

A Revised and Verified Proof of the Scalable Commutativity Rule

Lillian Tsai[†], Eddie Kohler^{*}, M. Frans Kaashoek[†], and Nikolai Zeldovich[†]
[†] MIT CSAIL ^{*} Harvard University

1 Introduction

This paper explains a flaw in the published proof of the Scalable Commutativity Rule (SCR) [1], presents a revised and formally verified proof of the SCR in the Coq proof assistant, and discusses the insights and open questions raised from our experience proving the SCR.

2 The Scalable Commutativity Rule

In order to explore the connection between commutativity and scalability in practical systems, Clements et al. [1] defined a new type of commutativity called *SIM commutativity*,¹ a property that can hold of certain interface specifications. This was used to state and prove the Scalable Commutativity Rule (SCR), which claims that every SIM-commutative interface has a conflict-free implementation—that is, on modern machines, a *scalable* implementation. Even if an interface is commutative only in a restricted context, there exists an implementation that scales in that context.

The rest of this section summarizes SIM commutativity and the precise statement of the rule, and describes the published proof of the rule.

2.1 Specifications

Specifications are represented using *actions*, where each action is either an *invocation* (representing an operation call with arguments) or a *response* (representing a return value). Each invocation is made by a specific thread, and the corresponding response is returned to the same thread. The division into invocations and responses models blocking interfaces and concurrent operations [2]. Invocations are written as “ $f(\text{args})_\alpha$ ” and responses are written as “ $\overline{\text{retval}}_\alpha$,” where overbars mark responses and Greek letters represent thread IDs.

A particular execution of a system is a *history* or *trace*, which is just a sequence of actions. For example,

$$H = [A_\alpha, B_\gamma, C_\beta, \overline{A}_\alpha, \overline{C}_\beta, \overline{B}_\gamma, D_\alpha, \overline{D}_\alpha, E_\beta, F_\gamma, G_\alpha, \overline{E}_\beta, \overline{G}_\alpha, \overline{F}_\gamma]$$

consists of seven invocations and seven corresponding responses across three different threads. In a *well-formed* history, each thread’s actions alternate invocations and responses, so each thread has at most one outstanding invocation at any point. H above is well-formed; for instance, in the thread-restricted subhistory $H|\alpha = [A_\alpha, \overline{A}_\alpha, D_\alpha, \overline{D}_\alpha, G_\alpha, \overline{G}_\alpha]$, which selects α ’s actions from H , invocations and responses alternate as required.

¹SIM stands for State-dependent, Interface-based, and Monotonic.

A *specification* models an interface’s behavior as a prefix-closed set of well-formed histories. A system execution is “correct” according to the specification if its trace is included in the specification. For instance, if \mathcal{S} corresponded to the POSIX specification, then $[\text{getpid}_\alpha, \overline{92}_\alpha] \in \mathcal{S}$ (a process may have PID 92) but $[\text{getpid}_\alpha, \overline{\text{ENOENT}}_\alpha] \notin \mathcal{S}$ (the `getpid()` system call may not return that error). A specification constrains both invocations and responses: $[\text{NtAddAtom}_\alpha]$ is not in the POSIX specification because `NtAddAtom` is not a POSIX system call.

An *implementation* is an abstract machine that takes invocations and calculates responses. The original proof of the SCR by Clements et al. [1] (also presented in Section 2.4) uses a class of machines on which conflict-freedom is defined; a good analogy is a Turing-type machine with a random-access tape, where conflict-freedom follows if the machine’s operations on behalf of different threads access disjoint portions of the tape.

An implementation M *exhibits* a history H if, when fed H ’s invocations at the appropriate times, M can produce H ’s responses (so that its external behavior equals H overall). An implementation M is *correct* for a specification \mathcal{S} if M ’s responses always obey the specification. This means that every history exhibited by M is either in \mathcal{S} or contains some invalid invocation.

2.2 Commutativity

SIM commutativity aims to capture state dependence at the interface level. State dependence means SIM commutativity must capture when operations commute in some states, even if those same operations do not commute in other states. SIM commutativity captures this contextually, without reference to any particular implementation’s state: to reason about *possible* implementations, SIM commutativity captures the scalability inherent in the interface itself. This in turn makes it possible to use the SCR early in software development, during interface design.

Commutativity states that actions may be reordered without affecting eventual results. A history H' is a *reordering* of H when $H|t = H'|t$ for every thread t . This allows actions to be reordered across threads, but not within them. For example, if $H = [A_\alpha, B_\beta, \overline{A}_\alpha, C_\alpha, \overline{B}_\beta, \overline{C}_\alpha]$, then $[B_\beta, \overline{B}_\beta, A_\alpha, \overline{A}_\alpha, C_\alpha, \overline{C}_\alpha]$ is a reordering of H , but $[B_\beta, C_\alpha, \overline{B}_\beta, \overline{C}_\alpha, A_\alpha, \overline{A}_\alpha]$ is not, since it doesn’t respect the order of actions in $H|\alpha$.

Now, consider a history $H = X ++ Y$ (where $++$ concatenates action sequences). Y *SI-commutes* in H when given any

reordering Y' of Y , and any action sequence Z ,

$$X ++ Y ++ Z \in \mathcal{S} \text{ if and only if } X ++ Y' ++ Z \in \mathcal{S}.$$

This definition captures state dependence at the interface level. The action sequence X puts the system into a specific state, without specifying a representation of that state (which would depend on an implementation). Switching regions Y and Y' requires that the exact responses in Y remain valid according to the specification even if Y is reordered. The presence of region Z in both histories requires that reorderings of actions in region Y cannot be distinguished by future operations, which is an interface-based way of saying that Y and Y' leave the system in the same state.

Unfortunately, SI commutativity is not sufficient to prove the SCR. To avoid certain degenerate cases,² the definition of commutativity must be strengthened to be *monotonic* (the M in SIM). An action sequence Y **SIM-commutes** in a history $H = X ++ Y$ when for any *prefix* P of any reordering of Y (including $P = Y$), P SI-commutes in $X ++ P$. Equivalently, Y SIM-commutes in H when, given any prefix P of any reordering of Y , any reordering P' of P , and any action sequence Z ,

$$X ++ P ++ Z \in \mathcal{S} \text{ if and only if } X ++ P' ++ Z \in \mathcal{S}.$$

Like SI commutativity, SIM commutativity captures state dependence and interface basis. Unlike SI commutativity, SIM commutativity excludes cases where the commutativity of a region changes depending on future operations. The SCR relies on SIM commutativity.

2.3 Rule

The Scalable Commutativity Rule formally states the following:

Assume an interface specification \mathcal{S} that has a correct implementation, and a history $H = X ++ Y$ exhibited by that implementation. Then whenever Y SIM-commutes in H , there exists a correct implementation of \mathcal{S} whose steps in Y are conflict-free. Since, given reasonable workload assumptions, conflict-free operations empirically scale on modern multicore hardware, this implementation is scalable in Y .

2.4 Proof

The published proof of the SCR proceeds by construction. We briefly describe how the proof proceeds (eliding certain details about interruptibility and thread switches).

²Consider an **undefinedbehavior** invocation that allows later invocations to return any value whatsoever. This resembles the specification governing C compilers. Many sequences containing **undefinedbehavior** SI-commute, even if the same sequences without **undefinedbehavior** require a non-scalable implementation; for instance, groups of sequentially-consistent **set** operations don't SI-commute, since future **gets** must return the latest value, but **sets** plus **undefinedbehavior** do SI-commute, since future **gets** are unconstrained. Most practical implementations cannot see the future, and so cannot choose unsafe scalable designs in the hope that **undefinedbehavior** will eventually occur.

```

1   $m_{\text{rule}}(s, i) \equiv$ 
2   $t \leftarrow \text{thread}(i)$ 
3  If  $\text{head}(s.h[t]) = \text{COMMUTE}$ : // enter conflict-free mode
4   $s.\text{commute}[t] \leftarrow \text{TRUE}$ ;  $s.h[t] \leftarrow \text{tail}(s.h[t])$ 
5  If  $\text{head}(s.h[t]) = i$ :
6   $s.h[t].\text{pop}()$ 
7  If  $\text{head}(s.h[t])$  is a response and  $\text{thread}(\text{head}(s.h[t])) = t$ :
8   $r \leftarrow \text{head}(s.h[t])$  // replay s.h
9  else if  $s.h[t] \neq \text{EMULATE}$ : // H complete/input diverged
10  $H' \leftarrow$  a witness consistent with  $s.h[t]$ 
11 For each invocation  $x$  in  $H'$ :
12  $\langle s.\text{refstate}, \_, \_ \rangle \leftarrow M(s.\text{refstate}, x)$ 
13  $s.h[u] \leftarrow \text{EMULATE}$  for each thread  $u$ 
14 If  $s.h[t] = \text{EMULATE}$ :
15  $\langle s.\text{refstate}, r, \_ \rangle \leftarrow M(s.\text{refstate}, i)$ 
16 else if  $s.\text{commute}[t]$ : // conflict-free mode
17  $s.h[t] \leftarrow \text{tail}(s.h[t])$ 
18 else: // replay mode
19  $s.h[u] \leftarrow \text{tail}(s.h[u])$  for each thread  $u$ 
20 Return  $\langle s, r \rangle$ 

```

Figure 1: Constructed scalable implementation m_{rule} for history H and reference implementation M .

Given a specification \mathcal{S} , an arbitrary *reference* implementation M satisfying \mathcal{S} , and a specific history $H = X ++ Y$ generated by M where Y SIM-commutes in H , the proof aims to construct an implementation m_{rule} that scales (is conflict-free) within the SIM-commutative region Y of H , but is correct (generates responses compatible with M) for any history.

m_{rule} operates in three modes: *replay*, *conflict-free*, and *emulation* modes. Its state consists of three parts:

1. $s.h[t]$, a per-thread history initialized as $X ++ [\text{COMMUTE}] ++ (Y|t)$ (where **COMMUTE** is a special marker);
2. $s.\text{commute}[t]$, a per-thread flag which is set to **TRUE** when **COMMUTE** is reached;
3. $s.\text{refstate}$, the state of reference implementation M .

Figure 1 shows pseudocode for m_{rule} 's implementation.

m_{rule} starts in replay mode. This mode replays the history H as long as each thread invokes operations in the same order as it did in H . While thread t has not yet reached the commutative region, t 's invocation matches that of $s.h[t]$ (line 5), and the next action in $s.h[t]$ is a response to t 's invocation (line 7), m_{rule} returns the matching response (line 8) and advances the stored history $s.h[u]$ of all threads (line 19).

If the next step in $s.h[t]$ is **COMMUTE**, then $s.\text{commute}[t]$ is set to **TRUE** (line 4) and m_{rule} enters conflict-free mode. In this mode, m_{rule} starts replaying steps in $Y|t$. Just like before, if the next action in $s.h[t]$ is a response to t 's invocation (line 7), m_{rule} returns the matching response (line 8). However, it advances only the stored history of t , namely $s.h[t]$ (line 17). This means that after m_{rule} enters conflict-free mode for a

thread t (i.e., $s.\text{commute}[t] = \text{TRUE}$), it accesses and modifies only those state components specific to t , and any steps in the conflict-free region Y are conflict-free.

Once $s.h[t]$ has fully replayed or if a thread t *diverges*—that is, t 's next invocation does not match the next invocation in $s.h[t]$ —then m_{rule} switches to emulate mode for all threads. In this mode, m_{rule} feeds the reference machine M invocations to determine the response to return. Before we can use the reference machine M in this way, however, the state of M must reflect the state of the execution history thus far.

We initialize M 's state by finding a *witness* of M that is consistent with the execution history (line 10). A *witness* of the execution history is a sequence of invocations that, when fed to M , generates the execution history. Once a witness is found, we know that feeding M the witness's sequence of invocations (line 12) will bring M to a valid state, where all future responses will be valid according to the spec.

Finding a witness is easy before m_{rule} reaches conflict-free mode: since m_{rule} generated the history $X ++ Y$, simply take all invocations in X (up to the current point) in order as the witness. However, if m_{rule} has entered the conflict-free mode and executed operations in the SIM-commutative region, the order in which operations were executed in this region may not equal the order in which operations were executed in Y . Here is where SIM commutativity comes in: we can reorder the operations in the commutative region of the execution history to achieve a witness. Because of SIM commutativity, **we can initialize M with a witness comprised of commutative actions in a different order than that in which they were executed, and all future responses will still be valid according to the specification.**

3 The flaw

This last, bold statement is where the proof goes wrong. To help build intuition, we now present a counterexample in which the proof-constructed machine fails: the construction cannot find a witness that will initialize M with a valid state.

3.1 Specification

Imagine a specification \mathcal{S} for opening and closing files with the following two operations:

1. int **open**(): returns a currently-unused file descriptor with value > 0 ;
2. int **close**(int fd): returns **OK** on success, **EBADFD** if fd was never opened, or **ECLOSEDFD** if fd was previously opened, but is currently closed.

To better understand \mathcal{S} , we look at a couple of examples of valid and invalid histories. Let α and β be thread IDs. The following histories would be valid:

$$\begin{aligned} H_1 &= [\text{open}_\alpha, \bar{\mathbf{1}}_\alpha, \text{open}_\beta, \bar{\mathbf{2}}_\beta, \text{close}(1)_\alpha, \overline{\mathbf{OK}}_\alpha] \\ H_2 &= [\text{open}_\alpha, \bar{\mathbf{1}}_\alpha, \text{close}(2)_\beta, \overline{\mathbf{EBADFD}}_\beta] \\ H_3 &= [\text{open}_\alpha, \bar{\mathbf{1}}_\alpha, \text{close}(1)_\alpha, \overline{\mathbf{OK}}_\alpha, \text{close}(1)_\alpha, \overline{\mathbf{ECLOSEDFD}}_\alpha] \end{aligned}$$

and the following histories would be invalid:

$$\begin{aligned} H'_1 &= [\text{open}_\alpha, \bar{\mathbf{1}}_\alpha, \text{open}_\beta, \bar{\mathbf{1}}_\beta] \\ &\quad (\text{returns used FD}) \\ H'_2 &= [\text{open}_\alpha, \bar{\mathbf{1}}_\alpha, \text{close}(2)_\alpha, \overline{\mathbf{OK}}_\alpha] \\ &\quad (\text{should return EBADFD}) \\ H'_3 &= [\text{open}_\alpha, \bar{\mathbf{1}}_\alpha, \text{close}(1)_\alpha, \overline{\mathbf{OK}}_\alpha, \text{close}(1)_\alpha, \overline{\mathbf{OK}}_\alpha] \\ &\quad (\text{should return ECLOSEDFD}) \end{aligned}$$

Note that all sequences of **open** operations are SIM-commutative regions: reordering any number of **open** operations satisfies \mathcal{S} , since the returned FDs are still unique and positive in value.

3.2 Reference implementation

We now choose a simple reference implementation M that implements \mathcal{S} . The implementation has two pieces of global state, namely a counter gfd initialized as 0 and a *closed* set initially empty. The two operations are implemented as follows:

1. int **open**: increment gfd and return the new value;
2. int **close**(int fd): if $0 < fd \leq gfd$ and $fd \notin \text{closed}$, then return **OK** and add fd to *closed*. Otherwise if $fd \in \text{closed}$, return **ECLOSEDFD**, else return **EBADFD**.

M satisfies \mathcal{S} : **open** returns only unused, positive file descriptors since the counter never (disregarding overflows) repeats values. **close**(fd) returns **OK** if the file has been opened, since all files below the current value of gfd must have been opened before, and **ECLOSEDFD** if the file has already been closed. Otherwise, fd is invalid and M returns **EBADFD**.

Although M satisfies \mathcal{S} , it is *not* scalable for SIM-commutative regions: all **open** and **close** operations access and write the shared, global counter and list.

We now have all the pieces to implement (and break) the proof construction from subsection 2.4.

3.3 Incorrect construction

We first choose a SIM-commutative region of \mathcal{S} . Since regions of **open** operations are SIM-commutative, we generate the following history using M :

$$H_{\text{commute}} = [\text{open}_\alpha, \bar{\mathbf{1}}_\alpha, \text{open}_\beta, \bar{\mathbf{2}}_\beta]$$

This history is used to set the state of m_{rule} , initializing $s.h[*]$ as

$$\begin{aligned} s.h[\alpha] &= [\text{COMMUTE}, \text{open}_\alpha, \bar{\mathbf{1}}_\alpha] \\ s.h[\beta] &= [\text{COMMUTE}, \text{open}_\beta, \bar{\mathbf{2}}_\beta] \end{aligned}$$

Now we execute m_{rule} on the following sequence of operations:

$$\text{open}_\beta, \text{close}(1)_\alpha$$

Following the proof construction, \mathbf{open}_β will first cause m_{rule} to switch to conflict-free mode for β (Figure 1, lines 3-4) because $s.h[\beta][0] = \text{COMMUTE}$. After line 4,

$$s.h[\beta] = [\mathbf{open}_\beta, \bar{\mathbf{2}}_\beta]$$

Next, since the invocation \mathbf{open} matches the first invocation by β in $s.h[\beta]$, m_{rule} will return the value $\mathbf{2}$ (lines 5-8). From line 17, the current state is now

$$\begin{aligned} s.h[\alpha] &= [\text{COMMUTE}, \mathbf{open}_\alpha, \bar{\mathbf{1}}_\alpha] \\ s.h[\beta] &= [] \end{aligned}$$

The next invocation m_{rule} receives is $\mathbf{close}(1)_\alpha$. This switches m_{rule} to conflict-free mode for α (lines 3-4), and sets the state to

$$s.h[\alpha] = [\mathbf{open}_\alpha, \bar{\mathbf{1}}_\alpha]$$

m_{rule} cannot, however, replay $s.h[\alpha]$ as it did for β . $\mathbf{close}(1)$ represents a divergence from H_{commute} , so m_{rule} enters emulate mode (lines 9-13). Our proof construction will now have to get M to a valid state consistent with the current history (line 10) so that we can feed M future invocations during emulation phase (as shown in line 15).

At this point, our recorded history is

$$H_{\text{current}} = [\mathbf{open}_\beta, \bar{\mathbf{2}}_\beta]$$

But even though this history is correct (it's in \mathcal{S}), the reference implementation *cannot exhibit this history*. In M , the first call to \mathbf{open} always returns file descriptor 1.

This causes a serious problem. The construction must initialize $s.\text{refstate}$ to some value that could have exhibited H_{current} , but no such state exists. As a result, any initialization sequence attempted by m_{rule} can cause errors on future invocations. Invoking \mathbf{open} only once sets $gfd = 1$; a following \mathbf{open} call would return 2, an invalid response (since 2 is still open). Invoking \mathbf{open} two or more times would avoid this problem, but would cause a following $\mathbf{close}(1)$ call to return the incorrect error (**ECLOSEDFD** instead of **EBADFD**). Thus, our proof construction fails both to find a witness consistent with the current execution history and to initialize M with a valid state.

3.4 When might the proof fail?

Note that M can reach a valid state if either *none* or *all* of the SIM-commutative region in question has been replayed. In other words, if m_{rule} diverges before or after the SIM-commutative region, then we can always get M to a valid state. For the former, when m_{rule} diverges before the SIM-commutative region, we can just feed M the recorded history's invocations sequentially in the order in which they occurred (that order is precisely known). For the latter, when m_{rule} diverges after the SIM-commutative region, we know there is at least one ordering of *all* operations in the SIM-commutative region that M can generate, namely H_{commute} .

Thus, feeding M all possible orderings of all the operations in the SIM-commutative region until M reaches a valid state must eventually terminate.

Furthermore, if a prefix of the operations in the SIM-commutative region *with the same order as in H_{commute}* occurs before divergence, then feeding M the operations in this order will also bring M to a valid state. This is because M generated H_{commute} by being fed operations sequentially in this order.

The potential for failure arises only when a *prefix of a reordering* of the SIM-commutative region occurs before divergence, as demonstrated in our example. In this scenario, we do not know if this reordered sequence of invocations and responses can ever be generated by M .

3.5 Why does the proof fail?

The key problem is that our reference implementation M cannot exhibit some histories required for SIM commutativity. M is correct, so every history H exhibited by M is in \mathcal{S} , but there may be some histories in \mathcal{S} that M cannot exhibit. That is, the reference implementation's *exhibited specification*—the set of histories that it can possibly exhibit—may lack some traces in \mathcal{S} . When this gap arises, and regions that are SIM-commutative in \mathcal{S} do not commute in the exhibited specification, then M may not be able to achieve a state consistent with the history our construction requires, as shown in our example.

On the other hand, the proof construction likely works whenever the given region SIM-commutes in the *exhibited* specification. Put another way, our construction likely works as long as “ M , the reference implementation, produces the same results for any reordering of the commutative region” (quoted from the published version of the proof in Clements et al. [1]). Any prefix of a reordering of the SIM-commutative region would still SIM-commute in the exhibited specification because M produces the same result for any ordering, and thus we can eliminate the failure case of our proof construction.

3.6 How to fix the proof?

We have seen that a bad reference implementation can prevent our proof technique from producing an implementation that scales within a given commutative region. We considered several fixes for this issue.

1. *Exhibited specification.* As we noted above, the proof construction will likely work if we restrict the rule to regions that are SIM-commutative in the exhibited specification, rather than those that are SIM-commutative in the defined specification.
2. *Specification oracle.* Alternately, we could remove the reference implementation from the proof entirely, and instead rely on a specification oracle that enumerates valid responses to invocations.

```

1   $m_{\text{oracle}}(s, i) \equiv$ 
2   $t \leftarrow \text{thread}(i)$ 
3   $h_{\text{copy}} \leftarrow []$ 
4   $h_{\text{perf}} \leftarrow []$ 
5  If  $s.\text{mode} \neq \text{ORACLE}$ :
6    If  $s.X_{\text{copy}} = []$ : // enter conflict-free mode
7       $s.\text{mode} \leftarrow \text{CONFLICT-FREE}$ 
8       $h_{\text{copy}} \leftarrow s.Y_{\text{copy}}[t]$ 
9       $h_{\text{perf}} \leftarrow s.Y_{\text{performed}}[t]$ 
10   else: // still in replay mode
11      $h_{\text{copy}} \leftarrow s.X_{\text{copy}}$ 
12      $h_{\text{perf}} \leftarrow s.X_{\text{performed}}$ 
13   If  $\text{head}(h_{\text{copy}}) = i$ :
14      $h_{\text{copy}}.\text{pop}()$ 
15   If  $\text{head}(h_{\text{copy}})$  is a response &  $\text{thread}(\text{head}(h_{\text{copy}})) = t$ :
16      $r \leftarrow \text{head}(h_{\text{copy}})$ 
17      $h_{\text{perf}}.\text{append}((i, r))$ 
18   else: // h.copy empty or input diverged
19      $s.\text{mode} \leftarrow \text{ORACLE}$ 
20   If  $s.\text{mode} = \text{ORACLE}$ :
21     for each possible response  $\text{resp}$  to invocation  $i$ :
22        $H' \leftarrow$  a history consistent with performed actions
23       If  $O_{\mathcal{S}}(H' ++ [(i, \text{resp})]) = \text{TRUE}$ :
24          $r \leftarrow \text{resp}$ 
25          $s.\text{oracle\_performed}.\text{append}((i, r))$ 
26       break
27   else if  $s.\text{mode} = \text{CONFLICT-FREE}$ :
28      $s.Y_{\text{copy}}[t] \leftarrow \text{tail}(s.Y_{\text{copy}}[t])$ 
29      $s.Y_{\text{performed}}[t] \leftarrow h_{\text{perf}}$ 
30   else: // replay mode
31      $s.X_{\text{copy}} \leftarrow \text{tail}(s.X_{\text{copy}})$ 
32      $s.X_{\text{performed}} \leftarrow h_{\text{perf}}$ 
33   Return  $\langle s, r \rangle$ 

```

Figure 2: Verified constructed scalable implementation m_{oracle} for history H and reference implementation M .

Exhibited specifications would preserve the somewhat “practical” feel of the flawed proof, and the useful intuition that a scalable implementation can be obtained from a non-scalable implementation by logging and reconciliation. However, this is stricter than SIM commutativity, which places requirements on the specification, not the implementation. This prevents the SCR from, for example, informing programmers about potential areas to increase the scalability of their implementations.

Specification oracles feel less realistic than reference implementations, but they have the advantage of completely avoiding the issue of whether a given specification can be implemented at all. They also fit nicely into Coq. Because of these reasons, our machine-verified proof uses specification oracles.

4 The Verified Proof

This section describes our machine-verified proof of the SCR based on a specification oracle. The pseudocode for our proof construction m_{oracle} is shown in Figure 2.

4.1 Oracle proof construction

An oracle $O_{\mathcal{S}}$ is a function from a history H to $\{\text{TRUE}, \text{FALSE}\}$ defined as

$$O_{\mathcal{S}}(H) = \begin{cases} \text{TRUE} & H \in \mathcal{S} \\ \text{FALSE} & H \notin \mathcal{S} \end{cases}$$

Given a specification \mathcal{S} , an oracle $O_{\mathcal{S}}$, and a specific history $H = X ++ Y$ where Y SIM-commutes in H , the proof constructs an implementation m_{oracle} that executes conflict-free within the SIM-commutative region Y .

m_{oracle} operates in three modes: *replay*, *conflict-free*, and *oracle* modes. Its state consists of three parts (with subparts):

1. Copies of H as histories to replay:

- $s.X_{\text{copy}}$, a global list of actions initialized as X
- $s.Y_{\text{copy}}[t]$, a per-thread list of actions initialized as $Y|t$

Note that these copies of H are equivalent to the per-thread copies $s.h[t]$ of H used in the original proof construction. The split here into X_{copy} and $Y_{\text{copy}}[t]$ allows us to cleanly separate the global state from the per-thread state in our Coq formulation. Instead of checking for a `COMMUTE` marker to indicate when the Y region has been reached, our construction will check if X_{copy} is empty.

2. Lists of performed actions:

- $s.X_{\text{performed}}$, a global list of performed actions of X initialized as $[]$
- $s.Y_{\text{performed}}[t]$, a per-thread list of performed actions of $Y|t$ initialized as $[]$
- $s.\text{oracle_performed}$, a global list of performed actions in oracle mode initialized as $[]$

3. $s.\text{mode}$, a global flag indicating the current mode of the machine

The replay and conflict-free modes act similar to the corresponding modes of m_{rule} from the prior proof. If m_{oracle} is not already in oracle mode, then m_{oracle} is in replay mode if $s.X_{\text{copy}}$ is nonempty, or in conflict-free mode if $s.X_{\text{copy}}$ is empty (line 6). If both $s.X_{\text{copy}}$ and $s.Y_{\text{copy}}[t]$ for all t are empty, m_{oracle} switches to oracle mode (line 19).

In replay mode, if the next requested invocation matches the next invocation in $s.X_{\text{copy}}$ and the next action in $s.X_{\text{copy}}$ is a response to that invocation, m_{oracle} pops the head off of $s.X_{\text{copy}}$, returns the response, and appends the

response to $s.X_performed$ (lines 13-17). Otherwise, m_{oracle} has diverged from $X ++ Y$ and switches to oracle mode.

In conflict-free mode (set up in lines 7-9), if the next requested invocation by t matches the next invocation in $s.Y_copy[t]$ and the next action in $s.Y_copy[t]$ is a response to that invocation, m_{oracle} returns the response and appends the response to $s.Y_performed[t]$ (lines 13-17). Otherwise, the execution has diverged and m_{oracle} switches to oracle mode.

In oracle mode, the next response is determined by querying the oracle function. m_{oracle} iterates through all possible responses r to the invocation i and calls $O_{\mathcal{F}}(H' ++ [r])$, where H' is a history consistent with the performed actions in the history ($s.X_performed$, $s.Y_performed[t]$ for all t , and $s.oracle_performed$). If the oracle returns *true*, then m_{oracle} stops iterating, returns r , and the chosen response is appended to $s.oracle_performed$ (lines 20-25).

More specifically, H' is constructed as

```
s.X_performed
++s.Y_performed[t0] ++ ... ++ s.Y_performed[t#threads]
++s.oracle_performed
```

Because of SIM commutativity, any ordering of operations in Y satisfies the spec. Thus, sequentially concatenating the $s.Y_performed[*]$ to construct H' generates a valid history.

Note that while in conflict-free mode, m_{oracle} executes in a scalable way: no thread accesses another's state. Thus, m_{oracle} should satisfy the SCR. In the next section, we describe how we verified this claim.

4.2 Coq formalization

Here we give an overview of how we formalized the proof construction and proved its correctness in Coq. The complete Coq source is available at https://github.com/tslilyai/coq_scr.

4.2.1 Definitions

Our Coq model includes definitions for action histories (see Section 2), the machine state and modes as described above, conflict-freedom, SIM commutativity, and machine execution. These definitions are presented in the Appendix (Figure 3). Actions are tuples of $\langle threadID, op, response \rangle$, and histories are (reversed) lists of actions. We create an enum for modes, and use a record to encode state, where the record contains either thread-specific or global histories. Per-thread state is represented as a function from tid to history.

Defining conflict-freedom requires defining conflict-free writes and reads. To define conflict-free writes, we use two constructions: **diff_histories_tid_set** takes two histories and returns the set of threads whose per-thread histories ($s.Y_performed[t]$ or $s.Y_copy[t]$) change from the first history to the second. **diff_states_tid_set** takes two states and uses **diff_histories_tid_set** to return the set of threads that have had their per-thread state changed between the two states. Note that per-thread state in our construction changes only if per-thread history changes.

With these constructions, we can define conflict-free writes (**conflict_free_writes**). All writes performed during a

step of the machine from state s_1 to s_2 on thread t are conflict-free writes if:

- the global state of s_1 and s_2 have equivalent values
- either the per-thread state of s_1 and s_2 have equivalent values, or the **diff_states_tid_set**(s_1, s_2) must contain at most the single calling thread t

These restrictions guarantee that any writes made during this step only modify t 's state.

We define conflict-free reads (**conflict_free_reads**) as follows: say the machine takes a step i on thread t from state s_1 . Then this step performs only conflict-free reads if, for any state s_2 with the same per-thread state of t and the same mode as s_1 ,

- the machine returns the same value when executing i from s_2 as it does executing i from s_1
- the state of the mode is unchanged when the machine executes i from either s_1 or s_2

Equivalent return values ensure that no read during the machine step depends on another thread's state or any global state apart from the mode; s_1 and s_2 are only restricted in their values of the per-thread state of t and the mode. Unchanging modes ensure that the only piece of global state (the mode) read by the machine in executing the step will not be modified during the step, which would cause the read to conflict with the modifying write.

With both conflict-free writes and conflict-free reads, we can prove a step is conflict-free by proving that the step's writes and reads are both conflict-free.

Note that this definition of conflict-freedom is quite different than the one presented in by Clements et al. in the original SCR paper, which reasoned about conflicts in terms of memory access sets. Instead, our Coq definition is specialized for the mechanics of our proof construction, and allows us to reason on the much higher level of our construction's abstract per-thread state (e.g., $s.Y_copy[t]$) rather than individual memory accesses.

4.2.2 Theorem statements

The final theorem and important lemmas are shown in the Appendix (Figure 4). Key helpers to prove these lemmas include determining the current mode and state of the machine, definitions for switching between modes when appropriate, and lemmas proving correctness of the machine when the machine is at each mode.

Our proof strategy was to first prove two lemmas, namely **machine_correct**, which states that the machine m_{oracle} generates only histories satisfying the spec, and **machine_conflict_free**, which states that the machine execution is conflict-free during the SIM-commutative region Y . We used these lemmas to prove our final statement of the SCR:

1. Correctness: All histories achievable by m_{oracle} satisfy the spec and m_{oracle} never returns an invalid response

- 2. Conflict-freedom: Any step m_{oracle} takes in the SIM-commutative region Y is conflict-free

4.3 Proof assumptions and evaluation

Our Coq proof makes several assumptions, encoded as parameters or as part of the definition:

- The oracle can enumerate all possible responses. For our proof, we assume something stronger, namely that the number of responses is finite; finiteness was necessary to convince Coq that the oracle function terminates.
- For every invocation that is valid to call, there exists a valid response to return.
- The oracle is correct: for all H , $O_{\mathcal{S}}(H) = \text{TRUE} \iff H \in \mathcal{S}$.

We made a number of decisions to ease the proof process. As we described earlier, we reason about conflict-freedom at a high level (without, for example, modeling memory arrays or low-level memory accesses as we did in our first attempt). This abstraction greatly simplified the proof process. We also switched to local reasoning, i.e., reasoning about steps of the machine, rather than proving facts about the machine’s entire history, and we found that non-inductive definitions (for example, for reordered histories) made the proofs easier to handle. Finally, we proved all lemmas using the reverse of histories: because histories are defined as a list, the generated inductive cases were more intuitive (new actions are added to the head, not the tail, of the list).

The entire Coq development is 2056 LOC, of which 562 lines are state definitions, lemmas, and theorems. The effort took approximately 3 person-months, including several weeks stuck on verifying the original (incorrect) proof. We believe that without attempting to verify the published proof in Coq, the flaw in the proof may have been difficult to find.

5 Discussion and Conclusion

We have presented an initial proof for the SCR, an example and brief discussion about how it is incorrect, and a new, verified proof for the SCR. However, the new proof is not wholly satisfactory. Most notably, the verified proof relies on the existence of a specification oracle, which causes it to stray even further from an imaginable construction than the original proof. We also assume enumerable responses, which may be practically true but fails to capture the intended semantics of any spec that semantically returns responses in, for example, \mathbb{R} . Furthermore, any construction that requires an oracle to iterate through an enumerable (but potentially infinite) number of responses is absurd in practice.

Problems in applying the SCR in practice are not restricted to the proof: the SCR statement itself may be unsuitable for conveying concrete information to implementers about how to design scalable systems in practice. For one, the SCR is a rule that applies only for a particular commutative region,

rather than all commutative regions of a spec. If a spec has 100 commutative regions, then we know each commutative region has a implementation scalable *for that region*; however, we do not know if all 100 implementations are different, or if one implementation exists that will scale for all (or even multiple) commutative regions.

Furthermore, as illustrated by our oracle-based proof construction, scalability may not always be optimal: an implementation may scale but not necessarily be more performant.

In its current form, the SCR serves best as a hint that certain implementations can be made more scalable and a suggestion for areas for potential implementation optimization. For example, this use of the SCR served Clements et al. well in constructing sv6 [1]. It remains an open question about whether there is a way to extend the SCR, or modify the proof construction to aid implementers in designing practical, scalable systems. Perhaps there is a way to formulate the SCR to apply not only to one particular commutative region, but rather a class of commutative regions. Or perhaps a proof construction exists for a more restricted class of commutative systems that lends itself toward efficient and practical systems.

Many other questions remain, such as whether there are specs for which all practical implementations of the spec will not commute for the SIM-commutative regions of the spec. In other words, are there specs for which the only fully scalable implementation *must* have the equivalent of an oracle? We see the space between the commutativity of implementations of a spec, and the commutativity of the spec itself as a fruitful area to explore in future work.

6 Acknowledgments

This research was supported by NSF award CNS-1302359.

References

- [1] A. T. Clements, M. F. Kaashoek, N. Zeldovich, R. T. Morris, and E. Kohler. The scalable commutativity rule: Designing scalable software for multicore processors. *ACM Transactions on Computer Systems (TOCS)*, 32(4):10, 2015.
- [2] M. P. Herlihy and J. M. Wing. Linearizability: A correctness condition for concurrent objects. *ACM Transactions on Programming Languages Systems*, 12(3):463–492, 1990.

```

Definition action : Type := tid * invocation * response.
Definition history : Type := list action.
Inductive mode : Type :=
  | ConflictFree : mode
  | Oracle : mode
  | Replay : mode.

Record state := mkState {
  X_copy : history;
  Y_copy : tid -> history;
  X_performed : history;
  Y_performed : tid -> history;
  oracle_performed : history;
  md : mode
}.

Parameter sim_commutates :
  forall hd tl tl' Z,
    reordered (hd ++ tl) Y ->
    reordered tl' tl ->
    spec (Z++tl++X) ->
    spec (Z++tl'++X).

Section Conflict.
  Definition diff_histories_tid_set {A : Type} (ts1 ts2 : tid -> A) : Ensemble tid :=
    fun tid => ts1 tid <> ts2 tid.
  Definition diff_states_tid_set (s1 s2 : state) : Ensemble tid :=
    Union tid
      (diff_histories_tid_set s1.(Y_performed) s2.(Y_performed))
      (diff_histories_tid_set s1.(Y_copy) s2.(Y_copy)).
  Definition conflict_free_writes (t :tid) (s1 s2 : state) :=
    diff_states_tid_set s1 s2 = Singleton tid t /\
    s1.(md) = s2.(md) /\
    s1.(X_copy) = s2.(X_copy) /\
    s1.(X_performed) = s2.(X_performed) /\
    s1.(oracle_performed) = s2.(oracle_performed).
  Definition conflict_free_reads t i s :=
    forall (s1 s2 s1' s2' : state) (a1 a2 : action),
      s1.(Y_copy) t = s.(Y_copy) t ->
      s2.(Y_copy) t = s.(Y_copy) t ->
      s1.(Y_performed) t = s.(Y_performed) t ->
      s2.(Y_performed) t = s.(Y_performed) t ->
      s1.(md) = s.(md) ->
      s2.(md) = s.(md) ->
      machine_act s1 t i = (s1', a1) ->
      machine_act s2 t i = (s2', a2) ->
      a1 = a2 /\ s1'.(md) = s.(md) /\ s2'.(md) = s.(md).
End Conflict.

Definition machine_act (s : state) (t : tid) (i : invocation) : (state * action) :=
  let mode := next_mode s t i in
  match mode with
  | Oracle => get_oracle_response (state_with_md s Oracle) t i
  | ConflictFree => get_conflictfree_response (state_with_md s ConflictFree) t i
  | Replay => match rev (s.(X_copy)) with
    | [hd] => get_replay_response (state_with_md s ConflictFree) t i
    | _ => get_replay_response (state_with_md s Replay) t i
  end
end.

```

Figure 3: Definitions for proving the SCR

```

Lemma machine_correct :
  forall s h,
    generated s h ->
      spec h.

Lemma machine_conflict_free :
  forall s s' h t i r,
    generated s (h ++ X) ->
      spec ((t,i,NoResp) :: h ++ X) ->
        (exists h', reordered (h' ++ (t,i,r) :: h) Y) ->
          machine_act s t i = (s', (t,i,r)) ->
            conflict_free_step t s s'.

Theorem scalable_commutativity_rule :
  (* All achievable histories satisfy the spec and *)
  (* the machine never returns an invalid response *)
  (forall s h t i r,
    current_state_history s h ->
      spec h /\
        (List.In (t,i,r) h -> exists rtyp, r = Resp rtyp))
  (* If the machine's next step is in the middle of a (reordering of) *)
  (* a SIM-commutative region Y, then the machine's execution of the *)
  (* step is conflict free. *)
  /\ (forall s s' h t i r,
    current_state_history s (h ++ X) ->
      spec ((t,i,NoResp) :: h ++ X) ->
        (exists h', reordered (h' ++ (t,i,r) :: h) Y) ->
          machine_act s t i = (s', (t,i,r)) ->
            conflict_free_writes t s s'
            /\ conflict_free_reads t i s).

```

Figure 4: Theorems proven about the SCR in Coq