

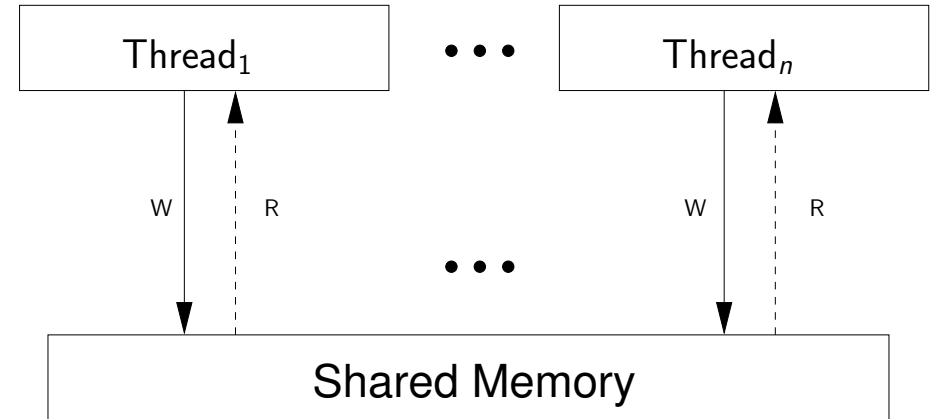
# Testing (and defining) Weak Memory Models

Luc Maranget

Inria Paris-Rocquencourt

1/74

## A simple shared memory computer



Threads execute programs as usual: instructions are executed completely and atomically (memory stores in particular).

2/74

## The Sequentially Consistent Model (SC)

Definition by L. Lamport:

“... the result of any execution is the same as if the operations of all the processors were executed in some sequential order, and the operations of each individual processor appear in this sequence in the order specified by its program”.

(One may add “stores take effect immediately” .)

**Interleaving semantics:** This is “interleaving semantics” as “some sequential order” results from interleaving “the order specified by the program of all individual processors” .

A first, one expect shared memory multiprocessors to behave that way, which of course they don't.

3/74

## Axiomatic or *post-mortem* semantics — Events

The effects of “operations executed by the processors” are represented by *events*. We define *memory events*  $(a):\mathbf{d}[\ell]v$ :

- ▶ Unique label typically  $(a)$ ,  $(b)$ , etc.
- ▶ Direction  $\mathbf{d}$ , that is read (R) or write (W)
- ▶ Memory location  $\ell$ , typically  $x$ ,  $y$ , etc.
- ▶ Value  $v$ , typically 0, 1 etc.
- ▶ Originating thread:  $T_0$ ,  $T_1$  (often omitted)

The program order  $\xrightarrow{po}$  (“order specified by program”) is a linear order amongst the events originating from the same thread.

Relation  $\xrightarrow{po}$  represents the sequential execution of events by one thread that follows the *uniprocessor model*: the usual processor execution model, where instructions are executed by following the order given in program.

4/74

## Example of program order

Despite its name, program order is a dynamic notion.

```
/* x,t and y are (shared) memory locations, t = { 2, 3, } */
int r1,r2=0 ; // non-shared locations (e.g. registers)
x = 1 ;
for (int k = 0 ; k < 2 ; k++) { r1 = t[k] ; r2 += r1 ; }
y = r2 ;
```

Events and program order :

(a):W[x]1  $\xrightarrow{po}$  (b):R[t + 0]2  $\xrightarrow{po}$  (c):R[t + 4]3  $\xrightarrow{po}$  (d):W[y]5

**Notice:** program order (and events) may depend on the values of the reads, *i.e.* on values written by other threads (**if**...).

In simple examples, program order is given by program text.

5/74

## Relating writes and reads

We define a first “communication relation” between events with the same location.

### Definition (Read-from $\xrightarrow{rf}$ )

Relates write events to read events that read the stored value (implicit initial writes).

1. Existence and unicity:

$$\forall r, \exists! w, w \xrightarrow{rf} r$$

2. Same location, same value:

$$\text{loc}(w) = \text{loc}(r) \wedge \text{val}(w) = \text{val}(r).$$

6/74

## Example of $\xrightarrow{rf}$

LB	
$T_0$	$T_1$
(a) $r_0 \leftarrow x$	(c) $r_1 \leftarrow y$
(b) $y \leftarrow 1$	(d) $x \leftarrow 1$

Observe:  $r_0; r_1;$

(Initial values of x and y are 0.)

$r_0=1; r_1=1;$

$r_0=1; r_1=0;$

a:  $Rx=1$

c:  $Ry=1$

a:  $Rx=1$

a:  $Rx=1$

c:  $Ry=0$

po ↓

po ↓

po ↓

po ↓

po ↓

po ↓

b:  $Wy=1$

d:  $Wx=1$

b:  $Wy=1$

b:  $Wy=1$

d:  $Wx=1$

$r_0=0; r_1=1;$

$r_0=0; r_1=0;$

rf a:  $Rx=0$

c:  $Ry=1$

po ↓

po ↓

b:  $Wy=1$

d:  $Wx=1$

rf

a:  $Rx=0$

rf c:  $Ry=0$

po ↓

po ↓

b:  $Wy=1$

d:  $Wx=1$

7/74

## A definition of SC

### Definition (SC 1)

An execution is SC when there exists a total order on events  $<$ , such that:

1. Order  $<$  is compatible with program order:

$$e_1 \xrightarrow{po} e_2 \implies e_1 < e_2.$$

2. A Read  $r$  reads from the recent write before  $r$  in  $<$ .

$$\xrightarrow{rf} \stackrel{\text{Def}}{=} \left\{ (w, r) \mid w = \max_{<}(w', \text{loc}(w') = \text{loc}(r) \wedge w' < r) \right\}.$$

8/74

## A question on SC

Program:

R	
$T_0$	$T_1$
(a) $x \leftarrow 1$	(c) $y \leftarrow 2$
(b) $y \leftarrow 1$	(d) $r0 \leftarrow x$
Observed? $y=2; r0=0$	

How can we know? Let us enumerate all interleavings.

$a, b, c, d$	$y=2; r0=1;$
$a, c, b, d$	$y=1; r0=1;$
$a, c, d, b$	$y=1; r0=1;$
$c, d, a, b$	$y=1; r0=0;$
$c, a, b, d$	$y=1; r0=1;$
$c, a, d, b$	$y=1; r0=1;$

**Conclusion:** No SC execution would ever yield the output " $y=2; r0=0;$ ".

9/74

## Modern machines are not SC, how can we know?

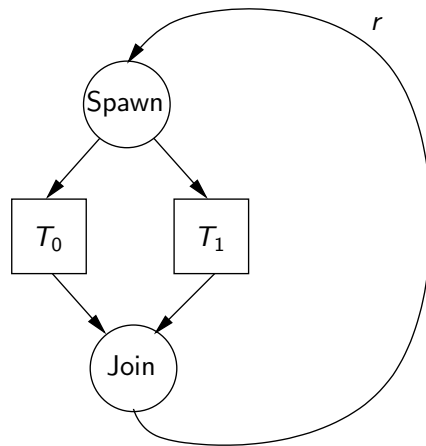
```
void *T0(void *_p) {
    ctx_t *p = _p ;
    common_t *q = p->common ;
    q->x = 1 ;
    q->y = 1 ;
    return NULL ;
}
```

```
void *T1(void *_p) {
    ctx_t *p = _p ;
    common_t *q = p->common ;
    q->y = 2 ;
    int r0 = q->x ;
    q->r0 = r0 ;
    return NULL ;
}
```

```
for ( ; ; ) {
    // Initialise
    common_t c ; c.x = c.y = 0 ; ctx_t a0,a1 ;
    // Run
    a0.id = 0 ; a0.common = &c ; create_thread(&th0,T0,&a0) ;
    a1.id = 1 ; a1.common = &c ; create_thread(&th1,T1,&a1) ;
    join_thread(&th0) ;
    join_thread(&th1) ;
    // Collect results
    ... c.y ... c.r0
}
```

10/74

## Naive testing, graphically



Let us run test **R** on this machine (demo/01/naive\_r.out)

11/74

## Minimizing the impact of thread creation on cost

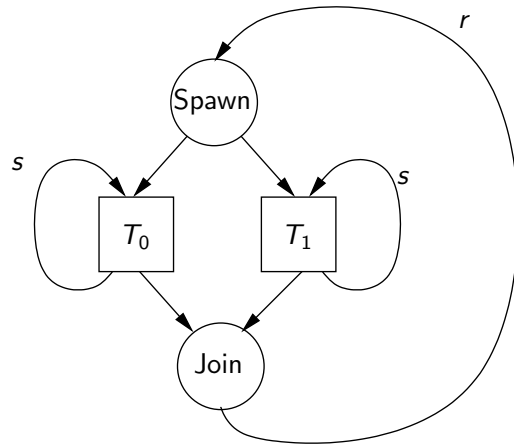
We perform size tests per thread creation, in arrays.

```
void *T0(void *_p) {
    ctx_t *p = _p ;
    common_t *q = p->common ;
    for
    (int k = 0 ;
     k < q->size ;
     k++) {
        q->x[k] = 1 ;
        q->y[k] = 1 ;
    }
    return NULL ;
}
```

```
void *T1(void *_p) {
    ctx_t *p = _p ;
    common_t *q = p->common ;
    for
    (int k = 0 ;
     k < q->size ;
     k++) {
        q->y[k] = 2 ;
        int r0 = q->x[k] ;
        q->r0[k] = r0 ;
    }
    return NULL ;
}
```

12/74

## Graphically



Let us run test **R** on this machine (demo/01/loop\_r.out)

13/74

## Let us synchronise iterations

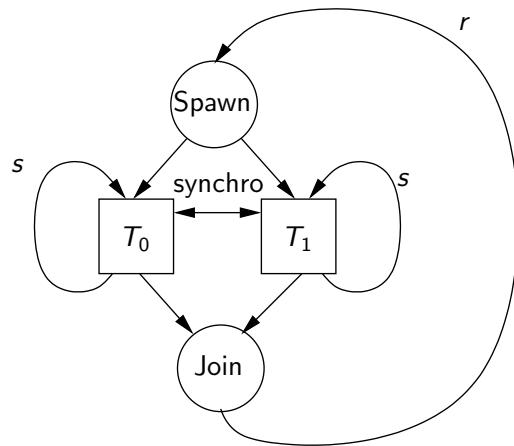
```
void *T0(void *_p) {
    ctx_t *p = _p ;
    common_t *q = p->common ;
    for
        (int k = 0 ;
         k < q->size ;
         k++) {
        wait_partner(&q->sync[k],k,0) ;
        q->x[k] = 1 ;
        q->y[k] = 1 ;
    }
    return NULL ;
}
```

```
void *T1(void *_p) {
    ctx_t *p = _p ;
    common_t *q = p->common ;
    for
        (int k = 0 ;
         k < q->size ;
         k++) {
        wait_partner(&q->sync[k],k,1) ;
        q->y[k] = 2 ;
        int r0 = q->x[k] ;
        q->r0[k] = r0 ;
    }
    return NULL ;
}
```

```
inline static void wait_partner(volatile int *p,int k,int id) {
    if (k % 2 == id) {
        *p = 1 ; __sync_synchronize() ; // Well...
    } else {
        while(*p == 0) ;
    }
}
```

14/74

## Graphically

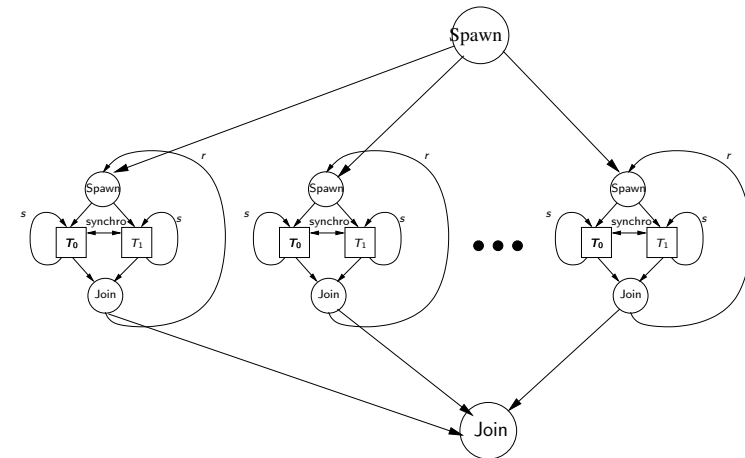


Let us run test **R** on this machine (demo/01/sync\_r.out)

15/74

## Running $n$ test instances at once

Why not run  $n = \lfloor a/2 \rfloor$  instances of **R** on a  $a$ -core machine?



- ▶ Much more outcomes on many-core machines per test run
- ▶ Important when one pays resources by chunks of say 32 cores.
- ▶ Makes noise and favors outcome variability.

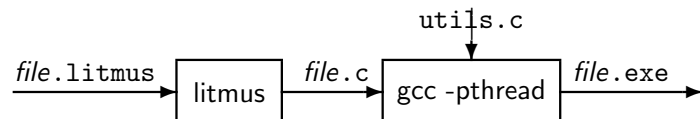
16/74

## The litmus tool

- ▶ Tests are written in assembler, for precision.

```
X86 R
{ }
P0      | P1      ;
MOV [x],$1 | MOV [y],$2 ;
MOV [y],$1 | MOV EAX,[x] ;
exists (y=2 /\ 1:EAX=0)
```

- ▶ Tests are compiled to C programs (with inline assembler):



Demo in demo/02:

```
% litmus -mach ./x86.cfg -o run R.litmus
% cd run
% make
% ./R.exe -v -v
```

17/74

## Some litmus settings

- ▶ Synchronisation
  - ▷ Loose synchronisation (we saw it).
  - ▷ Exact synchronisation with polling synchronisation barriers and time base.
  - ▷ Other, less useful, modes: POSIX thread barrier based and no synchronisation at all.
- ▶ Affinity: force threads to run on designated cores, or let the OS scheduler perform its job.
- ▶ Prefetching of flushing cache lines.
  - ▷ Automatic, depending on test.
  - ▷ Random.
  - ▷ Complete or none.
- ▶ Various scanning order of location arrays.
  - ▷ Random (by the means of a shuffled array of pointer).
  - ▷ Linear with a stride.

A complete testing campaign usually involves trying many settings (for instance, testing all strides from 1 to cache line size).

18/74

## Building significant tests

Perfect! We know how to run tests on hardware, but what tests do we run?

We study *relaxed* memory models, that is *relaxed* w.r.t SC.

Hence, we focus on programs that have non-SC behaviours.

The question is: how do we generate such programs.

Let us study SC in detail first.

19/74

## Back to our non-SC example

R	
$T_0$	$T_1$
(a) $x \leftarrow 1$	(c) $y \leftarrow 2$
(b) $y \leftarrow 1$	(d) $r0 \leftarrow x$
Observed? $y=2; r0=0$	

All interleavings.

$a, b, c, d$	$y=2; r0=1;$
$a, c, b, d$	$y=1; r0=1;$
$a, c, d, b$	$y=1; r0=1;$
$c, d, a, b$	$y=1; r0=0;$
$c, a, b, d$	$y=1; r0=1;$
$c, a, d, b$	$y=1; r0=1;$

We observe if  $b < c$  then  $y=2$ , if  $d < a$  then  $r0=0$ .

20/74

## Let us be a bit more clever

R	
$T_0$	$T_1$
(a) $x \leftarrow 1$	(c) $y \leftarrow 2$
(b) $y \leftarrow 1$	(d) $r0 \leftarrow x$
Observed? $y=2; r0=0$	

Collecting constraints on the scheduling order  $<$ :

We respect program order, thus  $a < b$  and  $c < d$ .

We observe  $r0=0$ , thus  $d < a$ .

We observe  $y=2$ , thus  $b < c$ .

Hence we have a cycle in  $<$ , which prevents it from being an order!

$$a < b < c < d < a \dots$$

**Conclusion:** No SC execution would ever yield the output “ $y=2; r0=0$ ”.

21/74

## Systematic approach

For a particular (partial) execution candidate (that is, for a set of events and a  $\xrightarrow{po}$  relation) we assume two additional relations:

- **Read-from** ( $\xrightarrow{rf}$ ): Relates write events to read events that read the stored value (implicit initial writes).

$$\forall r, \exists! w, w \xrightarrow{rf} r$$

- **Coherence** ( $\xrightarrow{co}$ ): Relates write events to the same location.

For any location  $\ell$ , the restriction of  $\xrightarrow{co}$  to write events to location  $\ell$  ( $\mathbf{W}_\ell$ ) is a total order.

**Notice:** To me, the very existence of  $\xrightarrow{co}$  stems from the existence of a shared, coherent, memory — Given location  $\ell$ , there is exactly one memory cell whose location is  $\ell$ .

22/74

## Example of $\xrightarrow{co}$

2+2W	
$T_0$	$T_1$
(a) $x \leftarrow 1$	(c) $y \leftarrow 1$
(b) $y \leftarrow 2$	(d) $x \leftarrow 2$
Observe: $x; y;$	

$x=1; y=2;$

a:  $W_x=2$       c:  $W_y=2$

po ↓      po ↓

b:  $W_y=1$       d:  $W_x=1$

$x=1; y=1;$

a:  $W_x=2$       c:  $W_y=2$

po ↓      po ↓

b:  $W_y=1$       d:  $W_x=1$

$x=2; y=2;$

a:  $W_x=2$       c:  $W_y=2$

po ↓      po ↓

b:  $W_y=1$       d:  $W_x=1$

$x=2; y=1;$

a:  $W_x=2$       c:  $W_y=2$

po ↓      po ↓

b:  $W_y=1$       d:  $W_x=1$

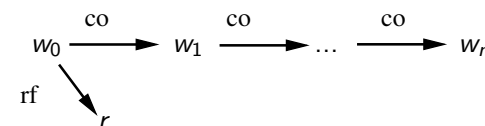
23/74

## One more relation: $\xrightarrow{fr}$

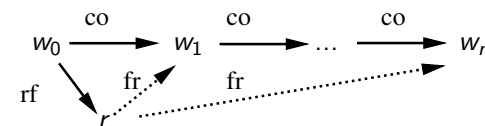
The new relation  $\xrightarrow{fr}$  (from read) relates reads to “younger writes” (younger w.r.t.  $\xrightarrow{co}$ ).

$$r \xrightarrow{fr} w \stackrel{\text{Def}}{=} w' \xrightarrow{rf} r \wedge w' \xrightarrow{co} w$$

This amounts to place a read into the coherence order of its location. Given



We have

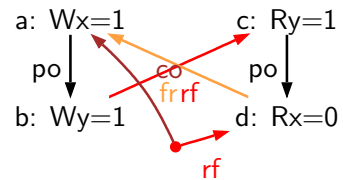


24/74

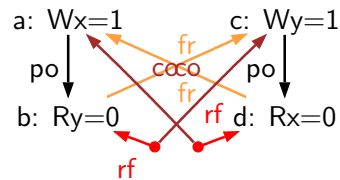
## Playing with $\xrightarrow{fr}$

Particular, easy, case: a read from the initial state is in  $\xrightarrow{fr}$  with writes by the program.

MP	
$T_0$	$T_1$
(a) $x \leftarrow 1$	(c) $r0 \leftarrow y$
(b) $y \leftarrow 1$	(d) $r1 \leftarrow x$
Observed? $r0=1; r1=0$	



SB	
$T_0$	$T_1$
(a) $x \leftarrow 1$	(c) $y \leftarrow 1$
(b) $r0 \leftarrow y$	(d) $r1 \leftarrow x$
Observed? $r0=0; r1=0$	



25/74

## Second definition of SC

### Definition (SC 2)

An execution is SC when:

$$\text{Acyclic} \left( \xrightarrow{rf} \cup \xrightarrow{co} \cup \xrightarrow{fr} \cup \xrightarrow{po} \right)$$

And of course:

### Theorem

The two definitions of SC are equivalent.

26/74

## Introducing herd a memory model simulator

A model sc.cat:

```
% cat sc.cat
"Sequential consistency"
let com = rf | co | fr
acyclic (po | com) as hb
```

Running **R** on SC (demo in demo/02):

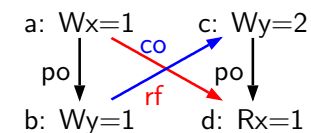
```
Test R Allowed
States 3
1:EAX=0; y=1;
1:EAX=1; y=1;
1:EAX=1; y=2;
No
Witnesses
Positive: 0 Negative: 3
Condition exists (y=2 /\ 1:EAX=0)
Observation R Never 0 3
```

**Notice:** Outcome 1:EAX=0; y=2; is forbidden by SC.

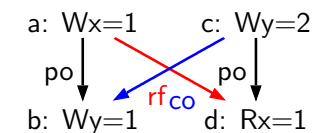
27/74

## Herd structure

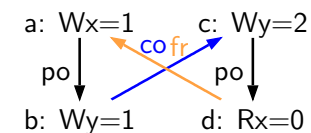
- Generate all candidate executions, i.e. all possible  $\xrightarrow{po}$ ,  $\xrightarrow{rf}$  and  $\xrightarrow{co}$  ( $\xrightarrow{fr}$  deduced):



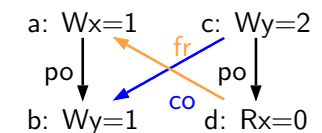
Ok



Ok



No



Ok

- Apply model checks to each candidate execution.

28/74

## Violations of SC

A cycle of  $\xrightarrow{po}$ ,  $\xrightarrow{rf}$ ,  $\xrightarrow{co}$ ,  $\xrightarrow{fr}$  describes a violation of SC.

From such a cycle, one may easily generate programs that potentially violate SC, and run them on an actual machine.

However, the cycle does not describe:

- ▶ How many threads are involved.
- ▶ How many memory locations are involved.

We now aim at:

- ▶ Extract a subset of *significant* cycles.
- ▶ Generate *one* program out of one cycle.

29/74

## Simplifying cycles: $\xrightarrow{po}$ and $\widehat{\xrightarrow{com}}$ steps alternate

A cycle in  $\xrightarrow{com} \cup \xrightarrow{po}$  is a cycle in  $(\xrightarrow{po}^+; \widehat{\xrightarrow{com}}^+)$  (group  $\xrightarrow{po}$  and  $\widehat{\xrightarrow{com}}$  steps together). Then:

- ▶  $\xrightarrow{po}$  is transitive  $\xrightarrow{po}^+ \subseteq \xrightarrow{po}$ .
- ▶  $\widehat{\xrightarrow{com}}^+$  is the union of the five following relations:

$$\widehat{\xrightarrow{com}} = \xrightarrow{rf} \cup \xrightarrow{co} \cup \xrightarrow{fr} \cup \left( \xrightarrow{co}; \xrightarrow{rf} \right) \cup \left( \xrightarrow{fr}; \xrightarrow{rf} \right).$$

Because  $(\xrightarrow{co}; \xrightarrow{co}) \subseteq \xrightarrow{co}$ ,  $(\xrightarrow{fr}; \xrightarrow{co}) \subseteq \xrightarrow{fr}$ , and  $(\xrightarrow{rf}; \xrightarrow{fr}) \subseteq \xrightarrow{co}$ .

**Conclusion:** Any cyclic  $\xrightarrow{com} \cup \xrightarrow{po}$  includes a cycle in  $(\xrightarrow{po}; \widehat{\xrightarrow{com}})$  — i.e. that alternates  $\xrightarrow{po}$  steps and  $\widehat{\xrightarrow{com}}$  steps.

30/74

## Simplifying cycles: all $\widehat{\xrightarrow{com}}$ steps are external

Given a cycle, we consider all  $\xrightarrow{com}$  and  $\widehat{\xrightarrow{com}}$  steps are external, that is source and target events are from pairwise distinct thread.

Given  $e_1 \widehat{\xrightarrow{com}} e_2$ , s.t.  $e_1$  and  $e_2$  are from the same thread:

- ▶ Either  $e_1 \xrightarrow{po} e_2$  and we consider this  $\xrightarrow{po}$  step in the cycle, in place of the  $\widehat{\xrightarrow{com}}$  step.
- ▶ Or  $e_2 \xrightarrow{po} e_1$  and we have a very simple cycle  $e_2 \xrightarrow{po} e_1 \widehat{\xrightarrow{com}} e_2$ . Such cycles are “*violations of coherence*” (more on them later).

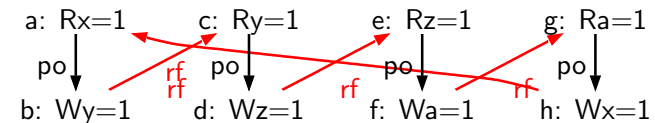
**Notice:** The same reasoning applies individual  $\xrightarrow{com}$  steps in composite  $\widehat{\xrightarrow{com}}$ .

31/74

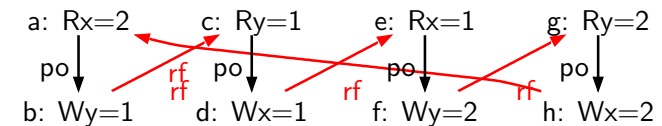
## Simplifying cycles — Locations

Cycle:  $W \xrightarrow{po} W \xrightarrow{rf} R \xrightarrow{po} R \xrightarrow{fr} W \xrightarrow{po} W \xrightarrow{rf} R \xrightarrow{po} R \xrightarrow{fr}$

- ▶ One interpretation (four locations):



- ▶ Another interpretation (two locations):

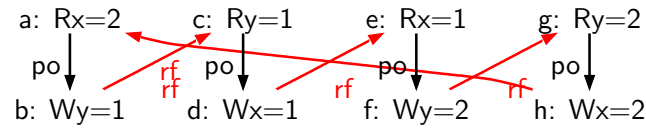


32/74



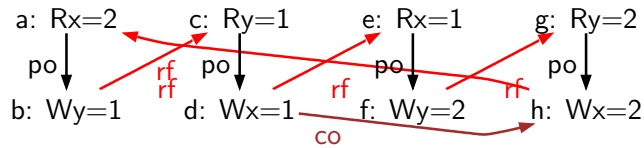
## The second interpretation is not “minimal”

Reminding the interpretation with two locations:



But, coherence  $\xrightarrow{co}$  totally orders write events to a given location.

Let us choose:  $Wx1 \xrightarrow{co} Wx2$ :



We have a smaller cycle:  $d \xrightarrow{co} h \xrightarrow{rf} a \xrightarrow{po} b \xrightarrow{rf} c \xrightarrow{po} d$ .

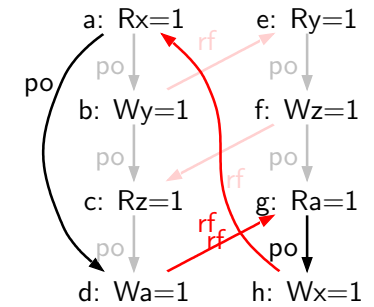
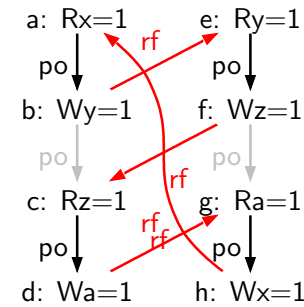
Choosing  $Wx2 \xrightarrow{co} Wx1$  would yield another smaller cycle.

Generally: do not repeat locations in cycles.

33/74

## Simplifying cycles — Threads

Cycle:  $W \xrightarrow{po} W \xrightarrow{rf} R \xrightarrow{po} R \xrightarrow{fr} W \xrightarrow{po} W \xrightarrow{rf} R \xrightarrow{po} R \xrightarrow{fr}$



Generally: one passage per thread

34/74

## ... Simplifying cycles

In a non SC execution we find:

- ▶ A *violation of coherence*, that is a cycle  $e_1 \xrightarrow{po} e_2 \xrightarrow{\widehat{com}} e_1$ .
- ▶ Or a *critical cycle* that is:
  - ▶ The cycle alternates  $\xrightarrow{po}$  steps and external  $\xrightarrow{\widehat{com}}$  steps, with at least four steps.
  - ▶ The cycle passes through a given thread at most once.
  - ▶ All  $\xrightarrow{\widehat{com}}$  steps have pairwise different locations.
  - ▶ The source and target of one given  $\xrightarrow{po}$  steps have different locations.

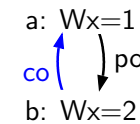
**Notice:** For a more formal presentation see D. Shasha and M. Snir Toplas 88 article, which introduced critical cycles.

35/74

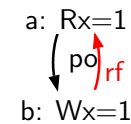
## Violations of coherence

There are five such cycles, which can occur as the following executions:

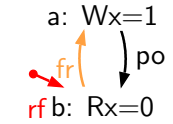
$\xrightarrow{po}$  contradicts  $\xrightarrow{co}$ ,  $\xrightarrow{rf}$ ,  $\xrightarrow{fr}$ , “ $\xrightarrow{co}$ ;  $\xrightarrow{rf}$ ”, “ $\xrightarrow{fr}$ ;  $\xrightarrow{rf}$ ”.



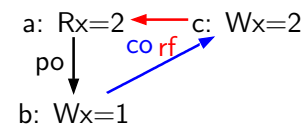
CoWW



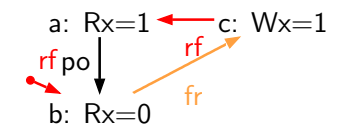
CoRW1



CoWR



CoRW2



CoRR

36/74

## Building a test from a cycle

The conditions on locations and threads allows generating one test (program + final condition) from one cycle.

Consider for instance:

$$\xrightarrow{\text{fre}} \xrightarrow{\text{rfe}} \xrightarrow{\text{po}} \xrightarrow{\text{rfe}} \xrightarrow{\text{po}}$$

- Set events (with directions, from communication relations)

$$\xrightarrow{\text{fre}} (a):W[x]? \xrightarrow{\text{rfe}} (b):R[x]? \xrightarrow{\text{po}} (c):W[y]? \xrightarrow{\text{rfe}} (d):R[y]? \xrightarrow{\text{po}} (e):W[x]?$$

- Set locations (one per  $\widehat{\text{com}}$ )

$$\xrightarrow{\text{fre}} (a):W[x]? \xrightarrow{\text{rfe}} (b):R[x]? \xrightarrow{\text{po}} (c):W[y]? \xrightarrow{\text{rfe}} (d):R[y]? \xrightarrow{\text{po}} (e):R[x]?$$

$$\xrightarrow{\text{fre}} (a):W[x]? \xrightarrow{\text{rfe}} (b):R[x]? \xrightarrow{\text{po}} (c):W[y]? \xrightarrow{\text{rfe}} (d):R[y]? \xrightarrow{\text{po}} (e):R[x]?$$

- Set values for writes (initial value 0, then follow cycle: 1, 2, etc.)

$$\xrightarrow{\text{fre}} (a):W[x]1 \xrightarrow{\text{rfe}} (b):R[x]? \xrightarrow{\text{po}} (c):W[y]1 \xrightarrow{\text{rfe}} (d):R[y]? \xrightarrow{\text{po}} (e):R[x]?$$

- Set values for reads (consider  $\xrightarrow{\text{rfe}} R$  or  $R \xrightarrow{\text{fre}}$ )

$$\xrightarrow{\text{fre}} (a):W[x]1 \xrightarrow{\text{rfe}} (b):R[x]1 \xrightarrow{\text{po}} (c):W[y]1 \xrightarrow{\text{rfe}} (d):R[y]1 \xrightarrow{\text{po}} (e):R[x]0$$

37/74

## Building a test from a cycle, continued

$$\xrightarrow{\text{fre}} (a):W[x]1 \xrightarrow{\text{rfe}} (b):R[x]1 \xrightarrow{\text{po}} (c):W[y]1 \xrightarrow{\text{rfe}} (d):R[y]1 \xrightarrow{\text{po}} (e):R[x]0$$

- Build program from events (change thread at every  $\widehat{\text{com}}$ ).

WRC		
$T_0$	$T_1$	$T_2$
(a) x ← 1	(b) r0 ← x	(d) r1 ← y
	(c) y ← 1	(e) r2 ← x

- Build condition from coherence sequences (here nothing) and from values read.

Observed? r0=1; r1=1; r2=0;

38/74

## Building another test, non-trivial coherence

Test R

$$\xrightarrow{\text{fre}} (a):W[x]1 \xrightarrow{\text{po}} (b):W[y]1 \xrightarrow{\text{coe}} (c):W[y]2 \xrightarrow{\text{po}} (x):R[x]0$$

R	
$T_0$	$T_1$
(a) x ← 1	(c) y ← 2
(b) y ← 1	(d) r0 ← x
Observed? y=2; r0=0	

Here, coherence order is  $0 \xrightarrow{\text{co}} 1 \xrightarrow{\text{co}} 2$ , it suffices to read the final value.

Longer coherence orders command other techniques, for instance adding an observer thread.

39/74

## A first tool: diyone

Generating WRC:

```
% diyone -arch X86 Rfe Pod** Rfe Pod** Fre
X86 A
P0          | P1          | P2          ;
MOV [x], $1 | MOV EAX, [x] | MOV EAX, [y] ;
              | MOV [y], $1  | MOV EBX, [x] ;
exists (1:EAX=1 /\ 2:EAX=1 /\ 2:EBX=0)
```

Doing the same for ARM is as simple as:

```
% diyone -arch ARM Rfe Pod** Rfe Pod** Fre
ARM A
{ %x0=x; %x1=x; %y1=y; %y2=y; %x2=x; }
P0          | P1          | P2          ;
MOV R0, #1  | LDR R0, [%x1] | LDR R0, [%y2] ;
STR R0, [%x0] | MOV R1, #1    | LDR R1, [%x2] ;
              | STR R1, [%y1] |              ;
exists (1:R0=1 /\ 2:R0=1 /\ 2:R1=0)
```

40/74

## Tool diyone, generating R

```
% diyone -arch X86 PodWW Wse PodWR Fre
X86 A
"PodWW Wse PodWR Fre"
{ }
  PO          | P1          ;
  MOV [x],$1 | MOV [y],$2 ;
  MOV [y],$1 | MOV EAX,[x] ;
exists (y=2 /\ 1:EAX=0)
```

**Notice:** We wrote PodWW, PodWR. The vocabulary of *Candidate Relaxations* is quite rich:

- ▶ Internal communications Rfi, Fri, Wsi.
- ▶  $\xrightarrow{po}$  edges with identical target and source locations: PosRR, etc.
- ▶ Dependencies (DpAddrW, etc.). fences (MFencedWR, etc.)

41/74

## Application, testing non-SC executions for two threads

Cycle → execution → program + final condition.

All (critical) cycles for two threads: six cycles.

<b>2+2W</b>	$\xrightarrow{po} \xrightarrow{co} \xrightarrow{po} \xrightarrow{co}$
<b>LB</b>	$\xrightarrow{po} \xrightarrow{rf} \xrightarrow{po} \xrightarrow{rf}$
<b>MP</b>	$\xrightarrow{po} \xrightarrow{rf} \xrightarrow{po} \xrightarrow{fr}$
<b>R</b>	$\xrightarrow{po} \xrightarrow{co} \xrightarrow{po} \xrightarrow{fr}$
<b>S</b>	$\xrightarrow{po} \xrightarrow{rf} \xrightarrow{po} \xrightarrow{co}$
<b>SB</b>	$\xrightarrow{po} \xrightarrow{fr} \xrightarrow{po} \xrightarrow{fr}$

Save of coherence violation, any non-SC execution on two threads includes one of the above six cycles.

Hence, testing the six tests built from the six cycles gives reasonable coverage of possible SC violation on two threads. (**Notice:** coherence violations neglected).

42/74

## Generating two-threads SC violations

The tool diy generates cycles (and tests) from a vocabulary of CR. It can be configured for the two threads case as follows:

```
-arch X86          # target architecture
-safe Pod**,Rfe,Fre,Wse # vocabulary
-nprocs 2          # 2 procs
-size 4            # max size of cycle (2 X nprocs)
-num false        # for naming tests
```

Demo in demo/03.

```
% diy -conf 2.conf
Generator produced 6 tests
% ls
2+2W.litmus 2.conf @all LB.litmus
MP.litmus R.litmus SB.litmus S.litmus
% diy -conf 4.conf
Generator produced 68 tests...
```

43/74

## Demo 03 continued, running the tests

Compiling:

```
% litmus -mach ./x86.cfg src/@all -o run
% make -C run -j 4
```

Running:

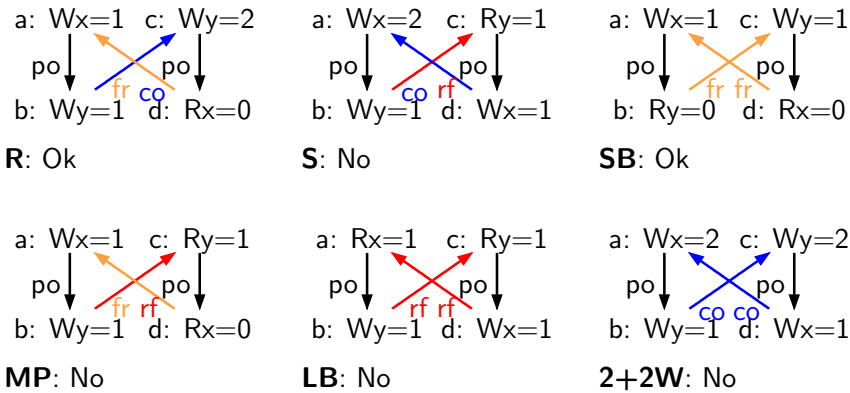
```
% cd run
% sh run.sh > X.00
```

Analysis:

```
% grep Observation X.00
Observation R Sometimes 79 1999921
Observation MP Never 0 2000000
Observation 2+2W Never 0 2000000
Observation S Never 0 2000000
Observation SB Sometimes 1194 1998806
Observation LB Never 0 2000000
```

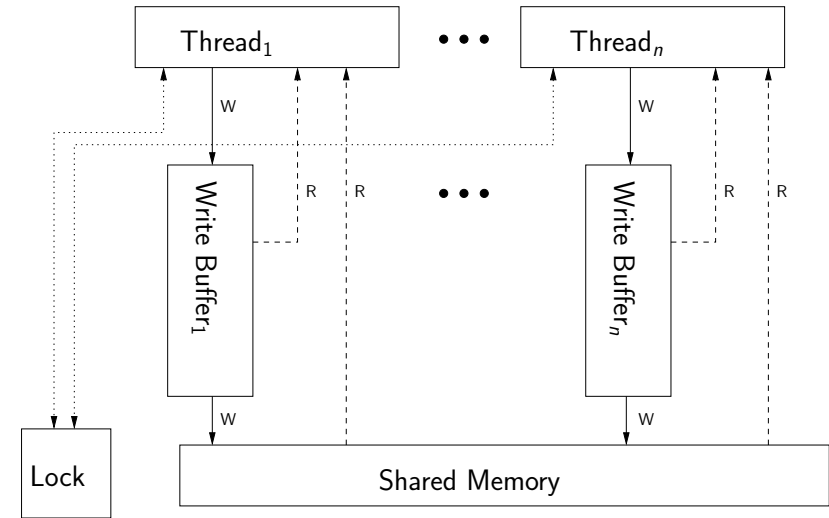
44/74

## Results for running the six tests on this machine



45/74

## TSO — The Model of X86 machines



The write buffer explains how “reads can pass over writes”.

46/74

## Axiomatic TSO

- Remember SC:

$$\text{Acyclic} \left( \xrightarrow{\text{rf}} \cup \xrightarrow{\text{co}} \cup \xrightarrow{\text{fr}} \cup \xrightarrow{\text{po}} \right)$$

A model for herd, our generic simulator:

```
let ppo = po # ppo stands for 'preserved program-order'
let com-hb = fr | rf | co # All communications create order
acyclic (ppo | com-hb)
```

- In TSO:

- Write-to-read does not create order:

```
let ppo = RM(po) | WW(po) # WR(po) omitted
```

- Local reads do not create order:

```
let com-hb = rfe | fr | co # rfi omitted
```

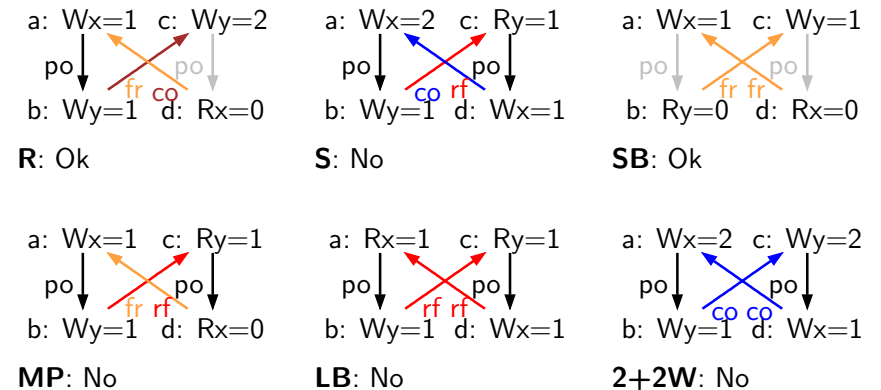
- TSO “happens-before” (HB) check:

```
acyclic (ppo | com-hb | mfence) as hb
```

**Notice:** Relations are between the points in time where a load binds its value and where a written value reaches memory.

47/74

## Results for running the six test on the model

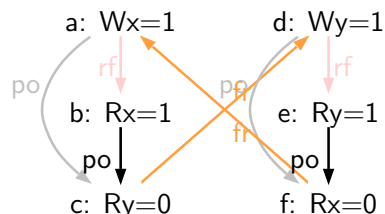


48/74

Internal  $\xrightarrow{rf}$  ( $\xrightarrow{rfi}$ ) are not in HB

SB+rfi-pos	
$T_0$	$T_1$
(a) $x \leftarrow 1$	(d) $y \leftarrow 1$
(b) $r0 \leftarrow x$	(e) $r2 \leftarrow y$
(c) $r1 \leftarrow y$	(f) $r3 \leftarrow x$

Observed?  $r0=1; r1=0; r2=1; r3=0;$

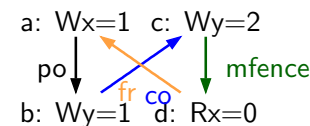


SB+rfi-pos: Ok

Restoring SC with mfence

R+po+mfence	
$T_0$	$T_1$
(a) $x \leftarrow 1$	(c) $y \leftarrow 2$
(b) $y \leftarrow 1$	mfence
	(d) $r0 \leftarrow x$

Observed?  $y=2; r0=0$

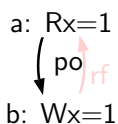


We are not done yet. . .

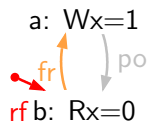
Our TSO model:

```
let ppo = RM(po) | WW(po) # WR(po) omitted
let com-hb = rfe | fr | co # rfi omitted
acyclic (ppo | com-hb)
show ppo | com-hb as hb
```

Allows two violations of coherence:



CoRW1



CoWR

$\xrightarrow{rfi}$  not in  $\xrightarrow{hb}$

$W \xrightarrow{po} R$  not in  $\xrightarrow{hb}$

Those behaviours *must* be rejected by our TSO model.

Correct axiomatic TSO

We add a specific UNIPROC check to rule out coherence violations:

$$\text{Irreflexive} \left( \xrightarrow{po-loc}; \widehat{\text{com}} \right)$$

Where  $\xrightarrow{po-loc}$  is  $\xrightarrow{po}$  between accesses to the same memory location.

irreflexive (po-loc; com+) as uniproc

...

In the TSO case we can "optimise":

irreflexive rf;RW(po-loc)

irreflexive fr;WR(po-loc)

because the other coherence violations are rejected by the HB check.

## A word on UNIPROC

From cycle analysis, we have the attractive definition (since relying on local action of the core and on the existence of coherence orders):

### Definition (Uniproc 1)

Program order  $\xrightarrow{po}$  does not contradict communication  $\xrightarrow{com^+}$ .

There is another definition “SC per location”. (Jason F. Cantin, Mikko H. Lipasti, James E. Smith ACM Symposium on Parallel Algorithms and Architectures 2004).

### Definition (Uniproc 2)

Relation  $\xrightarrow{po-loc} \cup \xrightarrow{com}$  is acyclic.

Definitions are equivalent.

It suffices to show that the existence of a cycle in  $\xrightarrow{po-loc} \cup \xrightarrow{com}$  implies the existence of a coherence violation (i.e. a cycle  $e_1 \xrightarrow{po} e_2 \xrightarrow{\widehat{com}} e_1$ ).

53/74

## Consequence of $\xrightarrow{co}$ ordering writes

### Lemma (Identical locations)

Let  $e_1, e_2$  be two different events with the same location,

1. either  $e_1 \xrightarrow{\widehat{com}} e_2$ ,
2. or  $e_2 \xrightarrow{\widehat{com}} e_1$ ,
3. or  $w \xrightarrow{rf} e_1$  and  $w \xrightarrow{rf} e_2$ .

Case analysis:

- ▶  $w_1, w_2$ , then either  $w_1 \xrightarrow{co} w_2$  or  $w_2 \xrightarrow{co} w_1$  (total order).
- ▶  $r_1, r_2$ , let  $w_1 \xrightarrow{rf} r_1$  and  $w_2 \xrightarrow{rf} r_2$ . Then, either  $w_1 = w_2$  and we are in case 3; or (for instance)  $w_1 \xrightarrow{co} w_2$  and we have  $r_1 \xrightarrow{fr} w_2 \xrightarrow{rf} r_2$ .
- ▶  $r_1, w_2$ , let  $w_1 \xrightarrow{rf} r_1$ . Then, either  $w_1 = w_2$  and  $w_2 \xrightarrow{rf} r_1$ ; or  $w_1 \xrightarrow{co} w_2$  and  $r_1 \xrightarrow{fr} w_2$ ; or  $w_2 \xrightarrow{co} w_1$  and  $w_2 \xrightarrow{co} \xrightarrow{rf} r_1$ .

**Corollary:**  $\xrightarrow{com}$  is acyclic.

54/74

## Equivalence of the two UNIPROC definitions

Proof is easy from “Identical locations” lemma.

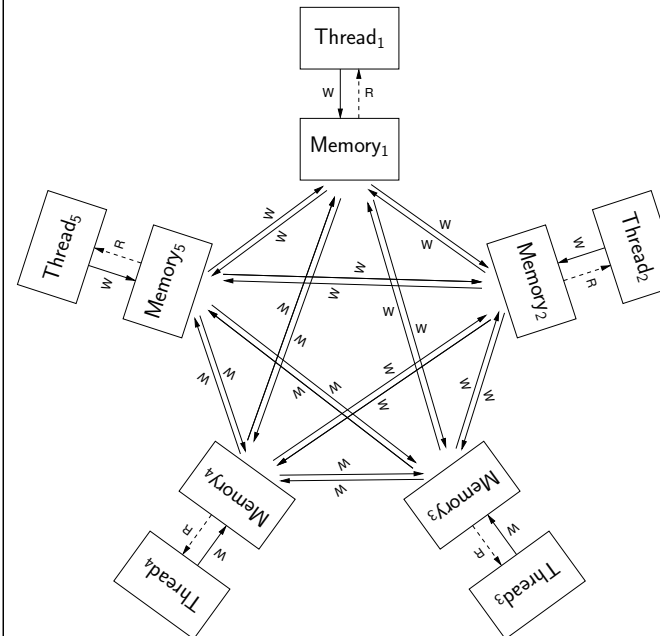
Consider a cycle in  $\xrightarrow{po} \cup \xrightarrow{com}$ .

- ▶ If there exists a  $e_1 \xrightarrow{po} e_2$  step s.t.  $e_2 \xrightarrow{\widehat{com}} e_1$ , then we are done.
- ▶ Otherwise, for each  $e_1 \xrightarrow{po} e_2$  step:
  - ▶ Either,  $r_1 \xrightarrow{po} r_2$ , with  $w \xrightarrow{rf} r_1$  and  $w \xrightarrow{rf} r_2$ . We short-circuit the  $\xrightarrow{po}$  step, replacing  $w \xrightarrow{rf} r_1 \xrightarrow{po} r_2$  by  $w \xrightarrow{rf} r_2$ .
  - ▶ Or,  $e_1 \xrightarrow{\widehat{com}} e_2$ . We replace the  $\xrightarrow{po}$  step by  $\xrightarrow{com}$  steps.

As a result we have a cycle in  $\xrightarrow{com}$ , which is impossible.

55/74

## A relaxed shared memory computer



More or less visible to user code:

- ▶ Cores:
  - ▶ Out of order execution
  - ▶ Branch speculation
  - ▶ Write buffers
- ▶ Memory
  - ▶ Physically distributed
  - ▶ Caches

56/74

## Situation of (our) ARM/Power models

- ▶ **Architecture public reference** Informal, cannot clearly explain how fences restore SC for instance.
- ▶ **Simple, global-time model:** (CAV'10) too relaxed. It remains useful as it supports simple reasoning on SC-violations (CAV'11).
- ▶ **Operational model:** (PLDI'11) more precise, developed with IBM experts. It is quite complex, and the simulator is very slow.
- ▶ **Multi-event axiomatic model:** (CAV'12) more precise (equivalent to PLDI'11), uses several events per access.
- ▶ **Single-event axiomatic model:** (...) more precise (proved to be more relaxed than PLDI'11, experimentally equivalent). A more simple axiomatic model.

Joint work with (in order of appearance) Jade Alglave, Susmit Sarkar, Peter Sewell, Derek Williams, Kayvan Memarian, Scott Owens, Mark Batty, Sela Mador-Haim, Rajeev Alur, Milo M. K. Martin and Michael Tautschnig.

57/74

## Some issues for ARM/Power

- ▶ No simple preserved-program-order. More precisely,  $\xrightarrow{ppo}$  will now account for core constraints, such as dependencies.
- ▶ Communication relations alone do not define happen-before steps.
- ▶ A variety of memory fences: lightweight (Power lwsync) and full (Power sync).

58/74

## Two-threads SC violation for ARM

Generating tests is as simple as:

```
% diy -conf 2.conf -arch ARM
```

With the same configuration file 2.conf as for X86.

Then, compile (in two steps, generate C locally, compile it on target machine), run and...

```
Observation R Sometimes 5722 1994278
Observation MP Sometimes 3571 1996429
Observation 2+2W Sometimes 17439 1982561
Observation S Sometimes 7270 1992730
Observation SB Sometimes 9788 1990212
Observation LB Sometimes 4782 1995218
```

All Non-SC behaviours observed!

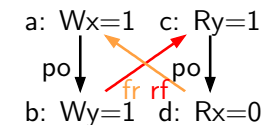
No hope to define  $\xrightarrow{ppo}$  as simply as for TSO.

59/74

## An experiment on ARM/Power

Consider test **MP**:

MP	
$T_0$	$T_1$
(a) $x \leftarrow 1$	(c) $r0 \leftarrow y$
(b) $y \leftarrow 1$	(d) $r1 \leftarrow x$
Observed?	$r0=1; r1=0$



We know that the test is Ok (observed, valid) on ARM/Power, what does it take (amongst fences, dependencies,) to make the test No (unobserved, invalid)?

- ▶ Fences: dsb, dmb, isb (ARM); sync, lwsync, isync (Power).
- ▶ Dependencies: address, data, control, control+isb/isync.

60/74

## Dependencies (Power)

Address dependency:

```

r1 ← x      lwz r1,0(r8) # r8 contains the address of 'x'
r2 ← t[r1]  slwi r7,r1,2 # sizeof(int) = 4
            lwzx r2,r7,r9 # r9 contains the address of 't'
    
```

Data dependency:

```

r1 ← x      lwz r1,0(r8) # r8 contains the address of 'x'
y ← r1+1    addi r2,r1,1
            stw r2,0(r9) # r9 contains the address of 'y'
    
```

Control dependency: (+isync)

```

            lwz r1,0(r8)
            cmpwi r1,0
            bne L1
            (isync)
            li r2,1
            stw r2,0(r9)
L1:
    
```

61/74

## Generating tests (ARM), yet another tool: diycross

Generating tests with diycross (demo in demo/04):

```

% diycross -arch ARM\
PodWW,DMBdWW,DSBdWW,ISBdWW\
Rfe\
PodRR,DpCtrlrDR,DpCtrlIsbDR,DpAddrDR,DMBdRR,DSBdRR,ISBdRR\
Fre
    
```

Generator produced 28 tests

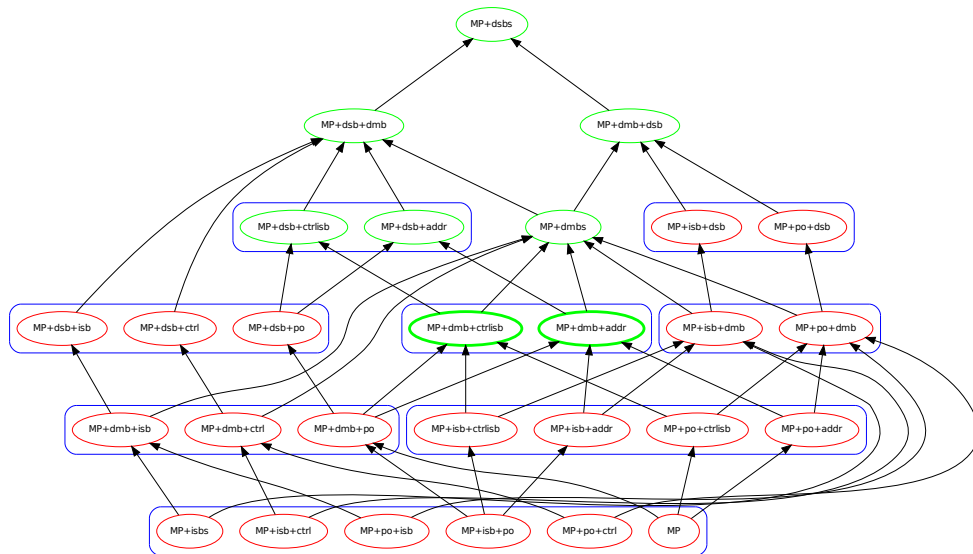
- ▶ One generates **MP** as diyone PodWW Rfe PodRR Fre
- ▶ diycross  $r_1^1, \dots, r_{N_1}^1 \dots r^M, \dots, r_{N_M}^M$ , generates the  $N_1 \times \dots \times N_M$  cycles  $r_{k_1}^1 \dots r_{k_\ell}^\ell \dots r_{k_M}^M$  by *cross-producting* the given CR list arguments.

This generates some variations in the **MP** family.

We then compile and run, and...

62/74

## Optimal fencing/dependencies for MP



63/74

## Optimal fencing for the 6 two-threads tests (Power)

a: Wx=1 c: Wy=2  
 sync ↓ cofr ↗  
 b: Wy=1 d: Rx=0

R+syncs

a: Wx=2 c: Ry=1  
 lwsync ↓ rfc ↗  
 b: Wy=1 d: Wx=1

S+lwsync+addr

a: Wx=1 c: Wy=1  
 sync ↓ frfr ↗  
 b: Ry=0 d: Rx=0

SB+syncs

a: Wx=1 c: Ry=1  
 lwsync ↓ rffr ↗  
 b: Wy=1 d: Rx=0

MP+lwsync+addr

a: Rx=1 c: Ry=1  
 addr ↓ rfrf ↗  
 b: Wy=1 d: Wx=1

LB+adds

a: Wx=2 c: Wy=2  
 lwsync ↓ cco ↗  
 b: Wy=1 d: Wx=1

2+2W+lwsyncs

64/74



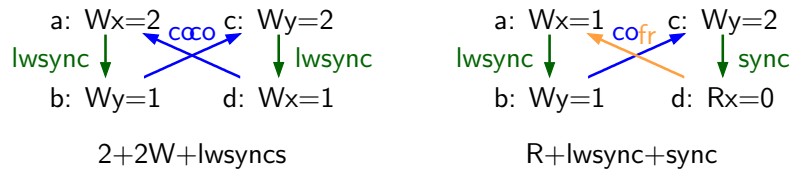
## Some observations

In the previous slide we considered increasing power (and cost):

$$addr < lwsync < sync$$

Then:

- Dependencies (address) are sufficient to restore order from reads to writes and reads in two-threads examples (but...)
- Fences restore order from writes to write and reads.
- Full fence (`sync`) is required from write to read.
- When to use the lightweight fence between writes is complex: **2+2W+lwsyncs** vs. **R+lwsync+sync**.

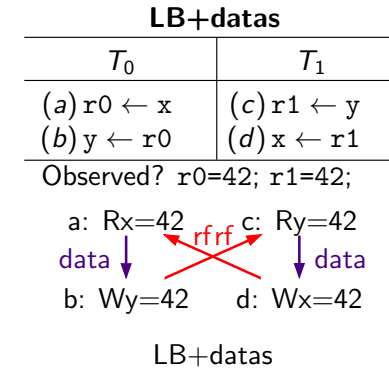


No

Ok

65/74

## Dependencies are enough



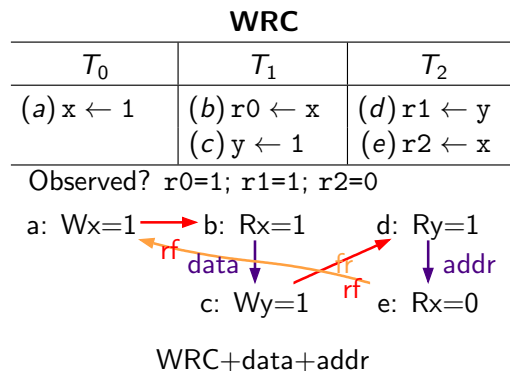
Of course we never observe this behaviour (values out of thin air) and any (hardware) model should forbid it.

**Happens-before** If we order: (1) stores: the point in time when the value is made available to other threads (2) loads: the point when the value is read by core.

66/74

## Dependencies from reads not always enough!

Consider test **WRC+data+addr**:



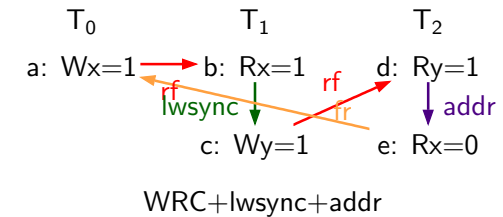
Behaviour observed on Power 6 and 7 (not on ARM, but documentation allows it).

**Stores are not “multi-copy atomic”**  $T_0$  and  $T_1$  share a private buffer/cache/memory (e.g. a cache in SMT context).  $T_2$  “does not see” the store by  $T_0$ , when  $T_1$  does.

67/74

## Restoring SC for WRC

Use a lightweight fence on  $T_1$ :



**Observation:** The fence orders the writes  $a$  (by  $T_0$ ) and  $c$  (by  $T_1$ ) for any observer (here  $T_2$ ).

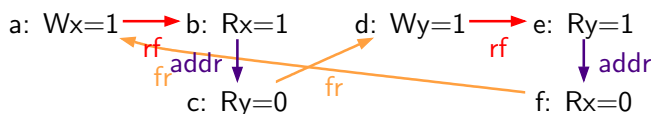
68/74

## Another case of insufficient dependencies

Consider test **IRIW+addr**:

IRIW			
$T_0$	$T_1$	$T_2$	$T_3$
(a) $x \leftarrow 1$	(b) $r0 \leftarrow x$ (c) $r1 \leftarrow y$	(d) $y \leftarrow 1$	(e) $r2 \leftarrow y$ (f) $r3 \leftarrow x$

Observed?  $r0=1; r1=0; r2=1; r3=0;$



IRIW+addr

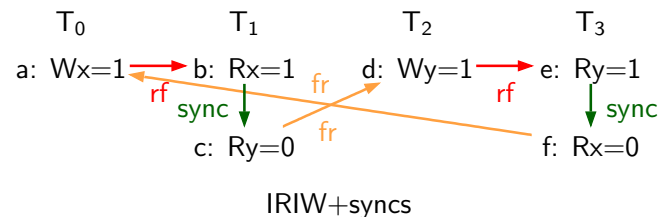
Behaviour observed on Power (not on ARM, but documentation allows it).

**Stores are not “multi-copy atomic”:**  $T_0$  and  $T_1$  have a private buffer/cache/memory,  $T_2$  and  $T_3$  also have one.

69/74

## Restoring SC for IRIW

Use a full fence on  $T_1$  and  $T_2$ :



IRIW+syncs

**Propagation:** Full fences order all communications.

70/74

## Relation summary

Communication relations:

- ▶ Read-from:  $w \xrightarrow{rf} r$ , with  $\text{loc}(w) = \text{loc}(r)$ ,  $\text{val}(w) = \text{val}(r)$ .
- ▶ Coherence:  $w \xrightarrow{co} w'$ , with  $\text{loc}(w) = \text{loc}(w') = x$ . Total order for given  $x$ : hence “coherence orders”.
- ▶ We deduce from-read:  $r \xrightarrow{fr} w$ , i.e.  $w' \xrightarrow{rf} r$  and  $w' \xrightarrow{co} w$ .
- ▶ We distinguish internal (same proc,  $\xrightarrow{rfi}$ ,  $\xrightarrow{coi}$ ,  $\xrightarrow{fri}$ ) and external (different procs,  $\xrightarrow{rfe}$ ,  $\xrightarrow{coe}$ ,  $\xrightarrow{fre}$ ) communications.

“Execution” relations

- ▶ Program order:  $e_1 \xrightarrow{po} e_2$ , with  $\text{proc}(e_1) = \text{proc}(e_2)$ .
- ▶ Same location program order:  $e_1 \xrightarrow{po-loc} e_2$ .
- ▶ Preserved program order:  $e_1 \xrightarrow{ppo} e_2$ , with  $\xrightarrow{ppo} \subseteq \xrightarrow{po}$ . Computed from other relations.
- ▶ Fences: effective strong and lightweight fences in between events  $\xrightarrow{\text{strong}}$  and  $\xrightarrow{\text{light}}$ . Effective means that for instance  $w \xrightarrow{\text{lwsync}} r$  does not implies  $w \xrightarrow{\text{light}} r$ .

71/74

## A model in four checks

UNIPROC

acyclic poloc | com as uniproc

HB

```
let fence = strong | light
let hb = ppo | fence | rfe
acyclic hb
```

OBSERVATION We now define the effect of fences (any fence) for ordering writes:

```
let propbase = (WW(fence)|(rfe;RW(fence)));hb*
irreflexive fre;propbase as observation
```

PROPAGATION Strong fences wait for all communications.

```
let propstrong = com*; propbase*; strong; hb*
let prop = WW(propbase)|(com*;propbase*;strong;hb*)
acyclic co | prop as propagation
```

72/74

## ARM/Power preserved program order

Rather complex, results from a two events per access analysis (cf. CAV'12).

(\* Utilities \*)

```
let dd = addr | data          let rdw = po-loc & (fre;rfe)
let detour = po-loc & (coe ; rfe) let addrpo = addr;po
```

(\* Initial value \*)

```
let ci0 = ctrllsync | detour
let ii0 = dd | rfi | rdw
let cc0 = dd | po-loc | ctrl | addrpo
let ic0 = 0
```

(\* Fixpoint from i -> c in instructions and transitivity \*)

```
let rec ci = ci0 | (ci;ii) | (cc;ci)
and ii = ii0 | ci | (ic;ci) | (ii;ii)
and cc = cc0 | ci | (ci;ic) | (cc;cc)
and ic = ic0 | ii | cc | (ic;cc) | (ii ; ic)
```

```
let ppo = RW(ic) | RR(ii)
```

Can be limited to dependencies...

73/74

## How good is our model?

Is it sound?

- ▶ A proof: any behaviour allowed is also allowed by the operational model of PLDI'11.
- ▶ Experiments
  - ▶ Soundness w.r.t. hardware (ARM being a bit problematic because of acknowledged read-after-read hazard).
  - ▶ Experimental equivalence with our previous models, saved from current debate on some subtle semantical point for lwsync.

In any case:

- ▶ Simulation is fast ( $\times 1000$  w.r.t. PLDI'11) ( $\times 10$  w.r.t. CAV'12).
- ▶ The existence of four checks UNIPROC, HB OBSERVATION and PROPAGATION stand on firm bases.
- ▶ The semantics of strong fences also does.
- ▶ The model and simulator (*i.e.* herd) are flexible, one easily change a few relations (*e.g.*  $\xrightarrow{ppo}$ , or the semantics of weak fences).

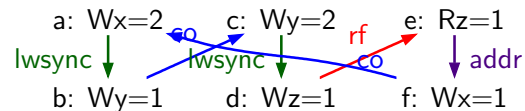
74/74

## Subtle point

### Z6.1

$T_0$	$T_1$	$T_2$
(a) $x \leftarrow 2$	(c) $y \leftarrow 2$	(d) $r0 \leftarrow z$
(d) $y \leftarrow 1$	(e) $z \leftarrow 1$	(f) $x \leftarrow 1$

Observed?  $x=2; y=2; r0=1$



Z61+lwsync+lwsync+addr

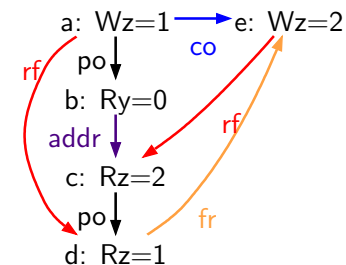
Unobserved and forbidden by model. May be allowed...

75/74

## A test of coherence violation

Our setting also finds bugs...

The following execution:



is observed on all (tested) ARM machines. It features a **CoRR**-style coherence violation (*i.e.*  $\xrightarrow{ppo}$  contradicts  $\xrightarrow{fr}$ ;  $\xrightarrow{fr}$ ).

76/74