# Relaxed Semantics of Concurrent Programs

G. Boudol, G. Petri and B. Serpette

EXPRESS-SOS 2012

# The Problem

- shared memory concurrency – aka multithreading = the concurrent programming model in mainstream programming languages (e.g. C/C++, Java), close to hardware multiprocessor architectures

- in shared memory multiprocessor (multicore) architectures, memory accesses may be reordered (also with optimizing compilers)

➡ there are behaviors of multithreaded programs that are not explained by the standard interleaving semantics

Formalization of such relaxed semantics?

# An Example

Initially $S(p) = \mathit{ff} = S(q)$

$$
\begin{array}{l|l}
p := \mathit{tt}; & q := \mathit{tt}; \\
r_0 := \,!q & r_1 := \,!p
\end{array}
$$

Possible outcomes with the interleaving semantics:

$$
\begin{aligned}
S(r_0) &= \mathit{ff} \quad \& \quad S(r_1) = \mathit{tt} \\
S(r_0) &= \mathit{tt} \quad \& \quad S(r_1) = \mathit{ff} \\
S(r_0) &= \mathit{tt} \quad \& \quad S(r_1) = \mathit{tt}
\end{aligned}
$$

On most multiprocessor machines, one may also observe

$$
S(r_0) = \mathit{ff} \quad \& \quad S(r_1) = \mathit{ff}
$$

# EXPLANATION

The two threads are executed on different processors

1. the write $p := tt$ is issued and put in a write buffer, but not yet performed on the shared memory
2. the read $!q$ is performed, returning the value $ff$ from the shared memory
3. the write $q := tt$ is performed, then
4. the read $!p$ is performed and returns $ff$ from the shared memory, since the second processor does not see the write buffer of the first

➥   the write $p := tt$ and the read $!q$ appear to be reordered

# ANOTHER (SIMILAR) EXAMPLE

Initially $S(p) = \mathit{ff} = S(q)$

$$
\begin{array}{c|c}
p := tt; & r_0 := !q; \\
q := tt & r_1 := !p
\end{array}
$$

On some multiprocessor machines, one may observe the outcome

$$
S(r_0) = \textcolor{red}{tt} \quad \& \quad S(r_1) = \textcolor{red}{\mathit{ff}}
$$

Reordering of the writes $p := tt$ and $q := tt$ (conceptually: one write buffer per memory location), or of the reads $!q$ and $!p$

# MEMORY BARRIERS

Low level instructions to prevent reorderings: full **fence**, or weaker **wr**, **ww**, **rr**, **rw**, for instance

$$
\begin{array}{c|c}
p := tt; & q := tt; \\
\text{wr}; & \text{wr}; \\
r_0 := !q \; (\mathit{ff}) & r_1 := !p \; (\mathit{ff})
\end{array}
$$

is forbidden. However

$$
\begin{array}{c|c}
p := tt; & r_0 := !q \; ; \; (\mathit{tt}) \\
\text{ww}; & \text{rr}; \\
q := tt & r_1 := !p \; (\mathit{ff})
\end{array}
$$

is still possible (on some machines)

# NON ATOMIC WRITES

The IRIW (Independent Reads of Independent Writes) example:

$$p := tt \quad \| \quad q := tt \quad \Big\| \quad \begin{array}{l} r_0 := !p \; ; \; (tt) \\ \mathsf{rr}; \\ r_1 := !q \; (\mathit{ff}) \end{array} \quad \Big\| \quad \begin{array}{l} r_2 := !q \; ; \; (tt) \\ \mathsf{rr}; \\ r_3 := !p \; (\mathit{ff}) \end{array}$$

A possible explanation: the write $p := tt$ is issued and made visible to the third thread, but not to the fourth (the previous example is similar). Notice: no reordering

# Relaxed Memory Models

According to Adve & Gharachorloo (Tutorial, 1996):

memory model = reordering, i.e. relaxation of program order, for
        memory accesses (write/read)
        + write visibility
        + means to maintain program order in some cases
        (memory barriers, acquire/release...)

Some examples:

▶ TSO: relaxing write/read + read-own-write-early

▶ PSO: TSO + write/write relaxation

▶ RMO: PSO + read/write + read/read relaxations

▶ PowerPC: RMO + read-others'-write-early

# SEMANTICS: APPROACHES

- axiomatic: an execution is a set of memory events, connected by various relations (po, rf, co... + various dependencies: data, addr, ctrl...), subject to a number of axioms. Given a set of events, one tentatively sets relations between them and checks if the axioms are satisfied

- operational: an abstract machine with transitions from configurations (shared memory, threads, interconnection structure – write buffers, caches...) to configurations. Rules to non-deterministically choose the next step

# A Language

Syntax – a call by value, imperative $\lambda$-calculus:

$$
\begin{array}{rcll}
p,\ q,\ r\ldots & \in & \mathcal{P} & pointers \\
b & \in & \mathcal{B} & barriers \\
v,\ w\ldots & ::= & p \mid x \mid \lambda x e \mid () \mid \cdots & values \\
e & ::= & v \mid (v e) & expressions \\
& \mid & (!v) \mid (v := w) \mid b \mid \cdots &
\end{array}
$$

Evaluation contexts

$$
\mathbf{E} \ ::= \ [] \ \mid \ (v\mathbf{E})
$$

Recall:

$$
e_0\,;\,e_1 \ =_{\mathrm{def}} \ (\lambda x e_1\, e_0) \ = \ (\lambda x e_1[e_0])
$$

(where $x$ is not free in $e_1$)

# CONFIGURATIONS

$(S, \sigma, T)$ where

- $S$, the shared memory, is a mapping from a finite set of pointers to closed values

- $T$, the thread system, is a mapping from a finite set of thread identifiers $\{t_1, \ldots, t_n\}$ to expressions, witten

$$(t_1, e_1) \parallel \cdots \parallel (t_n, e_n)$$

- $\sigma$ is the temporary store, a sequence $(t_{i_1}, \mu_1) \cdots (t_{i_k}, \mu_k)$ of pending memory operations $\mu_1, \ldots, \mu_k$ issued by the threads $t_{i_1}, \ldots, t_{i_k}$, but not yet "globally performed" (i.e. not yet touching the memory). An abtraction of the interconnection structure between processors and the shared memory

# MEMORY OPERATIONS

$$\mu ::= (v := w)^W \mid !^x v \mid b$$

- write operations $(v := w)^W$ where $v$ is the location to update (either a pointer or a variable, if not yet determined), $w$ the new value, $W$ is the visibility of the write, a set of thread identifiers

- read operations $!^x v$ where $v$ is the location to read and $x$ the place-holder for the value in the thread that reads (and in subsequent memory operations)

# TRANSITIONS

From the threads:

$$\frac{(\sigma, (t, e)) \to (\sigma', (t', e'))}{(S, \sigma, (t, e) \parallel T) \xrightarrow[\mathcal{M}]{} (S, \sigma', (t', e') \parallel T)}$$

where

$$
\begin{aligned}
(\sigma, (t, \mathbf{E}[(\lambda x e\, v)])) &\to (\sigma, (t, \mathbf{E}[\{x \mapsto v\} e])) \\
(\sigma, (t, \mathbf{E}[(!v)])) &\to (\sigma \cdot (t, !^x v), (t, \mathbf{E}[x])) \qquad x \text{ fresh} \\
(\sigma, (t, \mathbf{E}[(v := w)])) &\to (\sigma \cdot (t, (v := w)^\emptyset), (t, \mathbf{E}[()])) \\
(\sigma, (t, \mathbf{E}[b)])) &\to (\sigma \cdot (t, b), (t, \mathbf{E}[()]))
\end{aligned}
$$

# MEMORY MODEL

$\mathcal{M} = (\curvearrowleft, \mathcal{W})$ where

- $\curvearrowleft$, the commutability predicate, relates temporary stores $\sigma$ with pending memory operations $(t, \mu)$. If $\sigma \curvearrowleft (t, \mu)$ then $\mu$ may be immediately performed, overtaking the operations in $\sigma$ − allowing reorderings/relaxations of the program order

- $\mathcal{W}$, the write grain, is a set of sets of thread identifiers − the allowed visibilities

subject to some requirements, e.g. $\varepsilon \curvearrowleft (t, \mu)$ i.e. the empty temporary store allows any memory operation to be performed, $\emptyset \in \mathcal{W}, \dots$

# TRANSITIONS

From the temporary store:

$$(S, \sigma) \hookrightarrow (S', \sigma', \mathsf{Sub}) \;\Rightarrow\; (S, \sigma, T) \xrightarrow[\mathcal{M}]{} (S', \mathsf{Sub}(\sigma', T))$$

where

$$(S, \sigma) \;\hookrightarrow\; (S, \sigma_0 \cdot \sigma_1, \{x \mapsto v\}) \qquad\qquad read$$
$$if \;\; \sigma = \sigma_0 \cdot (t, !^x p) \cdot \sigma_1 \;\&$$
$$\sigma_0 \uparrow (t, !^x p) \;\&\; S(p) = v$$

$$(S, \sigma) \;\hookrightarrow\; (S[p := v], \sigma_0 \cdot \sigma_1, \emptyset) \qquad\qquad write$$
$$if \;\; \sigma = \sigma_0 \cdot (t, (p := v)^W) \cdot \sigma_1 \;\&$$
$$\sigma_0 \uparrow (t, (p := v)^W) \;\&\; v \text{ closed}$$

# TRANSITIONS

$$(S, \sigma) \hookrightarrow (S, \sigma_0 \cdot \sigma_1, \emptyset) \qquad \textit{barrier}$$
$$\textit{if} \quad \sigma = \sigma_0 \cdot (t, b) \cdot \sigma_1 \;\&\; \sigma_0 \curvearrowleft (t, b)$$

$$(S, \sigma) \hookrightarrow (S, \sigma_0 \cdot (t, (v := w)^{W'}) \cdot \sigma_1, \emptyset) \qquad \textit{write early}$$
$$\textit{if} \quad \sigma = \sigma_0 \cdot (t, (v := w)^{W}) \cdot \sigma_1 \;\&\;$$
$$t \in W' \;\&\; W \subset W' \in \mathcal{W}$$

$$(S, \sigma) \hookrightarrow (S, \sigma_0 \cdot \sigma_1, \{x \mapsto w\}) \qquad \textit{read early}$$
$$\textit{if} \quad \sigma = \sigma_0 \cdot (t, !^x v) \cdot \sigma_1 \;\&\;$$
$$\sigma_0 = \delta_0 \cdot (t', (v := w)^{W}) \cdot \delta_1 \;\&\;$$
$$t \in W \;\&\; \delta_1 \curvearrowleft (t, !^x v) \;\&\; \delta_0 \curvearrowleft^{\mathcal{B}} (t, !^x v)$$

# MEMORY MODELS: REQUIREMENTS

The commutability predicate should not be so relaxed that the semantics of sequential programs could be broken. Then $\curvearrowleft$ must satisfy

$$(t, \mu) \blacktriangleleft (t', \mu) \implies \forall \sigma, \sigma'. \ \neg\big(\sigma \cdot (t, \mu) \cdot \sigma' \curvearrowleft (t', \mu')\big)$$

where the precedence relation $\blacktriangleleft$ is inductively given by

$$\left. \begin{array}{c} v \approx v' \ \& \\ t' \in \{t\} \cup W \end{array} \right\} \implies \left\{ \begin{array}{l} (t, (v := w)^W) \blacktriangleleft (t', !^x v') \ \& \\ (t, (v := w)^W) \blacktriangleleft (t', (v' := w')^{W'}) \end{array} \right.$$

$$v \approx v' \implies (t, !^x v) \blacktriangleleft (t, (v' := w)^W)$$

where

$$v \approx v' \iff_{\mathrm{def}} v = v' \text{ or } v \in \mathcal{V}ar \text{ or } v' \in \mathcal{V}ar$$

# EXAMPLE 1

$$p := tt; \qquad \Big\| \quad q := tt;$$
$$r_0 := !q \ (\textit{ff}) \quad \Big\| \quad r_1 := !p \ (\textit{ff})$$

The initial configuration may evolve into

$$(S, \quad (t_0, (p := tt)^\emptyset) \cdot (t_0, !^x q), \quad (t_0, (r_0 := x)) \, \| (t_1, e_1))$$

With the relaxation of the write/read order:

$$(t_0, (p := tt)^\emptyset) \ \curvearrowleft \ (t_0, !^x q)$$

thus $!^x q$ can be performed, returning $\{x \mapsto \textit{ff}\}$. Then $e_1$ is executed: the write $(q := tt)^\emptyset$ commutes with $(p := tt)^\emptyset$, as well as the read $!^y p$ (which does not see this write), and finally $(p := tt)^\emptyset$ is performed

# MEMORY BARRIERS: SEMANTICS

By means of the commutability predicate:

$$(t, (v := w)^W) \quad \blacktriangleleft \quad (t, \mathsf{ww}) \quad \blacktriangleleft \quad (t, (v' := w')^{W'})$$
$$(t, (v := w)^W) \quad \blacktriangleleft \quad (t, \mathsf{wr}) \quad \blacktriangleleft \quad (t, (!^x v'))$$
$$\vdots$$

Notice: same thread

Global barrier: **sync** is a **ww**, **wr**, **rw** and **rr** (local) barrier, plus sees the writes from other threads:

$$t' \in W \quad \Rightarrow \quad (t, (v := w)^W) \quad \blacktriangleleft \quad (t', \mathsf{sync})$$

# EXAMPLE 2

$$p := tt; \qquad r_0 := !q \ ; \ (tt)$$
$$\text{ww}; \qquad \Big\| \qquad \text{rr};$$
$$q := tt \qquad r_1 := !p \ (f\!f)$$

From the temporary store

$$(t_0, (p := tt)^{\emptyset}) \cdot (t_0, \text{ww}) \cdot (t_0, (q := tt)^{\emptyset})$$

the visibility of the last write is extended to the second thread:

$$(t_0, (p := tt)^{\emptyset}) \cdot (t_0, \text{ww}) \cdot (t_0, (q := tt)^{\{t_0, t_1\}})$$

Then $t_1$ proceeds: $!^x q$ returns $\{x \mapsto tt\}$ from the temporary store, the barrier **rr** vanishes, and $!^y p$ returns $\{y \mapsto f\!f\}$ from the shared memory. Replacing **rr** with **sync** prevents this behavior

# In the PAPER

▶ a more refined abstract machine to deal with the (subtle) lwsync barrier (PowerPC)

$$\sigma \upharpoonleft (t, \mu) \;\neq\; \forall (t', \mu') \text{ in } \sigma. \; \neg\big((t', \mu') \blacktriangleleft (t, \mu)\big)$$

▶ extension with speculation – branch prediction: from a conditional branching (if $v$ then $e_0$ else $e_1$) one issues a prediction $[v = tt]$ or $[v = ff]$ in the temporary store. A correct prediction $[v = v]$ vanishes, while a prediction blocks memory updates:

$$(t, [v = w]) \;\blacktriangleleft\; (t, (v' := w')^W)$$

▶ a software simulator that allows us to run a (large) number of litmus tests, despite state explosion – trying to explore these without the simulator is highly error prone!