# Close Encounters of the Java Memory Model Kind
Aleksey Shipilev

Presented by Sankalp Gambhir for the LAMP PL Seminar, 14 Nov 2023

**This is a summary of a long-running discussion**

Thanks to:

**What is a memory model?**

- Spec describes an abstract machine model

## What is a memory model?

- Spec describes an abstract machine model
- The model has to subsume memory interactions

## What is a memory model?

- Spec describes an abstract machine model
- The model has to subsume memory interactions
- This description is the memory model

## What is a memory model?

- Spec describes an abstract machine model
- The model has to subsume memory interactions
- This description is the memory model
- Can be largely disconnected from the rest of the model

## What is a memory model?

- Spec describes an abstract machine model
- The model has to subsume memory interactions
- This description is the memory model
- Can be largely disconnected from the rest of the model
- But forms its foundations

## Sequential Code

```scala
// main thread
var x: Int = 0
val a0 = x
x = 1
val a1 = x
x = 2
val a2 = x
x = 3
val a3 = x
x = 4
val a4 = x
x = 5
val a5 = x
(a0, a1, a2, a3, a4, a5)
// (0, 1, 2, 3, 4, 5)
```

## Sequential Consistency

```
1 // t0
2 x = 1
3 y = 1
4 println(x)
```

```
1 // t1
2 x = 2
3 println(y)
4 y = 2
5 println(x)
```

## Sequential Consistency

```
1 // t0
2 x = 1
3 y = 1
4 println(x)
```

```
1 // t1
2 x = 2
3 println(y)
4 y = 2
5 println(x)
```

The result should be an interleaving of the events

## Sequential Consistency

```
1 // t0
2 x = 1
3 y = 1
4 println(x)
```

```
1 // t1
2 x = 2
3 println(y)
4 y = 2
5 println(x)
```

The result should be an interleaving of the events

As the name suggests, as close to emulating sequential behaviour as possible.

## Sequential Consistency

```
1 // t0              1 // t1
2 x = 1              2 x = 2
3 y = 1              3 println(y)
4 println(x)         4 y = 2
                     5 println(x)
```

The result should be an interleaving of the events

As the name suggests, as close to emulating sequential behaviour as possible.
Standard mental model of programmers too.

## Running a small test

```
1 var a: Boolean = false
2 var b: Boolean = false
3
```

```
1 // t0
2 a = true
3 if b then
4   0
5 else
6   1
7
```

```
1 // t1
2 b = true
3 if a then
4   0
5 else
6   1
7
```

Outputs?

## Running a small test

```
1 var a: Boolean = false
2 var b: Boolean = false
3
```

```
1 // t0
2 a = true
3 if b then
4   0
5 else
6   1
7
```

```
1 // t1
2 b = true
3 if a then
4   0
5 else
6   1
7
```

Outputs?

(0, 0)?

## Running a small test

```
1 var a: Boolean = false
2 var b: Boolean = false
3
```

```
1 // t0
2 a = true
3 if b then
4   0
5 else
6   1
7
```

```
1 // t1
2 b = true
3 if a then
4   0
5 else
6   1
7
```

Outputs?

(0, 0)? (1, 0)?

## Running a small test

```
1 var a: Boolean = false
2 var b: Boolean = false
3
```

```
1 // t0
2 a = true
3 if b then
4   0
5 else
6   1
7
```

```
1 // t1
2 b = true
3 if a then
4   0
5 else
6   1
7
```

Outputs?
(0, 0)? (1, 0)? (0, 1)?

## Running a small test

```
1 var a: Boolean = false
2 var b: Boolean = false
3
```

```
1 // t0
2 a = true
3 if b then
4    0
5 else
6    1
7
```

```
1 // t1
2 b = true
3 if a then
4    0
5 else
6    1
7
```

Outputs?
(0, 0)? (1, 0)? (0, 1)? (1, 1)?

```scala
 5   object NonSC extends TestHarness[Int]:
 6     class NonSCTest extends StressTest[Int]:
 7       // state
 8       var a: Boolean = false
 9       var b: Boolean = false
10
11       Test { () ⇒
12         a = true
13         if b then
14           0
15         else
16           1
17       }
18
19       Test { () ⇒
20         b = true
21         if a then
22           0
23         else
24           1
25       }
26     end NonSCTest
27
        ↑ testCount
28     val testCount: Int = 1000000
29
        run | debug
30     @main def runNonSC =
31       val res = test(NonSCTest(), 4)
32       println(res)
```

```
 5  object NonSC extends TestHarness[Int]:
 6      class NonSCTest extends StressTest[Int]:
 7          // state
 8          var a: Boolean = false
 9          var b: Boolean = false
10
11          Test { () =>
12              a = true
13              if b then
14                  0
15              else
16                  1
17          }
18
19          Test { () =>
20              b = true
21              if a then
22                  0
23              else
24                  1
25          }
26      end NonSCTest
27
        ↑ testCount
28      val testCount: Int = 1000000
29
        run | debug
30      @main def runNonSC =
31          val res = test(NonSCTest(), 4)
32          println(res)
```

```
Map(
    List(1,  0)  ->  499854,
    List(0,  1)  ->  500141,
    List(0,  0)  ->  1,
    List(1,  1)  ->  4
)
```

```scala
 5  object NonSC extends TestHarness[Int]:
 6      class NonSCTest extends StressTest[Int]:
 7          // state
 8          var a: Boolean = false
 9          var b: Boolean = false
10
11          Test { () ⇒
12              a = true
13              if b then
14                  0
15              else
16                  1
17          }
18
19          Test { () ⇒
20              b = true
21              if a then
22                  0
23              else
24                  1
25          }
26      end NonSCTest
27
         ↑ testCount
28      val testCount: Int = 1000000
29
     run | debug
30      @main def runNonSC =
31          val res = test(NonSCTest(), 4)
32          println(res)
```

```
Map(
    List(1,  0) -> 499854,
    List(0,  1) -> 500141,
    List(0,  0) -> 1,
    List(1,  1) -> 4
)
```

Only $\sim 0.0004\%$ cases!

*It ain't what you don't know that gets you into trouble. It's what you know for sure that just ain't so.*

*- Mark Twain*

## Just read the spec!

*The Java Memory Model is the most complicated part of Java spec that must be understood by at least library and runtime developers. Unfortunately, it is worded in such a way that it takes a few senior guys to decipher it for each other. Most developers, of course, are not using JMM rules as stated, and instead make a few constructions out of its rules, or worse, blindly copy the constructions from senior developers without understanding the limits of their applicability. [1]*

## Java Memory Model (compacted garbage summary)

- All hail happens-before (hb)

## Java Memory Model (compacted garbage summary)

- All hail happens-before (hb)
- Respect each thread's ordering — program order (po)

## Java Memory Model (compacted garbage summary)

- All hail happens-before (hb)
- Respect each thread's ordering — program order (po)
- Respect reads and writes — reads-from (rf)

## Java Memory Model (compacted garbage summary)

- All hail happens-before (hb)
- Respect each thread's ordering — program order (po)
- Respect reads and writes — reads-from (rf)
- Synchronization order (so) (the big one)

## Java Memory Model (compacted garbage summary)

- All hail happens-before (hb)
- Respect each thread's ordering — program order (po)
- Respect reads and writes — reads-from (rf)
- Synchronization order (so) (the big one)
  - volatile ordering (vl)

## Java Memory Model (compacted garbage summary)

- All hail happens-before (hb)
- Respect each thread's ordering — program order (po)
- Respect reads and writes — reads-from (rf)
- Synchronization order (so) (the big one)
    - volatile ordering (vl)
    - synchronizes-with (sw/so)
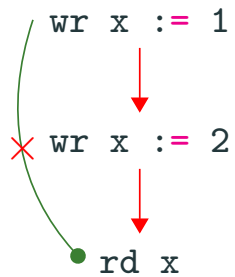
## Java Memory Model (compacted garbage summary)

- All hail happens-before (hb)
- Respect each thread's ordering — program order (po)
- Respect reads and writes — reads-from (rf)
- Synchronization order (so) (the big one)
    - volatile ordering (vl)
    - synchronizes-with (sw/so)
    - thread ordering (th)

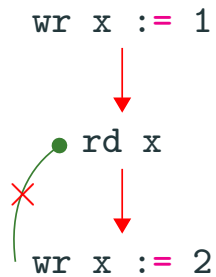## Java Memory Model (compacted garbage summary)

- All hail happens-before (hb)
- Respect each thread's ordering — program order (po)
- Respect reads and writes — reads-from (rf)
- Synchronization order (so) (the big one)
  - volatile ordering (vl)
  - synchronizes-with (sw/so)
  - thread ordering (th)

## Java Memory Model (compacted garbage summary)

- All hail happens-before (hb)
- Respect each thread's ordering — program order (po)
- Respect reads and writes — reads-from (rf)
- Synchronization order (so) (the big one)
    - volatile ordering (vl)
    - synchronizes-with (sw/so)
    - thread ordering (th)

Two conditions for a valid execution:

- no cycles
- no invalid reads

Legend

a —— po —— ▸ b    b is after a in the program order

a —— rf —— • b    b reads from a in an execution

a —— vl —— ◆ b    b is after a in an execution, and they are events on the same `volatile` variable

a —— so —— ▪ b    b is after a in an execution, and they are `synchronized` block start or end events

a —— th —— ▫ b    b is after a in an execution, and they are related as thread-start/first-event or as last-event/thread-join



https://lampepfl.github.io/
courses/cs206/
ex-08-solution-9475bbf99cf472ec48e4

https://lampepfl.github.io/
courses/cs206/ex-08

```
1 var a: Boolean = false
2 var b: Boolean = false
3 var x: Int = -1
4 var y: Int = -1
5
```

```
1 // t0
2 a = true
3 if b then
4    x = 0
5 else
6    x = 1
7
```

```
1 // t1
2 b = true
3 if a then
4    y = 0
5 else
6    y = 1
7
```

(wr y -1) (wr b f)  (wr a f) (wr x -1)

initialization

th0                         th1

(wr a t)                   (wr b t)

(rd b ??)                  (rd a ??)

(wr x ??)                  (wr y ??)

reads

(rd x ??) (rd y ??)

(wr y -1) (wr b f)  (wr a f) (wr x -1)

initialization

th0                          th1
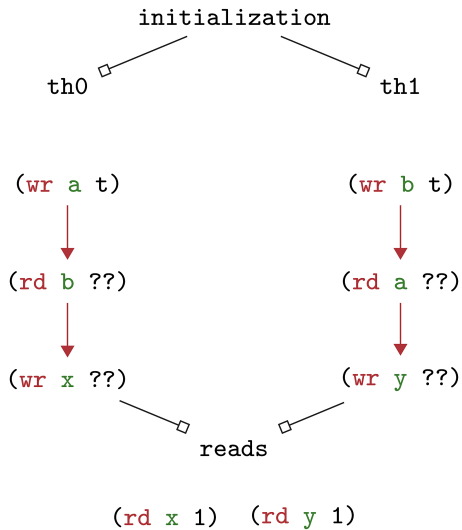

(wr a t)                    (wr b t)


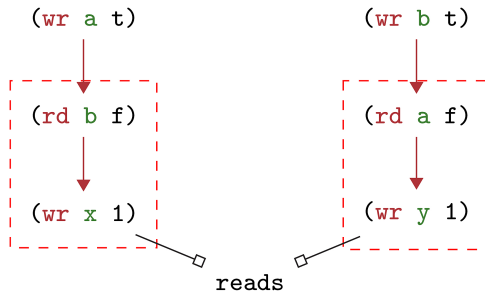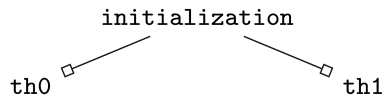(rd b ??)                   (rd a ??)


(wr x ??)                   (wr y ??)

reads
```
(rd x 1)  (rd y 1)
```

(wr y -1) (wr b f)  (wr a f) (wr x -1)

initialization

th0                                    th1

(wr a t)                          (wr b t)

(rd b f)                          (rd a f)

(wr x 1)                          (wr y 1)

reads

(rd x 1)   (rd y 1)
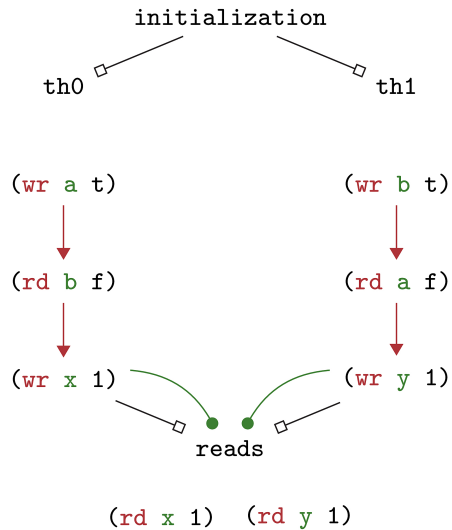
18

(wr y -1) (wr b f)  (wr a f) (wr x -1)



initialization

th0

th1

(wr a t)

(wr b t)

(rd b f)

(rd a f)

(wr x 1)

(wr y 1)

reads

(rd x 1)  (rd y 1)

**Myth: my program is compiled per the spec**

The spec is described on an abstract machine. The runtime only *emulates* it.

## Myth: my program is compiled per the spec

The spec is described on an abstract machine. The runtime only *emulates* it.

As long as no one can call your bluff, you're good.

## Myth: my program is compiled per the spec

The spec is described on an abstract machine. The runtime only *emulates* it.

As long as no one can call your bluff, you're good.

Swap statements? Remove entire blocks? As long as no one, *anywhere*, can differentiate the resulting effects, go ahead.

## Myth: my program is compiled per the spec

The spec is described on an abstract machine. The runtime only *emulates* it.

As long as no one can call your bluff, you're good.

Swap statements? Remove entire blocks? As long as no one, *anywhere*, can differentiate the resulting effects, go ahead.

Pain: students seem to *love* thinking about compiler reorderings. To the JMM spec, they don't exist. Students refuse to accept this.

```
1 var x: Int = 0
2 var y: Int = 0
```

```
1 // t0
2 synchronized {x = 1}
3 synchronized {y = 1}
```

```
1 // t1
2 val a = y
3 val b = x
4 println((a, b))
```

```
[OK] net.shipilev.jmm.LockCoarsening
(fork: #1, iteration #1, JVM args: [-server, -XX:+UnlockDiagnosticVMOptions, -XX:+StressLCM,
-XX:+StressGCM])
  Observed state    Occurrences           Expectation  Interpretation
            0, 0    43,558,372             ACCEPTABLE  All other cases are acceptable.
            0, 1        22,512             ACCEPTABLE  All other cases are acceptable.
            1, 0         1,565  ACCEPTABLE_INTERESTING  X and Y are visible in different order
            1, 1     1,372,341             ACCEPTABLE  All other cases are acceptable.
```

## "Commit to memory" as a mental model

```
1 var x: Int = 0
2 var y: Int = 0
3
```

```
1 // t0
2 x = 1
3
```

```
1 // t1
2 y = 1
3
```

```
1 // t2
2 val a = x
3 val b = y
4 (a, b)
5
```

```
1 // t3
2 val b = y
3 val a = x
4 (a, b)
5
```

## "Commit to memory" as a mental model

```
1 var x: Int = 0
2 var y: Int = 0
3
```

```
1 // t0
2 x = 1
3
```

```
1 // t1
2 y = 1
3
```

```
1 // t2
2 val a = x
3 val b = y
4 (a, b)
5
```

```
1 // t3
2 val b = y
3 val a = x
4 (a, b)
5
```

Possible to read t2: (x = 1, y = 0), t3: (x = 0, y = 1)!

22

## "Commit to memory" as a mental model

```
1 var x: Int = 0
2 var y: Int = 0
3
```

```
1 // t0              1 // t1              1 // t2              1 // t3
2 x = 1              2 y = 1              2 val a = x         2 val b = y
3                    3                    3 val b = y         3 val a = x
                                          4 (a, b)            4 (a, b)
                                          5                   5
```

Possible to read t2: (x = 1, y = 0), t3: (x = 0, y = 1)!

*Pain: lack of "multi-copy atomicity".* If many threads are reading a variable, their history of those reads do not need to be consistent in general. Hardware-dependent!

## "Fences" as a mental model

```
1 @volatile var greatBarrierReef: Int = 0
2 var x: Int = 0
3 var y: Int = 0
```

```
1 // t0
2 x = 1
3 y = 1
4 greatBarrierReef = 1
```

```
1 // t1
2 greatBarrierReef = 2
3 (x, y)
```

## "Fences" as a mental model

```
1 @volatile var greatBarrierReef: Int = 0
2 var x: Int = 0
3 var y: Int = 0
```

```
1 // t0                                    1 // t1
2 x = 1                                    2 greatBarrierReef = 2
3 y = 1                                    3 (x, y)
4 greatBarrierReef = 1
```

Surely it's not possible to read x = 0, y = 1...?

## "Fences" as a mental model

```scala
1 @volatile var greatBarrierReef: Int = 0
2 var x: Int = 0
3 var y: Int = 0
```

```scala
1 // t0                          1 // t1
2 x = 1                          2 greatBarrierReef = 2
3 y = 1                          3 (x, y)
4 greatBarrierReef = 1
```

Surely it's not possible to read $x = 0$, $y = 1$...? Exactly 1 out of 1M 😜

## "Separability" as a mental model

*If I write **my** code well, I'm good!*

- Separate memory locations you touch
- Carefully analyze critical sections

## "Separability" as a mental model

*If I write* **my** *code well, I'm good!*

- Separate memory locations you touch
- Carefully analyze critical sections

## "Separability" as a mental model

```
1 var x: Int = 0 // your territory
2 var y: Int = 0 // your nailbiting neighbour's
```

```
1 // t0
2 y = 1
3 x = 1
4 val a = y
5 val c = x
6
```

```
1 // t1
2 x = 2
3 y = 2
4 val b = y
5 val d = x
6
```

## "Separability" as a mental model

```
1 var x: Int = 0 // your territory
2 @volatile var y: Int = 0 // neighbour felt unsafe
```

```
1 // t0
2 y = 1
3 x = 1
4 val a = y
5 val c = x
6
```

```
1 // t1
2 x = 2
3 y = 2
4 val b = y
5 val d = x
6
```

## "Separability" as a mental model

```scala
1 var x: Int = 0 // your territory
2 @volatile var y: Int = 0 // neighbour felt unsafe
```

```scala
1 // t0
2 y = 1
3 x = 1
4 val a = y
5 val c = x
6
```

```scala
1 // t1
2 x = 2
3 y = 2
4 val b = y
5 val d = x
6
```

You are now implicitly attached to an abstract, and *very fragile*, global state!

```
1 var a: Long = 0L
2
```

```
1 // t0
2 a = 70000000000L
3
```

```
1 // t1
2 a = 80000000000L
3 println(a)
4
```

Outputs?

## Not everything is an `Int`

```
1 var a: Long = 0L
2
```

```
1 // t0
2 a = 70000000000L
3
```

```
1 // t1
2 a = 80000000000L
3 println(a)
4
```

Outputs? 0?

**Not everything is an** `Int`

```
1 var a: Long = 0L
2
```

```
1 // t0
2 a = 70000000000L
3
```

```
1 // t1
2 a = 80000000000L
3 println(a)
4
```

Outputs? 0? 70000000000L?

## Not everything is an `Int`

```
1 var a: Long = 0L
2
```

```
1 // t0
2 a = 70000000000L
3
```

```
1 // t1
2 a = 80000000000L
3 println(a)
4
```

Outputs? 0? 70000000000L? 80000000000L?

## Not everything is an `Int`

```
1 var a: Long = 0L
2
```

```
1 // t0
2 a = 70000000000L
3
```

```
1 // t1
2 a = 80000000000L
3 println(a)
4
```

Outputs? 0? 70000000000L? 80000000000L? 78589934592L . . . ?

## Not everything is an `Int`

```
1 var a: Long = 0L
2
```

```
1 // t0                              1 // t1
2 a = 70000000000L                   2 a = 80000000000L
3                                     3 println(a)
                                      4
```

Outputs? 0? 70000000000L? 80000000000L? 78589934592L ...?

*Other technical baggage: access atomicity v memory ordering and the burden of volatility*

## Horror Circus: sync on strings

```
1 // t0
2 "Lock".synchronized {
3     x = x + 1
4 }
```

```
1 // t1
2 "Lock".synchronized {
3     x = x + 1
4 }
```

More horrific examples on:

Aleksey Shipilev. **Close Encounters of The Java Memory Model Kind.** 2016. URL: https://shipilev.net/blog/2016/close-encounters-of-jmm-kind/

**Print everything!**

```
1036
1037   /**
1038    * Prints an Object and then terminate the line.  This method calls
1039    * at first String.valueOf(x) to get the printed object's string value
1040    * then behaves as
1041    * though it invokes {@link #print(String)} and then
1042    * {@link #println()}.
1043    *
1044    * @param x  The {@code Object} to be printed.
1045    */
1046   public void println(Object x) {
1047       String s = String.valueOf(x);
1048       if (getClass() == PrintStream.class) {
1049           // need to apply String.valueOf again since first invocation
1050           // might return null
1051           writeln(String.valueOf(s));
1052       } else {
1053           synchronized (this) {
1054               print(s);
1055               newLine();
1056           }
1057       }
1058   }
1059
```

## Conclusions

- Not much
- Take what you can
- Don't write concurrent code (yourself)

## References

[1] Aleksey Shipilev. **Java Memory Model Pragmatics (transcript).** 2014. URL: https://shipilev.net/blog/2014/jmm-pragmatics/.

[2] Aleksey Shipilev. **Close Encounters of The Java Memory Model Kind.** 2016. URL: https://shipilev.net/blog/2016/close-encounters-of-jmm-kind/.