# Concurrency in weak memory models

Andreas Lööw

Chalmers, 4th year PhD student, formal methods
(= mathematical reasoning about software and hardware)

# Aim of talk

Memory model related differences between programming in:

- "modelling languages" like Promela and pseudocode, and

- "real languages" like Java.

The talk is both Java specific and not Java specific:

- Java used as an example of a language with a "weak memory model",

- but at least (subsets of) C and C++ similar

# Talk in one slide

In "modelling languages", synchronization is used for:

- atomicity

In "real languages", synchronization is used for:

- atomicity, and
- **visibility**

# Outline

- What are memory models?

- Why weak memory models?

- Something about the Java memory model (as an example of a weak memory model)

- Programming in the Java memory model

# Outline

- **<span style="color:red">What are memory models?</span>**

- Why weak memory models?

- Something about the Java memory model (as an example of a weak memory model)

- Programming in the Java memory model

# Java has a "weak memory model"

- Memory model part of language semantics (what programs mean)

- Different memory models exist

- In pseudocode, sequential consistency (SC) often assumed -- one of the "strongest" memory models

- Java, instead, offers the Java memory model (JMM), one particular "weak" memory model

# OK… but what is a memory model?

- More or less: Semantics of shared variables (and synchronization)

- Consider the question: What values are variable reads allowed to return?

- ???

# Reading variables: Sequential programming

Obvious answer: The <span style="color:red">latest</span> value we wrote to the variable

```
int x = 0, y = 0;
x = 1;
y = 1;
print(y); // will obviously print 1
print(x); // again, prints 1
```

# Reading variables: Concurrent programming

```
int x = 0, y = 0;

t1 {
  x = 1;
  y = 1;
}


t2 {
  print(y);
  print(x);
}
```

Simple! Just consider the non-deterministic interleavings!

E.g., t1 completes before t2:

1

1


Or, other interleaving:

0

1

# Reading variables: Concurrent programming

```
int x = 0, y = 0;

t1 {
  x = 1;
  y = 1;
}


t2 {
  print(y);
  print(x);
}
```

But can we print the following?
```
1
0
```

Depends on memory model:

- Sequential consistency: No!

- Java memory model: Yes!

# Reading variables: Sequential consistency (SC)

```
int x = 0, y = 0;

t1 {
  x = 1;
  y = 1;
}

t2 {
  print(y);
  print(x);
}
```

"Program order" always maintained in CS

In particular, $x = 1$ always before $y = 1$ in any interleaving

Consequently, will not see

1

0

But the above program order guarantee not provided by some weak memory models!

# Reading variables: Weak memory models

```
int x = 0, y = 0;

t1 {
  x = 1;
  y = 1;
}


t2 {
  print(y);
  print(x);
}
```

"Interleaving-based semantics" in some sense the "obvious" semantics for concurrency

Why make things more difficult? Why give up program order and other nice things?

Because: SC costs too much

# Outline

- What are memory models?

- Why weak memory models?

- Something about the Java memory model (as an example of a weak memory model)

- Programming in the Java memory model

# SC cost 1: Prohibits (too many) compiler optimizations

- Aaaaah!!! Messiness! Real-world things! In pseudocode we do not have to consider ugliness such as compiler "details" etc.

- Example: For some compiler optimizations we want to reorder writes to variables. (For whatever reason: Might improve register allocation or anything.)

# SC cost 1: Prohibits (too many) compiler optimizations

- E.g., the transformation to the right "semantics preserving" in sequential setting if we only consider final state of program

- Not equivalent if we can inspect program under execution, which we can if x and y are shared variables in a concurrent setting

- Breaks illusion of "program order"!

Original program:

x = 1;

y = 2;

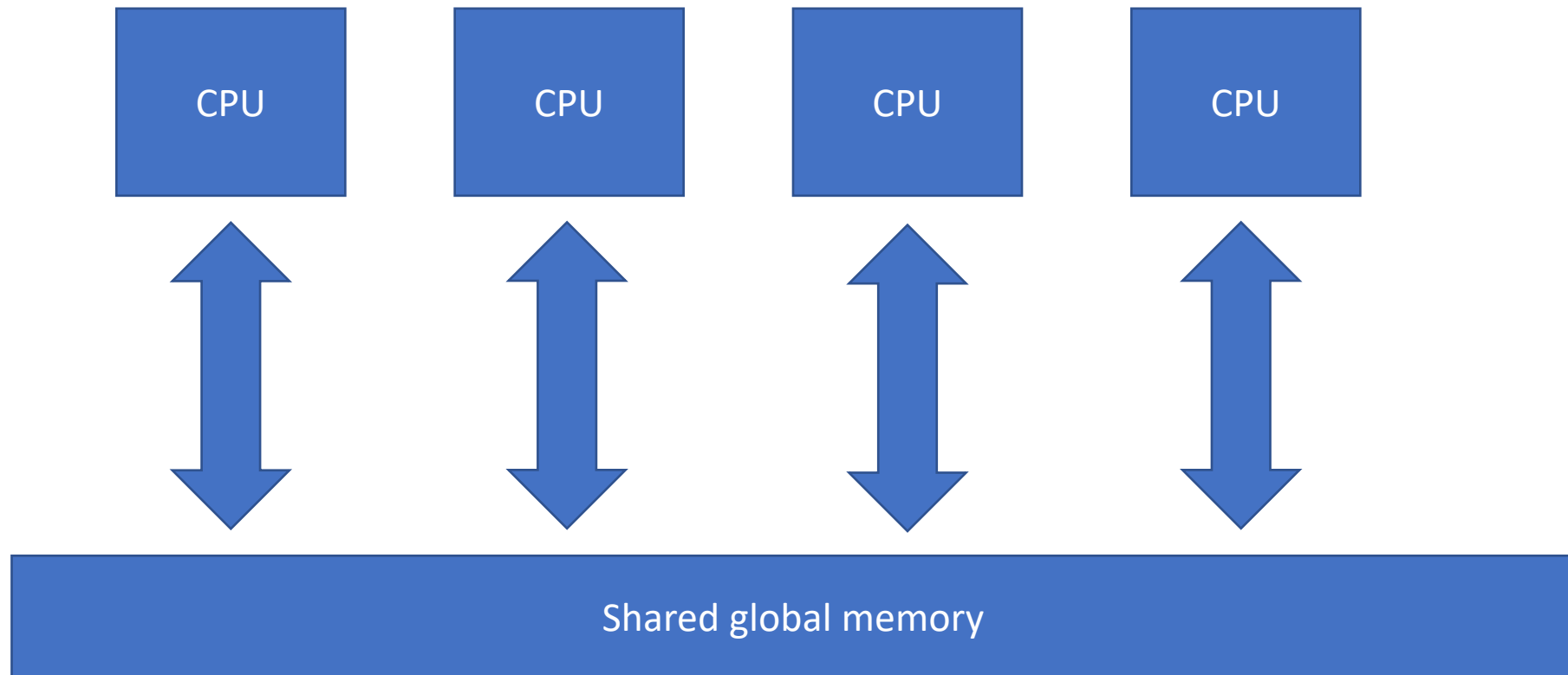z = x + y; // x = 1, y = 2, z = 3

Transformed program:

y = 2;

x = 1;

z = x + y; // x = 1, y = 2, z = 3

Write order swapped

# SC cost 2: Causes too much cache synchronization

Cost of SC not obvious with too simplified machine models:

# SC cost 2: Causes too much cache synchronization

More realistic model of today's computers:

Btw, modern CPUs execute instructions out-of-order and in parallel (which can also break illusion of program order)

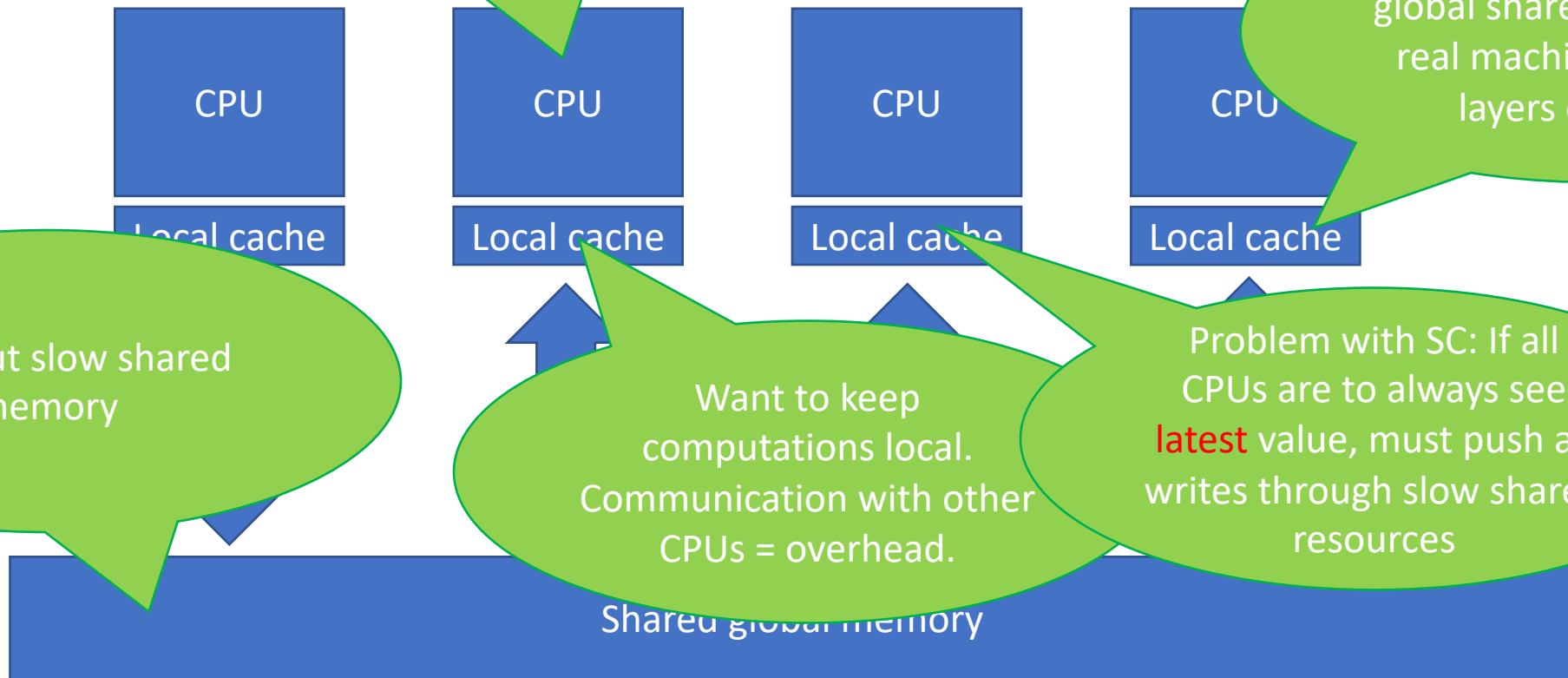Small but fast compared to global shared memory. (In real machines: multiple layers of cache.)

CPU
CPU
CPU
CPU

Local cache
Local cache
Local cache
Local cache

Large but slow shared memory

Want to keep computations local. Communication with other CPUs = overhead.

Problem with SC: If all CPUs are to always see latest value, must push all writes through slow shared resources

Shared global memory

# Why not SC: Summary

- Not a complete list of reasons, just two examples!

- Anyhow, in summary:
  <span style="color:red">SC too expensive in many situations</span>

- Solution to mentioned problems:
  Relax some guarantees offered by SC → we get weak memory models

- Weaker memory models (potentially) <span style="color:red">more performant</span>, but <span style="color:red">more difficult to program in</span>
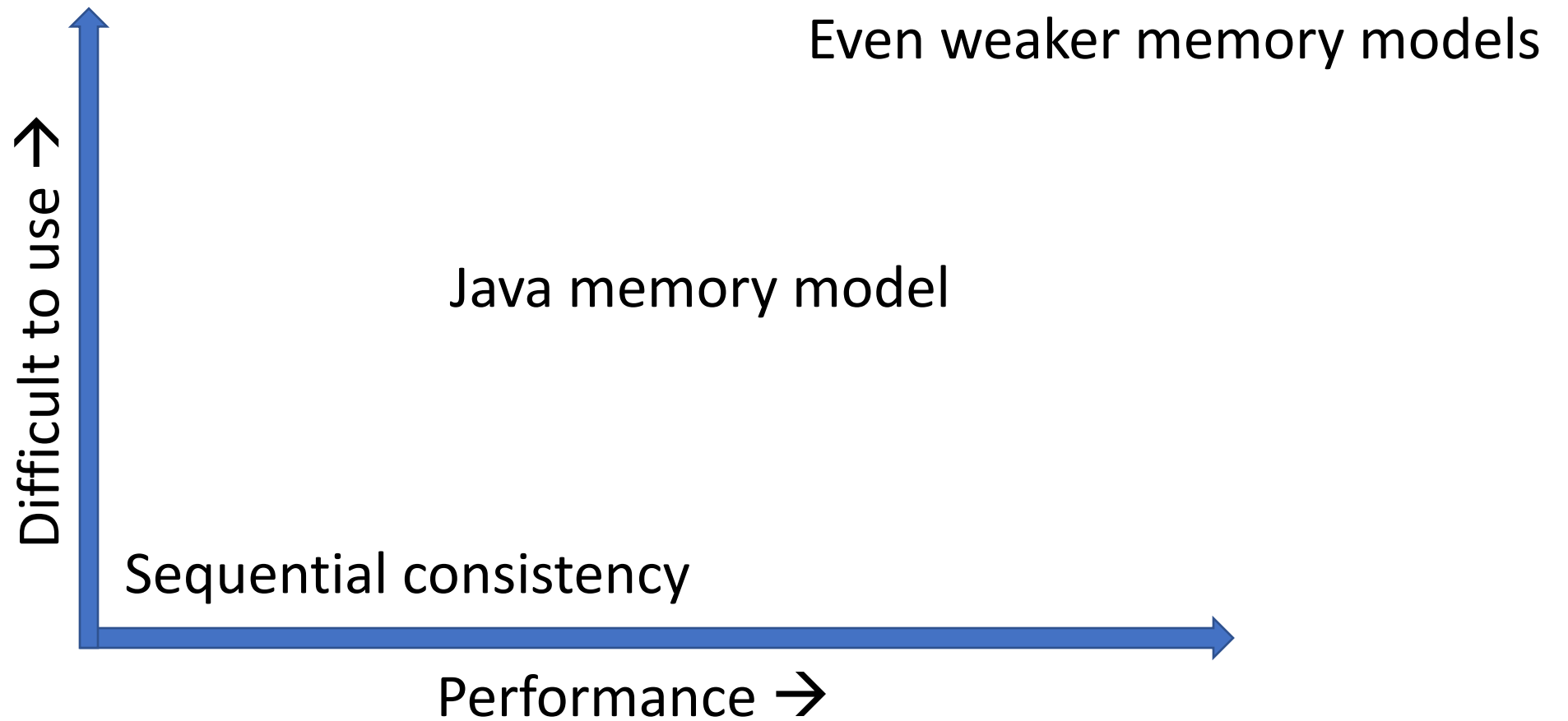
# Outline

- What are memory models?

- Why weak memory models?

- <span style="color:red">Something about the Java memory model (as an example of a weak memory model)</span>

- Programming in the Java memory model (as an example of programming in a weak memory model)
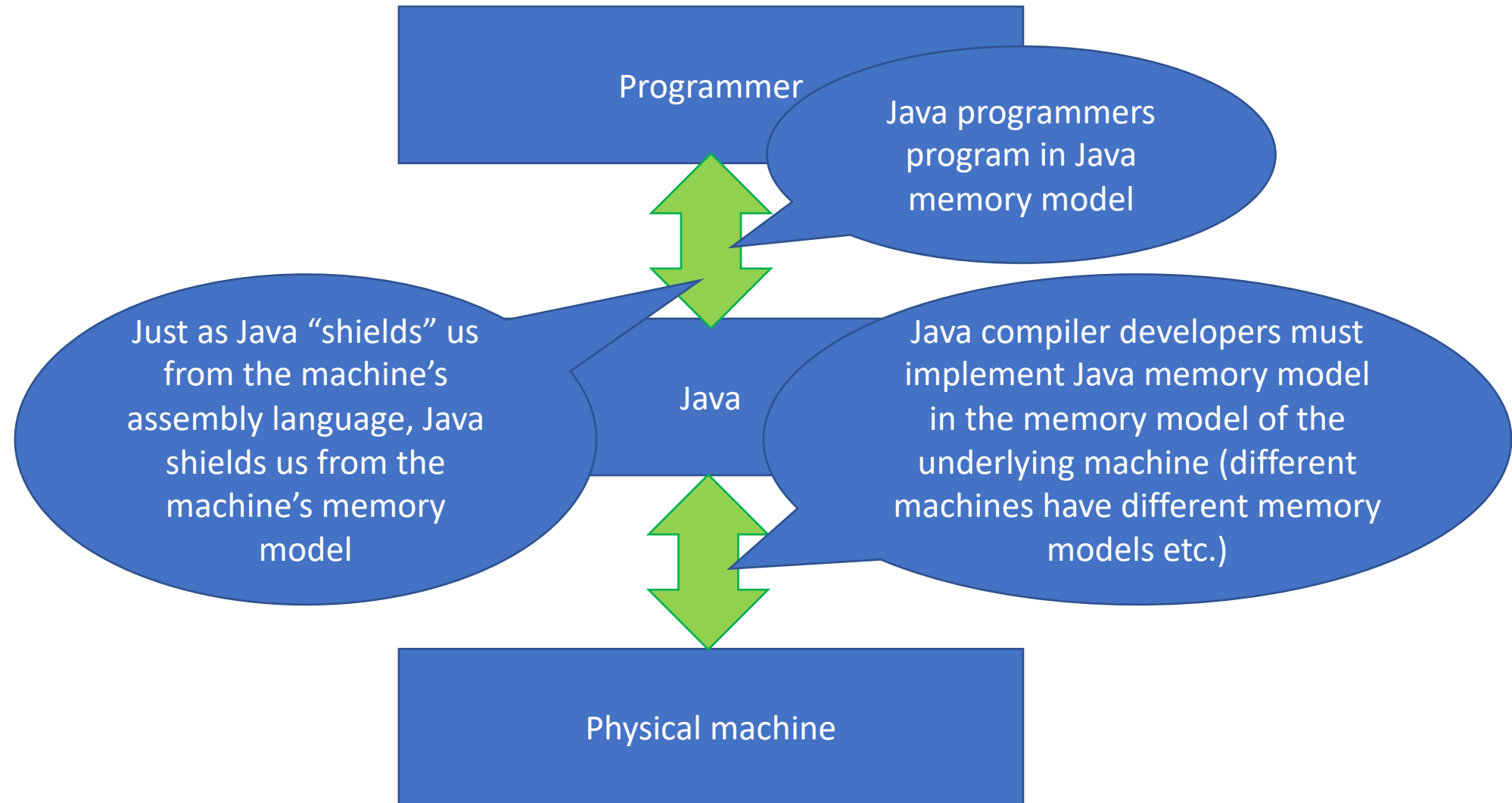
# The Java memory model

- Less convenient than SC, but implementable on modern machine architectures without too much performance loss

- Opinion: Memory model part of language design, and different coordinates in the design space have different tradeoffs. As with any other language feature: No "right" answer.

# Design tradeoff space

# More context: Some more machine details

Programmer

Java programmers program in Java memory model

Just as Java "shields" us from the machine's assembly language, Java shields us from the machine's memory model

Java

Java compiler developers must implement Java memory model in the memory model of the underlying machine (different machines have different memory models etc.)

Physical machine

# SC for data-race-free programs

- A few (C-like) languages have converged at "sequential consistency for data-race-free programs" memory models

- Java included in this family

- Reasoning principle: If there are no data races (under SC), we can assume SC when reasoning about our program

- Important to remember definition of data race (and difference with race conditions)

# Data races

Slight variation of previous definition you seen, to fit Java better:

**Def.** Two memory accesses are in a data race iff
- they access the same memory location simultaneously (they are interleaved next to each other),
- at least one access is a write,
- insufficient explicit synchronization used to protect the accesses

**Def.** A program is data-race-free iff no SC execution of the program contain a data race.

("Slight variation"? Note that we quantify over all SC executions in the second definition.)

Note that data-race-freedom is a "language-level" property!

# Definition of data race surprisingly subtle

E.g., does this program contain any data races?

```
bool x = false, y = false;

t1 {
    if (x) y = true;
}

t2 {
    if (y) x = true;
}
```

No!

# Race conditions

Definition from course slides:

**Def.** A *race condition* is a situation where the <span style="color:red">correctness</span> of a concurrent program depends on the specific execution.

Note that this is an "<span style="color:red">application-level</span>" property!

I.e., for a given program p, to answer the question "is p free from race conditions?" we must have access to the specification of p.

# SC for data-race-free programs, again

- For Java programs, we have SC for programs <span style="color:red">without data races</span>

- Presence of race conditions does not rob us of SC – important to know (the difference between) the two definitions

- What about the semantics of programs *with* data races?
  - Will not be considered here
  - In e.g. C++ data races result in undefined behavior (see C++ specification or https://en.cppreference.com/w/cpp/language/memory_model)
  - Java is supposed to be a "safe language", some guarantees (e.g. out-of-thin-air safety)

# Outline

- What are memory models?

- Why weak memory models?

- Something about the Java memory model (as an example of a weak memory model)

- Programming in the Java memory model (as an example of programming in a weak memory model)

# Practice?

- But what does this mean in practice?

- I.e: How does "weak memory models" affect my daily life as a programmer?

- Answer: You must "annotate" your program more (compared to CS). "Annotations" in the form of variable qualifiers, synchronization mechanisms etc.

- Essentially annotating which things are shared and which are not

# Simple example

Finally, an example!!!

```
bool done = false;

t1 {
    done = true;
}

t2 {
    if (done) print(33);
}
```

- Does this program contain
  - data races?
  - race conditions?

- Data race = yes, done is accessed without synchronization and one of the accesses is a write

- Race condition = depends on the specification we are to satisfy (what it means for the program to be correct)

- (Note: Difficult to reason about race conditions (correctness) because we cannot assume SC because we have data races!)

# Simple example

Finally, an example!!!

```
bool done = false;

t1 {
    done = true;
}

t2 {
    if (done) print(33);
}
```

- Wait a minute!

- Are you telling me there's a problem in this program?

- From a SC perspective, everything is fine!

- No atomicity problems or anything like that... but visibility problems!

# Simple example (fixed)

Finally, an example!!!

```
volatile bool done = false;

t1 {
    done = true;
}

t2 {
    if (done) print(33);
}
```

- Solution: Annotate your program. E.g., in Java `volatile` is considered synchronization.

- Does this program contain
  - data races?
  - race conditions?

- Data race = no, in Java `volatile` accesses are considered synchronized

- Race condition = ???, still depends on specification

# Simple example (fixed)

Finally, an example!!!

```
volatile bool done = false;

t1 {
    done = true;
}

t2 {
    if (done) print(33);
}
```

Example specification:

- Spec = "If the program outputs something, it must output 33"

- (In other words: Spec = "Output nothing or 33")

- Race conditions w.r.t. above specification?

- No race conditions! (As correct output does not depend on specific "execution"/ interleaving.)

# Simple example (fixed)

Finally, an example!!!

```
volatile bool done = false;

t1 {
    done = true;
}


t2 {
    if (done) print(33);
}
```

Example specification:

- Spec = "The program outputs 33"

- Race conditions w.r.t. above specification?

- Yes, have race condition. Some interleavings give us correct output, others do not.

# Similar example, with locks

```
lock lock = new lock();
int id = 0;

t1 {
    lock.lock();
    id++;
    lock.unlock();
}

t2 {
    print(id);
}
```
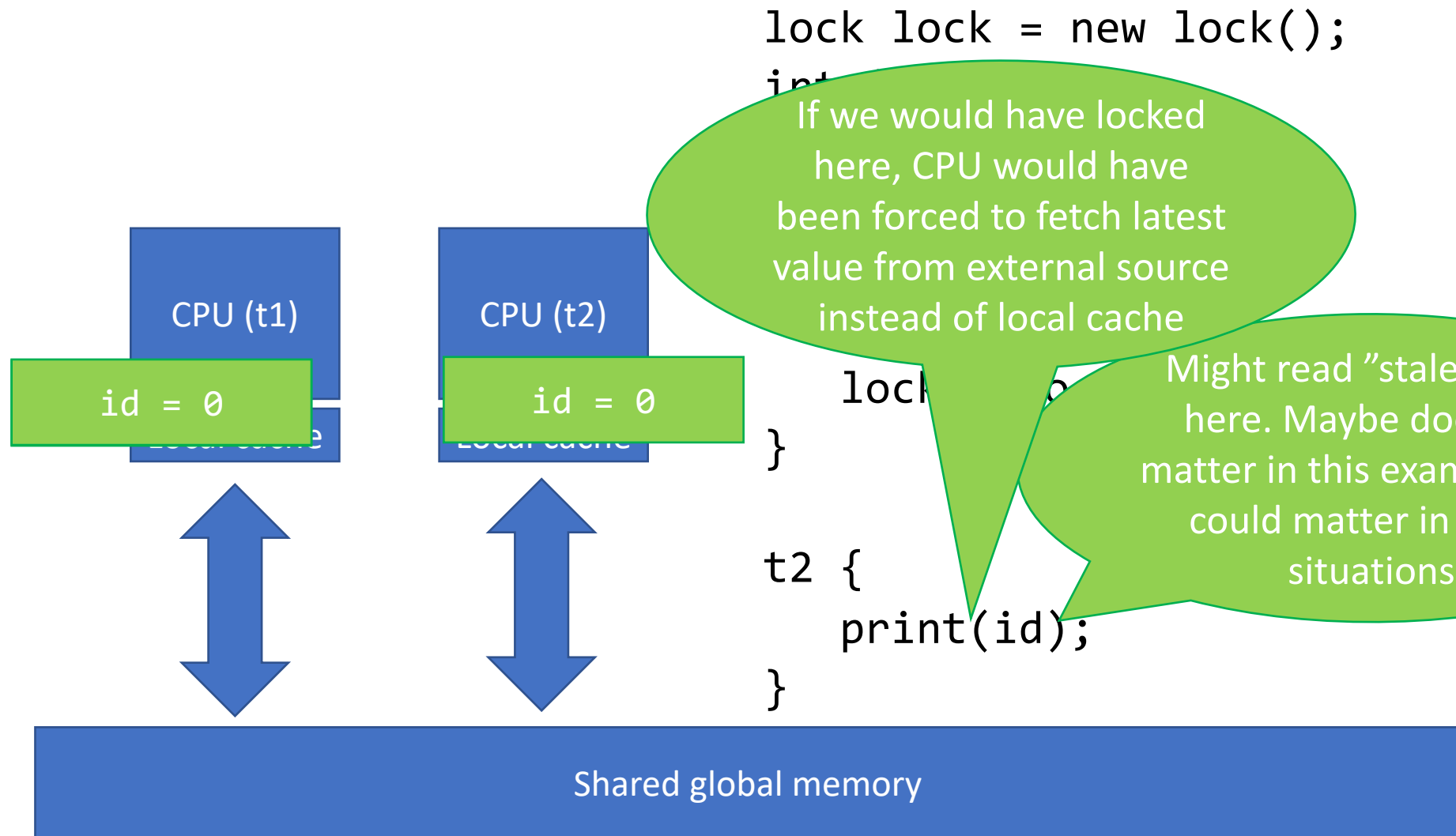
Data races?

We have a race! All accesses to the shared variable done must be synchronized!

Here we have (again) atomicity, but not: visibility

# `id` flag might exist as multiple copies...

```
lock lock = new lock();
```

CPU (t1)

CPU (t2)

id = 0

id = 0

If we would have locked here, CPU would have been forced to fetch latest value from external source instead of local cache

```
        lock
    }
```

Might read "stale" value here. Maybe does not matter in this example, but could matter in other situations

```
    t2 {
        print(id);
    }
```

Shared global memory

NOTE: Everything on this slide simplified, and makes unsound assumptions about JVM implementation details

# Similar example, with locks (fixed)

```
lock lock = new lock();
int id = 0;

t1 {
    lock.lock();
    id++;
    lock.unlock();
}

t2 {
    lock.lock(); // new
    print(id);
    lock.unlock(); // new
}
```

This is how the program would look like with proper annotations/synchronization

No data races in sight!

```
      Consumer c1 = new Consumer(q);
      Consumer c2 = new Consumer(q);
      new Thread(p).start();
      new Thread(c1).start();
      new Thread(c2).start();
   }
 }
```

Memory consistency effects: As with other concurrent collections, actions in a thread prior to placing an object into a `BlockingQueue` *happen-before* actions subsequent to the access or removal of that element from the `BlockingQueue` in another thread.

This interface is a member of the Java Collections Framework.

**Since:**

1.5

## Method Summary

| **All Methods** | **Instance Methods** | **Abstract Methods** |

| Modifier and Type | Method | Description |
| --- | --- | --- |
| boolean | **add**(**E** e) | Inserts the specified element into this queue if it is possible to do so immediately without violating capacity restrictions, returning `true` upon success and throwing an |

# Another example

```java
int x = 1;

x = 2;
// What can be printed?
Thread t = new Thread(() ->
 System.out.println(x));
t.start();
```

- Data race because t reads x without synchronization?

- (Could potentially argue read and write not overlapping in any CS execution.)

- More detailed reasoning principle: x write *happens-before* x read (remember screenshot on previous slide…)

# Reading suggestions

- See *Java Concurrency in Practice* (2006) if you want more of this. The book presents simplified rules you can follow to do concurrent programming in Java instead of having to learn the details of the Java memory model.

- E.g., the book provides useful "safe publication idioms"

- Also e.g.: Hans-J. Boehm, "Threads cannot be implemented as a library" (2005). (https://doi.org/10.1145/1065010.1065042)

- Also e.g.: Hans-J. Boehm and Sarita V. Adve, "You don't know jack about shared variables or memory models" (2012). (https://doi.org/10.1145/2076450.2076465)

# Summary?

- Make sure to not have data races in your Java programs

- One way to think about all of this: Atomicity *and* <span style="color:red">visibility</span>

- Visibility aspect new in weak memory models compared to SC!

# If you only will remember one thing, please:

In concurrent programming in Java, not only do we have to consider atomicity, we also must consider visibility!

visibility          visibility

visibility     visibility                    visibility

v i s i b i l i t y