# Semantics, languages and algorithms
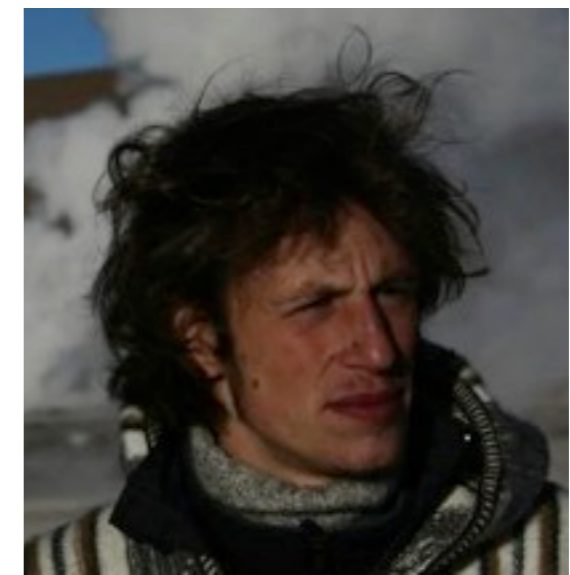## for multicore programming

Albert Cohen    Luc Maranget    Francesco Zappa Nardelli

# Vote: topics for my this lecture
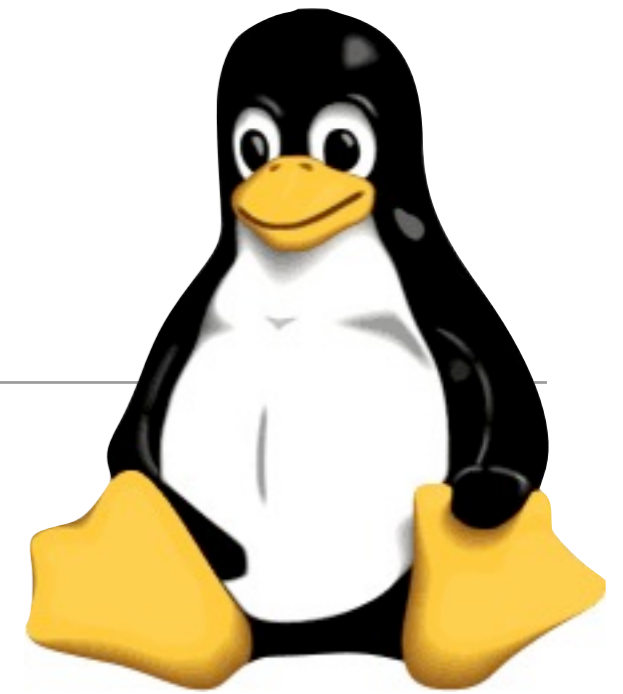


1. Operational and axiomatic formalisation of x86-TSO    (3)

2. The lwarx and stwcx Power instructions                (2)

3. Fence optimisations for x86-TSO                       (6)

4. The Java memory model                                 (3)

5. The C++11 memory model                                (10)

6. Static and dynamic techniques for data-race detection (5)

7. The Linux memory model (?!)                           (13)

8. Compiler correctness statements (compile non-determinism?) (5)

# 1. The Linux memory model  (ahem, kinda)
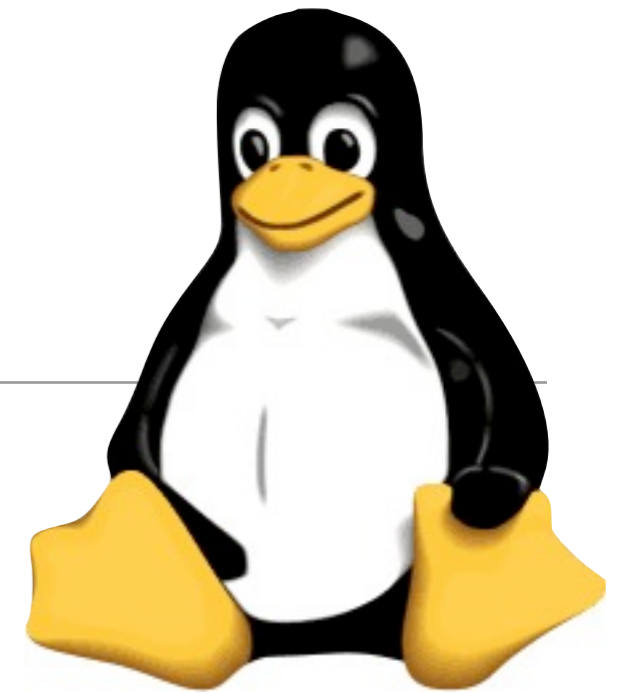
# The Linux memory model

*Facts*:

- abstraction layer over hardware and compilers

- relied upon by kernel developers to write "portable kernel code"

- documented by a text file:

`http://www.kernel.org/doc/Documentation/memory-barriers.txt`

# The Linux memory model

*Facts*:

- abstraction layer over hardware and compilers

- relied upon by kernel developers to write "portable kernel code"

- documented by a text file:

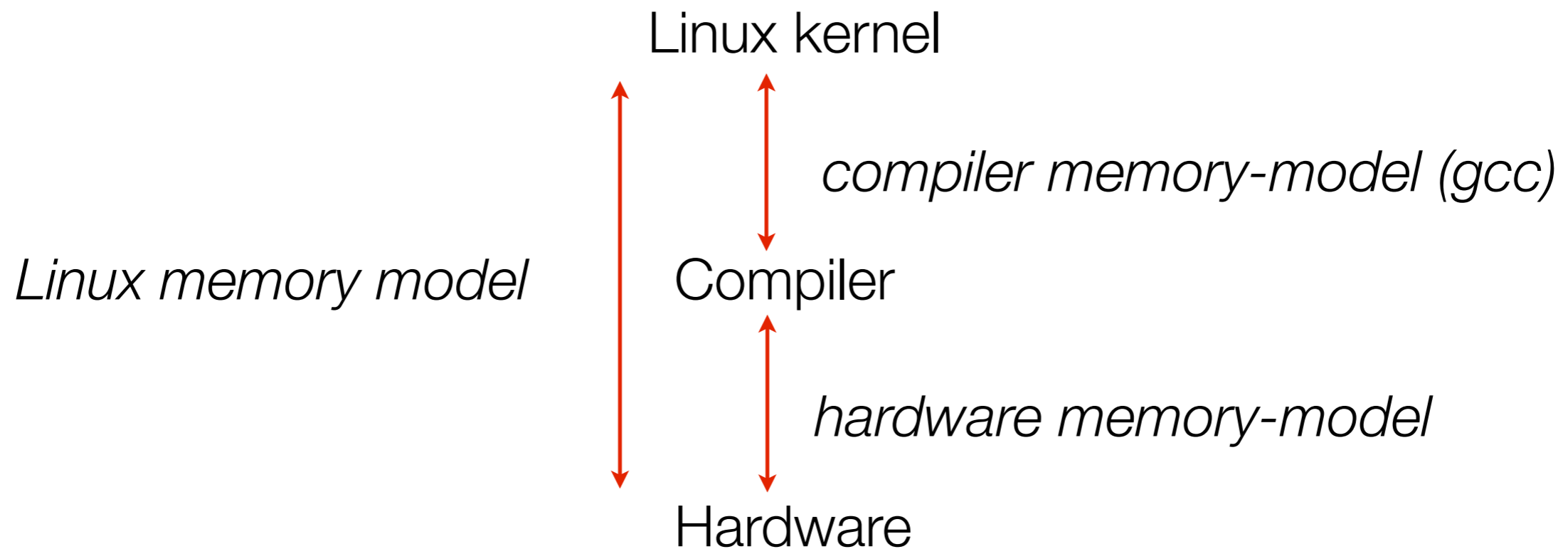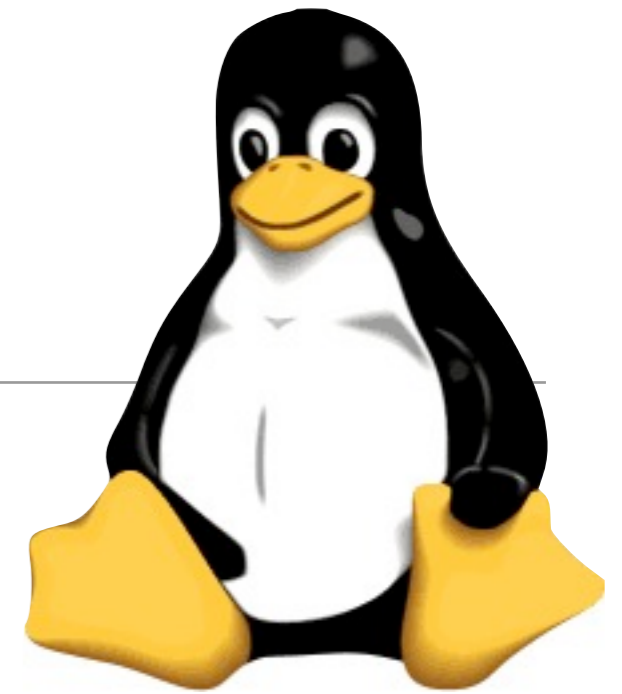`http://www.kernel.org/doc/Documentation/memory-barriers.txt`

*More facts*:

*I attempted to understand the doc, and exchanged a few email with Paul Mc Kenney.  However I don't understand much…*

*In the next hour, let's go over the documentation together and see if we can make sense of it…*

# The Linux memory model

Expected to account for all supported combinations of compiler and hardware memory model...

Linux kernel

*compiler memory-model (gcc)*

*Linux memory model*          Compiler

*hardware memory-model*

Hardware

*alpha*: Weak ordering.  No dependency ordering.  "Time does not go backwards" gives guarantees similar to Power/ARM A-cumulativity.  Possibly B-cumulativity as well.  I am not aware of formalization of this architecture's memory ordering other than Gharachorloo's PhD.

*arm*: You know at least as much as I do about this one.

*avr32*: Uniprocessor-only, kernel build failure for SMP.

*blackfin*: Uniprocessor-only to the best of my knowledge.  There are rumored to be some experimental SMP systems that lack cache coherence, and are thus outside of the Linux kernel's remit.  See for example:  https://docs.blackfin.uclinux.org/doku.php?id=linux-kernel:smp-like   The system.h file flushes cache when a memory barrier is encountered, which is consistent with an attempt to run the Linux kernel on a non-cache-coherent system…

*cris*: Uniprocessor-only to the best of my knowledge.  Though there appears to be recent addition of some SMP support. Its system.h file is consistent with full sequential  consistency.  Or extreme optimism on the part of the cris developers.

*frv*: Uniprocessor-only to the best of my knowledge.

*h8300*: Uniprocessor-only to the best of my knowledge. There is code in system.h that appears to be intended for SMP, but it looks to me like a (harmless) copy-paste error.  Either that or SMP h8300 systems are sequentially consistent.

*ia64*: Total order of all release operations, which include the "mf" (memory fence) instruction.  Memory fences cannot restore sequential consistency.

*m32r*: Uniprocessor-only to the best of my knowledge. However, there does appear to be some recent multiprocessor support.  This is quite strange -- atomic instructions flush cache, but memory barriers are no-ops.  Looks quite experimental.

*m68k*: Uniprocessor-only to the best of my knowledge.

*microblaze*: Uniprocessor-only to the best of my knowledge. At least one SMP attempt: http://microblazesmp.blogspot.com/ Its system.h file looks uniprocessor-only.

*mips*: Multiprocessor.  Old SGI MIPS systems were sequentially consistent.  Newer systems used for network infrastructure are rumored to have weak memory models similar to Power and ARM.  And its system.h file is consistent with a weak memory model.

*mn10300*: Recent SMP support which I know little about. The system.h file looks uniprocessor only, and contains comments on Intel, so copy-pasted from x86.

*parisc*: TSO, similar to x86.

*powerpc*: You know at least as much about this as I do.

*s390*: TSO, but with self-snooping of store buffer prohibited.

*score*: Uniprocessor-only to the best of my knowledge.

*sh*: Recent SMP support which I know little about. Its system.h file is consistent with weak memory ordering.

*sparc*: TSO, similar to x86.  There is documentation about weaker memory models (PSO and RMO), but in practice the hardware is TSO.

*tile*: Recent SMP CPU which I know little about.  Seems to be weakly ordered based on its system.h file.

*um*: Looks like an x86 knockoff judging by the system.h file.

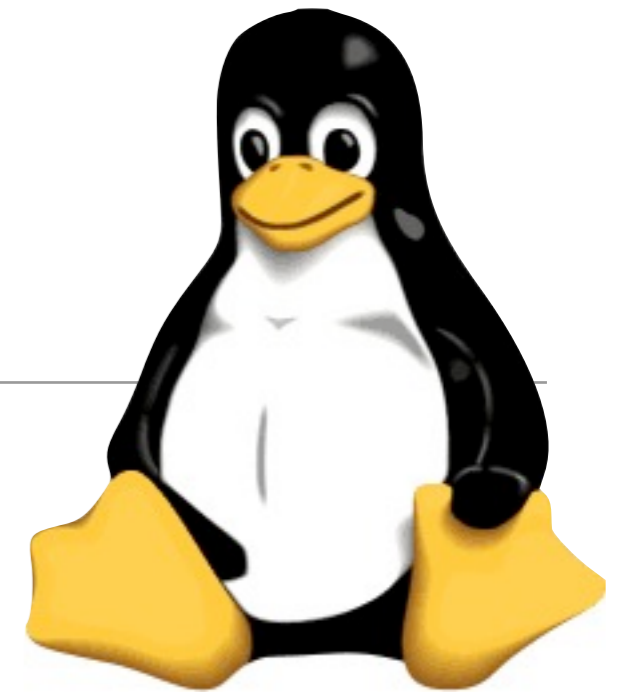*unicore32*: Uniprocessor-only to the best of my knowledge.

*x86*: You know this one at least as well as do I.

*xtensa*: Uniprocessor-only -- kernel build failure otherwise.

# The Linux memory model
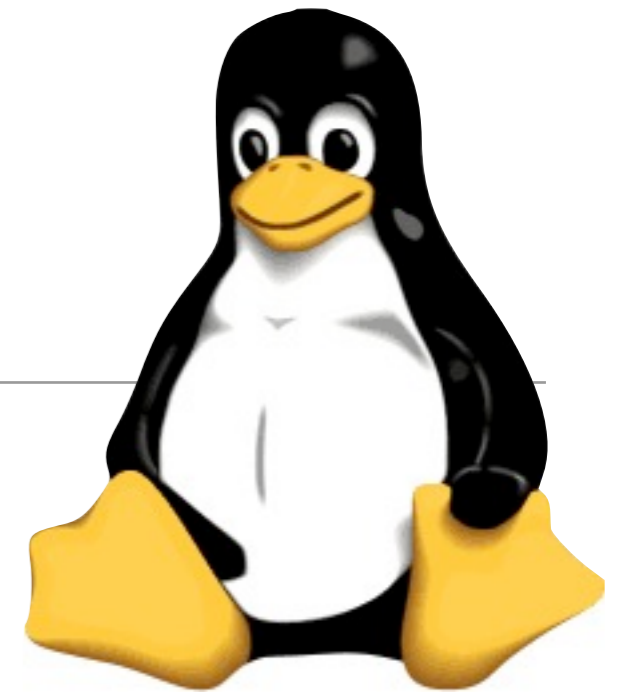
*My intuition:*

*Annoying facts:*

# The Linux memory model

*My intuition:*

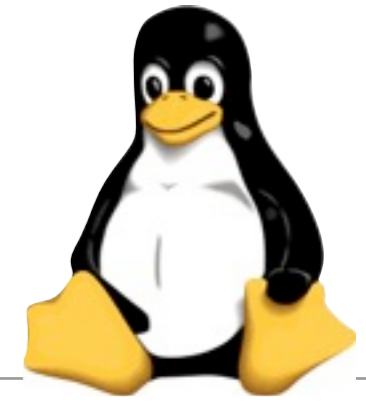*kinda of lowest common denominator between all hardware memory models of architectures Linux can be compiled to, taking into account also some common gcc optimisations, with some weirdnesses.*

*Annoying facts:*

*semantics of "read barriers" really weak, unclear how to formalise it*

*compilation of barriers on Itanium looks broken -- hardware might exhibit behaviours prohibited by the MM.*

...let's read the doc...

# The Linux memory model: macros

on x86:

```
#define mb()   asm volatile("mfence":::"memory")
#define rmb()  asm volatile("lfence":::"memory")
#define wmb()  asm volatile("sfence" ::: "memory")
```

as far as we know, lfence and sfence are noop in x86TSO

on Power:

```
#define mb()   __asm__ __volatile__ ("sync" : : : "memory")
#define rmb()  __asm__ __volatile__ ("sync" : : : "memory")
#define wmb()  __asm__ __volatile__ ("sync" : : : "memory")
#define read_barrier_depends()  do { } while(0)
```

Internship proposal on the fly...

## Sort out what the REAL Linux memory model is

Yes. Of course, if people come up with lots of situations where the more-complex programming model would help significantly, then it might be worth revisiting this.

*Pros*:

Challenging!

Can have a huge impact!

Collaboration with Paul Mc Kenney possible!

# 2. The C++11 memory model

a good example of an axiomatic memory model

# The C++11 memory model

1300 page prose specification defined by the ISO.

The design is a detailed compromise:

  hardware/compiler implementability

  useful abstractions

  broad spectrum of programmers

Welcome to the official home of

**ISO IEC JTC1/SC22/WG21 - The C++ Standards Committee**

2011-09-15: standards | projects | papers | mailings | internals | meetings | contacts

News 2011-09-11: The new C++ standard - C++11 - is published!

# The syntactic divide

```
// for regular programmers:
atomic_int x = 0;
x.store(1);
y = x.load();

// for experts:
x.store(2, memory_order);
y = x.load(memory_order);
atomic_thread_fence(memory_order);
```

where *memory_order* is one of the following:

```
mo_seq_cst   mo_release   mo_acquire
mo_acq_rel   mo_consume   mo_relaxed
```

# How may a program execute?

Two layer semantics:

1) an operational semantics processes programs, identifying memory actions, and constructs candidate executions ($E$opsem);

$$P \longrightarrow E_1, \ldots, E_n$$

2) an axiomatic memory model judges $E$opsem paired with a memory ordering $X$witness

$$E_i \longrightarrow X_{i1}, \ldots, X_{im}$$

3) searches the consistent executions for races and uncostrained reads

is there an $X_{ij}$ with a race?

# Relations

An $E_{\text{opsem}}$ part containing:

   *sb*     sequenced before, program order

   *asw*   additional synchronizes with, inter-thread ordering

An $X_{\text{witness}}$ part containing:

   *rf*     relates a write to any reads that take its value

   *sc*    a total order over mo_seq_cst and mutex actions

   *mo*   modification order, per location total order of writes

From these, compute synchronise-with (*sw*) and happens-before (*hb*).

We ignore *consume* atomics, which enables us to live in a simplified model.

Full details in Batty et al., POPL 11.

# Formally

```
cpp_memory_model_opsem (p : program) =
let pre_executions =
  {(Eopsem,Xwitness).  opsem p Eopsem ∧
    consistent execution (Eopsem, Xwitness)}
in
if ∃X ∈ pre executions.
    (indeterminate reads X = {}) ∨
    (unsequenced races X = {}) ∨
    (data races X = {})
then NONE
else SOME pre_executions
```

# A single-threaded example

1. sequenced before (sb) - given by opsem

```
int main() {
    int x = 2;
    int y = 0;
    y = (x==x);
    return 0;
}
```



$a{:}W_{na}\ x{=}2$

$sb$

$b{:}W_{na}\ y{=}0$

$sb$     $sb$

$c{:}R_{na}\ x{=}2$     $d{:}R_{na}\ x{=}2$

$sb$     $sb$

$e{:}W_{na}\ y{=}1$

# A single-threaded example

1. sequenced before (sb) - given by opsem
2. read-from (rf) - part of the witness

```
int main() {
    int x = 2;
    int y = 0;
    y = (x==x);
    return 0;
}
```

# A single-threaded ex. with undefined behaviour

An unsequenced race.

```
int main() {
    int x = 2;
    int y = 0;
    y = (x==(x=3));
    return 0;
}
```

# A simple concurrent program

```
int y, x = 2;
x = 3;                    | y = (x==3);
```



We will omit `asw` arrows whenever
we are not interested in the initialisation.

# Locks and unlocks

```
int x, r;
mutex m;

m.lock();          m.lock();
x = ...            r = x;
m.unlock();
```

1. the operational semantics defines the sb arrows

$c{:}L$ mutex

sb $\downarrow$

$d{:}W_{na}\ x{=}1$

sb $\downarrow$

$f{:}U$ mutex

$h{:}L$ mutex

sb $\downarrow$

$i{:}R_{na}\ x{=}1$

# Locks and unlocks

```
int x, r;
mutex m;

m.lock();          m.lock();
x = ...            r = x;
m.unlock();
```

1. the operational semantics defines the sb arrows

2. guess an sc order on Unlock/Lock actions (part of the witness)

$c{:}L$ mutex      $h{:}L$ mutex

sb      sb

$d{:}W_{na}\ x{=}1$      $i{:}R_{na}\ x{=}1$

sb    sc

$f{:}U$ mutex

# Locks and unlocks

```
int x, r;
mutex m;

m.lock();          m.lock();
x = ...            r = x;
m.unlock();
```

1. the operational semantics defines the sb arrows

2. guess an sc order on Unlock/Lock actions (part of the witness)

3. the sc order is included in the syncronised-with relation

c:L mutex                    h:L mutex

sb ↓                         sb ↓

d:$W_{na}$ x=1               i:$R_{na}$ x=1

sb ↓              sw

f:U mutex

# Locks and unlocks

$$\xrightarrow{\textit{simple-happens-before}} = \left( \xrightarrow{\textit{sequenced-before}} \cup \xrightarrow{\textit{synchronizes-with}} \right)^{+}$$
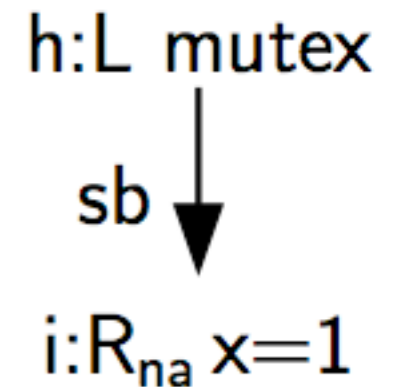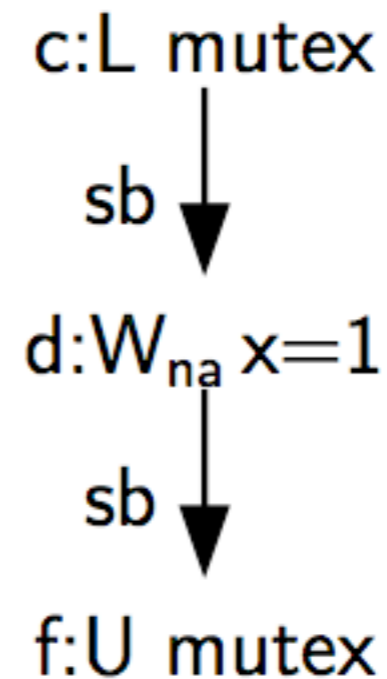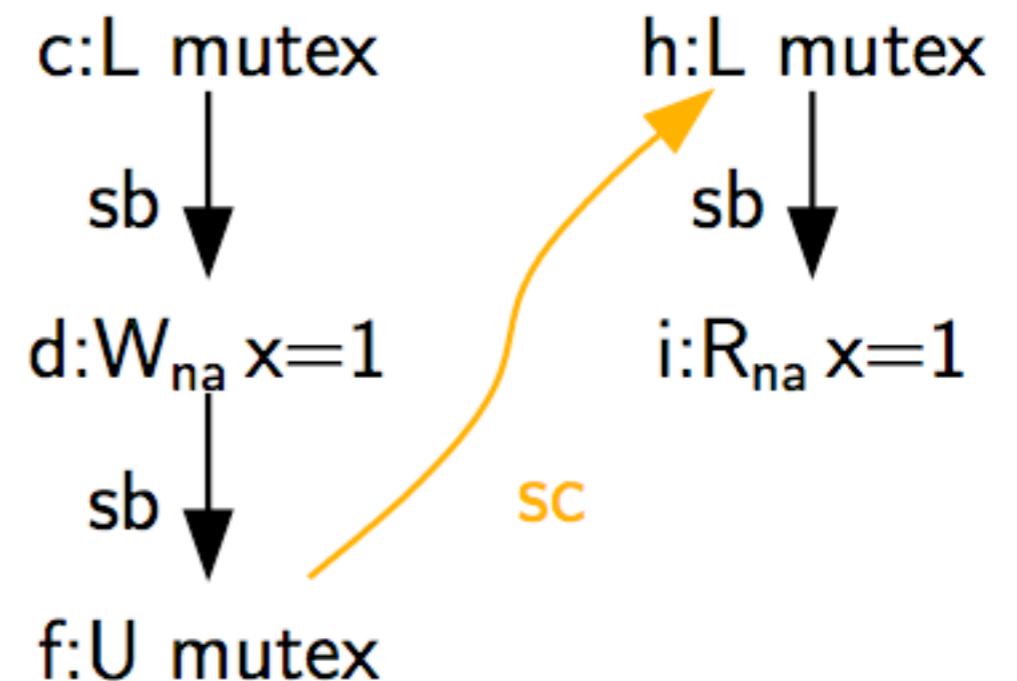
```
int x, r;
mutex m;

m.lock();          m.lock();
x = ...            r = x;
m.unlock();
```

1. the operational semantics defines the sb arrows

2. guess an sc order on Unlock/Lock actions (part of the witness)

3. the sc order is included in the syncronised-with relation

4. which in turn defines the happens-before relation...

# Happens before

The *happens before* relation is key to the model:

1. non-atomic loads read the most recent write in happens before.
   (This is unique in DRF programs)

2. the story is more complex for atomics, as we shall see.

3. data races are defined as an absence of happens before
   between conflicting actions.

$$\xrightarrow{\textit{simple-happens-before}} \;=\; (\xrightarrow{\textit{sequenced-before}} \cup \xrightarrow{\textit{synchronizes-with}})^+$$

# A data race

```
int y, x = 2;
x = 3;              | y = (x==3);
```



a:$W_{na}$ x=2

asw        asw,rf

b:$W_{na}$ x=3    c:$R_{na}$ x=2

sb

d:$W_{na}$ y=0

# A data race

```
int y, x = 2;
x = 3;                      | y = (x==3);
```

$a:W_{na}\ x=2$

asw    asw,rf

$b:W_{na}\ x=3$ —— $c:R_{na}\ x=2$
                dr

sb

$d:W_{na}\ y=0$

Here we have two conflicting accesses not related by happens-before.

# Data race definition

```
let data_races actions hb =
    { (a, b) | ∀ a∈actions b∈actions |
        ¬ (a = b) ∧
        same_location a b ∧
        (is_write a ∨ is_write b) ∧
        ¬ (same_thread a b) ∧
        ¬ (is_atomic_action a ∧ is_atomic_action b) ∧
        ¬ ((a, b) ∈ hb ∨ (b, a) ∈ hb) }
```

Programs with a data race have undefined behaviour (DRF model).

# Simple concurrency: Dekker's example and SC

```
atomic_int x = 0;
atomic_int y = 0;

x.store(1, seq_cst);  |  y.store(1, seq_cst);
y.load(seq_cst);      |  x.load(seq_cst);
```

$c:W_{sc}\,y{=}1$        $e:W_{sc}\,x{=}1$

FORBIDDEN

sb           sb

$d:R_{sc}\,x{=}0$        $f:R_{sc}\,y{=}0$

Why is this behaviour forbidden?

# Simple concurrency, Dekker's example and SC

```
atomic_int x = 0;
atomic_int y = 0;

x.store(1, seq_cst);  | y.store(1, seq_cst);
y.load(seq_cst);      | x.load(seq_cst);
```

$c{:}W_{sc}\,y{=}1$        $e{:}W_{sc}\,x{=}1$

sc       sc

sc

$d{:}R_{sc}\,x{=}0$        $f{:}R_{sc}\,y{=}1$

The `sc` relation must define a total order over unlocks/locks and `seq_cst` accesses… `sc` is included in `hb`, an `rf` must respect `hb`.

# Expert concurrency: the release-acquire idiom

```
// sender
x = ...
y.store(1, release);


// receiver
while (0 == y.load(acquire));
r = x;
```

Here we have an `rf` arrow beetwen a pair of release/acquire accesses.

$a\text{:}W_{na}\ x{=}1$

sb

$b\text{:}W_{rel}\ y{=}1$

rf

$c\text{:}R_{acq}\ y{=}1$

sb

$d\text{:}R_{na}\ x{=}1$

# Expert concurrency: the release-acquire idiom

```
// sender
x = ...
y.store(1, release);

// receiver
while (0 == y.load(acquire));
r = x;
```

W x=1

sb

$W_{REL}$ y=1

sw

$R_{ACQ}$ y=1

sb

R x=1

Here we have an `rf` arrow beetwen a pair of release/acquire accesses.

The `rf` arrow beetwen release/acquire accesses induces an `sw` arrow between those accesses.

# Expert concurrency: the release-acquire idiom

```
// sender
x = ...
y.store(1, release);


// receiver
while (0 == y.load(acquire));
r = x;
```

Here we have an `rf` arrow beetwen a pair of release/acquire accesses.

The `rf` arrow beetwen release/acquire accesses induces an `sw` arrow between those accesses.

And in turn defines an `hb` constraint.

$W\,x{=}1$

$\xrightarrow{\text{sb}}$

hb

$W_{REL}\,y{=}1$

sw

$R_{ACQ}\,y{=}1$

sb

$R\,x{=}1$

$$\xrightarrow{simple\text{-}happens\text{-}before} = (\xrightarrow{sequenced\text{-}before} \cup \xrightarrow{synchronizes\text{-}with})^+$$

# Relaxed writes

```
x.load(relaxed);           y.load(relaxed);
y.store(1, relaxed);       x.store(1, relaxed);
```

c:Rrlx x=1                 e:Rrlx y=1

sb ↓                       sb ↓

rf rf

d:Wrlx y=1                 f:Wrlx x=1

No data-races, no synchronisation cost, but weakly ordered.

# Relaxed writes, ctd.

```
atomic_int x = 0;
atomic_int y = 0;
x.store(1, relaxed); | y.store(2, relaxed); | x.load(relaxed); | y.load(relaxed);
                     |                      | y.load(relaxed); | x.load(relaxed);
```

c:Wrlx x=1    d:Wrlx y=1 — e:Rrlx x=1 → g:Rrlx y=1

rf    sb rf    sb

f:Rrlx y=0    h:Rrlx x=0

Again, no data-races, no synchronisation cost, but weakly ordered (IRIW).

# Expert concurrency: fences avoid excess sync.

```
// sender                          // receiver
x = ...                           while (0 == y.load(acquire));
y.store(1, release);              r = x;
```

```
// sender                          // receiver
x = ...                           while (0 == y.load(relaxed));
y.store(1, release);              fence(acquire);
                                  r = x;
```

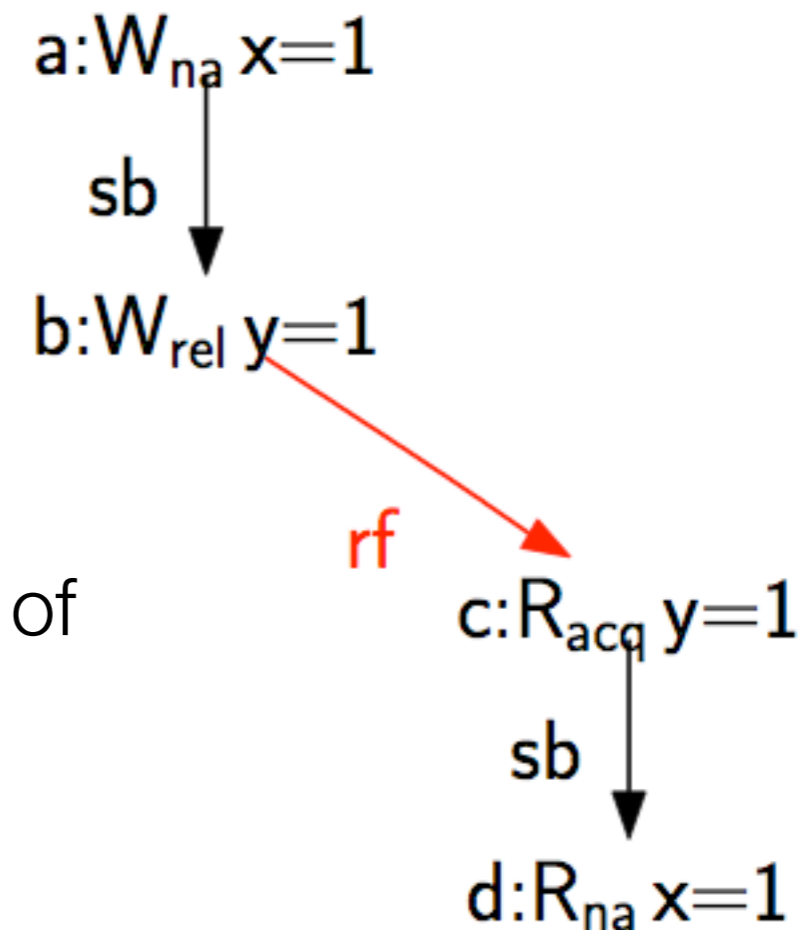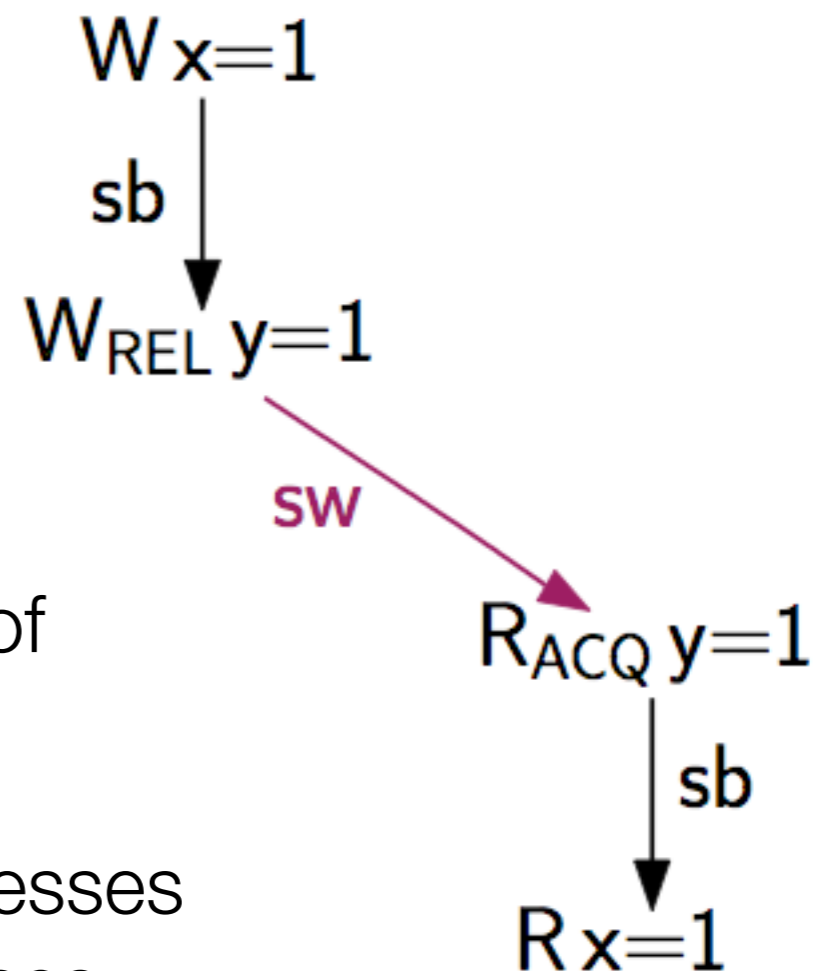# Expert concurrency: fences avoid excess sync.

```
// sender
x = ...
y.store(1, release);
```

```
// receiver
while (0 == y.load(relaxed));
fence(acquire);
r = x;
```

Here we have an **rf** arrow beetwen a release write and a relaxed write.

$$c:W_{na}\ x=1 \qquad\qquad e:R_{rlx}\ y=1$$

$$\text{sb}\downarrow \qquad\qquad\qquad \text{sb}\downarrow$$

$$d:W_{rel}\ y=1 \qquad\qquad f:F_{acq}$$

$$\text{rf}$$

$$\text{sb}\downarrow$$

$$g:R_{na}\ x=1$$

# Expert concurrency: fences avoid excess sync.

```
// sender
x = ...
y.store(1, release);
```

```
// receiver
while (0 == y.load(relaxed));
fence(acquire);
r = x;
```

Here we have an `rf` arrow beetwen a release write and a relaxed write.

The acquire fence follows the `sb`/`rf` relations looking for the corresponding release write, adding a `sw` arrow.

$c: W_{na}\ x=1$

$e: R_{rlx}\ y=1$

rf

sb

sb

$d: W_{rel}\ y=1$

$f: F_{acq}$

sw

sb

$g: R_{na}\ x=1$

# Expert concurrency: fences avoid excess sync.

```
// sender
x = ...
y.store(1, release);
```

```
// receiver
while (0 == y.load(relaxed));
fence(acquire);
r = x;
```

Here we have an `rf` arrow beetwen a release write and a relaxed write.

The acquire fence follows the `sb`/`rf` relations looking for the corresponding release write, adding a `sw` arrow.

Happens-before follows as usual...

$c: W_{na}\ x{=}1$

$e: R_{rlx}\ y{=}1$

rf

sb

hb

sb

$d: W_{rel}\ y{=}1$

$f: F_{acq}$

sw

sb

$g: R_{na}\ x{=}1$

# Modification order

```
   atomic_int x = 0;
x.store(1, relaxed);  ‖  x.load(relaxed);
x.store(2, relaxed);  ‖  x.load(relaxed);
```

$$W_{RLX}\, x{=}1 \xrightarrow{\ rf\ } R_{RLX}\, x{=}1$$

$$mo \downarrow \qquad\qquad sb \downarrow$$

$$W_{RLX}\, x{=}2 \xrightarrow{\ rf\ } R_{RLX}\, x{=}2$$

Modification order is a total order over atomic writes of any memory order.

# Coherence and atomic reads

All forbidden:



Idea: atomics cannot read from later writes in happens-before.

# Coherence and atomic reads

All forb

a:

b:

b:W x=2

CoWW

d:W x=2

CoRW

A pair $E_{\text{opsem}}$, $X_{\text{witness}}$ (a pre-execution)

defines a *consistent execution* when it satisfies

the constraints we have sketched

on `hb`/`rf`/`mo` and is race-free.

Idea: atomics cannot read from later writes in happens-before.

# The full model

# Is C++11 hopelessly complicated?

Programmers cannot be given this model.

However, with a formal definition, we can do proofs!

- Can we compile to x86?

| Operation | x86 Implementation |
|---|---|
| load(non-seq_cst) | mov |
| load(seq_cst) | lock xadd(0) |
| store(non-seq_cst) | mov |
| store(seq_cst) | lock xchg |
| fence(non-seq_cst) | no-op |

- Can we compile to Power?

| C++0x Operation | POWER Implementation |
|---|---|
| Non-atomic Load | ld |
| Load Relaxed | ld |
| Load Consume | ld (and preserve dependency) |
| Load Acquire | ld; cmp; bc; isync |
| Load Seq Cst | sync; ld; cmp; bc; isync |
| Non-atomic Store | st |
| Store Relaxed | st |
| Store Release | lwsync; st |
| Store Seq Cst | sync; st |

# Is C++11 hopelessly complicated?

Simplifications:

Full model: *visible sequences of side effects* are unneded (HOL4)

Derivative models:

- without consume, happens-before is transitive

- DRF programs using only `seq_cst` atomics are SC (false)

```
atomic_int x = 0;
atomic_int y = 0;
if (1 == x.load(seq_cst)) | if (1 == y.load(seq_cst))
    atomic_init(&y, 1);    |     atomic_init(&x, 1);
```

`atomic_init` is a non-atomic write, and in C++11 they race.

# The current state of the standard

**Fixed**:

- in some cases, happens-before was cyclic

- coherence

- `seq_cst` atomics were more broken

**Not fixed**:

- self satisfying conditional

```
r1 = x.load(mo_relaxed);        r2 = y.load(mo_relaxed);
if (r1 == 42)                   if (r2 == 42)
   y.store(r1, mo_relaxed);        x.store(42, mo_relaxed);
```

c:Rrlx x=1     e:Rrlx y=1
sb             sb
rf rf
d:Wrlx y=1     f:Wrlx x=1

- `seq_cst` atomics are still not SC

# 3. Sketch of an operational formalisation of x86-TSO

...starting with a formalisation of SC

# Separate language and memory semantics

```
1   class ArrayWrapper
2   {
3       public:
4           ArrayWrapper (int n)
5               : _p_vals( new int[ n ] )
6               , _size( n )
7           {}
8           // copy constructor
9           ArrayWrapper (const ArrayWrapper& other)
10              : _p_vals( new int[ other._size  ] )
11              , _size( other._size )
12          {
13              for ( int i = 0; i < _size; ++i )
14              {
15                  _p_vals[ i ] = other._p_vals[ i ];
16              }
17          }
18          ~ArrayWrapper ()
19          {
20              delete [] _p_vals;
21          }
22      private:
23      int *_p_vals;
24      int _size;
25  };
```



*program*
semantics defined via an LTS

*memory*
semantics defined via an LTS

$W_t[a]v$ : a write of value $v$ to address $a$ by thread $t$

*Labels for interaction:*   $R_t[a]v$  : a read of $v$ from $a$ by $t$ by thread $t$

+ other events for barriers and locked instructions

# Separate language and memory semantics

```
 1   class Arr
 2   {
 3       publi
 4          A
 5
 6
 7          {
 8          /
 9          A
10
11
12          {
13
14
15
16
17          }
18          ~
19          {
20
21          }
22      priva
23      int *
24      int _
25   };
```

Separate language and state semantics

proved to be a very good choice

in many (unrelated) projects I worked on!

semantics defined via an LTS                    semantics defined via an LTS

$W_t[a]v$ : a write of value $v$ to address $a$ by thread $t$

*Labels for interaction:*    $R_t[a]v$  : a read of $v$ from $a$ by $t$ by thread $t$

+ other events for barriers and locked instructions

# A tiny language

| | | | |
|---|---|---|---|
| $location,\ x,\ m$ | | address (or pointer value) | |
| $integer,\ n$ | | integer | |
| $thread\_id,\ t$ | | thread id | |
| $k,\ i,\ j$ | | | |

| $expression,\ e$ | $::=$ | | expression |
|---|---|---|---|
| | $\mid$ | $n$ | integer literal |
| | $\mid$ | $*x$ | read from pointer |
| | $\mid$ | $*x = e$ | write to pointer |
| | $\mid$ | $e;\ e'$ | sequential composition |
| | $\mid$ | $e + e'$ | plus |

| $process,\ p$ | $::=$ | | process |
|---|---|---|---|
| | $\mid$ | $t{:}e$ | thread |
| | $\mid$ | $p\vert p'$ | parallel composition |

# What can a thread do in isolation?

$$\boxed{e \xrightarrow{l} e'} \qquad e \text{ does } l \text{ to become } e'$$

$$\frac{}{*x \xrightarrow{R\ x=n} n} \quad \text{READ}$$

$$\frac{}{*x = n \xrightarrow{W\ x=n} n} \quad \text{WRITE}$$

$$\frac{e \xrightarrow{l} e'}{*x = e \xrightarrow{l} *x = e'} \quad \text{WRITE\_CONTEXT}$$

$$\frac{}{n; e \xrightarrow{\tau} e} \quad \text{SEQ}$$

$$\frac{e_1 \xrightarrow{l} e_1'}{e_1; e_2 \xrightarrow{l} e_1'; e_2} \quad \text{SEQ\_CONTEXT}$$

$$\frac{e_1 \xrightarrow{l} e_1'}{e_1 + e_2 \xrightarrow{l} e_1' + e_2} \quad \text{PLUS\_CONTEXT\_1}$$

$$\frac{e_2 \xrightarrow{l} e_2'}{n_1 + e_2 \xrightarrow{l} n_1 + e_2'} \quad \text{PLUS\_CONTEXT\_2}$$

$$\frac{n = n_1 + n_2}{n_1 + n_2 \xrightarrow{\tau} n} \quad \text{PLUS}$$

Observe that we can read an arbitrary value from the memory.

# Lifting to processes

$$\boxed{p \xrightarrow{l_t} p'} \quad p \text{ does } l_t \text{ to become } p'$$

$$\frac{e \xrightarrow{l} e'}{t{:}e \xrightarrow{l_t} t{:}e'} \quad \text{THREAD}$$

Actions are labelled by the thread that performed the action.

$$\frac{p_1 \xrightarrow{l_t} p_1'}{p_1 | p_2 \xrightarrow{l_t} p_1' | p_2} \quad \text{PAR\_CONTEXT\_LEFT}$$

*Free interleaving.*

$$\frac{p_2 \xrightarrow{l_t} p_2'}{p_1 | p_2 \xrightarrow{l_t} p_1 | p_2'} \quad \text{PAR\_CONTEXT\_RIGHT}$$

# A sequentially consistent memory

Take **M** to be a function from addresses to integers.

$$\boxed{M \xrightarrow{l} M'} \quad M \text{ does } l \text{ to become } M'$$

$$\frac{M(x) = n}{M \xrightarrow{\text{R } x=n} M} \quad \text{M{\scriptsize READ}}$$

$$\frac{}{M \xrightarrow{\text{W } x=n} M \oplus (x \mapsto n)} \quad \text{M{\scriptsize WRITE}}$$

# SC semantics: whole system transitions

$\boxed{s \xrightarrow{l_t} s'}$    $s$ **does** $l_t$ **to become** $s'$

$$\frac{p \xrightarrow{R_t\ x=n} p' \qquad M \xrightarrow{R\ x=n} M'}{\langle p,\ M \rangle \xrightarrow{R_t\ x=n} \langle p',\ M' \rangle} \ \textsc{Sread}$$

Synchronising between the processes and the memory.

$$\frac{p \xrightarrow{W_t\ x=n} p' \qquad M \xrightarrow{W\ x=n} M'}{\langle p,\ M \rangle \xrightarrow{W_t\ x=n} \langle p',\ M' \rangle} \ \textsc{Swrite}$$

$$\frac{p \xrightarrow{\tau_t} p'}{\langle p,\ M \rangle \xrightarrow{\tau_t} \langle p',\ M \rangle} \ \textsc{Stau}$$

# SC semantics, example

All threads read and write the shared memory. Threads execute asynchronously, the semantics allows any interleaving of the thread transitions.

$$\langle t_1{:}{*}x = 1 \,|\, t_2{:}{*}x = 2, \{x \mapsto 0\}\rangle$$

$W_{t_1}\ x{=}1$

$W_{t_2}\ x{=}2$

$$\langle t_1{:}1 \,|\, t_2{:}{*}x = 2, \{x \mapsto 1\}\rangle \qquad\qquad \langle t_1{:}{*}x = 1 \,|\, t_2{:}2, \{x \mapsto 2\}\rangle$$

$W_{t_2}\ x{=}2$

$W_{t_1}\ x{=}1$

$$\langle t_1{:}1 \,|\, t_2{:}2, \{x \mapsto 2\}\rangle \qquad\qquad \langle t_1{:}1 \,|\, t_2{:}2, \{x \mapsto 1\}\rangle$$

Each interleaving has a linear order of reads and writes to memory.

...now we just have to define a TSO memory...

# x86-TSO abstract machine



Thread ••• Thread

Wt[a]v    Rt[a]v

Write Buffer

Text •••

Wri

Lock

Events visible by each thread (aka. interface between each thread and the memory system):

$W_t[a]v$ : a write of value $v$ to address $a$ by thread $t$
$R_t[a]v$  : a read of $v$ from $a$ by $t$ by thread $t$
+ other events for barriers and locked instructions

# x86-tso: a formalisation using an LTS

The machine state `s` can be represented by a tuple `(M,B,L)`:

```
M : address -> value option
B : tid -> (address * value) list
L : tid option
```

where:

  `M` is the shared memory, mapping addresses to values

  `B` gives the store buffer for each thread

  `L` is the global machine lock indicating when a thread has exclusive access to memory (omitted in these slides)

# x86-tso abstract machine: selected transition rules

t is *not blocked* in machine state s = (M,B,L) if [… or] the lock is not held.

In buffer B(t) there are *no pending writes* for address x if there are no (x,v) elements in B(t).

**RM: Read from memory**

$$\frac{\text{not\_blocked}(s, t) \\ s.M(x) = v \\ \text{no\_pending}(s.B(t), x)}{s \quad \xrightarrow{\quad R_t \, x = v \quad} \quad s}$$

Thread $t$ can read $v$ from memory at address $x$ if $t$ is not blocked, the memory does contain $v$ at $x$, and there are no writes to $x$ in $t$'s store buffer.

# x86-tso abstract machine: selected transition rules

**RB: Read from write buffer**

$$\frac{\begin{array}{l} \mathrm{not\_blocked}(s, t) \\ \exists b_1 \, b_2. \; s.B(t) = b_1 \mathbin{+\!\!+} [(x, v)] \mathbin{+\!\!+} b_2 \\ \mathrm{no\_pending}(b_1, x) \end{array}}{s \xrightarrow{\;\;R_t \, x = v\;\;} s}$$

Thread $t$ can read $v$ from its store buffer for address $x$ if $t$ is not blocked and has $v$ as the newest write to $x$ in its buffer;

# x86-tso abstract machine: selected transition rules

**WB: Write to write buffer**

$$s \xrightarrow{\mathrm{W}_t\, x=v} s \oplus \langle\!| B := s.B \oplus (t \mapsto ([(x, v)] + \!\!+ s.B(t)))|\!\rangle$$

Thread $t$ can write $v$ to its store buffer for address $x$ at any time;

**WM: Write from write buffer to memory**

$$\frac{\mathrm{not\_blocked}(s, t) \qquad s.B(t) = b + \!\!+ [(x, v)]}{s \xrightarrow{\tau_t\, x=v}}$$

$$s \oplus \langle\!| M := s.M \oplus (x \mapsto v)|\!\rangle \oplus \langle\!| B := s.B \oplus (t \mapsto b)|\!\rangle$$

If $t$ is not blocked, it can silently dequeue the oldest write from its store buffer and place the value in memory at the given address, without coordinating with any hardware thread

# 4. Veryfing fence elimination optimisations

aka reasoning on the x86TSO operational memory model
and compiler correctness

# CompCertTSO



**ClightTSO**

simplify

**C#minor**

local vars

**Cstacked**

simplify

**Cminor**

instruction selection

**CminorSel**

CFG generation

**RTL**

const prop.

**RTL**

CSE

**RTL**

register allocation

**LTL**

branch tunnelling

**LTL**

linearize

**LTLin**

reload/spill

**Linear**

act.records

**Machabstr**

**Machconc** → **x86**

[POPL 2011]

# CompCertTSO + fence optimisations



**ClightTSO**

simplify

**C#minor**

local vars

**Cstacked**

simplify

**Cminor**

instruction selection

**CminorSel**

CFG generation

**RTL**

const prop.

**RTL**

CSE

**RTL**

FE1

**RTL**

PRE

**RTL**

FE2

**RTL**

register allocation

**LTL**

branch tunnelling

**LTL**

linearize

**LTLin**

reload/spill

**Linear**

act.records

**Machabstr**

**Machconc** → **x86**

[SAS 2011]

# Compilers are *ideal* for verification

```
┌─────────────────────────────┐        Compiler         ┌─────────────────────────────┐
│   source program (e.g., C)  │ ──────────────────────▶ │  target program (e.g., x86) │
└─────────────────────────────┘                         └─────────────────────────────┘
```

Compilers are:

— Basic computing infrastructure

— Generally reliable, but nevertheless contain many bugs
    e.g., Yang et al. [PLDI 2011] found 79 `gcc` & 202 `llvm` bugs

— "Specifiable": compiler correctness = preservation of behaviours

— Interesting: naturally higher-order, involve clever algorithms

— Big, but modular

# Language semantics

The semantics of all the CompCertTSO languages is defined by:

– a type of programs, $prg$

– a type of states, $states$

– a set of initial states for each program, $\text{init} \in prg \to \mathbb{P}(states)$

– a transition relation, $\to \in \mathbb{P}(states \times \boxed{event} \times states)$

$$\texttt{call}, \texttt{return}, \texttt{fail}, \texttt{oom}, \tau$$

The visible behaviour of a program is defined by the external function calls (`call`) and returns (`return`), errors (`fail`), and running out of memory (`oom`).

# Traces

– *Finite sequences* of call & return events ending with:

    **end**:    successful termination,
    **inftau**:  infinite execution that stops performing visible events
    **oom**:   execution runs out of memory

– *Infinite sequences* of call & return events;

$$
\begin{aligned}
\mathbf{traces}(p) \stackrel{\mathrm{def}}{=}\ & \{\ell \cdot \mathsf{end} \mid \exists s \in \mathsf{init}(p).\ \exists s'.\ s \stackrel{\ell}{\Rightarrow} s' \wedge s' \not\rightarrow\} \\
\cup\ & \{\ell \cdot tr \mid \exists s \in \mathsf{init}(p).\ \exists s'.\ s \xRightarrow{\ell \cdot \mathtt{fail}} s'\} \\
\cup\ & \{\ell \cdot \mathsf{inftau} \mid \exists s \in \mathsf{init}(p).\ \exists s'.\ s \stackrel{\ell}{\Rightarrow} s' \wedge \mathsf{inftau}(s')\} \\
\cup\ & \{\ell \cdot \mathsf{oom} \mid \exists s \in \mathsf{init}(p).\ \exists s'.\ s \stackrel{\ell}{\Rightarrow} s'\} \\
\cup\ & \{tr \mid \exists s \in \mathsf{init}(p).\ s \text{ can do the infinite trace } tr\}
\end{aligned}
$$

NB: Erroneous computations become undefined after the first error.

# Compiler correctness

source program (e.g., C)    — Compiler →    target program (e.g., x86)

$$\text{traces(source\_program)} \supseteq \text{traces(target\_program)}$$

print "a" || print "b"    → (green)    print "ab"

print "ab"    ✗ →    print "a" || print "b"

fail    → (green)    print "ab"

print "ab"    ✗ →    fail

# Store buffering

MOV [x] ← 1          MOV [y] ← 1

MOV EAX ← [y]        MOV EBX ← [x]

| | |
|---|---|
| Thread | ... Thread |

EAX : 32                                    EBX : 47

Write Buffer          Write Buffer

Shared Memory

x : 0     y : 0

# Store buffering

# Store buffering

MOV [x] ← 1

MOV EAX ← [y]

MOV [y] ← 1

MOV EBX ← [x]

EAX : 32

Thread

...

Thread

EBX : 47

Write
Buffer

x:1

Write
Buffer

y:1

Shared Memory

x : 0    y : 0

# Store buffering

```
MOV [x] ← 1              MOV [y] ← 1
MOV EAX ← [y]            MOV EBX ← [x]
```

EAX : 0

Thread          ...          Thread          EBX : 47

Write
Buffer

x:1

Write
Buffer

y:1

Shared Memory

x : 0      y : 0

# Store buffering

MOV [x] ← 1          MOV [y] ← 1

MOV EAX ← [y]        MOV EBX ← [x]

EAX : 0

| Thread | ... | Thread |

EBX : 0

| Write Buffer |   | Write Buffer |

x:1                 y:1

Shared Memory

x : 0      y : 0

# Store buffering

MOV [x] ← 1

MOV EAX ← [y]

MOV [y] ← 1

MOV EBX ← [x]

EAX : 0

Thread

...

Thread

EBX : 0

Write Buffer

Write Buffer

y:1

Shared Memory

x : 1    y : 0

# Store buffering

MOV [x] ← 1          MOV [y] ← 1

MOV EAX ← [y]        MOV EBX ← [x]

EAX : 0    Thread    ...    Thread    EBX : 0

Write Buffer                Write Buffer

Shared Memory

x : 1    y : 1

# Store buffering + fences

MOV [x] ← 1
MFENCE
MOV EAX ← [y]

MOV [y] ← 1
MFENCE
MOV EBX ← [x]

Thread        ...        Thread

EAX : 32                          EBX : 47

Write
Buffer

Write
Buffer

Shared Memory

x : 0    y : 0

# Store buffering + fences

```
MOV [x] ← 1                    MOV [y] ← 1
MFENCE                         MFENCE
MOV EAX ← [y]                  MOV EBX ← [x]
```

EAX : 32          Thread          ...          Thread          EBX : 47

Write
Buffer

x:1

Write
Buffer

Shared Memory

x : 0    y : 0

# Store buffering + fences

# Store buffering + fences

```
MOV [x] ← 1              MOV [y] ← 1
MFENCE                   MFENCE
MOV EAX ← [y]            MOV EBX ← [x]
```

Thread          ...         Thread

EAX : 32                                EBX : 47

Write Buffer                Write Buffer

                                        y:1

Shared Memory

x : 1    y : 0

# Who inserts fences?

1. The *programmer*, explicitly.  Example: Fraser's lockfree-lib:

```
/*
 * II. Memory barriers.
 *  MB():  All preceding memory accesses must commit before any later accesses.
 *
 *  If the compiler does not observe these barriers (but any sane compiler
 *  will!), then VOLATILE should be defined as 'volatile'.
 */
#define MB()  __asm__ __volatile__ ("lock; addl $0,0(%%esp)" : : : "memory")
```

2. The *compiler*, to implement a high-level memory model,
   e.g. `SEQ_CST` C++0x low-level atomics on x86:

   Load SEQ_CST:   `MFENCE; MOV`

   Store SEQ_CST:  `MOV; MFENCE`

# Fence instructions

1. *Fences are necessary*

   to implement locks & not fully-commutative linearizable objects (e.g., stacks, queues, sets, maps).

   [Attiya et al., POPL 2011]

2. *Fences can be expensive*

# Redundant fences (1)

If we have two consecutive fence instructions, we can remove the *latter*:

```
MFENCE                          MFENCE
MFENCE                          NOP
```

The *buffer is already empty* when the second fence is executed.

*Generalisation:*

```
MFENCE                          MFENCE
NON-WRITE INSTR                 NON-WRITE INSTR
…                               …
NON-WRITE INSTR                 NON-WRITE INSTR
MFENCE                          NOP
```

# FE1

A *forward* data-flow problem over the boolean domain $\{\bot, \top\}$

Associate to each program point:

$\bot$ : along all execution paths there is an atomic instruction *before* the current program point, with no intervening writes;

$\top$ : otherwise.

$$
\begin{aligned}
T_1(\mathbf{nop}, \mathcal{E}) &= \mathcal{E} \\
T_1(\mathbf{op}(op, \vec{r}, r), \mathcal{E}) &= \mathcal{E} \\
T_1(\mathbf{load}(\kappa, addr, \vec{r}, r), \mathcal{E}) &= \mathcal{E} \\
T_1(\mathbf{store}(\kappa, addr, \vec{r}, src), \mathcal{E}) &= \top \\
T_1(\mathbf{call}(sig, ros, args, res), \mathcal{E}) &= \top \\
T_1(\mathbf{cond}(cond, args), \mathcal{E}) &= \mathcal{E} \\
T_1(\mathbf{return}(optarg), \mathcal{E}) &= \top \\
T_1(\mathbf{threadcreate}(optarg), \mathcal{E}) &= \top \\
T_1(\mathbf{atomic}(aop, \vec{r}, r), \mathcal{E}) &= \bot \\
T_1(\mathbf{fence}, \mathcal{E}) &= \bot
\end{aligned}
$$

$$
\mathcal{FE}_1(n) = \begin{cases} \top & \text{if predecessors}(n) = \emptyset \\ \bigsqcup_{p \in \text{predecessors}(n)} T_1(instr(p), \mathcal{FE}_1(p)) & \text{otherwise} \end{cases}
$$

# FE1

A fence is redundant if it always follows a previous fence or locked instruction in program order, and no memory store instructions are in between.

A *forward* data-flow problem over the bo...

Assoc...

$\perp$ : alo...
   is a...
   cur...
   no...

$\top$ : oth...

*Implementation*:

1. Use CompCert implementation of Kildall algorithm to solve the data-flow equations.

2. Replace $\mathtt{MFENCE}$s for which the analysis returns $\perp$ with $\mathtt{NOP}$ instructions.

$$T_1(\mathbf{nop}, \mathcal{E}) = \mathcal{E}$$
$$= \mathcal{E}$$
$$= \mathcal{E}$$
$$= \top$$
$$) = \top$$
$$= \mathcal{E}$$
$$= \top$$
$$) = \top$$
$$= \perp$$
$$= \perp$$

$$\emptyset$$

$$\mathcal{FE}_1(n) = \begin{cases} & \\ \bigsqcup_{p \in \text{predecessors}(n)} T_1(instr(p), \mathcal{FE}_1(p)) & \text{otherwise} \end{cases}$$
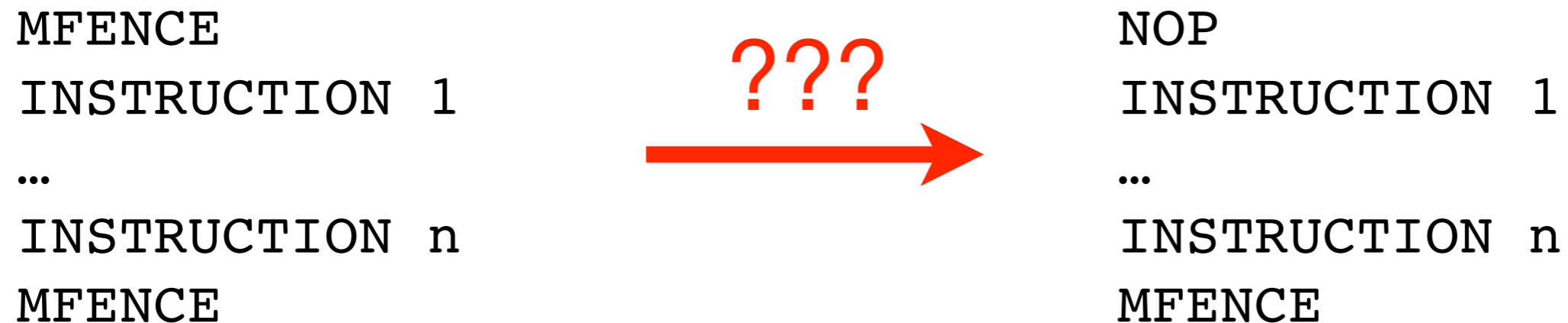
# Redundant fences (2)

If we have two consecutive fence instructions, we can remove the *former*:

```
MFENCE                                  NOP
MFENCE         ──────────▶              MFENCE
```

*Intuition:* the visible effects initially published by the former fence, are now published by the latter, and nobody can tell the difference.

*Generalisation:*

```
MFENCE                     ???          NOP
INSTRUCTION 1                           INSTRUCTION 1
…              ──────────▶              …
INSTRUCTION n                           INSTRUCTION n
MFENCE                                  MFENCE
```

# Redundant fences (2)

If there are reads in between the fences…

[x]=[y]=0

| Thread 0 | Thread 1 |
|---|---|
| MOV [x] ← 1<br>**MFENCE**<br>MOV EAX ← [y]<br>MFENCE | MOV [y] ← 1<br>MFENCE<br>MOV EBX ← [x] |

EAX = EBX = 0
forbidden

but

[x]=[y]=0

| Thread 0 | Thread 1 |
|---|---|
| MOV [x] ← 1<br>**NOP**<br>MOV EAX ← [y]<br>MFENCE | MOV [y] ← 1<br>MFENCE<br>MOV EBX ← [x] |

EAX = EBX = 0
allowed

# Redundant fences (2)

If there are reads in between the fences…

[x]=[y]=0

| Thread 0 | Thread 1 |
|---|---|
| MOV [x] ← 1<br>**MFENCE**<br>MOV EAX ← [y]<br>MF | MOV [y] ← 1<br>MFENCE |

EAX = EBX = 0
forbidden

but

If there are reads in between, the optimisation is unsound.

[x]=[y]=0

| Thread 0 | Thread 1 |
|---|---|
| MOV [x] ← 1<br>**NOP**<br>MOV EAX ← [y]<br>MFENCE | MOV [y] ← 1<br>MFENCE<br>MOV EBX ← [x] |

EAX = EBX = 0
allowed

# Redundant fences (2)

Swapping a `STORE` and a `MFENCE` is sound:

<div align="center">

`MFENCE; STORE`   ⟶   `STORE; MFENCE`

</div>

1. transformed program's behaviours ⊆ source program's behaviours

(source program might leave pending write in its buffer)

2. There is the new intermediate state if the buffer was initially non-empty, but this intermediate state *is not observable.*

(a local read is needed to access the local buffer)

*Intuition:* Iterate this swapping*...*

# FE2

A fence is redundant if it always precedes a later fence or locked instruction in program order, and no memory read instructions are in between.

A *backward* data-flow problem over the boolean domain $\{\bot, \top\}$

Associate to each program point:

$\bot$ : along all execution paths there is an atomic instruction *after* the current program point, with no intervening reads;

$\top$ : otherwise.

$$
\begin{aligned}
T_2(\mathbf{nop}, \mathcal{E}) &= \mathcal{E} \\
T_2(\mathbf{op}(op, \vec{r}, r), \mathcal{E}) &= \mathcal{E} \\
T_2(\mathbf{load}(\kappa, addr, \vec{r}, r), \mathcal{E}) &= \top \\
T_2(\mathbf{store}(\kappa, addr, \vec{r}, src), \mathcal{E}) &= \mathcal{E} \\
T_2(\mathbf{call}(sig, ros, args, res), \mathcal{E}) &= \top \\
T_2(\mathbf{cond}(cond, args), \mathcal{E}) &= \mathcal{E} \\
T_2(\mathbf{return}(optarg), \mathcal{E}) &= \top \\
T_2(\mathbf{threadcreate}(optarg), \mathcal{E}) &= \top \\
T_2(\mathbf{atomic}(aop, \vec{r}, r), \mathcal{E}) &= \bot \\
T_2(\mathbf{fence}, \mathcal{E}) &= \bot
\end{aligned}
$$

$$
\mathcal{FE}_2(n) = 
\begin{cases}
\top & \text{if } \mathrm{successors}(n) = \emptyset \\
\bigsqcup_{s \in \mathrm{successors}(n)} T_2(instr(s), \mathcal{FE}_2(s)) & \text{otherwise}
\end{cases}
$$

# FE1 and FE2 are both useful

Removed by FE1 but not FE2:

```
MFENCE
MOV EAX <- [y]
MFENCE
MOV EBX <- [y]
```

Removed by FE2 but not FE1:

```
MOV [x] <- 1
MFENCE
MOV [x] <- 2
MFENCE
```

# Informal correctness argument

*Intuition*: FE2 can be thought as iterating

```
MFENCE; STORE          ──▶          STORE; MFENCE

MFENCE; non-mem        ──▶          non-mem; MFENCE
```

and then applying

```
MFENCE; MFENCE         ──▶          NOP; MFENCE
```

This argument works for *finite traces*, but not for *infinite traces* as the later fence might never be executed:

```
MFENCE;                              NOP;
STORE;          ──▶                  STORE;
WHILE(1);                            WHILE(1);
MFENCE                               MFENCE
```

# Basic simulations

A pair of relations

$$\sim \in \mathbb{P}(src.states \times tgt.states) \qquad\qquad > \in \mathbb{P}(tgt.states \times tgt.states)$$

is a *basic simulation* for $\quad \textbf{compile} : src.prg \to tgt.prg \quad$ if:

$sim\_init : \forall p\, p'.\ \textsf{compile}(p) = p' \implies \forall t \in \textsf{init}(p').\ \exists s \in \textsf{init}(p).\ s \sim t$

$sim\_end : \forall s\, t.\ s \sim t \wedge t \not\to_{-} \implies s \not\to_{-}$

$sim\_step : \forall s\, t\, t'\, ev.\ s \sim t \wedge t \xrightarrow{ev} t' \wedge ev \neq \textsf{oom} \implies$

$\qquad\qquad (s \xrightarrow{\tau}{}^* \xrightarrow{\textbf{fail}} {}_{-}) \qquad\qquad\qquad — s \text{ reaches a failure}$

$\qquad\qquad \vee\, (\exists s'.\ s \xrightarrow{\tau}{}^* \xrightarrow{ev} s' \wedge s' \sim t') \quad — s \text{ does matching step sequence}$

$\qquad\qquad \vee\, (ev = \tau \wedge t > t' \wedge s \sim t').\quad — s \text{ stutters (only allowed if } t > t')$

Exhibiting a basic simulation implies:

$$\textsf{traces}(\textsf{compile}(p)) \setminus \{t \cdot \textsf{inftau} \mid t \text{ trace}\} \subseteq \textsf{traces}(p)$$

"simulation can stutter forever"

# Usual approach: measured simulations

**Definition 2 (Measured sim.).** *A* measured simulation *is any basic simulation* $(\sim, >)$ *such that* $>$ *is well-founded.*

**Theorem 1.** *If there exists a measured simulation for the compilation function* compile, *then for all programs* $p$, $\text{traces}(\text{compile}(p)) \subseteq \text{traces}(p)$.

# Simulation for FE2

$s \equiv_i t$  iff thread $i$ of $s$ and $t$ have identical pc, local states and buffers

$s \leadsto_i s'$ iff thread $i$ of $s$ can execute zero or more `NOP, OP, STORE` and `MFENCE` instructions and end in the state $s'$

$s \sim t$  iff
  – $t$'s CFG is the optimised version of $s$'s CFG; and
  – $s$ and $t$ have identical memories; and
  – $\forall$ thread $i,$ either $s \equiv_i t$ or
        the analysis for $i$'s pc returned $\bot$ and $\exists s', s \leadsto_i s'$ and $s' \equiv_i t$
          "$s$ is some instructions behind and can catch up"

*Stutter condition*:
  $t > t'$  iff  $t \rightarrow t'$ by a thread executing a `NOP, OP, STORE` or `MFENCE`
        (and $t$'s buffer being non-empty)

# Simulation for FE2

$s \equiv_i t$   iff thread $i$ of $s$ and $t$ have identical pc, local states and buffers

$s \leadsto_i s'$ iff th

    MFE

$s \sim t$   iff

  – $t$'s CFG

  – $s$ and $t$

  – $\forall$ thread

the analysis for $i$'s pc returned $\bot$ and $\exists s'$, $s \leadsto_i s'$ and $s' \equiv_i t$

"$s$ is some instructions behind and can catch up"

*Stutter condition*:

   $t > t'$  iff   $t \rightarrow t'$ by a thread executing a NOP, OP, STORE or MFENCE

(and $t$'s buffer being non-empty)

But if (1) all threads have non-empty buffers, and
            (2) are stuck executing infinite loops, and
            (3) no writes are ever propagated to memory,
then we can stutter forever.

(i.e., $>$ is not well-founded.)

# Simulation for FE2

$s \equiv_i t$    iff thread *i* of *s* and *t* have identical pc, local states and buffers

$s \leadsto_i s'$ iff th

MFE

$s \sim t$    iff
- *t*'s CFG
- *s* and *t*  (i.
- ∀ threa

But if (1) all threads have non-empty buffers, and
                  (2) are stuck executing infinite loops, and
                  (3) no writes are ever propagated to memory,
then we can stutter forever.

*Stutter conditio*
   $t > t'$   iff   $t \rightarrow$

(and

**Solution 1:** Assume this case never arises (*fairness*)

**Solution 2:** Do a case split.
— If this case does not arise, we are done.
— If it does, use a different (weaker) simulation to construct an infinite trace for the source

$\equiv_i t$

# Weaktau simulation

**Definition 3 (Weaktau sim.).** $A$ weaktau simulation *consists of a basic simulation $(\sim, >)$ with and an additional relation between source and target states, $\simeq \in \mathbb{P}(src.states \times tgt.states)$ satisfying the following properties:*

$$sim\_weaken : \forall s, t.\; s \sim t \implies s \simeq t$$
$$sim\_wstep : \forall s\, t\, t'.\; s \simeq t \wedge t \xrightarrow{\tau} t' \wedge t > t' \implies$$
$$(s \xrightarrow{\tau}{}^* \xrightarrow{\texttt{fail}} \_) \qquad\qquad — s \text{ reaches a failure}$$
$$\vee\; (\exists s'.\; s \xrightarrow{\tau}{}^* \xrightarrow{\tau} s' \wedge s' \simeq t') \quad — s \text{ does a matching step sequence.}$$

**Theorem 2.** *If there exists a weaktau-simulation $(\sim, >, \simeq)$ for the compilation function* compile, *then for all programs $p$,* traces(compile$(p)$) $\subseteq$ traces$(p)$.

Remarks:
— Once the simulation game moves from ~ to ≃, stuttering is forbidden;

— Can view difference between ~ and ≃ as a boolean prophecy variable.

# Weaktau simulation for FE2

$s \sim t$ , $t > t'$ as before.

$s \simeq t$ iff

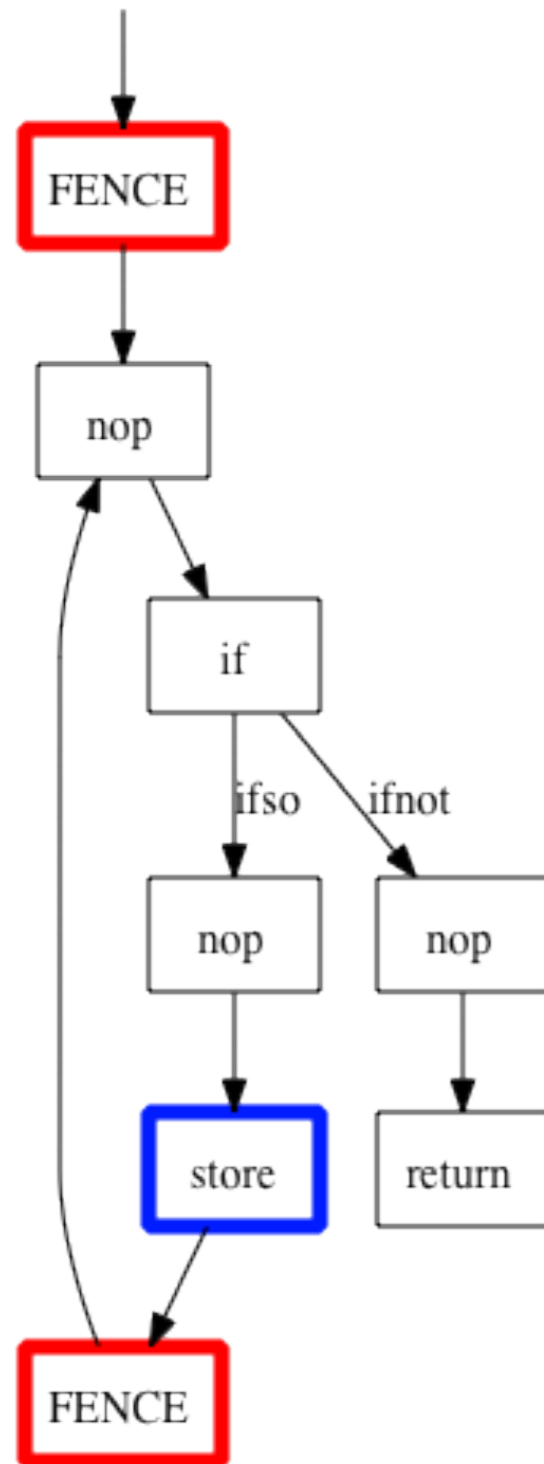- *t*'s CFG is the optimised version of *s*'s CFG; and
- $\forall i, \exists s'$ s.t. $s \leadsto_i s' \equiv_i t.$

(i.e., same as $s \sim t$ except that the memories memories are unrelated.)
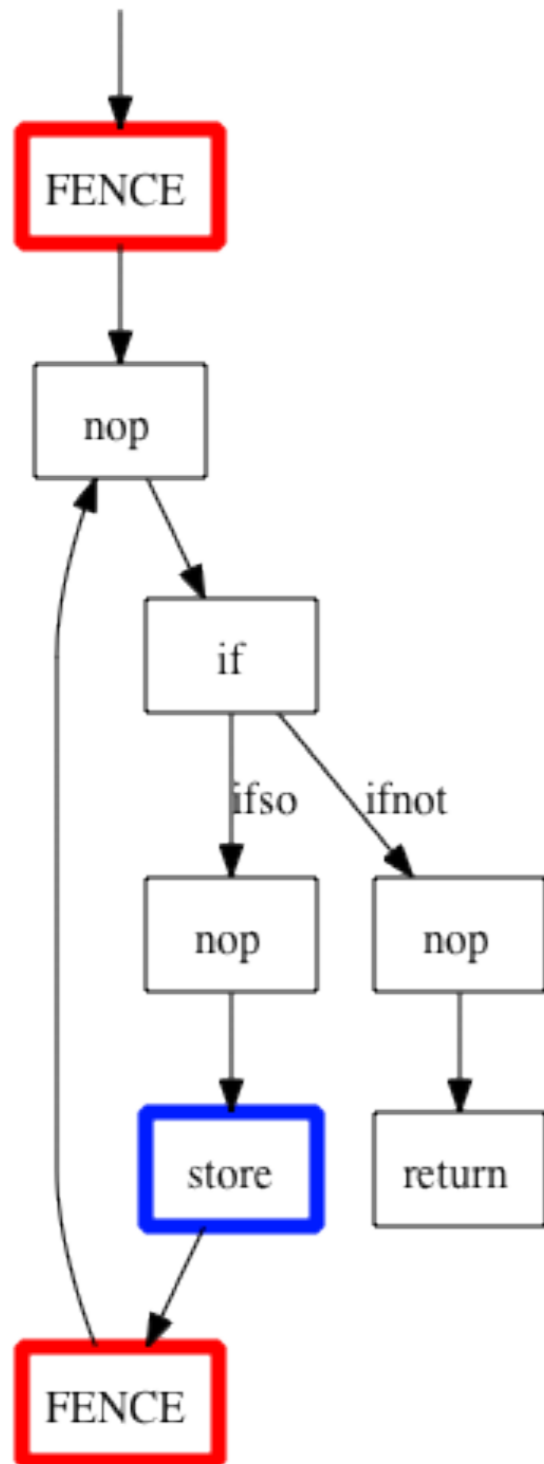
# A closer look at the RTL



Patterns like that on the left are common.

FE1 and FE2 do not optimise these patterns.

It would be nice to hoist those fences out of the loop.

# A closer look at the RTL
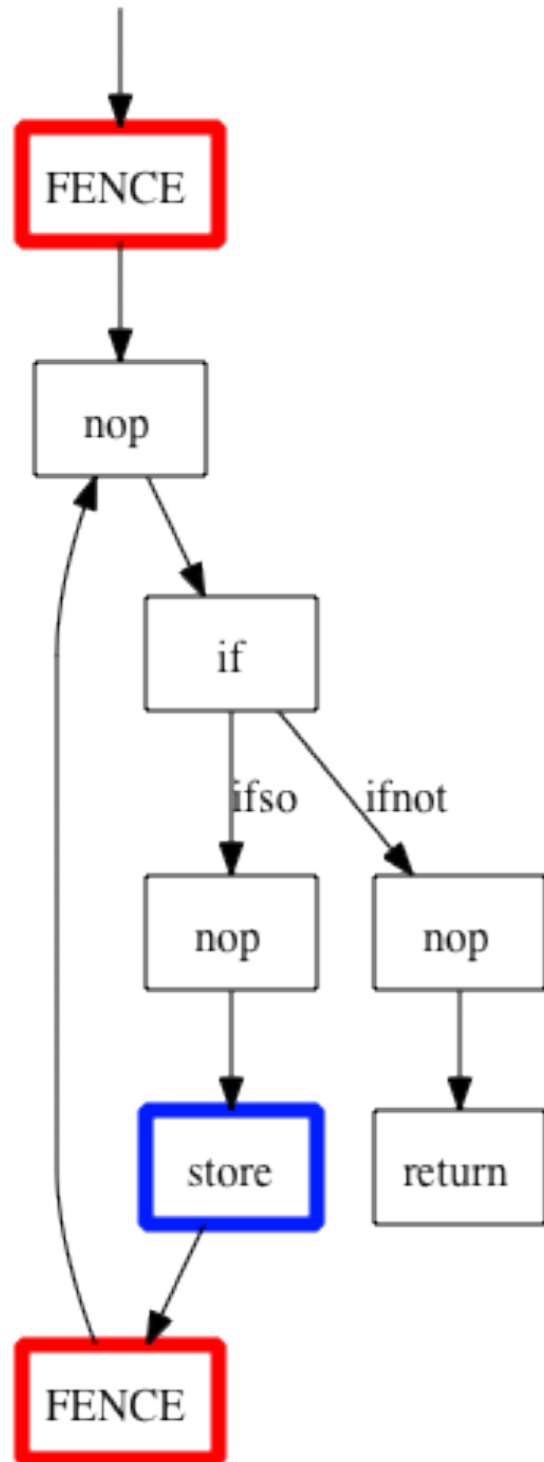
Patterns like that on the left are common.

FE1 and FE2 do not optimise these patterns.

It would be nice to hoist the fences out of the loop.

Do you perform PRE?

# A closer look at the RTL



Patterns like that on the left are common.
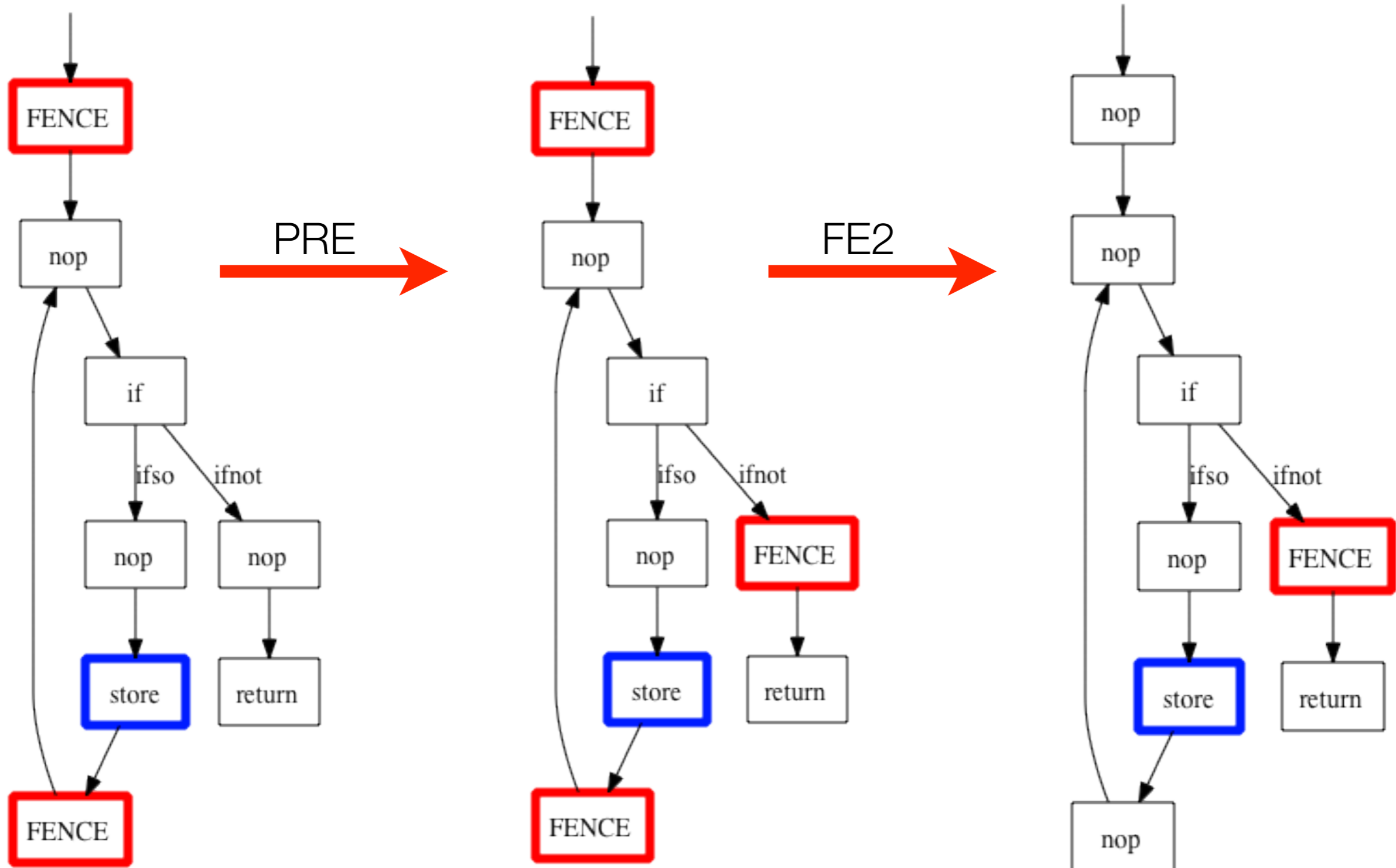
FE1 and FE2 do not optimise these patterns.

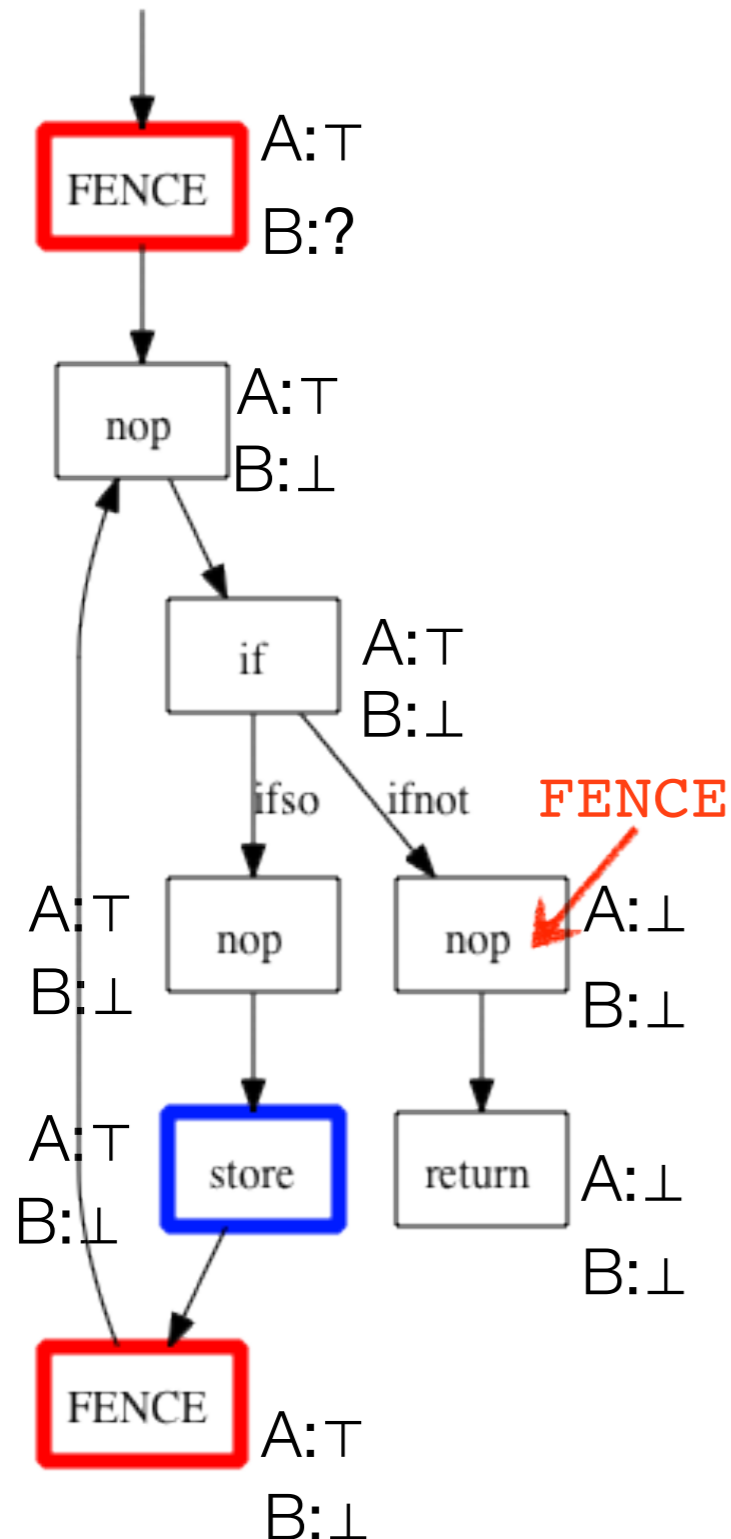It would be nice to hoist some of these out of the loop.

Do you perform PRE?

...adding a fence is always safe...

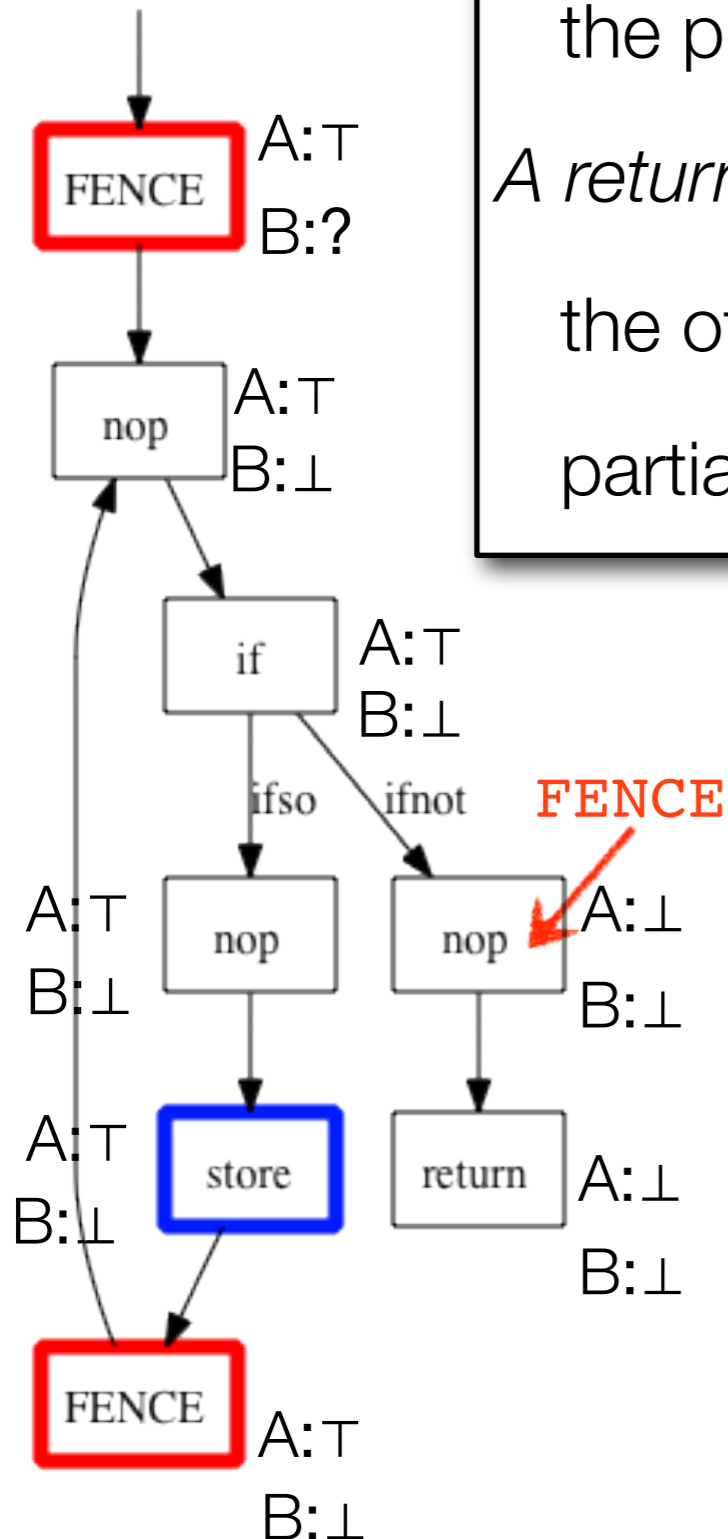# Partial redundancy elimination

# Partial redundancy elimination



*A*: a backward analysis returning ⊤ if along *some path* after the current program point there is an atomic instruction with no intervening reads;

*B*: a forward analysis returning ⊥ if along all paths to the current program point there is a fence with no later reads or atomic instructions.

*Replace* `NOP` *with* `FENCE` *after conditionals if:*
- B returns ⊥
- A returns ⊥
- A returns ⊤ on the other branch

# Partial redu...



A box overlay:

*B returns ⊥:*

a previous fence will be eliminated if we insert a fence at both branches of conditional nodes.

*A returns ⊥:*

the previous fence won't be removed by FE2.

*A returns ⊤ on the other branch:*

the other branch already makes the previous fence

partially redundant.

*B*: a forward analysis returning ⊥ if along all paths to the current program point there is a fence with no later reads or atomic instructions.

*Replace* `NOP` *with* `FENCE` *after conditionals if:*
- B returns ⊥
- A returns ⊥
- A returns ⊤ on the other branch

# Evaluation of the optimisations

– Insert `MFENCE`s *before every read* (br)*, or *after every write* (aw).

– Count the `MFENCE` instructions in the generated code.

|  | br | br+FE1 | aw | aw+FE2 | aw+PRE+FE2 |
|---|---|---|---|---|---|
| Dekker | 3 | 2 | 5 | 4 | 4 |
| Bakery | 10 | 2 | 4 | 3 | 3 |
| Treiber | 5 | 2 | 3 | 1 | 1 |
| Fraser | 32 | 18 | 19 | 12 | 11 |
| TL2 | 166 | 95 | 101 | 68 | 68 |
| Genome | 133 | 79 | 62 | 41 | 41 |
| Labyrinth | 231 | 98 | 63 | 42 | 42 |
| SSCA | 1264 | 490 | 420 | 367 | 367 |

# Evaluation of the optimisations

– Insert `MFENCE`s *before every read* (br)*, or *after every write* (aw).

– Cou[...]

| | | | | | 2 |
|---|---|---|---|---|---|

*Important remark for your future work:*

This is not a decent evaluation…  we know nothing about real code, and the number of fences is not a good measure.  But unclear how to do better.

*Evaluation should be taken seriously by CS scientists!*

`http://evaluate.inf.usi.ch/`

| | | | | | |
|---|---|---|---|---|---|
| Labyrinth | 231 | 98 | 63 | 42 | 42 |
| SSCA | 1264 | 490 | 420 | 367 | 367 |

# Conclusion

# Syllabus

In these lectures we have covered the hardware models of two modern computer architectures (x86 and Power/ARM - at least for a large subset of their instruction set).

We have seen how compiler optimisations can also break concurrent programs and the importance of defining the memory model of high-level programming languages (and we have seen in detail the C++11 memory model).

We have also introduced some proof methods to reason about concurrency.

*After these lectures, you might have the feeling that multicore programming is a mess and things can't just work.*

The memory models of modern hardware are better understood.

Programming languages attempt to specify and implement reasonable memory models.

Researchers and programmers are now interested in these problems.

The memory models of modern hardware are better understood.

attempt

*Still, many open problems...*

els.

mmers

se

problems.

The memory models of modern hardware are better understood.
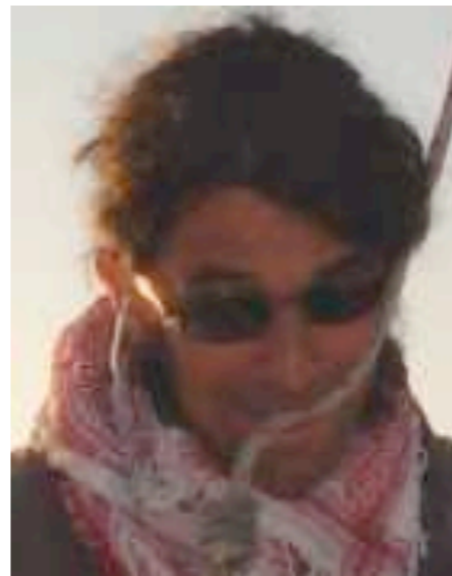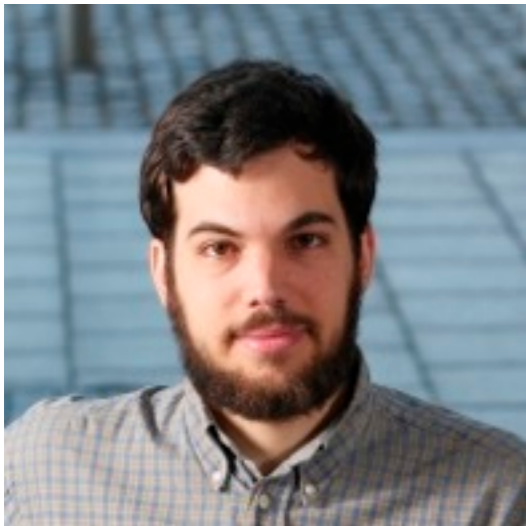
attempt

Still, many research opportunities! els.

mmers se

problems.

All these lectures are based on work done with/by my colleagues. Thank you!

# And thank you all for attending these lectures!

Please, fill the course evaluation form.
It is vital feedback to make a better course next year.