

Introduction to weak memory consistency

Ori Lahav

Viktor Vafeiadis

28 August 2017

Weak memory consistency

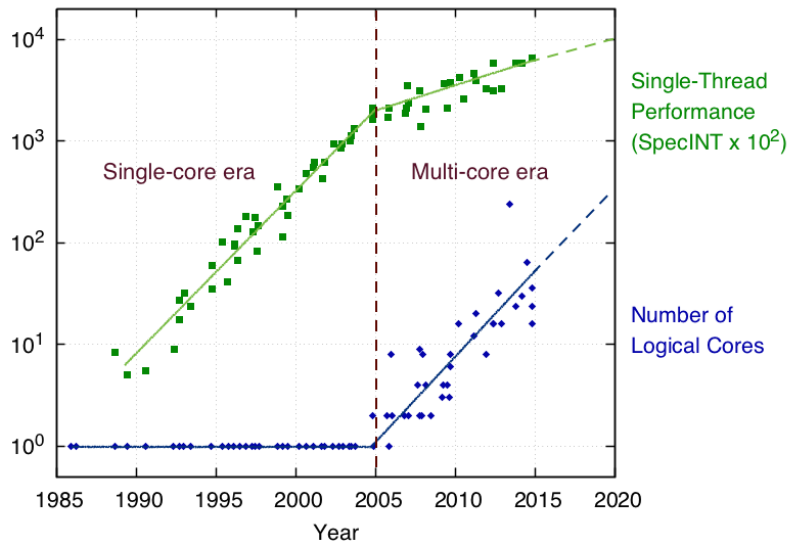
is about the

semantics of concurrent programs

taking into account the effects of:

- ▶ multicore **hardware** implementations
- ▶ and **compiler** optimizations.

CPU trends: Parallelism is here!



Concurrent programming is hard!

*If you can get away with it, avoid using threads.
Threads can be **difficult to use**, and they make
programs **harder to debug**.*

(Java documentation, \approx 15 years ago)

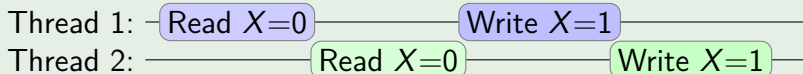
“Difficult to use”

- ▶ Requires a fundamentally different way of thinking.
- ▶ Interference among threads.

“Harder to debug”

- ▶ Huge non-determinism \rightsquigarrow testing is ineffective.

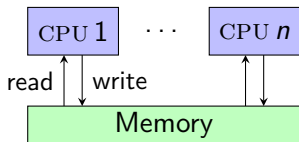
$X := X + 1 \parallel X := X + 1$ might increment X only once.



The illusion of sequential consistency

Sequential consistency (SC)

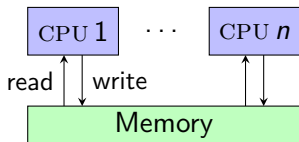
- ▶ The standard simplistic concurrency model.
- ▶ Threads access shared memory in an interleaved fashion.



The illusion of sequential consistency

Sequential consistency (SC)

- ▶ The standard simplistic concurrency model.
- ▶ Threads access shared memory in an interleaved fashion.



But...

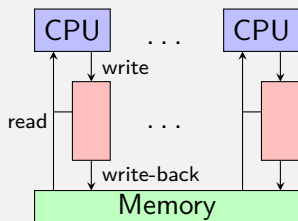
- ▶ No multicore processor implements SC.
- ▶ Compiler optimizations invalidate SC.

Weak consistency

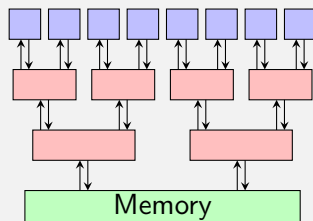
Hardware provides **weak consistency**.

- ▶ **Weak memory models** \rightsquigarrow semantics of shared memory.
- ▶ Every hardware architecture has its own WMM:
x86-TSO, ARM, Power, Itanium.

x86-TSO model (2010)



ARMv8 model (2016)



Weak consistency examples

Store buffering (SB)

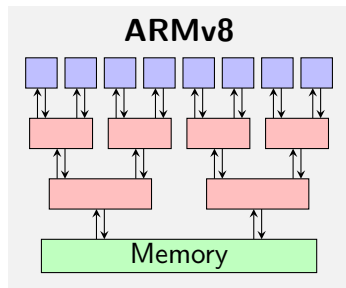
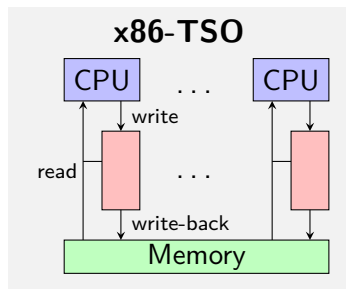
Initially, $x = y = 0$

$x := 1;$ $y := 1;$
 $a := y$ //0 $b := x$ //0

Load buffering (LB)

Initially, $x = y = 0$

$a := y;$ //1 $b := x;$ //1
 $x := 1$ $y := 1$

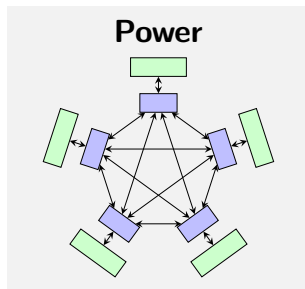


Independent reads of independent writes (IRIW)

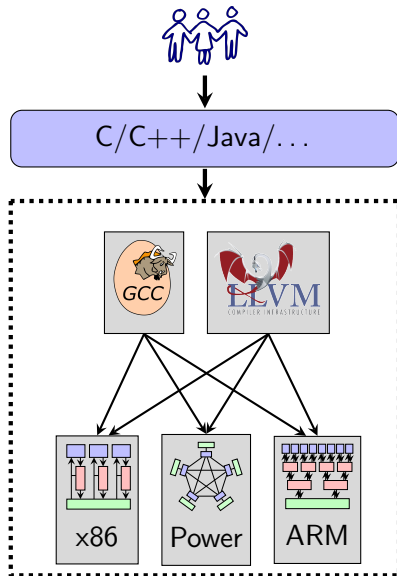
Initially, $x = y = 0$

$$x := 1 \quad \left\| \begin{array}{l} a := x; \text{ //1} \\ \text{lwsync;} \\ b := y \text{ //0} \end{array} \right\| \left\| \begin{array}{l} c := y; \text{ //1} \\ \text{lwsync;} \\ d := x \text{ //0} \end{array} \right\| \quad y := 1$$

- ▶ Thread II and III can observe the $x := 1$ and $y := 1$ writes happen in different orders.
- ▶ Because of the `lwsync` fences, no reorderings are possible!



WMC is not just about hardware



Quiz. Should these transformations be allowed?

1. CSE over acquiring a lock:

<code>a = x;</code>		<code>a = x;</code>
<code>lock();</code>	\rightsquigarrow	<code>lock();</code>
<code>b = x;</code>		<code>b = a;</code>

2. Load hoisting:

<code>if (c)</code>		<code>t = x;</code>
<code> a = x;</code>	\rightsquigarrow	<code>a = c ? t : a;</code>

[`x` is a global variable; `a`, `b`, `c` are local; `t` is a fresh temporary.]

Consider the transformation sequence:

if (c)		$t = x;$		$t = x;$
$a = x;$	hoist	$a = c ? t : a;$	CSE	$a = c ? t : a;$
lock();		lock();		lock();
$b = x;$		$b = x;$		$b = t;$

When c is false, x is moved out of the critical region!

So we have to forbid one transformation.

- ▶ C11 forbids load hoisting, allows CSE over lock().
- ▶ LLVM allows load hoisting, forbids CSE over lock().

Weak consistency in “real life”

- ▶ Messages may be delayed.


$$\begin{array}{l} \text{MsgX} := 1; \\ a := \text{MsgY}; \quad //0 \end{array} \parallel \begin{array}{l} \text{MsgY} := 1; \\ b := \text{MsgX}; \quad //0 \end{array}$$


- ▶ Messages may be sent/received out of order.


$$\begin{array}{l} \text{Email} := 1; \\ \text{Sms} := 1; \end{array} \parallel \begin{array}{l} a := \text{Sms}; \quad //1 \\ b := \text{Email}; \quad //0 \end{array}$$


Weak consistency is not a threat, but an opportunity.

- ▶ Can lead to more scalable concurrent algorithms.
- ▶ Several open research problems.
 - ▶ What is a good memory model?

Reasoning under WMC is often easier than under SC.

- ▶ Avoid thinking about thread interleavings.
- ▶ Many/most concurrent algorithms do not need SC!
- ▶ Positive vs negative knowledge.

Memory model definitions

- ▶ Operational memory models
- ▶ Axiomatic/declarative memory models
- ▶ Promising semantics

WMM metatheory

- ▶ Relating memory models
- ▶ Correctness of compilation and program transformations
- ▶ Programming guarantees: the DRF theorem

Verification techniques for WMM

- ▶ Program logics (relaxed separation logic, OGRA)
- ▶ Model checking