kernel-Bypass Networks, KB-Nets@SIGCOMM 2017, Los Angeles, CA, USA, August 21, 2017*, pages 25–30. ACM, 2017. DOI: 10.1145/3098583.3098588.
- [7] S. Owens et al. A better x86 memory model: x86-tso. In S. Berghofer et al., editors, *Theorem Proving in Higher Order Logics, 22nd International Conference, TPHOLs 2009, Munich, Germany, August 17-20, 2009. Proceedings*, volume 5674 of *Lecture Notes in Computer Science*, pages 391–407. Springer, 2009. DOI: 10.1007/978-3-642-03359-9_27.
- [8] C. Pulte et al. Simplifying ARM concurrency: multicopy-atomic axiomatic and operational models for ARMv8. *Proc. ACM Program. Lang.*, 2(POPL):19:1–19:29, 2018. DOI: 10.1145/3158107.
- [9] A. M. Dan et al. Modeling and analysis of remote memory access programming. In E. Visser and Y. Smaragdakis, editors, *Proceedings of the 2016 ACM SIGPLAN International Conference on Object-Oriented Programming, Systems, Languages, and Applications, OOPSLA 2016, part of SPLASH 2016, Amsterdam, The Netherlands, October 30 - November 4, 2016*, pages 129–144. ACM, 2016. DOI: 10.1145/2983990.2984033.
- [10] G. Ambal et al. Semantics of remote direct memory access: operational and declarative models of RDMA on TSO architectures. *Proc. ACM Program. Lang.*, 8(OOPSLA2):1982–2009, 2024. DOI: 10.1145/3689781.
- [11] M. Herlihy. Wait-free synchronization. *ACM Trans. Program. Lang. Syst.*, 13(1):124–149, 1991. DOI: 10.1145/114005.102808.
- [12] A. McMenamin. The end of dennard scaling, 2013. URL: <https://cartesianproduct.wordpress.com/2013/04/15/the-end-of-dennard-scaling/>.

- [13] J. Alglave et al. Frightening small children and disconcerting grown-ups: concurrency in the linux kernel. In X. Shen et al., editors, *Proceedings of the Twenty-Third International Conference on Architectural Support for Programming Languages and Operating Systems, ASPLOS 2018, Williamsburg, VA, USA, March 24-28, 2018*, pages 405–418. ACM, 2018. DOI: 10.1145/3173162.3177156.
- [14] Heisenbug - Wikipedia. URL: <https://en.wikipedia.org/wiki/Heisenbug>.
- [15] M. Herlihy and N. Shavit. *The Art of Multiprocessor Programming, Revised Reprint*. Morgan Kaufmann Publishers Inc., San Francisco, CA, USA, 1st edition, 2012. ISBN: 9780123973375.
- [16] P. Sewell et al. x86-TSO: a rigorous and usable programmer’s model for x86 multiprocessors. *Commun. ACM*, 53(7):89–97, 2010. DOI: 10.1145/1785414.1785443.
- [17] SPARC International, editor. *The SPARC architecture manual: version 8*. Prentice Hall, Englewood Cliffs, N.J, 1992. ISBN: 978-0-13-825001-0.
- [18] A. Godbole et al. Automated conversion of axiomatic to operational models: theory and practice. *CoRR*, abs/2208.06733, 2022. DOI: 10.48550/ARXIV.2208.06733. arXiv: 2208.06733.
- [19] G. D. Plotkin. A structural approach to operational semantics. *J. Log. Algebraic Methods Program.*, 60-61:17–139, 2004.
- [20] PCI-SIG. PCI Express Base Specification Revision 6.0 Version 1.0. Technical report, PCI-SIG, 2022. URL: <https://pcisig.com/pci-express-6.0-specification>.
- [21] *InfiniBand Architecture Specification Volume 1, Release 1.2.1*. Available upon request. InfiniBand Trade Association. Beaverton, OR, 2007.
- [22] J. Preshing. Atomic vs. Non-Atomic Operations, 2013-06-18. URL: <https://preshing.com/20130618/atomic-vs-non-atomic-operations/>.
- [23] F. Mantovani et al. Performance and energy consumption of HPC workloads on a cluster based on arm thunderx2 CPU. *CoRR*, abs/2007.04868, 2020. arXiv: 2007.04868. URL: <https://arxiv.org/abs/2007.04868>.
- [24] N. A. Simakov et al. Are we ready for broader adoption of ARM in the HPC community: performance and energy efficiency analysis of benchmarks and applications executed on high-end ARM systems. In *Proceedings of the HPC Asia 2023 Workshops, HPC Asia 2023, Singapore, 27 February 2023 - 2 March 2023*, pages 78–86. ACM, 2023. DOI: 10.1145/3581576.3581618.
- [25] A. Jackson et al. Evaluating the ARM ecosystem for high performance computing. *CoRR*, abs/1904.04250, 2019. arXiv: 1904.04250. URL: <http://arxiv.org/abs/1904.04250>.
- [26] V. Klimis et al. Taking back control in an intermediate representation for GPU computing. *Proc. ACM Program. Lang.*, 7(POPL), 2023-01. DOI: 10.1145/3571253.
- [27] D. Iorga et al. Simulating operational memory models using off-the-shelf program analysis tools. *IEEE Trans. Software Eng.*, 49(12):5084–5102, 2023. DOI: 10.1109/TSE.2023.3326056.
- [28] J. Wickerson et al. Automatically comparing memory consistency models. In G. Castagna and A. D. Gordon, editors, *Proceedings of the 44th ACM SIGPLAN Symposium on Principles of Programming Languages, POPL 2017, Paris, France, January 18-20, 2017*, pages 190–204. ACM, 2017. DOI: 10.1145/3009837.3009838.

A Declarations

A.1 Use of Generative AI

No generative AI models have been used in the production of this report.

A.2 Ethical Considerations

As with all research into the precise behaviour of computer systems, it is possible, however unlikely, that some behaviour identified may be abused by bad actors to achieve nefarious goals, such as unauthorised access to a system. However, RDMA is typically used on private network executing trusted code, so this is very unlikely in our case. As a primarily theoretical project, there are no other notable ethical considerations.

A.3 Sustainability

No resource intensive tasks were involved in this project. The computer used for the research and writing of this report was a Macbook with an energy efficient Apple M1 chip.

A.4 Availability of Materials

All materials are available within the text of this report, in particular the theoretical models are described in §4 and §5, and the Alloy encoding of the declarative model is provided in Appendix C.

B Annotated Semantics

B.1 Annotated Labels and Inference Rules

On top of the 12 labels presented in §5, we create six new labels: $\text{Put}(\bar{y}, x)$, $\text{Get}(x, \bar{y})$, $\text{nCAS}(z, \bar{x}, v, u)$, $\text{nFAA}(z, \bar{x}, u)$, $\text{nLEX}(\bar{n})$, and $\text{nrEX}(\bar{n})$. These labels can also be used to create events (when bundled with an event identifier and a thread identifier).

We note E^{ext} the extended set of all events, including the six new labels.

Recall that $\mathcal{R} = \text{1R} \cup \text{CAS} \cup \text{n1R} \cup \text{nrR} \cup \text{narR} \subseteq E^{\text{ext}}$ and $\mathcal{W} = \text{1W} \cup \text{CAS} \cup \text{n1W} \cup \text{nrW} \cup \text{narW} \subseteq E^{\text{ext}}$. We also note $\text{nEX} = \text{nLEX} \cup \text{nrEX}$ and $\text{nRMW} = \text{nCAS} \cup \text{nFAA}$.

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

We use $\text{type}(\lambda)$ to denote the type of the annotated label (1R , 1W , CAS , F , Push , NIC , n1R , nrR , n1W , nrW , CN , P , nF , B , \mathcal{E}). We use $r(\lambda)$, $w(\lambda)$, $u(\lambda)$, $a(\lambda)$, $f(\lambda)$, $p(\lambda)$, $e(\lambda)$, \dots to access the elements of a $\lambda \in \text{ALabel}$ where applicable. Also, we note $t(\lambda)$ for the thread of the first argument of λ .

The annotated program transitions (Fig. B.2) use an additional annotated label $\text{CASF}\langle r, w \rangle$ with $r \in \text{1R}$ 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 §B.2. Also, note that the annotated domains (e.g. the store buffers and the queue pairs) contain events, not annotated labels.

initialisation Given a program P , let

$$\begin{array}{ll}
 M_0 \in \text{AMem} & \text{s.t. } \forall x \in \text{Loc. } M_0(x) = \text{init}_x \text{ with } l(\text{init}_x) \triangleq \text{1W}(x, 0) \\
 b_0 \in \text{ASBuff} & b_0 \triangleq \varepsilon \\
 B_0 \in \text{ASBMap} & B_0 \triangleq \lambda t. b_0 \\
 A_0 \in \text{RAMap} & A_0 \triangleq \lambda t. \perp \\
 qp_0 \in \text{AQPair} & qp_0 \triangleq \langle \varepsilon, \varepsilon, \varepsilon \rangle \\
 QP_0 \in \text{AQPMap} & QP_0 \triangleq \lambda t. \lambda n. qp_0
 \end{array}$$

$\lambda \in \text{ALabel}$

$\lambda \triangleq$	$ \text{1R}\langle r, w \rangle$	where $r \in \text{1R}, w \in \mathcal{W}, \text{eq}_{\text{loc}\&\text{v}}(r, w)$
	$ \text{1W}\langle w \rangle$	where $w \in \text{1W}$
	$ \text{CAS}\langle u, w \rangle$	where $u \in \text{CAS}, w \in \mathcal{W}, \text{eq}_{\text{loc}\&\text{v}}(u, w)$
	$ \text{F}\langle f \rangle$	where $f \in \text{F}$
	$ \text{Push}\langle a \rangle$	where $a \in (\text{Put} \cup \text{Get} \cup \text{nCAS} \cup \text{nFAA} \cup \text{nF})$
	$ \text{NIC}\langle a \rangle$	where $a \in (\text{Put} \cup \text{Get} \cup \text{nCAS} \cup \text{nFAA} \cup \text{nF})$
	$ \text{n1R}\langle r, w, a, w' \rangle$	where $r \in \text{n1R}, w \in \mathcal{W}, a \in \text{Put}, w' \in \text{nrW}, \text{eq}_{\text{loc}\&\text{v}}(r, w),$ $\text{loc}_r(a) = \text{loc}(r), \text{loc}_w(a) = \text{loc}(w'), v_r(r) = v_w(w')$
	$ \text{nrR}\langle r, w, a, w' \rangle$	where $r \in \text{nrR}, w \in \mathcal{W}, a \in \text{Get}, w' \in \text{n1W}, \text{eq}_{\text{loc}\&\text{v}}(r, w),$ $\text{loc}_r(a) = \text{loc}(r), \text{loc}_w(a) = \text{loc}(w'), v_r(r) = v_w(w')$
	$ \text{narR}\langle r, w, a, w', w'' \rangle$	where $r \in \text{narR}, w \in \mathcal{W}, a \in \text{nRMW}, w' \in \text{n1W}, w'' \in \text{narW},$ $\text{eq}_{\text{loc}\&\text{v}}(r, w), \text{loc}_r(a) = \text{loc}(r) = \text{loc}(w''),$ $\text{loc}_w(a) = \text{loc}(w'), v_r(r) = v_w(w'),$ $a \in \text{nCAS} \implies v_r(r) = v_e(a) \wedge v_w(w'') = v_u(a)$ $a \in \text{nFAA} \implies v_w(w'') = v_r(r) + v(a)$
	$ \text{naF}\langle r, w, a, w' \rangle$	where $r \in \text{narR}, w \in \mathcal{W}, a \in \text{nRMW}, w' \in \text{n1W}, \text{eq}_{\text{loc}\&\text{v}}(r, w),$ $\text{loc}_r(a) = \text{loc}(r), \text{loc}_w(a) = \text{loc}(w'),$ $v_r(r) = v_w(w'), v_r(r) \neq v_e(a)$
	$ \text{n1W}\langle w, e \rangle$	where $w \in \text{n1W}, e \in \text{n1EX}, \text{sameqp}(w, e)$
	$ \text{nrW}\langle w, e \rangle$	where $w \in \text{nrW}, e \in \text{nrEX}, \text{sameqp}(w, e)$
	$ \text{narW}\langle w \rangle$	where $w \in \text{narW}$
	$ \text{CN}\langle e \rangle$	where $e \in \text{nrEX}$
	$ \text{P}\langle p, e \rangle$	where $p \in \text{P}, e \in \text{nEX}, \text{sameqp}(p, e)$
	$ \text{nF}\langle f \rangle$	where $f \in \text{nF}$
	$ \text{B}\langle w \rangle$	where $w \in \mathcal{W}$
	$ \mathcal{E}\langle t \rangle$	where $t \in \text{Tid}$

$$\begin{aligned} \text{eq}_{\text{loc}\&\text{v}}(r, w) &\triangleq \text{loc}(r) = \text{loc}(w) \wedge v_r(r) = v_w(w) \\ \text{sameqp}(e, e') &\triangleq t(e) = t(e') \wedge \bar{n}(e) = \bar{n}(e') \end{aligned}$$

Figure B.1: Annotated Labels

Program transitions: $\text{Prog} \xrightarrow{\text{ALabel} \uplus \{\text{CASF}\}} \text{Prog}$
Command transitions: $\text{Comm} \xrightarrow{\text{ALabel} \uplus \{\text{CASF}\}} \text{Comm}$

$$\begin{array}{c}
\frac{C_1 \xrightarrow{\lambda} C'_1}{C_1; C_2 \xrightarrow{\lambda} C'_1; C_2} \quad \frac{}{\text{skip}; C \xrightarrow{\mathcal{E}\langle t \rangle} C} \quad \frac{i \in \{1, 2\}}{C_1 + C_2 \xrightarrow{\mathcal{E}\langle t \rangle} C_i} \quad \frac{}{C^* \xrightarrow{\mathcal{E}\langle t \rangle} \text{skip}} \\
\\
\frac{}{C^* \xrightarrow{\mathcal{E}\langle t \rangle} C; C^*} \quad \frac{C \rightsquigarrow C'}{C \xrightarrow{\mathcal{E}\langle t \rangle} C'} \quad \frac{\text{elocs}(e) = \emptyset \quad w = (\iota, t, \text{lW}(x, \llbracket e \rrbracket))}{x := e \xrightarrow{\text{lW}\langle w \rangle} \text{skip}} \\
\\
\frac{\text{elocs}(e_{\text{old}}) = \text{elocs}(e_{\text{new}}) = \emptyset \quad v \neq \llbracket e_{\text{old}} \rrbracket \quad r = (\iota, t, \text{lR}(x, v))}{z := \text{CAS}(x, e_{\text{old}}, e_{\text{new}}) \xrightarrow{\text{CASF}\langle r, w \rangle} z := v} \\
\\
\frac{\text{elocs}(e_{\text{old}}) = \text{elocs}(e_{\text{new}}) = \emptyset \quad u = (\iota, t, \text{CAS}(x, \llbracket e_{\text{old}} \rrbracket, \llbracket e_{\text{new}} \rrbracket))}{z := \text{CAS}(x, e_{\text{old}}, e_{\text{new}}) \xrightarrow{\text{CAS}\langle u, w \rangle} z := \llbracket e_{\text{old}} \rrbracket} \quad \frac{f = (\iota, t, \text{F})}{\text{mfence} \xrightarrow{\text{F}\langle f \rangle} \text{skip}} \\
\\
\frac{a = (\iota, t, \text{Get}(x, \bar{y}))}{x := \bar{y} \xrightarrow{\text{Push}\langle a \rangle} \text{skip}} \quad \frac{a = (\iota, t, \text{Put}(\bar{y}, x))}{\bar{y} := x \xrightarrow{\text{Push}\langle a \rangle} \text{skip}} \quad \frac{a = (\iota, t, \text{nF}(\bar{n}))}{\text{rfence } \bar{n} \xrightarrow{\text{Push}\langle a \rangle} \text{skip}} \\
\\
\frac{\text{elocs}(e_{\text{old}}) = \text{elocs}(e_{\text{new}}) = \emptyset \quad v = \llbracket e_{\text{old}} \rrbracket \quad u = \llbracket e_{\text{new}} \rrbracket \quad a = (\iota, t, \text{nCAS}(z, \bar{x}, v, u))}{z := \text{nCAS}(x, e_{\text{old}}, e_{\text{new}}) \xrightarrow{\text{Push}\langle a \rangle} \text{skip}} \\
\\
\frac{\text{elocs}(e) = \emptyset \quad u = \llbracket e \rrbracket \quad a = (\iota, t, \text{nFAA}(z, \bar{x}, u))}{z := \text{nFAA}(x, e) \xrightarrow{\text{Push}\langle a \rangle} \text{skip}} \quad \frac{p = (\iota, t, \text{P}(n))}{\text{poll}(n) \xrightarrow{\text{P}\langle p, e \rangle} \text{skip}} \\
\\
\frac{r = (\iota, t, \text{lR}(x, v))}{\text{assume}(x = v) \xrightarrow{\text{lR}\langle r, w \rangle} \text{skip}} \quad \frac{v \neq v' \quad r = (\iota, t, \text{lR}(x, v'))}{\text{assume}(x \neq v) \xrightarrow{\text{lR}\langle r, w \rangle} \text{skip}} \quad \frac{\text{P}(t(\lambda)) \xrightarrow{\lambda} C}{\text{P} \xrightarrow{\lambda} \text{P}[t(\lambda) \mapsto C]}
\end{array}$$

Figure B.2: RDMA^{TSO} program and command transitions for the annotated semantics

$$\begin{aligned}
M \in \text{AMem} &\triangleq \{m \in \text{Loc} \rightarrow \mathcal{W} \mid \forall x \in \text{Loc}. \text{loc}(m[x]) = x\} & B \in \text{ASBMap} &\triangleq \text{Tid} \rightarrow \text{ASBuff} \\
A \in \text{RAMap} &\triangleq \lambda n. \{\perp, \top\} & \text{QP} \in \text{AQPMa} &\triangleq \text{Tid} \rightarrow (\text{Node} \rightarrow \text{AQPair}) \\
b \in \text{ASBuff} &\triangleq (1W \cup \text{Get} \cup \text{Put} \cup \text{nF} \cup \text{nCAS} \cup \text{nFAA})^* & \text{sqp} \in \text{AQPair} &\triangleq \text{APipe} \times \text{AWBR} \times \text{AWBL} \\
\text{pipe} \in \text{APipe} &\triangleq (\text{Get} \cup \text{Put} \cup \text{nF} \cup \text{nrW} \cup \text{narW} \cup \text{nrEX} \cup \text{n1W} \cup \text{nCAS} \cup \text{nFAA})^* \\
\text{wb}_R \in \text{AWBR} &\triangleq (\text{nrW}, \text{narW})^* & \text{wb}_L \in \text{AWBL} &\triangleq (\text{n1W} \cup \text{n1EX} \cup \text{nrEX})^*
\end{aligned}$$

$$\begin{aligned}
&\frac{B' = B[t(w) \mapsto w \cdot B(t(w))]}{M, B, A, \text{QP} \xrightarrow{1W\langle w \rangle} M, B', A, \text{QP}} & \frac{(M \blacktriangleleft B(t(r)))(\text{loc}(r)) = w \quad v_r(r) = v_w(w)}{M, B, A, \text{QP} \xrightarrow{1R\langle r, w \rangle} M, B, A, \text{QP}} \\
&\frac{B(t(u)) = \varepsilon \quad M(\text{loc}(u)) = w \quad v_r(u) = v_w(w)}{M, B, A, \text{QP} \xrightarrow{\text{CAS}\langle u, w \rangle} M[x \mapsto u], B, A, \text{QP}} & \frac{B(t(f)) = \varepsilon}{M, B, A, \text{QP} \xrightarrow{F\langle f \rangle} M, B, A, \text{QP}} \\
&\frac{B' = B[t(a) \mapsto a \cdot B(t(a))]}{M, B, A, \text{QP} \xrightarrow{\text{Push}\langle a \rangle} M, B', A, \text{QP}} & \frac{B(t(w)) = b \cdot w \quad w \in 1W}{M, B, A, \text{QP} \xrightarrow{B\langle w \rangle} M[x \mapsto w], B[t(w) \mapsto b], A, \text{QP}} \\
&\frac{B(t(a)) = b \cdot a \quad a \notin 1W \quad \text{QP}(t(a))(n(a)) = \text{qp} \quad \text{qp}' = \text{qp}[\text{pipe} \mapsto a \cdot \text{qp.pipe}]}{M, B, A, \text{QP} \xrightarrow{\text{NIC}\langle a \rangle} M, B[t(a) \mapsto b], A, \text{QP}[t(a) \mapsto \text{QP}(t(a))[n(a) \mapsto \text{qp}']]} \\
&\frac{\text{QP}(t(p))(n(p)) = \text{qp} \quad \text{qp}.\text{wb}_L = \alpha \cdot e \quad e \in \text{nEX} \quad \text{qp}' = \text{qp}[\text{wb}_L \mapsto \alpha]}{M, B, A, \text{QP} \xrightarrow{P\langle p, e \rangle} M, B, A, \text{QP}[t(p) \mapsto \text{QP}(t(p))[n(p) \mapsto \text{qp}']]} \\
&\frac{M, A, \text{QP}(t(\lambda))(\bar{n}) \xrightarrow{\lambda}_{\text{sqp}} M', A', \text{qp}}{M, B, A, \text{QP} \xrightarrow{\lambda} M', B, A', \text{QP}[t(\lambda) \mapsto \text{QP}(t(\lambda))][\bar{n} \mapsto \text{qp}]}
\end{aligned}$$

$$\text{with } (M \blacktriangleleft \alpha)(x) = \begin{cases} M[x] & \alpha = \varepsilon \\ w & \alpha = w \cdot \beta \wedge w \in \mathcal{W} \wedge \text{loc}(w) = x \\ (M \blacktriangleleft \beta)(x) & \alpha = e \cdot \beta \wedge (e \notin \mathcal{W} \vee \text{loc}(e) \neq x) \end{cases}$$

Figure B.3: RDMA^{TSO} hardware domains and hardware transitions for the annotated semantics

$$\begin{array}{c}
\frac{\text{pipe} = \alpha \cdot f \quad f = (\iota, t, \text{nF}(n))}{\text{M}, \text{A}, \langle \text{pipe}, \text{wb}_R, \text{wb}_L \rangle \xrightarrow{\text{nF}\langle f \rangle}_{\text{sqp}} \text{M}, \text{A}, \langle \alpha, \text{wb}_R, \text{wb}_L \rangle} \\
\\
\frac{\begin{array}{l} \text{pipe} = \alpha \cdot a \cdot \beta \quad a = (\iota_a, t, \text{Put}(\bar{y}, x)) \quad \text{M}(x) = w \quad r = (\iota_r, t, \text{nLR}(x, v_w(w), n(\bar{y}))) \\ w' = (\iota_{w'}, t, \text{nrW}(\bar{y}, v_w(w))) \quad \beta \in (\text{nrW} \cup \text{narW} \cup \text{Get} \cup \text{nLW} \cup \text{nCAS} \cup \text{nFAA} \cup \text{nrEX})^* \quad \text{wb}_L \in \text{nEX}^* \end{array}}{\text{M}, \text{A}, \langle \text{pipe}, \text{wb}_R, \text{wb}_L \rangle \xrightarrow{\text{nLR}\langle r, w, a, w' \rangle}_{\text{sqp}} \text{M}, \text{A}, \langle \alpha \cdot w' \cdot \beta, \text{wb}_R, \text{wb}_L \rangle} \\
\\
\frac{\text{pipe} = \alpha \cdot w \cdot \beta \quad w = (\iota_w, t, \text{nrW}(\bar{y}, v)) \quad e = (\iota_e, t, \text{nrEX}(n(\bar{y}))) \quad \beta \in (\text{Get} \cup \text{nLW} \cup \text{nrEX})^*}{\text{M}, \text{A}, \langle \text{pipe}, \text{wb}_R, \text{wb}_L \rangle \xrightarrow{\text{nrW}\langle w, e \rangle}_{\text{sqp}} \text{M}, \text{A}, \langle \alpha \cdot e \cdot \beta, w \cdot \text{wb}_R, \text{wb}_L \rangle} \\
\\
\frac{\text{wb}_R = \alpha \cdot w \quad w \in \text{nrW}}{\text{M}, \text{A}, \langle \text{pipe}, \text{wb}_R, \text{wb}_L \rangle \xrightarrow{\text{B}\langle w \rangle}_{\text{sqp}} \text{M}[\text{loc}(w) \mapsto w], \text{A}, \langle \text{pipe}, \alpha, \text{wb}_L \rangle} \\
\\
\frac{\text{pipe} = \alpha \cdot e \quad e \in \text{nrEX}}{\text{M}, \text{A}, \langle \text{pipe}, \text{wb}_R, \text{wb}_L \rangle \xrightarrow{\text{CN}\langle e \rangle}_{\text{sqp}} \text{M}, \text{A}, \langle \alpha, \text{wb}_R, e \cdot \text{wb}_L \rangle} \\
\\
\frac{\begin{array}{l} \text{pipe} = \alpha \cdot a \cdot \beta \quad a = (\iota_a, t, \text{Get}(x, \bar{y})) \quad \text{M}(\bar{y}) = w \quad r = (\iota_r, t, \text{nrR}(\bar{y}, v_w(w))) \\ w' = (\iota_{w'}, t, \text{nLW}(x, v_w(w), n(\bar{y}))) \quad \beta \in (\text{Get} \cup \text{nLW} \cup \text{nrEX})^* \quad \text{wb}_R = \varepsilon \end{array}}{\text{M}, \text{A}, \langle \text{pipe}, \text{wb}_R, \text{wb}_L \rangle \xrightarrow{\text{nrR}\langle r, w, a, w' \rangle}_{\text{sqp}} \text{M}, \text{A}, \langle \alpha \cdot w' \cdot \beta, \text{wb}_R, \text{wb}_L \rangle} \\
\\
\frac{\begin{array}{l} \text{pipe} = \alpha \cdot w \quad w = (\iota_w, t, \text{nLW}(x, v, n)) \quad e = (\iota_e, t, \text{nLEX}(n)) \\ \text{M}, \text{A}, \langle \text{pipe}, \text{wb}_R, \text{wb}_L \rangle \xrightarrow{\text{nLW}\langle w, e \rangle}_{\text{sqp}} \text{M}, \text{A}, \langle \alpha, \text{wb}_R, e \cdot w \cdot \text{wb}_L \rangle \end{array}}{\text{M}, \text{A}, \langle \text{pipe}, \text{wb}_R, \text{wb}_L \rangle \xrightarrow{\text{B}\langle w \rangle}_{\text{sqp}} \text{M}[\text{loc}(w) \mapsto w], \text{A}, \langle \text{pipe}, \text{wb}_R, \alpha \cdot \beta \rangle} \\
\\
\frac{\text{wb}_L = \alpha \cdot w \cdot \beta \quad w \in \text{nLW} \quad \beta \in \text{nEX}^*}{\text{M}, \text{A}, \langle \text{pipe}, \text{wb}_R, \text{wb}_L \rangle \xrightarrow{\text{B}\langle w \rangle}_{\text{sqp}} \text{M}[\text{loc}(w) \mapsto w], \text{A}, \langle \text{pipe}, \text{wb}_R, \alpha \cdot \beta \rangle} \\
\\
\frac{\begin{array}{l} \text{pipe} = \alpha \cdot a \cdot \beta \quad \text{wb}_R = \varepsilon \quad \text{M}(\bar{x}) = w \quad v_w(w) \neq v \quad \text{A}(n(\bar{x})) = \perp \quad a = (\iota_a, t, \text{nCAS}(z, \bar{x}, v, u)) \\ r = (\iota_r, t, \text{narR}(\bar{x}, v_w(w))) \quad w' = (\iota_{w'}, t, \text{nLW}(z, v_w(w), n(\bar{x}))) \quad \beta \in (\text{Get} \cup \text{nLW} \cup \text{nrEX})^* \end{array}}{\text{M}, \text{A}, \langle \text{pipe}, \text{wb}_R, \text{wb}_L \rangle \xrightarrow{\text{narF}\langle r, w, a, w' \rangle}_{\text{sqp}} \text{M}, \text{A}, \langle \alpha \cdot w' \cdot \beta, \text{wb}_R, \text{wb}_L \rangle} \\
\\
\frac{\begin{array}{l} \text{pipe} = \alpha \cdot a \cdot \beta \quad \text{wb}_R = \varepsilon \quad \text{M}(\bar{x}) = w \\ v_w(w) = v \quad \text{A}(n(\bar{x})) = \perp \quad a = (\iota_a, t, \text{nCAS}(z, \bar{x}, v, u)) \quad r = (\iota_r, t, \text{narR}(\bar{x}, v_w(w))) \\ w'' = (\iota_{w''}, t, \text{narW}(\bar{x}, u)) \quad w' = (\iota_{w'}, t, \text{nLW}(z, v_w(w), n(\bar{x}))) \quad \beta \in (\text{Get} \cup \text{nLW} \cup \text{nrEX})^* \end{array}}{\text{M}, \text{A}, \langle \text{pipe}, \text{wb}_R, \text{wb}_L \rangle \xrightarrow{\text{narR}\langle r, w, a, w', w'' \rangle}_{\text{sqp}} \text{M}, \text{A}[n(\bar{x}) \mapsto \top], \langle \alpha \cdot w' \cdot w'' \cdot \beta, \text{wb}_R, \text{wb}_L \rangle} \\
\\
\frac{\begin{array}{l} \text{pipe} = \alpha \cdot a \cdot \beta \quad \text{wb}_R = \varepsilon \quad \text{M}(\bar{x}) = w \\ v_w(w) + v = u \quad \text{A}(n(\bar{x})) = \perp \quad a = (\iota_a, t, \text{nFAA}(z, \bar{x}, v)) \quad r = (\iota_r, t, \text{narR}(\bar{x}, v_w(w))) \\ w'' = (\iota_{w''}, t, \text{narW}(\bar{x}, u)) \quad w' = (\iota_{w'}, t, \text{nLW}(z, v_w(w))) \quad \beta \in (\text{Get} \cup \text{nLW} \cup \text{nrEX})^* \end{array}}{\text{M}, \text{A}, \langle \text{pipe}, \text{wb}_R, \text{wb}_L \rangle \xrightarrow{\text{narR}\langle r, w, a, w', w'' \rangle}_{\text{sqp}} \text{M}, \text{A}[n(\bar{x}) \mapsto \top], \langle \alpha \cdot w' \cdot w'' \cdot \beta, \text{wb}_R, \text{wb}_L \rangle} \\
\\
\frac{\text{pipe} = \alpha \cdot w \cdot \beta \quad \beta \in (\text{Get} \cup \text{nLW} \cup \text{nrEX})^* \quad w = (\iota_w, t, \text{narW}(\bar{x}, v))}{\text{M}, \text{A}, \langle \text{pipe}, \text{wb}_R, \text{wb}_L \rangle \xrightarrow{\text{narW}\langle w \rangle}_{\text{sqp}} \text{M}, \text{A}, \langle \alpha \cdot \beta, w \cdot \text{wb}_R, \text{wb}_L \rangle} \\
\\
\frac{\text{wb}_R = \alpha \cdot w \quad w = (\iota_w, t, \text{narW}(\bar{x}, v))}{\text{M}, \text{A}, \langle \text{pipe}, \text{wb}_R, \text{wb}_L \rangle \xrightarrow{\text{B}\langle w \rangle}_{\text{sqp}} \text{M}[\text{loc}(w) \mapsto w], \text{A}[n(\bar{x}) \mapsto \perp], \langle \text{pipe}, \alpha, \text{wb}_L \rangle}
\end{array}$$

Figure B.4: Annotated 3 Buffers NIC Semantics

B.2 Paths, Gluing, and Other Definitions

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

We define Annotated Operational Semantics Gluing with the following rules.

$$\begin{array}{c}
 \frac{P \xrightarrow{\mathcal{E}\langle t \rangle} P'}{P, M, B, A, QP, \pi \Rightarrow P', M, B, A, QP, \pi} \\
 \\
 \frac{P \xrightarrow{\lambda} P' \quad M, B, A, QP \xrightarrow{\lambda} M', B', A, QP' \quad \lambda \in (1R \cup 1W \cup \text{CAS} \cup F \cup \text{Push} \cup P) \quad \text{fresh}(\lambda, \pi)}{P, M, B, A, QP, \pi \Rightarrow P', M', B', A, QP', \lambda \cdot \pi} \\
 \\
 \frac{M, B, A, QP \xrightarrow{\lambda} M', B', A', QP' \quad \lambda \in (\text{NIC} \cup \text{n1R} \cup \text{nrR} \cup \text{n1W} \cup \text{nrW} \cup \text{naF} \cup \text{narR} \cup \text{narW} \cup \text{CN} \cup \text{nF} \cup B) \quad \text{fresh}(\lambda, \pi)}{P, M, B, A, QP, \pi \Rightarrow P, M', B', A', QP', \lambda \cdot \pi} \\
 \\
 \frac{P \xrightarrow{\text{CASF}\langle r, w \rangle} P' \quad \lambda_1 = F\langle (l, t(r), F) \rangle \quad \lambda_2 = 1R\langle r, w \rangle \quad M, B, A, QP \xrightarrow{\lambda_1} M, B, A, QP \xrightarrow{\lambda_2} M, B, A, QP \quad \text{fresh}(\lambda_1, \pi) \quad \text{fresh}(\lambda_2, \pi)}{P, M, B, A, QP, \pi \Rightarrow P', M, B, A, QP, \lambda_2 \cdot \lambda_1 \cdot \pi}
 \end{array}$$

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

$$\mathbf{Relevant} : \text{ALabel} \rightarrow 2^{E^{\text{ext}}}$$

$$\begin{array}{ll}
 \mathbf{Relevant}(1R\langle r, _ \rangle) \triangleq \{r\} & \mathbf{Relevant}(\text{narR}\langle r, _, a, w'', w' \rangle) \triangleq \{r, a, w', w''\} \\
 \mathbf{Relevant}(1W\langle w \rangle) \triangleq \{w\} & \mathbf{Relevant}(\text{narW}\langle w \rangle) \triangleq \{w\} \\
 \mathbf{Relevant}(\text{CAS}\langle u, _ \rangle) \triangleq \{u\} & \mathbf{Relevant}(\text{n1W}\langle w, e \rangle) \triangleq \{w, e\} \\
 \mathbf{Relevant}(F\langle f \rangle) \triangleq \{f\} & \mathbf{Relevant}(\text{nrW}\langle w, e \rangle) \triangleq \{w, e\} \\
 \mathbf{Relevant}(\text{Push}\langle a \rangle) \triangleq \{a\} & \mathbf{Relevant}(\text{CN}\langle e \rangle) \triangleq \{e\} \\
 \mathbf{Relevant}(\text{NIC}\langle a \rangle) \triangleq \{a\} & \mathbf{Relevant}(P\langle p, e \rangle) \triangleq \{p, e\} \\
 \mathbf{Relevant}(\text{n1R}\langle r, _, a, w' \rangle) \triangleq \{r, a, w'\} & \mathbf{Relevant}(\text{nF}\langle f \rangle) \triangleq \{f\} \\
 \mathbf{Relevant}(\text{nrR}\langle r, _, a, w' \rangle) \triangleq \{r, a, w'\} & \mathbf{Relevant}(B\langle w \rangle) \triangleq \{w\} \\
 \mathbf{Relevant}(\text{naF}\langle r, _, a, w' \rangle) \triangleq \{r, a, w'\} & \mathbf{Relevant}(\mathcal{E}\langle _ \rangle) \triangleq \{\}
 \end{array}$$

$$\begin{aligned}
 \lambda_1 \bowtie \lambda_2 &\triangleq \text{type}(\lambda_1) \neq \text{type}(\lambda_2) \vee \mathbf{Relevant}(\lambda_1) \cap \mathbf{Relevant}(\lambda_2) = \emptyset \\
 \text{fresh}(\lambda, \pi) &\triangleq \forall \lambda' \in \pi, \lambda \bowtie \lambda' \\
 \text{nodup}(\pi) &\triangleq \forall \pi_2, \lambda, \pi_1. \pi = \pi_2 \cdot \lambda \cdot \pi_1 \implies \text{fresh}(\lambda, \pi_1)
 \end{aligned}$$

Relevant(λ) 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 $\text{1R}\langle r_1, w \rangle$ and $\text{1R}\langle r_2, w \rangle$ during an execution is fine, since w can be looked up any number of time.
- Having both $\text{n1R}\langle r_1, w_1, a, e_1 \rangle$ and $\text{n1R}\langle r_2, w_2, a, e_2 \rangle$ during an execution is problematic, since it means the put operation a is being run twice.

Completeness

$\text{complete}(\pi) \triangleq \forall a, w', e, r, w, f, w''.$

$$\begin{aligned}
& \text{1W}\langle w \rangle \in \pi \implies \text{B}\langle w \rangle \in \pi \\
& \wedge \text{Push}\langle a \rangle \in \pi \implies \text{NIC}\langle a \rangle \in \pi \\
& \wedge \text{NIC}\langle f \rangle \in \pi \wedge f \in \text{nF} \implies \text{nF}\langle f \rangle \in \pi \\
& \wedge \text{NIC}\langle a \rangle \in \pi \wedge a \in \text{Put} \implies \exists r, w, w'. \text{n1R}\langle r, w, a, w' \rangle \in \pi \\
& \wedge \text{NIC}\langle a \rangle \in \pi \wedge a \in \text{Get} \implies \exists r, w, w'. \text{nrR}\langle r, w, a, w' \rangle \in \pi \\
& \wedge \text{NIC}\langle a \rangle \in \pi \wedge a \in \text{nFAA} \implies \exists r, w, w'. \text{narR}\langle r, w, a, w', w'' \rangle \in \pi \\
& \wedge \text{NIC}\langle a \rangle \in \pi \wedge a \in \text{nCAS} \implies \left(\begin{array}{l} \exists r, w, w'. \text{naF}\langle r, w, a, w' \rangle \in \pi \\ \vee \exists r, w, w', w''. \text{narR}\langle r, w, a, w', w'' \rangle \in \pi \end{array} \right) \\
& \wedge \text{n1R}\langle r, w, a, w' \rangle \in \pi \implies \exists e. \text{nrW}\langle w', e \rangle \in \pi \\
& \wedge \text{nrR}\langle r, w, a, w' \rangle \in \pi \implies \exists e. \text{n1W}\langle w', e \rangle \in \pi \\
& \wedge \text{narR}\langle r, w, a, w', w'' \rangle \in \pi \implies \text{narW}\langle w'' \rangle \in \pi \\
& \wedge \text{narR}\langle r, w, a, w', w'' \rangle \in \pi \implies \exists e. \text{n1W}\langle w', e \rangle \in \pi \\
& \wedge \text{naF}\langle r, w, a, w' \rangle \in \pi \implies \exists e. \text{n1W}\langle w', e \rangle \in \pi \\
& \wedge \text{n1W}\langle w, e \rangle \in \pi \implies \text{B}\langle w \rangle \in \pi \\
& \wedge \text{nrW}\langle w, e \rangle \in \pi \implies \text{B}\langle w \rangle \in \pi \wedge \text{CN}\langle e \rangle \in \pi \\
& \wedge \text{narW}\langle w \rangle \in \pi \implies \text{B}\langle w \rangle \in \pi
\end{aligned}$$

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 π without duplicate (e.g. if **nodup**(π) 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)

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

$$\begin{aligned}
& \mathbf{B}\langle w \rangle \in \pi \implies \left(\begin{array}{l} \mathbf{1W}\langle w \rangle \prec_{\pi} \mathbf{B}\langle w \rangle \\ \vee \exists e. \mathbf{n1W}\langle w, e \rangle \prec_{\pi} \mathbf{B}\langle w \rangle \\ \vee \exists e. \mathbf{nrW}\langle w, e \rangle \prec_{\pi} \mathbf{B}\langle w \rangle \\ \vee \exists e. \mathbf{narW}\langle w \rangle \prec_{\pi} \mathbf{B}\langle w \rangle \end{array} \right) \\
& \wedge \mathbf{NIC}\langle a \rangle \in \pi \implies \mathbf{Push}\langle a \rangle \prec_{\pi} \mathbf{NIC}\langle a \rangle \\
& \wedge \mathbf{nF}\langle f \rangle \in \pi \implies \mathbf{NIC}\langle f \rangle \prec_{\pi} \mathbf{nF}\langle f \rangle \\
& \wedge \mathbf{n1R}\langle r, w, a, w' \rangle \in \pi \implies \mathbf{NIC}\langle a \rangle \prec_{\pi} \mathbf{n1R}\langle r, w, a, w' \rangle \\
& \wedge \mathbf{nrR}\langle r, w, a, w' \rangle \in \pi \implies \mathbf{NIC}\langle a \rangle \prec_{\pi} \mathbf{nrR}\langle r, w, a, w' \rangle \\
& \wedge \mathbf{naF}\langle r, w, a, w' \rangle \in \pi \implies \mathbf{NIC}\langle a \rangle \prec_{\pi} \mathbf{naF}\langle r, w, a, w' \rangle \\
& \wedge \mathbf{narR}\langle r, w, a, w', w'' \rangle \in \pi \implies \mathbf{NIC}\langle a \rangle \prec_{\pi} \mathbf{narR}\langle r, w, a, w', w'' \rangle \\
& \wedge \mathbf{nrW}\langle w', e \rangle \in \pi \implies \exists r, w, a. \mathbf{n1R}\langle r, w, a, w' \rangle \prec_{\pi} \mathbf{nrW}\langle w', e \rangle \\
& \wedge \mathbf{n1W}\langle w', e \rangle \in \pi \implies \left(\begin{array}{l} \exists r, w, a. \mathbf{nrR}\langle r, w, a, w' \rangle \prec_{\pi} \mathbf{n1W}\langle w', e \rangle \\ \vee \exists r, w, a. \mathbf{naF}\langle r, w, a, w' \rangle \prec_{\pi} \mathbf{n1W}\langle w', e \rangle \\ \vee \exists r, w, a, w''. \left(\begin{array}{l} \mathbf{narR}\langle r, w, a, w', w'' \rangle \\ \prec_{\pi} \mathbf{narW}\langle w'' \rangle \prec_{\pi} \mathbf{n1W}\langle w', e \rangle \end{array} \right) \end{array} \right) \\
& \wedge \mathbf{narW}\langle w' \rangle \in \pi \implies \exists r, w, a, w''. \mathbf{narR}\langle r, w, a, w', w'' \rangle \prec_{\pi} \mathbf{narW}\langle w' \rangle \\
& \wedge \mathbf{CN}\langle e \rangle \in \pi \implies \exists w. \mathbf{nrW}\langle w, e \rangle \prec_{\pi} \mathbf{CN}\langle e \rangle \\
& \wedge \mathbf{P}\langle p, e \rangle \in \pi \implies \left(\begin{array}{l} \exists w. \mathbf{n1W}\langle w, e \rangle \prec_{\pi} \mathbf{B}\langle w \rangle \prec_{\pi} \mathbf{P}\langle p, e \rangle \\ \vee \mathbf{CN}\langle e \rangle \prec_{\pi} \mathbf{P}\langle p, e \rangle \end{array} \right)
\end{aligned}$$

Poll Order

$$\text{pollOrder}(\pi) \triangleq \forall e_1, e_2. \left(\begin{array}{l} \text{sameqp}(e_1, e_2) \\ \wedge \lambda_1 \in \{\mathbf{n1W}\langle _, e_1 \rangle, \mathbf{CN}\langle e_1 \rangle\} \\ \wedge \lambda_2 \in \{\mathbf{n1W}\langle _, e_2 \rangle, \mathbf{CN}\langle e_2 \rangle\} \\ \wedge \lambda_1 \prec_{\pi} \lambda_2 \\ \wedge \mathbf{P}\langle _, e_2 \rangle \in \pi \end{array} \right) \implies \mathbf{P}\langle _, e_1 \rangle \prec_{\pi} \mathbf{P}\langle _, e_2 \rangle$$

Flush Order

$\text{bufFlushOrd}(\pi) \triangleq$

$$\begin{aligned}
& \forall w_1, w_2 \in \text{1W}. \left(t(w_1) = t(w_2) \implies \right. \\
& \quad \left. (\text{B}\langle w_2 \rangle \in \pi \wedge \text{1W}\langle w_1 \rangle \prec_\pi \text{1W}\langle w_2 \rangle) \iff \text{B}\langle w_1 \rangle \prec_\pi \text{B}\langle w_2 \rangle \right) \\
& \wedge \forall a_1, a_2 \in (\text{Get} \cup \text{Put} \cup \text{nF} \cup \text{nCAS} \cup \text{nFAA}). \\
& \quad \left(t(a_1) = t(a_2) \implies \right. \\
& \quad \left. (\text{NIC}\langle a_2 \rangle \in \pi \wedge \text{Push}\langle a_1 \rangle \prec_\pi \text{Push}\langle a_2 \rangle) \iff \text{NIC}\langle a_1 \rangle \prec_\pi \text{NIC}\langle a_2 \rangle \right) \\
& \wedge \forall a_1 \in (\text{Get} \cup \text{Put} \cup \text{nF} \cup \text{nCAS} \cup \text{nFAA}), w_2 \in \text{1W}. \\
& \quad \left(t(a_1) = t(w_2) \implies \right. \\
& \quad \left. (\text{B}\langle w_2 \rangle \in \pi \wedge \text{Push}\langle a_1 \rangle \prec_\pi \text{1W}\langle w_2 \rangle) \iff \text{NIC}\langle a_1 \rangle \prec_\pi \text{B}\langle w_2 \rangle \right. \\
& \quad \left. \wedge (\text{NIC}\langle a_1 \rangle \in \pi \wedge \text{1W}\langle w_2 \rangle \prec_\pi \text{Push}\langle a_1 \rangle) \iff \text{B}\langle w_2 \rangle \prec_\pi \text{NIC}\langle a_1 \rangle \right) \\
& \wedge \forall w_1, w_2 \in \text{n1W}. \left(\text{sameqp}(w_1, w_2) \implies \right. \\
& \quad \left. (\text{B}\langle w_2 \rangle \in \pi \wedge \text{n1W}\langle w_1 \rangle \prec_\pi \text{n1W}\langle w_2 \rangle) \iff \text{B}\langle w_1 \rangle \prec_\pi \text{B}\langle w_2 \rangle \right) \\
& \wedge \forall w_1, w_2 \in \text{nrW}. \left(\text{sameqp}(w_1, w_2) \implies \right. \\
& \quad \left. (\text{B}\langle w_2 \rangle \in \pi \wedge \text{nrW}\langle w_1 \rangle \prec_\pi \text{nrW}\langle w_2 \rangle) \iff \text{B}\langle w_1 \rangle \prec_\pi \text{B}\langle w_2 \rangle \right) \\
& \wedge \forall w_1, w_2 \in \text{narW}. \left(\text{sameqp}(w_1, w_2) \implies \right. \\
& \quad \left. (\text{B}\langle w_2 \rangle \in \pi \wedge \text{narW}\langle w_1 \rangle \prec_\pi \text{narW}\langle w_2 \rangle) \iff \text{B}\langle w_1 \rangle \prec_\pi \text{B}\langle w_2 \rangle \right) \\
& \wedge \forall w_1 \in \text{nrW}, w_2 \in \text{narW}. \left(\text{sameqp}(w_1, w_2) \implies \right. \\
& \quad \left. (\text{B}\langle w_2 \rangle \in \pi \wedge \text{nrW}\langle w_1 \rangle \prec_\pi \text{narW}\langle w_2 \rangle) \iff \text{B}\langle w_1 \rangle \prec_\pi \text{B}\langle w_2 \rangle \right) \\
& \wedge \forall w_1 \in \text{narW}, w_2 \in \text{nrW}. \left(\text{sameqp}(w_1, w_2) \implies \right. \\
& \quad \left. (\text{B}\langle w_2 \rangle \in \pi \wedge \text{narW}\langle w_1 \rangle \prec_\pi \text{nrW}\langle w_2 \rangle) \iff \text{B}\langle w_1 \rangle \prec_\pi \text{B}\langle w_2 \rangle \right) \\
& \wedge \forall w \in \text{1W}, f \in \text{F}. \text{1W}\langle w \rangle \prec_\pi \text{F}\langle f \rangle \wedge t(w) = t(f) \implies \text{B}\langle w \rangle \prec_\pi \text{F}\langle f \rangle \\
& \wedge \forall w \in \text{1W}, u \in \text{CAS}. \text{1W}\langle w \rangle \prec_\pi \text{CAS}\langle u, _ \rangle \wedge t(w) = t(u) \implies \text{B}\langle w \rangle \prec_\pi \text{CAS}\langle u, _ \rangle \\
& \wedge \forall w \in \text{n1W}, r \in \text{n1R}. \\
& \quad (\text{n1W}\langle w, _ \rangle \prec_\pi \text{n1R}\langle r, _, _, _ \rangle \wedge \text{sameqp}(w, r)) \implies \text{B}\langle w \rangle \prec_\pi \text{n1R}\langle r, _, _, _ \rangle \\
& \wedge \forall w \in \text{nrW}, r \in \text{nrR}. \\
& \quad (\text{nrW}\langle w, _ \rangle \prec_\pi \text{nrR}\langle r, _, _, _ \rangle \wedge \text{sameqp}(w, r)) \implies \text{B}\langle w \rangle \prec_\pi \text{nrR}\langle r, _, _, _ \rangle \\
& \wedge \forall w \in \text{nrW}, r \in \text{narR}. \\
& \quad (\text{nrW}\langle w, _ \rangle \prec_\pi \text{naF}\langle r, _, _, _ \rangle \wedge \text{sameqp}(w, r)) \implies \text{B}\langle w \rangle \prec_\pi \text{naF}\langle r, _, _, _ \rangle \\
& \wedge \forall w \in \text{nrW}, r \in \text{narR}. \\
& \quad (\text{nrW}\langle w, _ \rangle \prec_\pi \text{narR}\langle r, _, _, _, _ \rangle \wedge \text{sameqp}(w, r)) \implies \text{B}\langle w \rangle \prec_\pi \text{narR}\langle r, _, _, _, _ \rangle \\
& \wedge \forall w \in \text{narW}, r \in \text{nrR}. \\
& \quad (\text{narW}\langle w \rangle \prec_\pi \text{nrR}\langle r, _, _, _ \rangle \wedge \text{sameqp}(w, r)) \implies \text{B}\langle w \rangle \prec_\pi \text{nrR}\langle r, _, _, _ \rangle \\
& \wedge \forall w \in \text{narW}, r \in \text{narR}. \\
& \quad (\text{narW}\langle w \rangle \prec_\pi \text{naF}\langle r, _, _, _ \rangle \wedge \text{sameqp}(w, r)) \implies \text{B}\langle w \rangle \prec_\pi \text{naF}\langle r, _, _, _ \rangle \\
& \wedge \forall w \in \text{narW}, r \in \text{narR}. \\
& \quad (\text{narW}\langle w \rangle \prec_\pi \text{narR}\langle r, _, _, _, _ \rangle \wedge \text{sameqp}(w, r)) \implies \text{B}\langle w \rangle \prec_\pi \text{narR}\langle r, _, _, _, _ \rangle
\end{aligned}$$

[illegible]

$$\begin{aligned}
& \wedge \left(\begin{array}{l} a_1 \in (\text{nCAS} \cup \text{nFAA}) \wedge a_2 \in \text{Get} \wedge \text{nrR}\langle _, _, a_2, _ \rangle \in \pi \\ \implies \left(\begin{array}{l} a_1 \in \text{nCAS} \wedge \text{naF}\langle _, _, a_1, _ \rangle \prec_\pi \text{nrR}\langle _, _, a_2, w_2 \rangle \\ \vee \text{narR}\langle _, _, a_1, w_1, _ \rangle \prec_\pi \text{narW}\langle w_1 \rangle \prec_\pi \text{nrR}\langle _, _, a_2, w_2 \rangle \end{array} \right) \end{array} \right) \\
& \wedge \left(\begin{array}{l} a_1 \in (\text{nCAS} \cup \text{nFAA}) \wedge a_2 \in \text{Get} \wedge \text{nrR}\langle _, _, a_2, w_2 \rangle \prec_\pi \text{nlW}\langle w_2, _ \rangle \\ \implies \left(\begin{array}{l} a_1 \in \text{nCAS} \wedge \text{naF}\langle _, _, a_1, w_1 \rangle \prec_\pi \text{nlW}\langle w_1, _ \rangle \prec_\pi \text{nlW}\langle w_2, _ \rangle \\ \vee \text{narR}\langle _, _, a_1, w_1, _ \rangle \prec_\pi \text{nlW}\langle w_1, _ \rangle \prec_\pi \text{nlW}\langle w_2, _ \rangle \end{array} \right) \end{array} \right) \\
& \wedge \left(\begin{array}{l} a_1 \in (\text{nCAS} \cup \text{nFAA}) \wedge a_2 \in \text{Put} \wedge \text{nlR}\langle _, _, a_2, w_2 \rangle \prec_\pi \text{nrW}\langle w_2, _ \rangle \\ \implies \left(\begin{array}{l} a_1 \in \text{nCAS} \wedge \text{naF}\langle _, _, a_1, w_1 \rangle \prec_\pi \text{nrW}\langle w_2, _ \rangle \\ \vee \text{narR}\langle _, _, a_1, w_1, _ \rangle \prec_\pi \text{narW}\langle w_1 \rangle \prec_\pi \text{nrW}\langle w_2, _ \rangle \end{array} \right) \end{array} \right) \\
& \wedge \left(\begin{array}{l} a_1 \in (\text{nCAS} \cup \text{nFAA}) \wedge a_2 \in \text{Put} \wedge \text{nlR}\langle _, _, a_2, w_2 \rangle \prec_\pi \text{nrW}\langle w_2, e_2 \rangle \prec_\pi \text{CN}\langle e_2 \rangle \\ \implies \left(\begin{array}{l} a_1 \in \text{nCAS} \wedge \text{naF}\langle _, _, a_1, w_1 \rangle \prec_\pi \text{nlW}\langle w_1, _ \rangle \prec_\pi \text{CN}\langle e_2 \rangle \\ \vee \text{narR}\langle _, _, a_1, w_1, _ \rangle \prec_\pi \text{nlW}\langle w_1, _ \rangle \prec_\pi \text{CN}\langle e_2 \rangle \end{array} \right) \end{array} \right) \\
& \wedge \left(\begin{array}{l} a_1 \in (\text{nCAS} \cup \text{nFAA}) \wedge a_2 \in \text{nCAS} \wedge \text{naF}\langle _, _, a_2, _ \rangle \in \pi \\ \implies \left(\begin{array}{l} a_1 \in \text{nCAS} \wedge \text{naF}\langle _, _, a_1, _ \rangle \prec_\pi \text{naF}\langle _, _, a_2, _ \rangle \\ \vee \text{narR}\langle _, _, a_1, _, w_1 \rangle \prec_\pi \text{narW}\langle w_1 \rangle \prec_\pi \text{naF}\langle _, _, a_2, _ \rangle \end{array} \right) \end{array} \right) \\
& \wedge \left(\begin{array}{l} a_1, a_2 \in (\text{nCAS} \cup \text{nFAA}) \wedge \text{narR}\langle _, _, a_2, _, _ \rangle \in \pi \\ \implies \left(\begin{array}{l} a_1 \in \text{nCAS} \wedge \text{naF}\langle _, _, a_1, w_1 \rangle \prec_\pi \text{narR}\langle _, _, a_2, _, _ \rangle \\ \vee \text{narR}\langle _, _, a_1, _, w_1 \rangle \prec_\pi \text{narW}\langle w_1 \rangle \prec_\pi \text{narR}\langle _, _, a_2, _, _ \rangle \end{array} \right) \end{array} \right) \\
& \wedge \left(\begin{array}{l} a_1, a_2 \in (\text{nCAS} \cup \text{nFAA}) \wedge \left(\begin{array}{l} a_1 \in \text{nCAS} \wedge \text{naF}\langle _, _, a_2, w_2 \rangle \prec_\pi \text{nlW}\langle w_2, _ \rangle \\ \vee \text{narR}\langle _, _, a_2, w_2, _ \rangle \prec_\pi \text{nlW}\langle w_2, _ \rangle \end{array} \right) \\ \implies \left(\begin{array}{l} a_1 \in \text{nCAS} \wedge \text{naF}\langle _, _, a_1, w_1 \rangle \prec_\pi \text{nlW}\langle w_1, _ \rangle \prec_\pi \text{nlW}\langle w_2, _ \rangle \\ \vee \text{narR}\langle _, _, a_1, w_1, _ \rangle \prec_\pi \text{nlW}\langle w_1, _ \rangle \prec_\pi \text{nlW}\langle w_2, _ \rangle \end{array} \right) \end{array} \right)
\end{aligned}$$

NIC Atomicity

$$\text{nicAtomicity}(\pi) \triangleq \forall a_1, a_2, r, w.$$

$$\left(\begin{array}{l} \lambda_1 = \text{narR}\langle r_1, _, a_1, _, w \rangle \\ \wedge \lambda_2 \in \{ \text{naF}\langle _, _, a_2, _ \rangle, \text{narR}\langle _, _, a_2, _, _ \rangle \} \\ \wedge a_1, a_2 \in \text{nRMW} \wedge \bar{n}(a_1) = \bar{n}(a_2) \wedge \lambda_1 \prec_\pi \lambda_2 \end{array} \right) \implies B\langle w \rangle \prec_\pi \lambda_2$$

Read Order

$$\begin{aligned}
\text{wfrd}(\pi) & \triangleq \forall \pi_2, r, w, \pi_1. \pi = \pi_2 \cdot \text{lR}\langle r, w \rangle \cdot \pi_1 \implies \text{wfrdCPU}(r, w, \pi_1) \\
& \wedge \forall \pi_2, u, w, \pi_1. \pi = \pi_2 \cdot \text{CAS}\langle u, w \rangle \cdot \pi_1 \implies \text{wfrdCPU}(u, w, \pi_1) \\
& \wedge \forall \pi_2, r, w, \pi_1. \pi = \pi_2 \cdot \text{nlR}\langle r, w, _, _ \rangle \cdot \pi_1 \implies \text{wfrdNIC}(r, w, \pi_1) \\
& \wedge \forall \pi_2, r, w, \pi_1. \pi = \pi_2 \cdot \text{nrR}\langle r, w, _, _ \rangle \cdot \pi_1 \implies \text{wfrdNIC}(r, w, \pi_1) \\
& \wedge \forall \pi_2, r, w, \pi_1. \pi = \pi_2 \cdot \text{naF}\langle r, w, _, _ \rangle \cdot \pi_1 \implies \text{wfrdNIC}(r, w, \pi_1) \\
& \wedge \forall \pi_2, r, w, \pi_1. \pi = \pi_2 \cdot \text{narR}\langle r, w, _, _, _ \rangle \cdot \pi_1 \implies \text{wfrdNIC}(r, w, \pi_1)
\end{aligned}$$

$$\begin{aligned}
\text{wfrdCPU}(r, w, \pi) &\triangleq \left(\begin{array}{l} \exists \pi_2, \lambda, \pi_1. \pi = \pi_2 \cdot \lambda \cdot \pi_1 \\ \wedge \lambda \in \{\mathbf{B}\langle w \rangle, \mathbf{CAS}\langle w, _ \rangle\} \\ \wedge \{\mathbf{B}\langle w' \rangle, \mathbf{CAS}\langle w', _ \rangle \in \pi_2 \mid \text{loc}(w') = \text{loc}(r)\} = \emptyset \\ \wedge \left\{ w' \mid \begin{array}{l} \mathbf{1W}\langle w' \rangle \in \pi \wedge \mathbf{B}\langle w' \rangle \notin \pi \wedge \\ \text{loc}(w') = \text{loc}(r) \wedge t(w') = t(r) \end{array} \right\} = \emptyset \end{array} \right) \\
&\vee \left(\begin{array}{l} \exists \pi_2, \lambda, \pi_1. \pi = \pi_2 \cdot \lambda \cdot \pi_1 \\ \wedge \lambda = \mathbf{1W}\langle w \rangle \wedge t(w) = t(r) \wedge \mathbf{B}\langle w \rangle \notin \pi_2 \\ \wedge \{\mathbf{1W}\langle w' \rangle \in \pi_2 \mid \text{loc}(w') = \text{loc}(r) \wedge t(w') = t(r)\} = \emptyset \end{array} \right) \\
&\vee \left(\begin{array}{l} w = \text{init}_{\text{loc}(w)} \wedge \\ \left\{ \begin{array}{l} \mathbf{B}\langle w' \rangle, \mathbf{CAS}\langle w', _ \rangle \in \pi, \\ \mathbf{1W}\langle w'' \rangle \in \pi \end{array} \mid \begin{array}{l} \text{loc}(w') = \text{loc}(r) \wedge \\ \text{loc}(w'') = \text{loc}(r) \wedge t(w'') = t(r) \end{array} \right\} = \emptyset \end{array} \right) \\
\text{wfrdNIC}(r, w, \pi) &\triangleq \left(\begin{array}{l} \exists \pi_2, \lambda, \pi_1. \pi = \pi_2 \cdot \lambda \cdot \pi_1 \\ \wedge \lambda \in \{\mathbf{B}\langle w \rangle, \mathbf{CAS}\langle w, _ \rangle\} \\ \wedge \{\mathbf{B}\langle w' \rangle, \mathbf{CAS}\langle w', _ \rangle \in \pi_2 \mid \text{loc}(w') = \text{loc}(r)\} = \emptyset \end{array} \right) \\
&\vee \left(\begin{array}{l} w = \text{init}_{\text{loc}(w)} \wedge \\ \{\mathbf{B}\langle w' \rangle, \mathbf{CAS}\langle w', _ \rangle \in \pi \mid \text{loc}(w') = \text{loc}(r)\} = \emptyset \end{array} \right)
\end{aligned}$$

Well-formed path

$$\begin{aligned}
\text{wfp}(\pi) &\triangleq \quad \text{nodup}(\pi) \\
&\quad \wedge \text{backComp}(\pi) \\
&\quad \wedge \text{bufFlushOrd}(\pi) \\
&\quad \wedge \text{pollOrder}(\pi) \\
&\quad \wedge \text{nicActOrder}(\pi) \\
&\quad \wedge \text{nicAtomicity}(\pi) \\
&\quad \wedge \text{wfrd}(\pi)
\end{aligned}$$

Definition 5.

$$\begin{aligned}
\text{wf}(\mathbf{M}, \mathbf{B}, \mathbf{A}, \mathbf{QP}, \pi) &\triangleq \quad \text{wfp}(\pi) \\
&\quad \wedge \forall x \in \text{Loc}. \mathbf{M}(x) = \text{read}(\pi, x) \\
&\quad \wedge \forall t \in \text{Tid}. \mathbf{B}(t) = \text{mksbuff}(\varepsilon, t, \pi) \\
&\quad \wedge \forall n \in \text{Node}. \mathbf{A}(n) = \text{chkatm}(n, \pi) \\
&\quad \wedge \forall t \in \text{Tid}. \forall \bar{n} \in (\text{Node} \setminus \{n(t)\}). \left(\begin{array}{l} \mathbf{QP}(t)(\bar{n}).\text{pipe} = \text{mkpipe}(\varepsilon, t, \bar{n}, \pi) \\ \mathbf{QP}(t)(\bar{n}).\mathbf{wb}_R = \text{mkwbR}(\varepsilon, t, \bar{n}, \pi) \\ \mathbf{QP}(t)(\bar{n}).\mathbf{wb}_L = \text{mkwbL}(\varepsilon, t, \bar{n}, \pi) \end{array} \right)
\end{aligned}$$

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

$$\begin{aligned}
\text{read}(\lambda \cdot \pi, x) &\triangleq \begin{cases} w & \lambda \in \{\mathbf{B}\langle w \rangle, \mathbf{CAS}\langle w, _ \rangle\} \wedge \text{loc}(w) = x \\ \text{read}(\pi, x) & \text{otherwise} \end{cases} \\
\text{read}(\varepsilon, x) &\triangleq \text{init}_x
\end{aligned}$$

$$\text{mksbuff}(\mathbf{b}, t, \varepsilon) \triangleq \mathbf{b}$$

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

$$\text{chkatm}(n, \pi) \triangleq \begin{cases} \perp & \forall w. \left(\begin{array}{l} \mathbf{narR}\langle _, _, a, _, w \rangle \in \pi \\ \wedge \bar{n}(a) = n \end{array} \right) \implies \mathbf{B}\langle w \rangle \in \pi \\ \top & \text{otherwise} \end{cases}$$

$$\text{mkpipe}(\mathbf{pipe}, t, \bar{n}, \varepsilon) \triangleq \mathbf{pipe}$$

$$\text{mkpipe}(\mathbf{pipe}, t, \bar{n}, \pi \cdot \lambda) \triangleq \begin{cases} \text{mkpipe}(a \cdot \mathbf{pipe}, t, \bar{n}, \pi) & \text{if } \left(\begin{array}{l} t(\lambda) = t \wedge \bar{n}(\lambda) = \bar{n} \wedge \lambda = \mathbf{NIC}\langle a \rangle \\ \wedge \mathbf{n1R}\langle _, _, a, _ \rangle \notin \pi \wedge \mathbf{nrR}\langle _, _, a, _ \rangle \notin \pi \\ \wedge \mathbf{nF}\langle a \rangle \notin \pi \wedge \mathbf{naF}\langle _, _, a, _ \rangle \notin \pi \\ \wedge \mathbf{narR}\langle _, _, a, _, _ \rangle \notin \pi \end{array} \right) \\ \text{mkpipe}(w \cdot \mathbf{pipe}, t, \bar{n}, \pi) & \text{if } \left(\begin{array}{l} t(\lambda) = t \wedge \bar{n}(\lambda) = \bar{n} \wedge \lambda = \mathbf{NIC}\langle a \rangle \\ \wedge \mathbf{n1R}\langle _, _, a, w \rangle \in \pi \wedge \mathbf{nrW}\langle w, _ \rangle \notin \pi \end{array} \right) \\ \text{mkpipe}(e \cdot \mathbf{pipe}, t, \bar{n}, \pi) & \text{if } \left(\begin{array}{l} t(\lambda) = t \wedge \bar{n}(\lambda) = \bar{n} \wedge \lambda = \mathbf{NIC}\langle a \rangle \\ \wedge \mathbf{n1R}\langle _, _, a, w \rangle \in \pi \wedge \mathbf{nrW}\langle w, e \rangle \in \pi \\ \wedge \mathbf{CN}\langle e \rangle \notin \pi \end{array} \right) \\ \text{mkpipe}(w \cdot \mathbf{pipe}, t, \bar{n}, \pi) & \text{if } \left(\begin{array}{l} t(\lambda) = t \wedge \bar{n}(\lambda) = \bar{n} \wedge \lambda = \mathbf{NIC}\langle a \rangle \\ \wedge \mathbf{nrR}\langle _, _, a, w \rangle \in \pi \wedge \mathbf{n1W}\langle w, _ \rangle \notin \pi \end{array} \right) \\ \text{mkpipe}(w \cdot \mathbf{pipe}, t, \bar{n}, \pi) & \text{if } \left(\begin{array}{l} t(\lambda) = t \wedge \bar{n}(\lambda) = \bar{n} \wedge \lambda = \mathbf{NIC}\langle a \rangle \\ \wedge \mathbf{naF}\langle _, _, a, w \rangle \in \pi \wedge \mathbf{n1W}\langle w, _ \rangle \notin \pi \end{array} \right) \\ \text{mkpipe}(w \cdot w' \cdot \mathbf{pipe}, t, \bar{n}, \pi) & \text{if } \left(\begin{array}{l} t(\lambda) = t \wedge \bar{n}(\lambda) = \bar{n} \wedge \lambda = \mathbf{NIC}\langle a \rangle \\ \wedge \mathbf{narR}\langle _, _, a, w, w' \rangle \in \pi \wedge \mathbf{narW}\langle w' \rangle \notin \pi \end{array} \right) \\ \text{mkpipe}(\mathbf{pipe}, t, \bar{n}, \pi) & \text{otherwise} \end{cases}$$

$$\text{mkwbR}(\mathbf{wb}_R, t, \bar{n}, \varepsilon) \triangleq \mathbf{wb}_R$$

$$\text{mkwbR}(\mathbf{wb}_R, t, \bar{n}, \pi \cdot \lambda) \triangleq \begin{cases} \text{mkwbR}(w \cdot \mathbf{wb}_R, t, \bar{n}, \pi) & \text{if } \left(\begin{array}{l} t(\lambda) = t \wedge \bar{n}(\lambda) = \bar{n} \wedge \mathbf{B}\langle w \rangle \notin \pi \\ \wedge \lambda \in \{\mathbf{nrW}\langle w, _ \rangle, \mathbf{narW}\langle w \rangle\} \end{array} \right) \\ \text{mkwbR}(\mathbf{wb}_R, t, \bar{n}, \pi) & \text{otherwise} \end{cases}$$

$$\begin{aligned}
& \text{mkwbL}(\mathbf{wb}_L, t, \bar{n}, \varepsilon) \triangleq \mathbf{wb}_L \\
& \text{mkwbL}(\mathbf{wb}_L, t, \bar{n}, \pi \cdot \lambda) \triangleq \begin{cases} \text{mkwbL}(e \cdot w \cdot \mathbf{wb}_L, t, \bar{n}, \pi) & \text{if } \left(\begin{array}{l} t(\lambda) = t \wedge \bar{n}(\lambda) = \bar{n} \wedge \lambda = \mathbf{n1W}\langle w, e \rangle \\ \wedge \mathbf{B}\langle w \rangle \notin \pi \wedge \mathbf{P}\langle _, e \rangle \notin \pi \end{array} \right) \\ \text{mkwbL}(e \cdot \mathbf{wb}_L, t, \bar{n}, \pi) & \text{if } \left(\begin{array}{l} t(\lambda) = t \wedge \bar{n}(\lambda) = \bar{n} \wedge \lambda = \mathbf{n1W}\langle w, e \rangle \\ \wedge \mathbf{B}\langle w \rangle \in \pi \wedge \mathbf{P}\langle _, e \rangle \notin \pi \end{array} \right) \\ \text{mkwbL}(e \cdot \mathbf{wb}_L, t, \bar{n}, \pi) & \text{if } \left(\begin{array}{l} t(\lambda) = t \wedge \bar{n}(\lambda) = \bar{n} \wedge \lambda = \mathbf{CN}\langle e \rangle \\ \wedge \mathbf{P}\langle _, e \rangle \notin \pi \end{array} \right) \\ \text{mkwbL}(\mathbf{wb}_L, t, \bar{n}, \pi) & \text{otherwise} \end{cases}
\end{aligned}$$

Theorem 2. For all $P, P', M, M', B, B', A, A', QP, QP', \pi, \pi'$:

- $\text{wf}(M_0, B_0, A_0, QP_0, \varepsilon)$;
- if $P, M, B, A, QP, \pi \Rightarrow P', M', B', A', QP', \pi'$ and $\text{wf}(M, B, A, QP, \pi)$ then $\text{wf}(M', B', A', QP', \pi')$;
- if $P, M_0, B_0, A_0, QP_0, \varepsilon \Rightarrow^* (\lambda t.\text{skip}), M, B_0, A_0, QP, \pi$ such that forall t, \bar{n} we have $QP(t)(\bar{n}) = \langle \varepsilon, \varepsilon, \mathbf{nEX}^* \rangle$, then $\text{wf}(M, B_0, A_0, QP, \pi)$ and $\text{complete}(\pi)$.

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 \mathbf{nEX} in \mathbf{wb}_L) implies $\text{complete}(\pi)$.

B.3 From Annotated Semantics to Declarative Semantics

We define

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

with

$$\text{Event} \triangleq \text{Event}_0 \cup \{\text{getA}(\lambda) \mid \lambda \in \pi\}$$

Recall that Event_0 is the set of initialisation events $\{\text{init}_x \mid x \in \text{Loc}\}$, where $l(\text{init}_x) = 1W(x, 0)$

$$\text{getA} : \text{ALabel} \rightarrow \text{Event}$$

$$\begin{array}{ll} \text{getA}(1R\langle r, _ \rangle) \triangleq r & \text{getA}(\text{narR}\langle r, _, _, _, _ \rangle) \triangleq r \\ \text{getA}(1W\langle w \rangle) \triangleq w & \text{getA}(\text{narW}\langle w \rangle) \triangleq w \\ \text{getA}(\text{CAS}\langle u, _ \rangle) \triangleq u & \text{getA}(P\langle p, _ \rangle) \triangleq p \\ \text{getA}(F\langle f \rangle) \triangleq f & \text{getA}(\text{nF}\langle f \rangle) \triangleq f \\ \text{getA}(\text{n1R}\langle r, _, _, _ \rangle) \triangleq r & \text{getA}(B\langle w \rangle) \triangleq w \\ \text{getA}(\text{nrW}\langle w \rangle) \triangleq w & \text{getA}(\text{Push}\langle _ \rangle) \text{ is undefined} \\ \text{getA}(\text{nrR}\langle r, _, _, _ \rangle) \triangleq r & \text{getA}(\text{NIC}\langle _ \rangle) \text{ is undefined} \\ \text{getA}(\text{n1W}\langle w, _ \rangle) \triangleq w & \text{getA}(\text{CN}\langle _ \rangle) \text{ is undefined} \\ \text{getA}(\text{naF}\langle r, _, _, _ \rangle) \triangleq r & \text{getA}(\mathcal{E}\langle _ \rangle) \text{ is undefined} \end{array}$$

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

$$\text{getI}\lambda(_, \pi), \text{getO}\lambda(_, \pi) : \{\text{getA}(\lambda) \mid \lambda \in \pi\} \rightarrow \text{ALabel}$$

For all $\lambda \in \pi$:

- if $\text{type}(\lambda) \in \{1R, \text{CAS}, F, P, \text{n1R}, \text{nrR}, \text{narR}, \text{naF}, \text{nF}\}$,
then $\text{getI}\lambda(\text{getA}(\lambda), \pi) \triangleq \text{getO}\lambda(\text{getA}(\lambda), \pi) \triangleq \lambda$;
- if $\text{type}(\lambda) \in \{1W, \text{n1W}, \text{nrW}, \text{narW}\}$,
then $\text{getI}\lambda(\text{getA}(\lambda), \pi) \triangleq \lambda$ while $\text{getO}\lambda(\text{getA}(\lambda), \pi) \triangleq B\langle \lambda \rangle$.
- if $\lambda = B\langle w \rangle$, then from $\text{backComp}(\pi)$ there is $\lambda' \prec_\pi \lambda$ such that $\text{type}(\lambda') \in \{1W, \text{n1W}, \text{nrW}, \text{narW}\}$ and $\text{getA}(\lambda') = \text{getA}(\lambda) = w$. From the previous case, we have $\text{getI}\lambda(w, \pi) \triangleq \lambda'$ and $\text{getO}\lambda(w, \pi) \triangleq \lambda$.

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

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

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

From $\text{wfp}(\pi)$, **IB** and **OB** are transitive and irreflexive. Note: we could make **IB** and **OB** total by adding an arbitrary total order on Event_0 .

$$\mathbf{rf} \triangleq \left\{ (w, r) \mid \begin{array}{l} \text{IR}\langle r, w \rangle \in \pi \vee \text{nIR}\langle r, w, _, _ \rangle \in \pi \vee \text{nrR}\langle r, w, _, _ \rangle \in \pi \vee \text{CAS}\langle r, w \rangle \in \pi \\ \vee \text{narR}\langle r, w, _, _, _ \rangle \in \pi \vee \text{naF}\langle r, w, _, _ \rangle \in \pi \end{array} \right\}$$

$$\mathbf{pf} \triangleq \left\{ (w, p) \mid \begin{array}{l} \text{nIW}\langle w, e \rangle \prec_{\pi} \text{P}\langle p, e \rangle \\ \vee \text{nrW}\langle w, e \rangle \prec_{\pi} \text{P}\langle p, e \rangle \end{array} \right\}$$

$$\lambda \text{ generates } e \text{ in } \pi \triangleq \left(\begin{array}{l} \lambda \in \{\text{IR}\langle e, _ \rangle, \text{IW}\langle e \rangle, \text{CAS}\langle e, _ \rangle, \text{Push}\langle e \rangle, \text{P}\langle e, _ \rangle, \text{F}\langle e \rangle\} \\ \vee \lambda = \text{Push}\langle a \rangle \wedge \left(\begin{array}{l} \lambda \prec_{\pi} \text{nIR}\langle e, _, a, _ \rangle \\ \vee \lambda \prec_{\pi} \text{nIR}\langle _, _, a, e \rangle \\ \vee \lambda \prec_{\pi} \text{nrR}\langle e, _, a, _ \rangle \\ \vee \lambda \prec_{\pi} \text{nrR}\langle _, _, a, e \rangle \\ \vee \lambda \prec_{\pi} \text{naF}\langle e, _, a, _ \rangle \\ \vee \lambda \prec_{\pi} \text{naF}\langle _, _, a, e \rangle \\ \vee \lambda \prec_{\pi} \text{narR}\langle e, _, a, _, _ \rangle \\ \vee \lambda \prec_{\pi} \text{narR}\langle _, _, a, e, _ \rangle \\ \vee \lambda \prec_{\pi} \text{narR}\langle _, _, a, _, e \rangle \end{array} \right) \end{array} \right)$$

$$\mathbf{po} \triangleq \left(\begin{array}{l} \text{Event}_0 \times (\text{Event} \setminus \text{Event}_0) \\ \cup \left\{ (e_1, e_2) \mid \begin{array}{l} \lambda_1 \prec_{\pi} \lambda_2 \wedge t(\lambda_1) = t(\lambda_2) \\ \wedge \lambda_1 \text{ generates } e_1 \text{ in } \pi \\ \wedge \lambda_2 \text{ generates } e_2 \text{ in } \pi \end{array} \right\} \\ \cup \left\{ (r, w) \mid \begin{array}{l} \text{nIR}\langle r, _, _, w \rangle \in \pi \\ \vee \text{nrR}\langle r, _, _, w \rangle \in \pi \\ \vee \text{naF}\langle r, _, _, w \rangle \in \pi \\ \vee \text{narR}\langle r, _, _, _, w \rangle \in \pi \end{array} \right\} \\ \cup \{(w_1, w_2) \mid \text{narR}\langle _, _, _, w_2, w_1 \rangle \in \pi\} \end{array} \right)$$

$$\mathbf{mo} \triangleq \left\{ (w_1, w_2) \mid \begin{array}{l} w_1 = \text{init}_x \\ \wedge (\text{B}\langle w_2 \rangle \in \pi \vee \text{CAS}\langle w_2, _ \rangle \in \pi) \\ \wedge \text{loc}(w_1) = x = \text{loc}(w_2) \end{array} \right\} \cup \left\{ (w_1, w_2) \mid \begin{array}{l} \lambda_1 \prec_{\pi} \lambda_2 \\ \wedge \lambda_1 \in \{\text{B}\langle w_1 \rangle, \text{CAS}\langle w_1, _ \rangle\} \\ \wedge \lambda_2 \in \{\text{B}\langle w_2 \rangle, \text{CAS}\langle w_2, _ \rangle\} \\ \wedge \text{loc}(w_1) = \text{loc}(w_2) \end{array} \right\}$$

$$\mathbf{nfo} \triangleq \left(\begin{array}{l} \{(r, w) \mid \text{sameqp}(r, w) \wedge \text{nIR}\langle r, _, _, _ \rangle \prec_{\pi} \text{nIW}\langle w, _ \rangle \prec_{\pi} \text{B}\langle w \rangle\} \\ \cup \{(r, w) \mid \exists \lambda_r, \lambda_w. \text{sameqp}(r, w) \wedge \lambda_r \prec_{\pi} \lambda_w \prec_{\pi} \text{B}\langle w \rangle\} \\ \cup \{(w, r) \mid \text{sameqp}(w, r) \wedge \text{nIW}\langle w, _ \rangle \prec_{\pi} \text{B}\langle w \rangle \prec_{\pi} \text{nIR}\langle r, _, _, _ \rangle\} \\ \cup \{(w, r) \mid \exists \lambda_r, \lambda_w. \text{sameqp}(w, r) \wedge \lambda_w \prec_{\pi} \text{B}\langle w \rangle \prec_{\pi} \lambda_r\} \end{array} \right)$$

where $\lambda_r \in \{\text{nrR}\langle r', \dots \rangle, \text{naF}\langle r', \dots \rangle, \text{narR}\langle r', \dots \rangle\}$
 $\lambda_w \in \{\text{nrW}\langle w, _ \rangle, \text{narW}\langle w \rangle\}$

$$\text{rao} \triangleq \left(\left\{ (r_1, r_2) \mid \bar{n}(a_1) = \bar{n}(a_2) \wedge \left(\begin{array}{l} \lambda_1 \prec_{\pi} \lambda_2 \\ \wedge \lambda_1 \in \{\text{naF}\langle r_1, a_1, \dots \rangle, \text{narR}\langle r_1, a_1, \dots \rangle\} \\ \wedge \lambda_2 \in \{\text{naF}\langle r_2, a_2, \dots \rangle, \text{narR}\langle r_2, a_2, \dots \rangle\} \end{array} \right) \right\} \right)$$

From an execution graph $E = \text{getEG}(\pi)$, we use the definitions of the paper to define **oppo**, **ippo**, **rf_b**, **rf_{b̄}**, **rb**, **rb_b**, **ar**, **ob**, and **ib**.

Lemma 1. $w \in \text{nIW} \implies \exists r. \left(\begin{array}{l} r \in \text{nrR} \wedge (r, w) \in \text{po}|_{\text{imm}} \\ \vee r \in \text{narR} \wedge (r, w) \in \text{po}|_{\text{imm}} \cup (\text{po}|_{\text{imm}})^2 \end{array} \right)$

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

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

□

Theorem 3. $\text{getEG}(\pi)$ is well-formed.

Proof. We need to check the conditions of a pre-execution (Def. 2) and of well-formedness (Def. 3). For the pre-execution conditions:

- Checking $\text{Event}^0 \times (\text{Event} \setminus \text{Event}^0) \subseteq \text{po}$:
by definition.
- Checking **po** is a union of strict partial orders each on one thread:
If $t(e_1) \neq t(e_2)$, then $(e_1, e_2) \notin \text{po}$ and $(e_2, e_1) \notin \text{po}$ by definition. If $t(e_1) = t(e_2)$, then either $(e_1, e_2) \in \text{po}$ or $(e_2, e_1) \in \text{po}$. This comes from the second case of the definition of **po**: if there is λ_1 and λ_2 such that λ_i generates e_i in π , 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 $\text{lR}\langle r, _ \rangle$, $\text{nIR}\langle r, _, _, _ \rangle$, $\text{nrR}\langle r, _, _, _ \rangle$, $\text{naF}\langle r, _, _, _ \rangle$, or $\text{narR}\langle r, _, _, _ \rangle$ in π , and r have at least one antecedent.
If $(w, r) \in \text{rf}$, let us assume $r \in \text{nIR}$, then by definition $\text{nIR}\langle r, w, _, _ \rangle \in \pi$. Since $\text{nodup}(\pi)$, for all $w' \neq w$, we have $\text{nIR}\langle r, w', _, _ \rangle \notin \pi$, and syntactically we cannot write $\text{lR}\langle r, _ \rangle$ or $\text{nrR}\langle r, _, _, _ \rangle$, so $(w', r) \notin \text{rf}$. Similarly, $r \in \text{lR}$, $r \in \text{nrR}$, $r \in \text{naF}$ or $r \in \text{narR}$ only have one antecedent.
- Checking that **rf** relates events on the same location with matching values:
By syntactic definition of the annotated labels **lR**, **nIR**, **nrR**, **naF** and **narR**, e.g., $\text{lR}\langle r, w \rangle \implies \text{eq}_{\text{loc}\&\text{v}}(r, w)$.
- Checking that **mo** is a union of strict total orders for writes on each variables:
By definition of **mo**, given that we have $\text{complete}(\pi)$, e.g., if $\text{lW}\langle w \rangle \in \pi$ then $\text{B}\langle w \rangle \in \pi$.

- Checking that $\mathbf{pf} \subseteq \mathbf{po} \cap \mathbf{sqp}$:

If $(w, p) \in \mathbf{pf}$ with $w \in \mathbf{nlW}$ (resp. \mathbf{nrW}), then we have $\mathbf{nlW}\langle w, e \rangle \prec_{\pi} \mathbf{P}\langle p, e \rangle$. There is λ such that λ generates w in π , and we have $\lambda \prec_{\pi} \mathbf{nlW}\langle w, e \rangle \prec_{\pi} \mathbf{P}\langle p, e \rangle$. Also, $t(p) = t(w)$ and $\bar{n}(p) = \bar{n}(e) = \bar{n}(w)$, so we have $(w, p) \in \mathbf{po}$ and $(w, p) \in \mathbf{sqp}$.

- Checking that \mathbf{pf} is functional on its domain:

If $(w, p) \in \mathbf{pf}$ with $w \in \mathbf{nlW}$ (resp. \mathbf{nrW}), then we have $\mathbf{nlW}\langle w, e \rangle \prec_{\pi} \mathbf{P}\langle p, e \rangle$. From $\mathbf{nodup}(\pi)$, for all $p' \neq p$ we have $\mathbf{P}\langle p', e \rangle \notin \pi$, so w has at most one image.

- Checking that \mathbf{pf} is total and functional on its range:

If $p \in \mathbf{Event}$, then there is $e \in \mathbf{nlex}$ (resp. \mathbf{nrEX}) such that $\mathbf{P}\langle p, e \rangle \in \pi$. From $\mathbf{backComp}(\pi)$ there is $w \in \mathbf{nlW}$ (resp. \mathbf{nrW}) such that $\mathbf{nlW}\langle w, e \rangle \prec_{\pi} \mathbf{P}\langle p, e \rangle$, and so $(w, p) \in \mathbf{pf}$. From $\mathbf{nodup}(\pi)$, e cannot be used in another \mathbf{nlW} (resp. \mathbf{nrW}) annotated label, and p has exactly one antecedent.

- Checking that for all $(a, b) \in \mathbf{sqp}$, $a \in \mathbf{nrR} \cup \mathbf{naF} \cup \mathbf{narR}$, $b \in \mathbf{nrW} \cup \mathbf{narW}$, (resp. $\mathbf{nLR}/\mathbf{nLW}$) then $(a, b) \in \mathbf{nfo} \cup \mathbf{nfo}^{-1}$:

By definition of \mathbf{nfo} , given that $\mathbf{bufFlushOrd}(\pi)$ forbids such interleavings as $\mathbf{nrW}\langle w, _ \rangle \prec_{\pi} \mathbf{nrR}\langle r, _, _, _ \rangle \prec_{\pi} \mathbf{B}\langle w \rangle$ (resp. \mathbf{nLW} and \mathbf{nLR}) when $\mathbf{sameqp}(r, w)$.

- Checking that \mathbf{rao} is a union of strict total orders for remote atomic reads:

By definition of \mathbf{rao} .

For the well-formedness conditions:

- (1) Let us assume $(w_1, w_2) \in \mathbf{po} \cap \mathbf{sqp}$ and $(w_2, p_2) \in \mathbf{pf}$. The three events are on the same thread and queue pair.

If $w_1 \in \mathbf{nLW}$, then by $\mathbf{complete}(\pi)$ there is a chain $\mathbf{Push}\langle a_1 \rangle \prec_{\pi} \mathbf{NIC}\langle a_1 \rangle \prec_{\pi} \mathbf{nR} \prec_{\pi} \mathbf{nLW}\langle w_1, e_1 \rangle$ for some $\mathbf{nR} \in \{\mathbf{nrR}\langle _, _, a_1, w_1 \rangle, \mathbf{naF}\langle _, _, a_1, w_1 \rangle, \mathbf{narR}\langle _, _, a_1, w_1, _ \rangle\}$; if $w_1 \in \mathbf{nrW}$, there is instead a chain $\mathbf{Push}\langle a_1 \rangle \prec_{\pi} \mathbf{NIC}\langle a_1 \rangle \prec_{\pi} \mathbf{nLR}\langle _, _, a_1, w_1 \rangle \prec_{\pi} \mathbf{nrW}\langle w_1, e_1 \rangle \prec_{\pi} \mathbf{CN}\langle e_1 \rangle$. Similarly there is a chain for w_2 . By $(w_1, w_2) \in \mathbf{po}$ we have $\mathbf{Push}\langle a_1 \rangle \prec_{\pi} \mathbf{Push}\langle a_2 \rangle$, and by $\mathbf{bufFlushOrd}(\pi)$ we have $\mathbf{NIC}\langle a_1 \rangle \prec_{\pi} \mathbf{NIC}\langle a_2 \rangle$.

Let us call λ_1 the last annotated label on the chain for w_1 , i.e., either $\mathbf{nLW}\langle w_1, e_1 \rangle$ or $\mathbf{CN}\langle e_1 \rangle$. Similarly, λ_2 is the last annotated label on the chain for w_2 . There are four cases to consider, but in all four $\mathbf{nicActOrder}(\pi)$ implies $\lambda_1 \prec_{\pi} \lambda_2$.

Then, from $\mathbf{pollOrder}(\pi)$, there is p_1 such that $\mathbf{P}\langle p_1, e_1 \rangle \prec_{\pi} \mathbf{P}\langle p_2, e_2 \rangle$. By definitions, we have both $(w_1, p_1) \in \mathbf{pf}$ and $(p_1, p_2) \in \mathbf{po}$.

- (2) If $r \in \mathbf{nLR}$, then there is $w \in \mathbf{nrW}$ (taken from $\mathbf{nLR}\langle r, _, _, w \rangle$) such that $(r, w) \in \mathbf{po}|_{\mathbf{imm}}$. This is by the last case of definition of \mathbf{po} , since there is λ_a such that we have both λ_a generates r in π and λ_a generates w in π . Similarly for $\mathbf{nrR}/\mathbf{nLW}$ and $\mathbf{nrW}/\mathbf{nLR}$.

- (3) If $(r, w) \in \mathbf{po}|_{\mathbf{imm}}$, $\mathbf{type}(r) \in \{\mathbf{nLR}, \mathbf{nrR}\}$, and $\mathbf{type}(w) \in \{\mathbf{nLW}, \mathbf{nrW}\}$, then $(r, w) \in \mathbf{po}$ comes from the third case of the definition of \mathbf{po} , and we have either $\mathbf{nLR}\langle r, _, _, w \rangle$ or $\mathbf{nrR}\langle r, _, _, w \rangle$ in π . In both cases, we have $v_r(r) = v_w(w)$ by syntactic definition of the annotated labels.

- (4) (a) If $r \in \text{narR}$, then either: There is $\text{naF}\langle r, _, _, w \rangle \in \pi$ such that $w \in \text{nlW}$ and $(r, w) \in \text{po}|_{\text{imm}}$. This follows from the second case definition of po . There is $\text{narR}\langle r, _, _, w_2, w_1 \rangle \in \pi$ such that $w_1 \in \text{narW}$, $w_2 \in \text{nlW}$, and $(r, w_1), (w_1, w_2) \in \text{po}|_{\text{imm}}$. This follows from the second and third cases of the definition of po since there is λ_a which generates r , w_1 and w_2 in π . (b) If $w \in \text{narW}$ then $(r, w), (w, w') \in \text{po}|_{\text{imm}}$ with $r \in \text{narR}$ and $w' \in \text{nlW}$ comes from the second case definition of po .
- (5) If $(r, w) \in G.\text{po}|_{\text{imm}}$, $\text{type}(r) = \text{narR}$ and $\text{type}(w) = \text{nlW}$, then (r, w) comes from the second case definition of po and we have $\text{naF}\langle r, _, _, w \rangle \in \pi$. Then $v_r(r) = v_w(w)$ by the syntax of annotated labels. If $(r, w_1), (w_1, w_2) \in G.\text{po}|_{\text{imm}}$, $\text{type}(r) = \text{narR}$, $\text{type}(w_1) = \text{narW}$ and $\text{type}(w_2) = \text{nlW}$, then (r, w_1) comes from the second case definition of po and (w_1, w_2) from the third case, so we have $\text{narR}\langle r, _, _, w_2, w_1 \rangle \in \pi$. Then $v_r(r) = v_w(w_2)$ by the syntax of annotated labels.
- (6) Comes from lemma 1.

□

Lemma 2. $\text{OB}; [\text{Inst}] \subseteq \text{IB}$ and $[\text{Inst}]; \text{IB} \subseteq \text{OB}$.

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

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

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

□

Theorem 4. $\text{getEG}(\pi)$ is consistent.

Proof. From Definition 4, we need to check that both ib and ob are irreflexive. Since IB and OB are irreflexive, it is enough to show that $\text{ib} \subseteq \text{IB}$ and $\text{ob} \subseteq \text{OB}$.

The explicit definition using limits is the following (where $\text{rf}_{\bar{b}} \triangleq (\text{rf} \setminus \text{rf}_b)$ includes $(\text{rf} \cap \text{spp})$ since we assume the PCIe guarantees hold):

$$\begin{aligned}
 \text{ib}^0 &\triangleq (\text{ippo} \cup \text{rf} \cup \text{pf} \cup \text{rb}_b \cup \text{nfo})^+ \\
 \text{ob}^0 &\triangleq (\text{oppo} \cup \text{rf}_{\bar{b}} \cup [\text{nlW}]; \text{pf} \cup \text{rb} \cup \text{nfo} \cup \text{mo} \cup \text{ar}; \text{rao})^+ \\
 \text{ib}^{n+1} &\triangleq (\text{ib}^n \cup \text{ob}^n; [\text{Inst}])^+ \\
 \text{ob}^{n+1} &\triangleq (\text{ob}^n \cup [\text{Inst}]; \text{ib}^n)^+ \\
 \text{ib} &\triangleq \lim_{n \rightarrow \infty} \text{ib}^n \\
 \text{ob} &\triangleq \lim_{n \rightarrow \infty} \text{ob}^n
 \end{aligned}$$

It is then enough to show that $\text{ib}^0 \subseteq \text{IB}$ and $\text{ob}^0 \subseteq \text{OB}$. Using Lemma 2 above, we can check the induction case:

$$\begin{aligned}
 \text{ib}^{n+1} &= (\text{ib}^n \cup \text{ob}^n; [\text{Inst}])^+ \subseteq (\text{ib}^n \cup \text{OB}; [\text{Inst}])^+ \subseteq (\text{IB} \cup \text{IB})^+ = \text{IB} \\
 \text{ob}^{n+1} &= (\text{ob}^n \cup [\text{Inst}]; \text{ib}^n)^+ \subseteq (\text{ob}^n \cup [\text{Inst}]; \text{IB})^+ \subseteq (\text{OB} \cup \text{OB})^+ = \text{OB}
 \end{aligned}$$

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 $\text{ippo} \subseteq \text{IB}$.

Let $E^{\text{cpu}} = \{1R, 1W, \text{CAS}, F, P\}$ and $E^{\text{nic}} = \{n1R, nrR, narR, naF, n1W, nrW, narW, nF\}$. $[E^{\text{cpu}}]; \text{po} \subseteq \text{IB}$ by definition of po and IB : E^{cpu} are the events for which the same annotated label is used to define po and IB , i.e., $\forall e \in E^{\text{cpu}}, \text{get}\lambda(e, \pi)$ generates e in π . 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 $\text{oppo} \subseteq \text{OB}$.

From above we have $[\text{Inst}]; \text{oppo} \subseteq [\text{Inst}]; \text{ippo} \subseteq [\text{Inst}]; \text{IB} \subseteq \text{OB}$.

$[1W]; \text{po}; [\text{Event} \setminus (1R \cup P)] \subseteq \text{OB}$ by using $\text{bufFlushOrd}(\pi)$.

For the remaining cases:

- (G7) $[\text{nrW}]; (\text{po} \cap \text{sqp}); [\text{nrW}] \subseteq \text{OB}$ comes from $\text{nicActOrder}(\pi)$ (i.e., $\text{nrW}\langle \dots \rangle \prec_{\pi} \text{nrW}\langle \dots \rangle$) and $\text{bufFlushOrd}(\pi)$ (i.e., $B\langle \dots \rangle \prec_{\pi} B\langle \dots \rangle$).
- (G8) $[\text{nrW}]; (\text{po} \cap \text{sqp}); [\text{narR}] \subseteq \text{OB}$ comes from $\text{nicActOrder}(\pi)$ (i.e., $\text{nrW}\langle \dots \rangle \prec_{\pi} \text{narR}\langle \dots \rangle$) and $\text{bufFlushOrd}(\pi)$ (i.e., $B\langle \dots \rangle \prec_{\pi} \text{narR}\langle \dots \rangle$).
- (G9) If $e_1 \in \text{nrW}$, $e_3 \in \text{narW}$, and $(e_1, e_3) \in (\text{po} \cap \text{sqp})$, then from Def. 3 there is $e_2 \in \text{narR}$ such that $(e_2, e_3) \in \text{po}|_{\text{imm}}$ and thus $(e_1, e_2) \in (\text{po} \cap \text{sqp})$. From case G8 above, we have $(e_1, e_2) \in \text{OB}$. From $\text{backComp}(\pi)$, we have $(e_2, e_3) \in [\text{Inst}]; \text{IB} \subseteq \text{OB}$. Thus $[\text{nrW}]; (\text{po} \cap \text{sqp}); [\text{narW}] \subseteq \text{OB}$.
- (G10) $[\text{nrW}]; (\text{po} \cap \text{sqp}); [\text{nrR}] \subseteq \text{OB}$ comes from $\text{nicActOrder}(\pi)$ (i.e., $\text{nrW}\langle \dots \rangle \prec_{\pi} \text{nrR}\langle \dots \rangle$) and $\text{bufFlushOrd}(\pi)$ (i.e., $\text{nrW}\langle \dots \rangle \prec_{\pi} B\langle \dots \rangle \prec_{\pi} \text{nrR}\langle \dots \rangle$).
- (G11) If $e_1 \in \text{nrW}$, $e_3 \in \text{n1W}$, and $(e_1, e_3) \in (\text{po} \cap \text{sqp})$, then from Def. 3 there is $e_2 \in (\text{narR} \cup \text{nrR})$ such that $(e_2, e_3) \in \text{po}|_{\text{imm}}^{\{1,2\}}$ and thus $(e_1, e_2) \in (\text{po} \cap \text{sqp})$. Then $(e_1, e_2) \subseteq \text{OB}$ comes from cases G9 and G10 respectively. From $\text{backComp}(\pi)$, we have $(e_2, e_3) \in [\text{Inst}]; \text{IB} \subseteq \text{OB}$. Thus $[\text{nrW}]; (\text{po} \cap \text{sqp}); [\text{n1W}] \subseteq \text{OB}$.
- (I7) $[\text{narW}]; (\text{po} \cap \text{sqp}); [\text{nrW}] \subseteq \text{OB}$ comes from $\text{nicActOrder}(\pi)$ (i.e., $\text{narW}\langle \dots \rangle \prec_{\pi} \text{nrW}\langle \dots \rangle$) and $\text{bufFlushOrd}(\pi)$ (i.e., $B\langle \dots \rangle \prec_{\pi} B\langle \dots \rangle$).
- (I8) $[\text{narW}]; (\text{po} \cap \text{sqp}); [\text{narR}] \subseteq \text{OB}$ follows from Def. 3, $\text{nicActOrder}(\pi)$ and $\text{bufFlushOrd}(\pi)$ by similar reasoning to I7.
- (I9) $[\text{narW}]; (\text{po} \cap \text{sqp}); [\text{narW}] \subseteq \text{OB}$ follows similarly to I7.
- (I10) $[\text{narW}]; (\text{po} \cap \text{sqp}); [\text{nrR}] \subseteq \text{OB}$ follows similarly to I7.
- (K11) $[\text{n1W}]; (\text{po} \cap \text{sqp}); [\text{n1W}] \subseteq \text{OB}$ comes from $\text{nicActOrder}(\pi)$ (i.e., $\text{n1W}\langle \dots \rangle \prec_{\pi} \text{n1W}\langle \dots \rangle$) and $\text{bufFlushOrd}(\pi)$ (i.e., $B\langle \dots \rangle \prec_{\pi} B\langle \dots \rangle$).

- Checking $\text{rf}_{\bar{b}} \subseteq \text{OB}$.

If $(w, r) \in \text{rf}_{\bar{b}}$, there is π_1 and π_2 such that $\pi = \pi_2 \cdot \text{getO}\lambda(r, \pi) \cdot \pi_1$, and we use $\text{wfrd}(\pi)$.

- If $r \in 1R$, we have $\text{wfrdCPU}(r, w, \pi_1)$. The definition allow for three different cases. In the first case, $\lambda \in \{B\langle w \rangle, \text{CAS}\langle w, _ \rangle\}$ is in π_1 ; we have $\lambda = \text{getO}\lambda(w, \pi) \prec_{\pi} \text{getO}\lambda(r, \pi)$ and so $(w, r) \in \text{OB}$. In the second case, we have $\lambda = 1W\langle w \rangle$ and $t(w) = t(r)$; so $(w, r) \in [1W]; (\text{rf} \cap \text{sthd}); [1R] = \text{rf}_{\bar{b}}$, which contradicts $(w, r) \in \text{rf}_{\bar{b}} = \text{rf} \setminus \text{rf}_{\bar{b}}$. In the third case, $w = \text{init}_x$ for some location x , so $(w, r) \in \text{Event}_0 \times (\text{Event} \setminus \text{Event}_0) \subseteq \text{OB}$.

- If $r \in \mathbf{CAS}$, similarly to above, except the second case of $\mathbf{wfrdCPU}(r, w, \pi_1)$ is not possible because of $\mathbf{bufFlushOrd}(\pi)$: $\mathbf{B}\langle w \rangle \notin \pi_1$ while \mathbf{CAS} acts as a memory fence.
- If $r \in \mathbf{n1R}$, we have $\mathbf{wfrdNIC}(r, w, \pi_1)$, with two possibilities. In the first case, $\lambda \in \{\mathbf{B}\langle w \rangle, \mathbf{CAS}\langle w, _ \rangle\}$ is in π_1 ; we have $\lambda = \mathbf{getOL}(w, \pi) \prec_\pi \mathbf{getOL}(r, \pi)$ and so $(w, r) \in \mathbf{OB}$. In the second case, $w = \mathit{init}_x$ for some location x , so $(w, r) \in \mathbf{Event}_0 \times (\mathbf{Event} \setminus \mathbf{Event}_0) \subseteq \mathbf{OB}$.
- If $r \in \mathbf{nrR}$ or \mathbf{narR} , similarly to above.
- Checking $\mathbf{rf} \subseteq \mathbf{IB}$.
 From above we have $\mathbf{rf}_b = \mathbf{rf}_b; [\mathbf{Inst}] \subseteq \mathbf{OB}; [\mathbf{Inst}] \subseteq \mathbf{IB}$.
 If $(w, r) \in \mathbf{rf}_b \subseteq [\mathbf{1W}]; \mathbf{rf}; [\mathbf{1R}]$, then there is $\mathbf{1R}\langle r, w \rangle \in \pi$. There is π_1 and π_2 such that $\pi = \pi_2 \cdot \mathbf{1R}\langle r, w \rangle \cdot \pi_1$. So by $\mathbf{wfrd}(\pi)$ we have $\mathbf{wfrdCPU}(r, w, \pi_1)$ which implies $\mathbf{1W}\langle w \rangle \prec_\pi \mathbf{1R}\langle r, w \rangle$ and $(w, r) \in \mathbf{IB}$.
- Checking $[\mathbf{n1W}]; \mathbf{pf} \subseteq \mathbf{OB}$.
 If $(w, p) \in \mathbf{pf}$ with $w \in \mathbf{n1W}$, then there exists e such that $\mathbf{n1W}\langle w, e \rangle \prec_\pi \mathbf{P}\langle p, e \rangle$. From $\mathbf{backComp}(\pi)$, we have $\mathbf{n1W}\langle w, e \rangle \prec_\pi \mathbf{B}\langle w \rangle \prec_\pi \mathbf{P}\langle p, e \rangle$ and so $(w, p) \in \mathbf{OB}$.
- Checking $\mathbf{pf} \subseteq \mathbf{IB}$.
 If $(w, p) \in \mathbf{pf}$, then there exists e such that either $\mathbf{n1W}\langle w, e \rangle \prec_\pi \mathbf{P}\langle p, e \rangle$ or $\mathbf{nrW}\langle w, e \rangle \prec_\pi \mathbf{P}\langle p, e \rangle$. In both cases we immediately have $(w, p) \in \mathbf{IB}$.
- Checking $\mathbf{rb}_b \subseteq \mathbf{IB}$.
 If $(r, w') \in \mathbf{rb}_b$ then $r \in \mathbf{1R}$, $w' \in \mathbf{1W}$, $t(r) = t(w')$, and there exists w such that $(w, r) \in \mathbf{rf}$ and $(w, w') \in \mathbf{mo}$. There is π_4 and π_3 such that $\pi = \pi_4 \cdot \mathbf{1R}\langle r, w \rangle \cdot \pi_3$. So by $\mathbf{wfrd}(\pi)$ we have $\mathbf{wfrdCPU}(r, w, \pi_3)$, and there are three cases to consider.
 - In the first case, $\pi_3 = \pi_2 \cdot \lambda_w \cdot \pi_1$, with $\lambda_w \in \{\mathbf{B}\langle w \rangle, \mathbf{CAS}\langle w, _ \rangle\}$, and $\mathbf{B}\langle w' \rangle \notin \pi_2$. Since $(w, w') \in \mathbf{mo}$ we have $\mathbf{B}\langle w' \rangle \notin \pi_1$, and so $\mathbf{B}\langle w' \rangle \notin \pi_3$. The last condition of the first case then gives us $\mathbf{1W}\langle w' \rangle \notin \pi_3$, which implies $(r, w') \in \mathbf{IB}$.
 - In the second case, $\pi_3 = \pi_2 \cdot \lambda_w \cdot \pi_1$, with $\lambda_w = \mathbf{1W}\langle w \rangle$, $\mathbf{thread}(w) = \mathbf{thread}(r)$, and $\mathbf{B}\langle w \rangle \notin \pi_3$. Then w and w' are on the same thread, and by $\mathbf{bufFlushOrd}(\pi)$ and $(w, w') \in \mathbf{mo}$ we have $\mathbf{1W}\langle w \rangle \prec_\pi \mathbf{1W}\langle w' \rangle$ and $\mathbf{1W}\langle w' \rangle \notin \pi_1$. The last condition of the second case gives us $\mathbf{1W}\langle w' \rangle \notin \pi_2$, so $\mathbf{1W}\langle w' \rangle \notin \pi_3$ and $(r, w') \in \mathbf{IB}$.
 - In the last case, $w = \mathit{init}_x$ for some location x , and we immediately get $\mathbf{1W}\langle w' \rangle \notin \pi_3$, which implies $(r, w') \in \mathbf{IB}$.
- Checking $\mathbf{rb} \subseteq \mathbf{OB}$.
 If $(r, w') \in \mathbf{rb}$, then there exists w such that $(w, r) \in \mathbf{rf}$ and $(w, w') \in \mathbf{mo}$. By definition of \mathbf{rf} , there is π_4 and π_3 such that $\pi = \pi_4 \cdot \lambda_r \cdot \pi_3$, with $\lambda_r \in \{\mathbf{1R}\langle r, w \rangle, \mathbf{CAS}\langle r, w \rangle, \mathbf{n1R}\langle r, w, _, _ \rangle, \mathbf{nrR}\langle r, w, _, _ \rangle, \mathbf{naF}\langle r, w, _, _ \rangle, \mathbf{narR}\langle r, w, _, _, _ \rangle\}$. So by $\mathbf{wfrd}(\pi)$ we have either $\mathbf{wfrdNIC}(r, w, \pi_3)$ or $\mathbf{wfrdCPU}(r, w, \pi_3)$, and there are five cases to consider.
 - In the first case of $\mathbf{wfrdNIC}(r, w, \pi_3)$, $\pi_3 = \pi_2 \cdot \mathbf{getOL}(w, \pi) \cdot \pi_1$, and $\mathbf{getOL}(w', \pi) \notin \pi_2$. Since $(w, w') \in \mathbf{mo}$ we have $\mathbf{getOL}(w', \pi) \notin \pi_1$, and thus $\mathbf{getOL}(w', \pi) \notin \pi_3$. So $\mathbf{getOL}(w', \pi) \in \pi_4$ and $(r, w') \in \mathbf{OB}$.

- In the last case $\text{wfrdNIC}(r, w, \pi_3)$, $w = \text{init}_x$ for some location x , and we immediately have $\text{getO}\lambda(w', \pi) \notin \pi_3$, which implies $(r, w') \in \text{OB}$.
 - For the first case of $\text{wfrdCPU}(r, w, \pi_3)$, same reasoning as for the first case of wfrdNIC .
 - For the second case of $\text{wfrdCPU}(r, w, \pi_3)$, $\pi_3 = \pi_2 \cdot \text{getl}\lambda(w, \pi) \cdot \pi_1$, with $\text{thread}(w) = \text{thread}(r)$, and $\text{getO}\lambda(w, \pi) \notin \pi_3$. So $\text{getO}\lambda(w, \pi) \in \pi_4$, and since $(w, w') \in \text{mo}$ we have $\text{getO}\lambda(w', \pi) \in \pi_4$ as well, and $(r, w') \in \text{OB}$.
 - For the last case of $\text{wfrdCPU}(r, w, \pi_3)$, same reasoning as for the last case of wfrdNIC .
- Checking $\text{nfo} \subseteq \text{IB}$.
By definition of nfo .
 - Checking $\text{nfo} \subseteq \text{OB}$.
By definition of nfo .
 - Checking $\text{mo} \subseteq \text{OB}$.
By definition of mo , as what matters are the init_x , $\text{B}\langle w \rangle$, and $\text{CAS}\langle w, _ \rangle$ events.
 - Checking $\text{rao} \subseteq \text{IB}$.
By definition of rao .
 - Checking $\text{ar}; \text{rao} \subseteq \text{OB}$.
If $(w, r_1) \in \text{ar}$ then $\text{narR}\langle r_1, a_1, _, _, w \rangle \in \pi$ for some $a_1 \in \text{nRMW}$, and if $(r_1, r_2) \in \text{rao}$ then $\text{narR}\langle r_1, _, a_1, _, w \rangle \prec_\pi \lambda_r$ for some $\lambda_r \in \{\text{naF}\langle r_2, _, a_2, _ \rangle, \text{narR}\langle r_2, _, a_2, _, _ \rangle\}$, with $\bar{n}(a_1) = \bar{n}(a_2)$. Then using $\text{nicAtomicity}(\pi)$ we have that $\text{B}\langle w \rangle \prec_\pi \lambda_r$.

□

B.4 From Declarative Semantics to Annotated Semantics

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

Theorem 5. ib and ob can be extended into total relations IB and OB on Event such that:

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

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 $(a, b) \in \text{Event}^2$ such that $(a, b) \notin \text{ib}$ and $(b, a) \notin \text{ib}$. We arbitrarily decide to include (a, b) in our relation and we define $\text{ib}' = (\text{ib} \cup \{(a, b)\})^+$ and $\text{ob}' = (\text{ob} \cup [\text{Inst}]; \text{ib}')^+$.

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

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

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

Once ib is a total relation on Event , we can similarly freely extend ob into a total relation. \square

We use Theorem 5 above to extend ib and ob into total relations IB and OB .

Since $(\text{Event}, \text{po})$ is derived from P , by §5.2 we have that for all $t \in \text{Tid}$ there are s_t and G_t such that $G_t \in G^t(s_t)$, $P(t) \mapsto s_t$ and $(\text{Event}, \text{po}) = G_{\text{init}}; (\|_{t \in \text{Tid}} G_t)$. We consider each premise of the form $C \mapsto s$, where C is a primitive command, to generate new events and annotated labels.

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

- If $s = f \in \mathbf{F}$, we create annotated labels $\mathbf{F}\langle f \rangle$.
- If $s = r, w$ where $r \in \mathbf{nlR}$ and $w \in \mathbf{nrW}$, we create two events $a \in \mathbf{Put}$ and $e \in \mathbf{nrEX}$, and the annotated labels $\mathbf{Push}\langle a \rangle$, $\mathbf{NIC}\langle a \rangle$, $\mathbf{nlR}\langle r, w', a, w \rangle$ (where $(w', r) \in \mathbf{rf}$), $\mathbf{nrW}\langle w, e \rangle$, $\mathbf{B}\langle w \rangle$, and $\mathbf{CN}\langle e \rangle$. If there is p such that $(w, p) \in \mathbf{pf}$, we also create an annotated label $\mathbf{P}\langle p, e \rangle$. To simplify later definition, we also extend \mathbf{po} such that the event a is placed just before r , and e just after w . I.e., let $\mathbf{po}' = \mathbf{po} \cup \{(e', a) \mid (e', r) \in \mathbf{po}\} \cup \{(a, e') \mid (r, e') \in \mathbf{po}^*\}$ and redefine $\mathbf{po} = \mathbf{po}' \cup \{(e', e) \mid (e', w) \in \mathbf{po}'^*\} \cup \{(e, e') \mid (w, e') \in \mathbf{po}'\}$.
Note: from well-formedness conditions, every \mathbf{nlR} and every \mathbf{nrW} are part of such a pair.
- If $s = r, w$ where $r \in \mathbf{nrR}$ and $w \in \mathbf{nlW}$, we similarly create $a \in \mathbf{Get}$, $e \in \mathbf{nlEX}$, $\mathbf{Push}\langle a \rangle$, $\mathbf{NIC}\langle a \rangle$, $\mathbf{nrR}\langle \dots \rangle$, $\mathbf{nlW}\langle \dots \rangle$, $\mathbf{B}\langle \dots \rangle$, and potentially $\mathbf{P}\langle \dots \rangle$.
- If $s = r, w$ where $r \in \mathbf{narR}$ and $w \in \mathbf{nlW}$, we have C of the form $z := \mathbf{nCAS}(\bar{x}, e, e')$, so we use the values $\llbracket e \rrbracket$ and $\llbracket e' \rrbracket$ to create $a \in \mathbf{nCAS}$, $\mathbf{Push}\langle a \rangle$, $\mathbf{NIC}\langle a \rangle$, $\mathbf{naF}\langle \dots \rangle$, $\mathbf{nlW}\langle \dots \rangle$, $\mathbf{B}\langle \dots \rangle$, and potentially $\mathbf{P}\langle \dots \rangle$.
- If $s = r, w_1, w_2$ where $r \in \mathbf{narR}$, $w_1 \in \mathbf{narW}$, $w_2 \in \mathbf{nlW}$, we have C either of the form $z := \mathbf{nFAA}(\bar{x}, e)$ or $z := \mathbf{nCAS}(\bar{x}, e_1, e_2)$, so we create $a \in \mathbf{nFAA}$ or $a \in \mathbf{nCAS}$ accordingly, and $\mathbf{Push}\langle a \rangle$, $\mathbf{NIC}\langle a \rangle$, $\mathbf{narR}\langle \dots \rangle$, $\mathbf{narW}\langle w_1 \rangle$, $\mathbf{nlW}\langle w_2, \dots \rangle$, $\mathbf{B}\langle w_1 \rangle$, $\mathbf{B}\langle w_2 \rangle$ and potentially $\mathbf{P}\langle \dots \rangle$.
- If $s = f \in \mathbf{nF}$, we create the annotated labels $\mathbf{Push}\langle f \rangle$, $\mathbf{NIC}\langle f \rangle$, and $\mathbf{nF}\langle f \rangle$.
- We ignore $s = p \in \mathbf{P}$, 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 π_0 such that:

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

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 $1W\langle w \rangle / \mathbf{nlW}\langle w, _ \rangle / \mathbf{nrW}\langle w, _ \rangle / \mathbf{narW}\langle w \rangle$ first and $\mathbf{B}\langle w \rangle$ last. Note that the annotated labels $\mathbf{Push}\langle \dots \rangle$, $\mathbf{NIC}\langle \dots \rangle$, and $\mathbf{CN}\langle \dots \rangle$ not covered by **IB**/**OB** are not yet integrated in π_0 .

Then we extend π_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.

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

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

$$E^{\text{ext}} \triangleq \mathbf{Event} \cup (\mathbf{Get} \cup \mathbf{Put} \cup \mathbf{nCAS} \cup \mathbf{nFAA} \cup \mathbf{nlEX} \cup \mathbf{nrEX})$$

$$\begin{aligned} \text{getCPU} &: E^{\text{ext}} \rightarrow \text{ALabel} \\ \text{getCPU}(e) &\triangleq \begin{cases} \text{getl}\lambda(e, \pi_0) & \text{if } e \in E^{\text{cpu}} = \{1R, 1W, \text{CAS}, F, P\} \\ \text{Push}\langle e \rangle & \text{if } e \in \{\text{Put}, \text{Get}, \text{nCAS}, \text{nFAA}, \text{nF}\} \\ \text{undefined} & \text{otherwise} \end{cases} \end{aligned}$$

And a similar function for events emptying a CPU buffer:

$$\begin{aligned} \text{getTSO} &: E^{\text{ext}} \rightarrow \text{ALabel} \\ \text{getTSO}(e) &\triangleq \begin{cases} B\langle e \rangle & \text{if } e \in 1W \\ \text{NIC}\langle e \rangle & \text{if } e \in \{\text{Put}, \text{Get}, \text{nCAS}, \text{nFAA}, \text{nF}\} \\ \text{undefined} & \text{otherwise} \end{cases} \end{aligned}$$

Let us consider $(a_1, \dots, a_n) = \text{Event} \cap \{\text{Put}, \text{Get}, \text{nCAS}, \text{nFAA}, \text{nF}\}$ in po order, i.e., if $i < j$ then $(a_j, a_i) \notin \text{po}$. We extend π_0 successively until we get π_n :

- We introduce **Push** as early as possible:
Let $\pi' = \text{extend}(\pi_{i-1}, \text{Push}\langle a_i \rangle, \{\text{getCPU}(e) \mid (e, a_i) \in \text{po}\})$
- We introduce **NIC** as early as possible:
Let $\pi'' = \text{extend}(\pi', \text{NIC}\langle a_i \rangle, \{\text{Push}\langle a_i \rangle\} \cup \{\text{getTSO}(e) \mid (e, a_i) \in \text{po}\})$
- If $a_i \in \text{Put}$, there is $e_i \in \text{nrEX}$ such that $\text{n1R}\langle _, _, a_i, w \rangle \prec_{\pi_0} \text{nrW}\langle w, e_i \rangle$. We also introduce **CN**: Let $S = \{\text{nrW}\langle w, e_i \rangle\} \cup \{\text{n1W}\langle _, e \rangle \mid (e, e_i) \in \text{po} \cap \text{sqp}\} \cup \{\text{CN}\langle e \rangle \mid (e, e_i) \in \text{po} \cap \text{sqp}\}$, we pose $\pi_i = \text{extend}(\pi'', \text{CN}\langle e_i \rangle, S)$.
Otherwise, i.e. $a_i \notin \text{Put}$, we simply have $\pi_i = \pi''$

Finally, $\pi = \pi_n$ is our path for an annotated semantics reduction. We clearly have $\text{complete}(\pi)$ by definition. Our goal is then to prove that $\text{wfp}(\pi)$ 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

$\text{nodup}(\pi)$ 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 3. For all $a \in \{\text{Put}, \text{Get}, \text{nCAS}, \text{nFAA}, \text{nF}\}$ 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 $e \in \{\text{Put}, \text{Get}, \text{nCAS}, \text{nFAA}, \text{nF}\}$ 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 $\text{getCPU}(e)$ with $(e, a) \in \text{po}$. If $e \in \{\text{Put}, \text{Get}, \text{nCAS}, \text{nFAA}, \text{nF}\}$, 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 $e \in E^{\text{cpu}} = \{1R, 1W, \text{CAS}, F, P\}$, we have $\text{getl}\lambda(e, \pi_0) \prec_{\pi''} \text{Push}\langle a \rangle \prec_{\pi''} \text{getl}\lambda(b, \pi_0)$ since $(e, b) \in \text{ippo} \subseteq \text{IB}$. \square

Lemma 4. $\forall a \in \{\text{Put}, \text{Get}, \text{nCAS}, \text{nFAA}, \text{nF}\}, \forall b \in \{\text{nF}, \text{nrR}, \text{nLR}, \text{narR}, \text{1W}\},$ 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 $b \in \{\text{nF}, \text{nrR}, \text{nLR}, \text{narR}\}$, and proceed for a in po order, i.e., we can assume it holds for $e \in \{\text{Put}, \text{Get}, \text{nCAS}, \text{nFAA}, \text{nF}\}$ such that $(e, a) \in \text{po}$. By definition, $\text{NIC}\langle a \rangle$ comes from an extension $\pi'' = \text{extend}(\pi', \text{NIC}\langle a \rangle, S)$, with $S = \{\text{Push}\langle a \rangle\} \cup \{\text{getTSO}(e) \mid (e, a) \in \text{po}\}$, 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 3 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 $e \in \{\text{Put}, \text{Get}, \text{nCAS}, \text{nFAA}, \text{nF}\}$, then we have $\lambda \prec_{\pi''} \text{NIC}\langle a \rangle \prec_{\pi''} \text{getO}\lambda(b, \pi_0)$ by induction hypothesis.
- If $\lambda = \text{getTSO}\langle e \rangle = \text{B}\langle e \rangle$ for some $e \in \text{1W}$, then we have $\text{B}\langle e \rangle \prec_{\pi''} \text{NIC}\langle a \rangle \prec_{\pi''} \text{getO}\lambda(b, \pi_0)$ since $(e, b) \in \text{oppo} \subseteq \text{OB}$.

□

Lemma 5. For all w, e, p , if $\text{nrW}\langle w, e \rangle \in \pi$ and $\text{P}\langle p, e \rangle \in \pi$, then $\text{CN}\langle e \rangle \prec_\pi \text{P}\langle p, e \rangle$.

Proof. Once again, we proceed for e in po order, i.e., we can assume the result holds for $e' \in \text{nrEX}$ such that $(e', e) \in \text{po}$. $\text{CN}\langle e \rangle$ is inserted in some operation $\pi'' = \text{extend}(\pi', \text{CN}\langle e \rangle, S)$, with $S = \{\text{nrW}\langle w, e \rangle\} \cup \{\text{n1W}\langle _, e' \rangle \mid (e', e) \in \text{po} \cap \text{sqp}\} \cup \{\text{CN}\langle e' \rangle \mid (e', e) \in \text{po} \cap \text{sqp}\}$. It is then placed just after some label $\lambda \in S$.

- If $\lambda = \text{nrW}\langle w, e \rangle$, we have $\lambda \prec_{\pi''} \text{CN}\langle e \rangle \prec_{\pi''} \text{P}\langle p, e \rangle$ because $(w, p) \in \text{pf} \subseteq \text{IB}$.
- If $\lambda = \text{CN}\langle e' \rangle$ with $(e', e) \in \text{po} \cap \text{sqp}$, then there is some w' such that $(w', w) \in \text{po} \cap \text{sqp}$ and $\text{nrW}\langle w', e' \rangle \in \pi'$. From well-formedness condition number 1 (see Definition 3), there is some p' such that $(w', p') \in \text{pf}$ and $(p', p) \in \text{po}$. By induction hypothesis, we have $\text{CN}\langle e' \rangle \prec_{\pi'} \text{P}\langle p', e' \rangle$, and from $(p', p) \in \text{IB}$ we have $\text{P}\langle p', e' \rangle \prec_{\pi'} \text{P}\langle p, e \rangle$. In the end, we have the result $\text{CN}\langle e' \rangle \prec_{\pi''} \text{CN}\langle e \rangle \prec_{\pi''} \text{P}\langle p, e \rangle$.
- If $\lambda = \text{n1W}\langle w', e' \rangle$ with $(e', e) \in \text{po} \cap \text{sqp}$, then we also have $(w', w) \in \text{po} \cap \text{sqp}$, so from well-formedness condition number 1 (see Definition 3), there is some p' such that $(w', p') \in \text{pf}$ and $(p', p) \in \text{po}$. We have $\text{n1W}\langle w', e' \rangle \prec_{\pi''} \text{CN}\langle e \rangle \prec_{\pi''} \text{P}\langle p', e' \rangle \prec_{\pi''} \text{P}\langle p, e \rangle$.

□

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

- $\text{1W}\langle w \rangle \prec_\pi \text{B}\langle w \rangle$ comes from the third property when defining π_0 ; similarly for n1W , nrW and narW .
- $\text{Push}\langle a \rangle \prec_\pi \text{NIC}\langle a \rangle$ comes from the extension process.
- $\text{NIC}\langle f \rangle \prec_\pi \text{nF}\langle f \rangle$ comes from Lemma 4; similarly for $\text{NIC}\langle a \rangle \prec_\pi \text{nLR}/\text{nrR}/\text{naF}/\text{narR}\langle \dots \rangle$.
- $\text{nLR}\langle r, w, a, w' \rangle \prec_\pi \text{nrW}\langle w', e \rangle$ comes from $(r, w') \in \text{ippo} \subseteq \text{IB}$; similarly for $\text{nrR}/\text{n1W}$, $\text{naF}/\text{n1W}$, $\text{narR}/\text{n1W}$ and narR/narW .
- $\text{nrW}\langle w, e \rangle \prec_\pi \text{CN}\langle e \rangle$ comes from the extension process

- $\text{n1W}\langle w, e \rangle \prec_{\pi} \text{B}\langle w \rangle \prec_{\pi} \text{P}\langle p, e \rangle$ comes from $(w, p) \in [\text{n1W}]; \text{pf} \subseteq \text{OB}$.
- $\text{CN}\langle e \rangle \prec_{\pi} \text{P}\langle p, e \rangle$ comes from Lemma 5.

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

bufFlushOrd

- $\text{1W}\langle w_1 \rangle \prec_{\pi} \text{1W}\langle w_2 \rangle \iff \text{B}\langle w_1 \rangle \prec_{\pi} \text{B}\langle w_2 \rangle$ when $t(w_1) = t(w_2)$ comes the fact that $[\text{1W}]; \text{po}; [\text{1W}] \subseteq (\text{IB} \cup \text{OB})$, so both sides are true if and only if $(w_1, w_2) \in \text{po}$; similarly for n1W and nrW/narW on the same queue pair.
- When $t(a_1) = t(a_2)$, $\text{Push}\langle a_1 \rangle \prec_{\pi} \text{Push}\langle a_2 \rangle \iff \text{NIC}\langle a_1 \rangle \prec_{\pi} \text{NIC}\langle a_2 \rangle \iff (a_1, a_2) \in \text{po}$ from the definition of the extension process (to define π_n).
- For $a \in \{\text{Put}, \text{Get}, \text{nF}, \text{nCAS}, \text{nFAA}\}$, $w \in \text{1W}$, such that $t(a) = t(w)$:
 - If $(w, a) \in \text{po}$, then $\text{1W}\langle w \rangle \prec_{\pi} \text{Push}\langle a \rangle$ and $\text{B}\langle w \rangle \prec_{\pi} \text{NIC}\langle a \rangle$ from the definition of the extension process.
 - If $(a, w) \in \text{po}$, then $\text{Push}\langle a \rangle \prec_{\pi} \text{1W}\langle w \rangle$ and $\text{NIC}\langle a \rangle \prec_{\pi} \text{B}\langle w \rangle$ from Lemmas 3 and 4.
- When $t(w) = t(f)$, $\text{1W}\langle w \rangle \prec_{\pi} \text{F}\langle f \rangle$ implies $(w, f) \in \text{po}$ (since $[\text{F}]; \text{po}; [\text{1W}] \subseteq \text{ippo} \subseteq \text{IB}$), which implies $\text{B}\langle w \rangle \prec_{\pi} \text{F}\langle f \rangle$ (since $[\text{1W}]; \text{po}; [\text{F}] \subseteq \text{oppo} \subseteq \text{OB}$); similarly for CAS.
- If $w \in \text{n1W}$, $r \in \text{n1R}$, and $\text{sameqp}(w, r)$, then from the definition of pre-executions (see condition 6 of Definition 2), either $(w, r) \in \text{nfo}$ or $(r, w) \in \text{nfo}$. If $\text{n1W}\langle w, _ \rangle \prec_{\pi} \text{n1R}\langle r, _, _, _ \rangle$, then $(r, w) \notin \text{nfo}$ (since $\text{nfo} \subseteq \text{IB}$) and $(w, r) \in \text{nfo}$. Thus, $\text{B}\langle w \rangle \prec_{\pi} \text{n1R}\langle r, _, _, _ \rangle$ (since $\text{nfo} \subseteq \text{OB}$); similarly for $w \in \{\text{nrW}, \text{narW}\}$ and $r \in \{\text{nrR}, \text{narR}\}$.

Thus we have $\text{bufFlushOrd}(\pi)$.

pollOrder

Lemma 6. For all $e_1, e_2 \in \{\text{n1EX}, \text{nrEX}\}$, such that $\text{sameqp}(e_1, e_2)$, let $\lambda_1 \in \{\text{n1W}\langle _, e_1 \rangle, \text{CN}\langle e_1 \rangle\}$, $\lambda_2 \in \{\text{n1W}\langle _, e_2 \rangle, \text{CN}\langle e_2 \rangle\}$, then $(e_1, e_2) \in \text{po} \iff \lambda_1 \prec_{\pi} \lambda_2$.

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

- If $\lambda_1 = \text{n1W}\langle w_1, e_1 \rangle$ and $\lambda_2 = \text{n1W}\langle w_2, e_2 \rangle$, then $(e_1, e_2) \in \text{po}$ implies $(w_1, w_2) \in (\text{po} \cap \text{sqp})$, so $(w_1, w_2) \in \text{ippo} \subseteq \text{IB}$ and $\lambda_1 \prec_{\pi} \lambda_2$.
- If $\lambda_1 = \text{n1W}\langle w_1, e_1 \rangle$ and $\lambda_2 = \text{CN}\langle e_2 \rangle$, then by definition of the extension process we have $\lambda_1 \prec_{\pi} \lambda_2$.
- If $\lambda_1 = \text{CN}\langle e_1 \rangle$ and $\lambda_2 = \text{n1W}\langle w_2, e_2 \rangle$, then λ_1 is inserted in some operation $\pi'' = \text{extend}(\pi', \text{CN}\langle e_1 \rangle, S)$, with $S = \{\text{nrW}\langle _, e_1 \rangle\} \cup \{\text{n1W}\langle _, e' \rangle \mid (e', e_1) \in \text{po} \cap \text{sqp}\} \cup \{\text{CN}\langle e' \rangle \mid (e', e_1) \in \text{po} \cap \text{sqp}\}$. It is then placed just after some label $\lambda \in S$.
 - If $\lambda = \text{nrW}\langle w_1, e_1 \rangle$, we have $\lambda \prec_{\pi''} \lambda_1 \prec_{\pi''} \lambda_2$ because $(w_1, w_2) \in \text{ippo} \subseteq \text{IB}$.

- If $\lambda = \text{CN}\langle e' \rangle$ or $\lambda = \text{nIW}\langle _, e' \rangle$, with $(e', e_1) \in \text{po} \cap \text{sqp}$, then by induction hypothesis $\lambda \prec_{\pi''} \lambda_1 \prec_{\pi''} \lambda_2$.
- If $\lambda_1 = \text{CN}\langle e_1 \rangle$ and $\lambda_2 = \text{CN}\langle e_2 \rangle$, 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 $\text{sameqp}(e_1, e_2)$, $\lambda_1 \in \{\text{nIW}\langle _, e_1 \rangle, \text{CN}\langle e_1 \rangle\}$, $\lambda_2 \in \{\text{nIW}\langle _, e_2 \rangle, \text{CN}\langle e_2 \rangle\}$, $\lambda_1 \prec_{\pi} \lambda_2$, and $P\langle p_2, e_2 \rangle \in \pi$.

From the creation of the events e_1 and e_2 , there is some $w_1, w_2 \in \{\text{nIW}, \text{nrW}\}$ such that $(w_i, e_i) \in \text{po}|_{\text{imm}}$. From Lemma 6, we have $(e_1, e_2) \in \text{po}$ and thus $(w_1, w_2) \in (\text{po} \cap \text{sqp})$. By definition, we also have $(w_2, p_2) \in \text{pf}$. From well-formedness condition number 1 (see Definition 3), there is some p_1 such that $(w_1, p_1) \in \text{pf}$ and $(p_1, p_2) \in \text{po}$. Thus we have $P\langle p_1, e_1 \rangle \prec_{\pi} P\langle p_2, e_2 \rangle$ as required to prove $\text{pollOrder}(\pi)$.

nicActOrder

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

- If $a_1 \in \text{nF}$ or $a_2 \in \text{nF}$, then most of the required results hold by definition of **ippo**. The only exception is $\text{CN}\langle e \rangle \prec_{\pi} \text{nF}\langle a_2 \rangle$ which holds (by induction on e in po order) because all the dependencies of $\text{CN}\langle e \rangle$ are before $\text{nF}\langle a_2 \rangle$ by **ippo**.
- If $(a_1 \in \text{Get} \wedge a_2 \in \text{Get})$, the result holds by **ippo**.
- If $(a_1 \in \text{Get} \wedge a_2 \in \text{Put})$, the result holds by Lemma 6.
- If $(a_1 \in \text{Get} \wedge a_2 \in \text{nCAS} \cup \text{nFAA})$, the results hold by **ippo**.
- If $(a_1 \in \text{Put} \wedge a_2 \in \text{Get})$, the first result holds by **ippo**, the second by Lemma 6.
- If $(a_1 \in \text{Put} \wedge a_2 \in \text{Put})$, the first two results hold by **ippo**, the last one by Lemma 6.
- If $(a_1 \in \text{Put} \wedge a_2 \in \text{nCAS} \cup \text{nFAA})$, the results hold by **ippo**.
- If $(a_1 \in \text{nCAS} \cup \text{nFAA} \wedge a_2 \in \text{Get})$, the results hold by **ippo**.
- If $(a_1 \in \text{nCAS} \cup \text{nFAA} \wedge a_2 \in \text{Put})$, the first result holds by **ippo**, the latter by Lemma 6.
- If $(a_1, a_2 \in \text{nCAS} \cup \text{nFAA})$, the first result holds by **ippo**, the latter by Lemma 6.

Thus we have $\text{nicActOrder}(\pi)$.

nicAtomicity

For every $a_1, a_2 \in \text{nRMW}$ where $\bar{n}(a_1) = \bar{n}(a_2)$, 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_n is a total order on $\{e \in \text{narR} \mid \bar{n}(e) = n\}$. 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 $(r_2, r_1) \in \text{rao} \subseteq \text{IB}$. However, since $\text{narR}\langle r_1, \dots \rangle \prec_{\pi} \lambda_r$, we have $(r_1, r_2) \in \text{IB}$, which is a contradiction, as **IB** is irreflexive. Therefore reject our original assumption. Thus $(r_1, r_2) \in \text{rao}$, then we have $(w, r_2) \in \text{ar}; \text{rao} \subseteq \text{OB}$, so $\text{B}\langle w \rangle \prec_{\pi} \lambda_r$. Thus we have $\text{nicAtomicity}(\pi)$.

wfrd

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

- If $\lambda_r = \mathbf{1R}\langle r, w \rangle$, we need to show $\mathbf{wfrdCPU}(r, w, \pi_3)$.
 - If $w = \mathit{init}_{\mathbf{loc}(w)}$, then we need to check that $\{\mathbf{B}\langle w' \rangle, \mathbf{CAS}\langle w', _ \rangle \in \pi_3 \mid \mathbf{loc}(w') = \mathbf{loc}(r)\} = \emptyset$ and $\{\mathbf{1W}\langle w'' \rangle \in \pi_3 \mid \mathbf{loc}(w'') = \mathbf{loc}(r) \wedge t(w'') = t(r)\} = \emptyset$. For the first, such a w' would imply $(r, w') \in \mathbf{rb} \subseteq \mathbf{OB}$, which contradicts the ordering with λ_r . For the second, such an w'' would imply $(r, w'') \in \mathbf{rb}_b \subseteq \mathbf{IB}$, and $\lambda_r \prec_\pi \mathbf{1W}\langle w'' \rangle$ which similarly contradicts the ordering with λ_r .
 - If $w \in \mathbf{1W}$, $t(w) = t(r)$, and $\mathbf{B}\langle w \rangle \notin \pi_3$. From $(w, r) \in \mathbf{rf}_b \subseteq \mathbf{IB}$, we have $\lambda_w = \mathbf{1W}\langle w \rangle \prec_\pi \lambda_r$, i.e., $\pi_3 = \pi_2 \cdot \lambda_w \cdot \pi_1$. We need to show that $\{\mathbf{1W}\langle w' \rangle \in \pi_2 \mid \mathbf{loc}(w') = \mathbf{loc}(r) \wedge t(w') = t(r)\} = \emptyset$. Such a w' would imply $(w, w') \in \mathbf{po}$ (from $[\mathbf{1W}]; \mathbf{po}; [\mathbf{1W}] \subseteq \mathbf{ippo} \subseteq \mathbf{IB}$, and the execution graph forcing either $(w, w') \in \mathbf{po}$ or $(w', w) \in \mathbf{po}$), $(w, w') \in \mathbf{mo}$ (from $[\mathbf{1W}]; \mathbf{po}; [\mathbf{1W}] \subseteq \mathbf{oppo} \subseteq \mathbf{OB}$, and well-formedness conditions forcing either $(w, w') \in \mathbf{mo}$ or $(w', w) \in \mathbf{mo}$), and $(r, w') \in \mathbf{rb}_b \subseteq \mathbf{IB}$ would contradict the ordering with λ_r .
 - Else we have $\lambda_w \in \pi_3$, with $\lambda_w \in \{\mathbf{B}\langle w \rangle, \mathbf{CAS}\langle w, _ \rangle\}$. If $w \in \mathbf{1W}$ and $t(w) = t(r)$, this is the remaining subcase, else it comes from $(w, r) \in \mathbf{rf}_b \subseteq \mathbf{OB}$. Thus we have $\pi_3 = \pi_2 \cdot \lambda_w \cdot \pi_1$, and we need to check two properties. First, we check that $\{\mathbf{B}\langle w' \rangle, \mathbf{CAS}\langle w', _ \rangle \in \pi_2 \mid \mathbf{loc}(w') = \mathbf{loc}(r)\} = \emptyset$. It holds because such a w' would again imply $(r, w') \in \mathbf{rb} \subseteq \mathbf{OB}$ and contradict the ordering with λ_r . Second, we check that $\left\{ w' \mid \begin{array}{l} \mathbf{1W}\langle w' \rangle \in \pi_3 \wedge \mathbf{B}\langle w' \rangle \notin \pi_3 \wedge \\ \mathbf{loc}(w') = \mathbf{loc}(r) \wedge t(w') = t(r) \end{array} \right\} = \emptyset$. It holds because such a w' would again imply $(w, w') \in \mathbf{mo}$, $(r, w') \in \mathbf{rb}_b \subseteq \mathbf{IB}$ and contradict the ordering with λ_r .
- If $\lambda_r = \mathbf{CAS}\langle r, w \rangle$, we similarly check that $\mathbf{wfrdCPU}(r, w, \pi_3)$ holds. The difference is that cases that previously contradicted ($\mathbf{rb}_b \subseteq \mathbf{IB}$) now contradict $\mathbf{bufFlushOrd}(\pi)$ that forces the buffer of $t(r)$ to be empty when performing λ_r .
- If $\lambda_r = \mathbf{n1R}\langle r, w, _, _ \rangle$, we need to show $\mathbf{wfrdNIC}(r, w, \pi_3)$.
 - If $w = \mathit{init}_{\mathbf{loc}(w)}$, then we need to check that $\{\mathbf{B}\langle w' \rangle, \mathbf{CAS}\langle w', _ \rangle \in \pi_3 \mid \mathbf{loc}(w') = \mathbf{loc}(r)\} = \emptyset$. Such a w' would imply $(r, w') \in \mathbf{rb} \subseteq \mathbf{OB}$, which contradicts the ordering with λ_r .
 - Else we have $\lambda_w \in \pi_3$, with $\lambda_w \in \{\mathbf{B}\langle w \rangle, \mathbf{CAS}\langle w, _ \rangle\}$. This comes from $(w, r) \in \mathbf{rf}_b \subseteq \mathbf{OB}$. Thus we have $\pi_3 = \pi_2 \cdot \lambda_w \cdot \pi_1$, and we need to check that $\{\mathbf{B}\langle w' \rangle, \mathbf{CAS}\langle w', _ \rangle \in \pi_2 \mid \mathbf{loc}(w') = \mathbf{loc}(r)\} = \emptyset$. It holds because such a w' would again imply $(r, w') \in \mathbf{rb} \subseteq \mathbf{OB}$ and contradict the ordering with λ_r .
- If $\lambda_r = \mathbf{nrR}\langle r, w, _, _ \rangle$, $\mathbf{naF}\langle r, w, _, _ \rangle$ or $\mathbf{narR}\langle r, w, _, _, _ \rangle$, we similarly check that $\mathbf{wfrdNIC}(r, w, \pi_3)$ for the same reasons.

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

Theorem 6. Let G be a well-formed consistent execution graph generated from a program P . Let π be the path obtained from G by the process defined above. Then there is M' , QP' (such that for all t, \bar{n} we have $QP'(t)(\bar{n}) = \langle \varepsilon, \varepsilon, \mathbf{nEX}^* \rangle$), and an equivalent path π' (producing the same outcome as π) such that $P, M_0, B_0, A_0, QP_0, \varepsilon \Rightarrow^* (\lambda t. \mathbf{skip}), M', B_0, A_0, QP', \pi'$.

Proof. From above, we have $\mathbf{wf}(\pi)$. This shows that the program configuration can perform the events described by the annotated labels of π . The remaining part of the proof is simply to check that the command rewritings used when deriving the execution graph from P (see Fig. 5.3) can be used as \mathcal{E} transitions in the annotated semantics for P , which follows from the definitions. \square

B.5 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.

$$\begin{aligned}
& \llbracket \cdot \rrbracket_M : \text{AMem} \rightarrow \text{Mem} \\
& \llbracket M \rrbracket_M \triangleq \lambda x.v_w(M(x)) \\
& \llbracket \cdot \rrbracket_{op} : E^{\text{ext}} \rightarrow \left\{ \begin{array}{l} y^{\bar{n}} := x^n, y^{\bar{n}} := v, \text{ack}_p, x^n := y^{\bar{n}}, x^n := v, \\ x := \text{nCAS}(y^n, v, v'), x := \text{nFAA}(y^n, v), \text{cn}, \text{rfence } \bar{n} \end{array} \right\} \\
& \begin{array}{ll} \llbracket \text{lW}(x, v_w) \rrbracket_{op} \triangleq x := v_w & \llbracket F \rrbracket_{op} \text{ is undefined} \\ \llbracket \text{nrW}(\bar{y}, v_r) \rrbracket_{op} \triangleq \bar{y} := v_r & \llbracket P(\dots) \rrbracket_{op} \text{ is undefined} \\ \llbracket \text{nlW}(x, v_w, \bar{n}) \rrbracket_{op} \triangleq x := v_w & \llbracket \text{lR}(\dots) \rrbracket_{op} \text{ is undefined} \\ \llbracket \text{nF}(\bar{n}) \rrbracket_{op} \triangleq \text{rfence } \bar{n} & \llbracket \text{CAS}(\dots) \rrbracket_{op} \text{ is undefined} \\ \llbracket \text{Put}(\bar{y}, x) \rrbracket_{op} \triangleq \bar{y} := x & \llbracket \text{nlR}(\dots) \rrbracket_{op} \text{ is undefined} \\ \llbracket \text{Get}(x, \bar{y}) \rrbracket_{op} \triangleq x := \bar{y} & \llbracket \text{nrR}(\dots) \rrbracket_{op} \text{ is undefined} \\ \llbracket \text{nCAS}(z, \bar{x}, v, v') \rrbracket_{op} \triangleq z := \text{nCAS}(\bar{x}, v, v') & \llbracket \text{naF}(\dots) \rrbracket_{op} \text{ is undefined} \\ \llbracket \text{nFAA}(z, \bar{x}, v) \rrbracket_{op} \triangleq z := \text{nFAA}(\bar{x}, v) & \llbracket \text{narR}(\dots) \rrbracket_{op} \text{ is undefined} \\ \llbracket \text{nlEX}(\bar{n}) \rrbracket_{op} \triangleq \text{cn} & \llbracket \text{narW}(\dots) \rrbracket_{op} \text{ is undefined} \\ \llbracket \text{nrEX}(\bar{n}) \rrbracket_{op} \triangleq \text{ack}_p & \end{array} \\
& \llbracket \cdot \rrbracket_{opl} : E^{\text{ext}} \rightarrow \left\{ \begin{array}{l} y^{\bar{n}} := x^n, y^{\bar{n}} := v, x^n := y^{\bar{n}}, x^n := v, \\ x := \text{nCAS}(y^n, v, v'), x := \text{nFAA}(y^n, v), \text{cn}, \text{rfence } \bar{n} \end{array} \right\} \\
& \llbracket l \rrbracket_{opl} = \begin{cases} \text{cn} & \text{if } l = \text{nrEX}(\bar{n}) \\ \llbracket l \rrbracket_{op} & \text{otherwise} \end{cases}
\end{aligned}$$

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 .

$\llbracket \cdot \rrbracket_{op}$ and $\llbracket \cdot \rrbracket_{opl}$ are extended to lists in an obvious way.

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

For $\text{qp} = \langle \text{pipe}, \text{wb}_R, \text{wb}_L \rangle \in \text{AQPair}$, we define $\llbracket \text{qp} \rrbracket \triangleq \langle \llbracket \text{pipe} \rrbracket_{op}, \llbracket \text{wb}_R \rrbracket_{op}, \llbracket \text{wb}_L \rrbracket_{opl} \rangle$.

For $\text{QP} \in \text{AQPMAP}$, we define $\llbracket \text{QP} \rrbracket \triangleq \lambda t. \lambda \bar{n}. \llbracket \text{QP}(t)(\bar{n}) \rrbracket$.

For $B \in \text{ASBMap}$, we define $\llbracket B \rrbracket \triangleq \lambda t. \llbracket B(t) \rrbracket_{op}$.

Theorem 7. For all $P, P' \in \text{Prog}$, $M, M' \in \text{AMem}$, $B, B' \in \text{ASBMap}$, $A, A' \in \text{RAMap}$, $\text{QP}, \text{QP}' \in \text{AQPMAP}$, $\pi, \pi' \in \text{Path}$, if $P, M, B, A, \text{QP}, \pi \Rightarrow P', M', B', A', \text{QP}', \pi'$ and $\text{wf}(M, B, A, \text{QP}, \pi)$, then $P, \llbracket M \rrbracket_M, \llbracket B \rrbracket, A, \llbracket \text{QP} \rrbracket \Rightarrow P', \llbracket M' \rrbracket_M, \llbracket B' \rrbracket, A', \llbracket \text{QP}' \rrbracket$.

Proof. By straightforward induction on \Rightarrow . \square

Theorem 8. For all $M \in \text{AMem}$, $M'' \in \text{Mem}$, $B \in \text{ASBMap}$, $B'' \in \text{SBMap}$, $A, A' \in \text{RAMap}$, $QP \in \text{AQPMMap}$, $QP'' \in \text{QPMMap}$, and $\pi \in \text{Path}$, if $P, \llbracket M \rrbracket_M, \llbracket B \rrbracket, A, \llbracket QP \rrbracket \Rightarrow P', M'', B'', A', QP''$ and $\text{wf}(M, B, A, QP, \pi)$, then there exists $M' \in \text{AMem}$, $B' \in \text{ASBMap}$, $QP' \in \text{AQPMMap}$, and $\pi' \in \text{Path}$ such that $\llbracket M' \rrbracket_M = M''$, $\llbracket B' \rrbracket = B''$, $\llbracket QP' \rrbracket = QP''$, 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 λ and we have $\pi' = \lambda \cdot \pi$; we then need $\text{wf}(M, B, A, QP, \pi)$ to check that λ is fresh enough for π . \square

Theorem 9 (Operational and Annotated Semantics Equivalence). For all program P .

- $\llbracket M_0 \rrbracket_M$, $\llbracket B_0 \rrbracket$, A_0 , and $\llbracket QP_0 \rrbracket$ 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 $P, \llbracket M_0 \rrbracket_M, \llbracket B_0 \rrbracket, A_0, \llbracket QP_0 \rrbracket \Rightarrow^* P', M'', B'', A', QP''$ then there exists $M' \in \text{AMem}$, $B' \in \text{ASBMap}$, $QP' \in \text{AQPMMap}$, and $\pi' \in \text{Path}$ such that $\llbracket M' \rrbracket_M = M''$, $\llbracket B' \rrbracket = B''$, $\llbracket QP' \rrbracket = QP''$, 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 7 and 8. The condition $\text{wf}(M, B, A, QP, \pi)$ is obtained by applying Theorem 2 when needed. \square

C Encoding the Declarative Model in Alloy

These models are written for Alloy Analyzer 6.2.0, which is available for download at <https://alloytools.org>.

C.1 Prototype Encoding

```

module rdma_tso

fun memoryLoc[e: Event]: lone Loc {
  { l: Loc | e in MemEvent and loc[e & MemEvent] = l }
}

fun reads: set Event {
  lR + CAS + nlR + nrR + narR
}

fun writes: set Event {
  lW + CAS + nlW + nrW + narW
}

fun nicWrites: set Event {
  nlW + nrW
}

fun nicEvents: set Event {
  nlR + nrW + narR + narW + nrR + nlW + nF
}

fun nicMemEvents: set Event {
  nlR + nrW + narR + narW + nrR + nlW
}

fun instEvents: set Event {
  Event - Write
}

fun ippo: Event -> Event {
  {e1, e2: Event |
    e1 -> e2 in po and
    (
      e1 in lR + lW + CAS + F + P
      or (e1 in nlR + nF and e2 in nicEvents and e1 -> e2 in sqp)
      or (e1 in nrW + narR + narW and e2 in nicEvents - nlR
          and e1 -> e2 in sqp)
      or (e1 in nrR + nlW and e2 in nlW + nF and e1 -> e2 in sqp)
    )
  }
}

```

```

fun oppo: Event -> Event {
  {e1, e2: Event |
    e1 -> e2 in po and
    (
      e1 in lR + CAS + F + P
      or (e1 in lW and e2 in lW + CAS + F + nicEvents)
      or (e1 in nlR + nF and e2 in nicEvents and e1 -> e2 in sqp)
      or (e1 in nrW and e2 in nicEvents - (nlR + nF)
        and e1 -> e2 in sqp)
      or (e1 in narR and e2 in nicEvents - nlR and e1 -> e2 in sqp)
      or (e1 in narW and e2 in nicEvents - (nlR + nlW + nF)
        and e1 -> e2 in sqp)
      or (e1 in nrR + nlW and e2 in nlW + nF and e1 -> e2 in sqp)
    )
  }
}

fun sloc: Event -> Event {
  {e1, e2: Event | memoryLoc[e1] = memoryLoc[e2] }
}

fun sthd: Event -> Event {
  {e1, e2: Event | thread[e1] = thread[e2] }
}

fun sqp: Event -> Event {
  { e1, e2: Event | thread[e1] = thread[e2]
    and node[memoryLoc[e1]] = node[memoryLoc[e2]] }
}

pred totalOrder[s: set Event, r: Event -> Event] {
  all a: s | not a in r[a]
  all a, b, c: s |
    (a -> b in r and b -> c in r) implies a -> c in r
  all a, b: s | a != b implies (a -> b in r or b -> a in r)
}

abstract sig Execution {
  events: set Event,
  po: Event -> Event,
  rf: Event -> Event,
  mo: Event -> Event,
  pf: Event -> Event,
  nfo: Event -> Event,
  rao: Event -> Event
}

one sig E extends Execution {}

fun Events: set Event { E.events }
fun po: Event -> Event { E.po }
fun rf: Event -> Event { E.rf }

```



```

fun mo: Event -> Event { E.mo }
fun pf: Event -> Event { E.pf }
fun nfo: Event -> Event { E.nfo }
fun rao: Event -> Event { E.rao }

fact po {
  all e1: Init | all e2: Event - Init |
    e1 -> e2 in po

  all t: Thread |
    let es = { e: Event | e.thread = t } | totalOrder[es, po]
}

fact rf {
  all r: Event | some w: Event |
    r in reads implies (w in writes and w -> r in rf
                        and w.loc = r.loc
                        and valueWritten[w] = valueRead[r])
}

fact mo {
  all ex: Execution | all x: Loc {
    let wx = { w: writes | w.loc = x } | totalOrder[wx, ex.mo]
  }
}

fact pf {
  all ex: Execution | all p: P | some w: Event |
    w in nicWrites and w -> p in ex.pf
}

fact nfo {
  all ex: Execution | all e1, e2: Event |
    (e1 in nicMemEvents and e2 in nicMemEvents and e1 != e2)
    implies (e1 -> e2 in ex.nfo or e2 -> e1 in ex.nfo)
}

fact rao {
  all ex: Execution | all n: Node {
    let rn = { r: narR | r.loc.node = n } | totalOrder[rn, ex.rao]
  }
}

fun rb: Event -> Event {
  (~rf . mo) - (iden & (Event -> Event))
}

fun rfb: Event -> Event {
  (iden & (lW -> lW)) . (rf & sthd) . (iden & (lR -> lR))
}

fun rfbc: Event -> Event {
  rf - rfb
}

```

```

fun rbb: Event -> Event {
  (iden & (lR -> lR)) . (rb & sthd) . (iden & (lW -> lW))
}

fun imm(r: Event -> Event): Event -> Event {
  { a, b: Event |
    a -> b in r
    and no c: Event |
      a -> c in r and c -> b in r and a != c and b != c }
}

fun ar: Event -> Event {
  (iden & (narW -> narW)) . imm[~po]
}

fun ib: Event -> Event {
  ~(ippo + rf + pf + nfo + rbb + rao)
}

fun ob: Event -> Event {
  ~(oppo + rfb + ((iden & (nlW -> nlW)) . pf) + nfo + mo + (ar . rao))
}

pred preExecution[X:Execution] {
  Event in X.events
}

pred wellFormedExecution[X: Execution] {
  preExecution[X]
  all w1, w2: nlW + nrW + narW | all p2: P |
    ((w1 -> w2) in (po & sqp) and (w2 -> p2) in pf)
    implies (some p1: P | (w1 -> p1) in pf and (p1 -> p2) in po)
  all r: nlR | some w: nrW |
    (r -> w) in imm[po] and valueRead[r] = valueWritten[w]
  all r: nrR | some w: nlW |
    r -> w in imm[po] and valueRead[r] = valueWritten[w]
  all w: nrW | some r: nlR |
    r -> w in imm[po] and valueRead[r] = valueWritten[w]
  all w: nlW | some r: nrR + narR |
    r -> w in imm[po] and valueRead[r] = valueWritten[w]
  all r: narR | some w1: nlW |
    valueRead[r] = valueWritten[w1]
    and (r -> w1 in imm[po]
      or some wr: narW |
        r -> wr in imm[po] and wr -> w1 in imm[po])
  all w: narW | some r: narR, w2: nlW |
    r -> w in imm[po] and w -> w2 in imm[po]
    and valueWritten[w2] = valueRead[r]
}

pred consistentExecution[X: Execution] {
  wellFormedExecution[X]
}

```

```
all e: Event | e -> e not in ib
all e: Event | e -> e not in ob
all e: Event |
  e -> e not in ^((iden & (instEvents -> instEvents)) . ib . ob)
}
```

C.2 Example Litmus Test

Alloy encoding of the litmus test Fig. 3.1a.

```

open rdma_tso

one sig x extends Loc {}
one sig y extends Loc {}

one sig n1 extends Node {}
one sig n2 extends Node {}

one sig t1 extends Thread {}
one sig t2 extends Thread {}

fact { t1.node = n1 }
fact { t2.node = n2 }

pred RDMA_seq_a[X:Execution] {
  consistentExecution[X]
  some disj e1, e2, e3, e4, e5: Event {
    e1 in Init and loc[e1] = x and thread[e1] = t1
    and e2 in Init and loc[e2] = y and thread[e2] = t2
    and e3 in lW and loc[e3] = x and valueWritten[e3] = 1 and thread[e3] = t1
    and e4 in nlR and loc[e4] = x and thread[e4] = t1
    and e5 in nrW and loc[e5] = y and thread[e5] = t1
    and e3 -> e4 in X.po
    and e4 -> e5 in X.po
  }
}

run RDMA_seq_a
for exactly 3 Loc, 2 Node, 2 Thread, 5 Event

```
