# Sequential Reasoning for Optimizing Compilers Under Weak Memory Concurrency

Minki Cho, Sung Hwan Lee, Dongjae Lee, Chung Kil Hur, Ori Lahav

Presented by
Akshay Gopalakrishnan

April 10, 2024

## Introduction

- Performing/designing thread-local optimizations in a concurrent context would require understanding the underlying memory consistency model.
- Often, these models are complex, as they are weak consistency based.
- Thus, it becomes difficult to design optimizations which are safe.
- Most optimizations that are performed on concurrent programs mainly involve reordering or eliminating non-atomics.
- Is there perhaps, then a simpler semantics that an optimization designer can rely on?

## Main Idea

- This paper proposes a Sequential Model (SEQ) to design compiler optimizations.
- The optimizations addressed are solely regarding non-atomics, or reordering non-atomics across atomic parts of code.
- SEQ is shown to preserve contextual refinement for non-atomic optimizations in the Promising Semantics (PS).
- Contextual refinement in PS makes sure that optimization developers can solely on SEQ to design optimizations.

# Optimizations on Non-atomics

- Reordering non-atomics (RR/WR/WW/RW).
- Reordering non-atomics outside loops (Loop Invariant Code motion).
- Elimination of non-atomics (RR/WR/WW/RW).
- Redundant non-atomic load introduction/elimination.
- Reordering atomics with non-atomics (RR/WR/WW/RW).
- Elimination of non-atomics across atomics.

2024-04-10

Optimizations on Non-atomics

- Reordering non-atomics (RR/WR/WW/RW).
- Reordering non-atomics outside loops (Loop Invariant Code motion).
- Elimination of non-atomics (RR/WR/WW/RW).
- Redundant non-atomic load introduction/elimination.
- Reordering atomics with non-atomics (RR/WR/WW/RW).
- Elimination of non-atomics across atomics.

The non-adjacent optimization involving elimination of non-atomics across atomics does not seem to be required in my opinion. This is because the proof for that can be justified by decomposing the optimizations into the above elementary adjacent components.

To add, Write introduction is not allowed by the model. This is a good design choice to preserve security guarantees, however it still limits from introducing redundant writes within dead-code blocks, which may be used by optimizing compilers to easily perform partial redundancy elimination.

- Atomic events of Type RLX/REL/ACQ.
- Non atomic events of type NA.

Each state $S = \langle \alpha, P, F, M \rangle$ keeps track of

- $\alpha$ - State transition.
- $P$ - Permission set, has set of non-atomic locations safe to access.
- $F$ - Written set, has set of non-atomic locations written to.
- $M$ - Shared memory set, consisting both atomic and non-atomic locations.

Final state is either *return v* or $\perp$ (which includes undefined behavior).

SEQ Machine: Memory Events and State

- Atomic events of Type RLX/REL/ACQ.
- Non atomic events of type NA.

Each state $S = (\alpha, P, F, M)$ keeps track of

- $\alpha$ - State transition.
- $P$ - Permission set, has set of non-atomic locations safe to access.
- $F$ - Written set, has set of non-atomic locations written to.
- $M$ - Shared memory set, consisting both atomic and non-atomic locations.

Final state is either $return\ v$ or $\perp$ (which includes undefined behavior).

It is unclear whether $P$ at the beginning of the execution includes all possible non-atomic locations (is it $\top$) or it is empty (is it $\perp$). My understanding is that it should be $\perp$, and that acquire semantics would start including in the set locations we are safe to read/write to, whereas release semantics would do the reverse. The transitions seem to make sense, the non-deterministic change (inclusion in acq and exclusion in rel) is meant to reflect any random context (other threads) that would be running concurrently.

**Figure 1.** Transitions of SEQ

A Behavior is defined as $\langle tr, r \rangle$ where

- $tr$ is the sequence of labels denoting an in which a trace of events were executed.
- $r$ is either
  - $trm(v, F, M)$ denoting return value $v$, updated writes $F$ and memory state $M$.
  - $\perp$ denoting erroneous termination, which includes Undefined behavior.

1. Transition labels:

$$\frac{}{e \sqsubseteq e}$$

$$\frac{v_{\text{tgt}} \sqsubseteq v_{\text{src}}}{\text{W}^{\text{rlx}}(x, v_{\text{tgt}}) \sqsubseteq \text{W}^{\text{rlx}}(x, v_{\text{src}})}$$

$$\frac{F_{\text{tgt}} \subseteq F_{\text{src}}}{\text{R}^{\text{acq}}(x, v, P, P', F_{\text{tgt}}, V) \sqsubseteq \text{R}^{\text{acq}}(x, v, P, P', F_{\text{src}}, V)}$$

$$\frac{v_{\text{tgt}} \sqsubseteq v_{\text{src}} \quad F_{\text{tgt}} \subseteq F_{\text{src}} \quad V_{\text{tgt}} \sqsubseteq V_{\text{src}}}{\text{W}^{\text{rel}}(x, v_{\text{tgt}}, P, P', F_{\text{tgt}}, V_{\text{tgt}}) \sqsubseteq \text{W}^{\text{rel}}(x, v_{\text{src}}, P, P', F_{\text{src}}, V_{\text{src}})}$$

2. Traces: $\quad e_{\text{tgt}}^1 \cdot \ldots \cdot e_{\text{tgt}}^n \sqsubseteq e_{\text{src}}^1 \cdot \ldots \cdot e_{\text{src}}^n \Leftrightarrow \forall k. \; e_{\text{tgt}}^k \sqsubseteq e_{\text{src}}^k$

3. Behaviors:

$$\frac{tr_{\text{tgt}} \sqsubseteq tr_{\text{src}} \quad v_{\text{tgt}} \sqsubseteq v_{\text{src}} \quad F_{\text{tgt}} \subseteq F_{\text{src}} \quad M_{\text{tgt}} \sqsubseteq M_{\text{src}}}{\langle tr_{\text{tgt}}, \text{trm}(v_{\text{tgt}}, F_{\text{tgt}}, M_{\text{tgt}}) \rangle \sqsubseteq \langle tr_{\text{src}}, \text{trm}(v_{\text{src}}, F_{\text{src}}, M_{\text{src}}) \rangle}$$

$$\frac{tr_{\text{tgt}} \sqsubseteq tr_{\text{src}} \quad F_{\text{tgt}} \subseteq F_{\text{src}}}{\langle tr_{\text{tgt}}, \text{prt}(F_{\text{tgt}}) \rangle \sqsubseteq \langle tr_{\text{src}}, \text{prt}(F_{\text{src}}) \rangle} \qquad \frac{tr_{\text{tgt}} \sqsubseteq tr_{\text{src}}}{\langle tr_{\text{tgt}} \cdot tr, r \rangle \sqsubseteq \langle tr_{\text{src}}, \bot \rangle}$$

Sequential Reasoning for Optimizing Compilers Under Weak Memory Concurrency

└─Refinement of Transition Steps



What baffles me is the last transition, which says the normal termination is a refinement of erroenous termination. I guess this means it is okay to make something that invokes undefined behavior to one that does not. In the context of data races however, this would clearly be incorrect, since removing races, even if desirable is not the original behavior of the program. The compiler is allowed to do anything yes, but I guess from a debugging standpoint this would not be ideal. Perhaps mail one of the authors about this.

- Reordering of non-atomics across atomics in certain cases are unsound.
- Most of these cases reduce to *Roach Motel Reordering*.
- Reason is the final set of labels are different, keeping the termination state as $\perp$.
- Example: $a = y^{rlx}; b = x^{na} \mapsto b = x^{na}; a = y^{rlx}$.

## Fixes: High Level

- Relax the definition of $\sqsubseteq$ on trace labels.
- Handle rlx-na to na-rlx: Assume no fixed read value returned for an atomic read.
- Handle acq-na to na-acq: Restrict *tr* of source such that the sequence of labels contain no acquire reads that lead to UB.
- Handle rel-na to na-rel: Commitment set of non-atomic writes to be fulfilled before an acquire read.

Fixes: High Level

- Relax the definition of ⊑ on trace labels.
- Handle rlx-na to na-rlx: Assume no fixed read value returned for an atomic read.
- Handle acq-na to na-acq: Restrict $tr$ of source such that the sequence of labels contain no acquire reads that lead to UB.
- Handle rel-na to na-rel: Commitment set of non-atomic writes to be fulfilled before an acquire read.

The solution to address the disallowing of non-atomic and atomic reorder-
ing is certainly confusing. While some of it is intuitive, it does not give the
confidence that it actually would hold in all cases. The crux is to relax the
refinement definition to incorporate more reorderings, noting that some
information of trace labels are irrelevant.

- The first idea is to note that if you anyways reach UB irrespective of
  reordering, then you can make do by not ensuring all trace labels are
  the same.

- To ensure we do not mess this up, we need to ensure that we do not
  introduce UB by reordering across acquire reads (as hb relation
  would make sure no data race in some sense). This translates to
  ensuring the trace labels have no trailing acquire reads.

- Next, to ensure we do not introduce UB again by optimizing, we
  should not assume the atomic reads have a fixed read value,
  assuming no fixed value of shared memory.

The main claim for SEQ was that a compiler developer can solely rely on sequential semantics (SEQ in this case) to design sound optimizations (that which satisfy contextual refinement).

- Adequacy of SEQ is shown by claiming any optimization (refinement) sound in SEQ is sound in Promising Semantics.
- Adequacy is contingent on program being deterministic, that concrete execution path always leads to the same final outcome.

## Assumptions/Limitations

- Do not tackle optimizations only involving atomics, claiming this is rarely done by compilers.
- LICM tackled only involves removing read-events outside the loop, what about writes? Usually writes are not involved or considered as a Loop invariant.
- SEQ thus, cannot be used for atomic optimizations, and thus atomic based optimizations would require understanding PS at its full.
- SEQ relies on the practical implementation and usage of Promising Semantics as a programming language model.
- Paper claims it addresses optimizations without catch-fire semantics, but still uses a notion of UB from LLVM semantics to prove soundness of optimizations.
- Paper's claims on optimizations are additive, with minor warnings on allowing false optimizations. However, it is not yet clear if it allows only what is desired.

- Paper seemed to falsify the rhetoric. The intention being to allow non-atomic based optimizations in the PS.
- This was done by designing a simpler sequential model SEQ, taking advantage of UB and being conservative while proving contextual refinement.
- SEQ was then showed COMPLETE w.r.t. PS for the intended optimizations.
- Promising Semantics was extended by adding non-atomic events (did not discuss here in this presentation).