

Script generated by TTT

Title: Petter: Programmiersprachenh
(24.10.2018)

Date: Wed Oct 24 14:09:04 CEST 2018

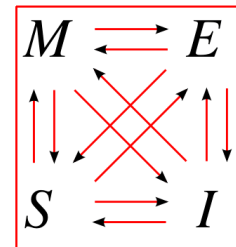
Duration: 94:44 min

Pages: 16

The MESI Protocol: States [?]

Processors use caches to avoid a costly round-trip to RAM for every memory access.

- programs often access the same memory area repeatedly (e.g. stack)
- keeping a local mirror image of certain memory regions requires bookkeeping about who has the latest copy



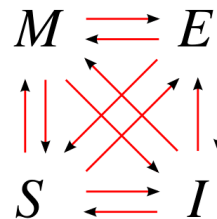
Each cache line is in one of the states M, E, S, I :

The MESI Protocol: Messages



Moving data between caches is coordinated by sending messages [McK10]:

- **Read**: sent if CPU needs to read from an address
- **Read Response**: when in state E or S, response to a **Read** message, carries the data for the requested address
- **Invalidate**: asks others to evict a cache line
- **Invalidate Acknowledge**: reply indicating that a cache line has been evicted
- **Read Invalidate**: like **Read** + **Invalidate** (also called “read with intend to modify”)
- **Writeback**: **Read Response** when in state M, as a side effect noticing main memory about modifications to the cacheline, changing sender’s state to S



We mostly consider messages between processors. Upon **Read Invalidate**, a processor replies with **Read Response/Writeback** before the **Invalidate Acknowledge** is sent.

MESI Example



Consider how the following code might execute:

Thread A

```
a = 1; // A.1
b = 1; // A.2
```

Thread B

```
while (b == 0) {}; // B.1
assert(a == 1); // B.2
```

- in all examples, the initial values of variables are assumed to be 0
- suppose that a and b reside in different cache lines
- assume that a cache line is larger than the variable itself
- we write the content of a cache line as
 - ▶ M_x : modified, with value x
 - ▶ E_x : exclusive, with value x
 - ▶ S_x : shared, with value x
 - ▶ I: invalid

MESI Example (I)



```

Thread A
a = 1; // A.1
b = 1; // A.2
    
```

```

Thread B
while (b == 0) {}; // B.1
assert(a == 1); // B.2
    
```

state- ment	CPU A		CPU B		RAM		message
	a	b	a	b	a	b	
A.1					0	0	<i>read invalidate</i> of a from CPU A
					0	0	<i>invalidate ack.</i> of a from CPU B
					0	0	<i>read response</i> of a=0 from RAM
B.1	M1				0	0	<i>read</i> of b from CPU B
	M1				0	0	<i>read response</i> with b=0 from RAM
B.1	M1			E0	0	0	
A.2	M1			E0	0	0	<i>read invalidate</i> of b from CPU A
	M1			E0	0	0	<i>read response</i> of b=0 from CPU B
	M1	S0		S0	0	0	<i>invalidate ack.</i> of b from CPU B
	M1	M1			0	0	

MESI Example (II)



```

Thread A
a = 1; // A.1
b = 1; // A.2
    
```

```

Thread B
while (b == 0) {}; // B.1
assert(a == 1); // B.2
    
```

state- ment	CPU A		CPU B		RAM		message
	a	b	a	b	a	b	
B.1	M1	M1			0	0	<i>read</i> of b from CPU B
	M1	M1			0	0	<i>write back</i> of b=1 from CPU A
B.2	M1	S1		S1	0	1	<i>read</i> of a from CPU B
	S1	S1	S1	S1	1	1	<i>write back</i> of a=1 from CPU A
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
A.1	S1	S1	S1	S1	1	1	<i>invalidate</i> of a from CPU A
	S1	S1		S1	1	1	<i>invalidate ack.</i> of a from CPU B
	M1	S1		S1	1	1	

MESI Example (II)



```

Thread A
a = 1; // A.1
b = 1; // A.2
    
```

```

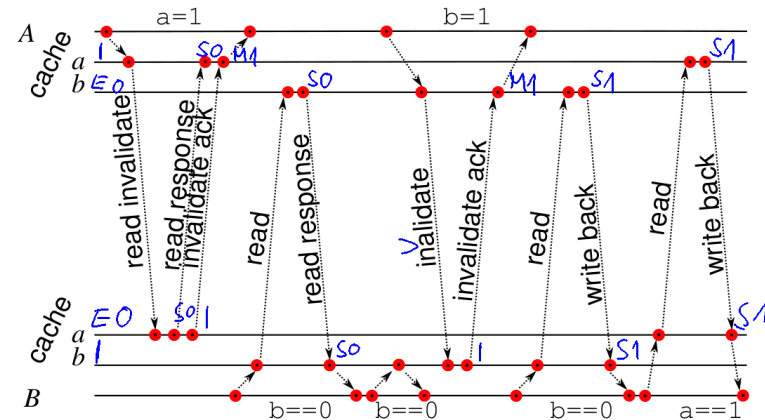
Thread B
while (b == 0) {}; // B.1
assert(a == 1); // B.2
    
```

state- ment	CPU A		CPU B		RAM		message
	a	b	a	b	a	b	
B.1	M1	M1			0	0	<i>read</i> of b from CPU B
	M1	M1			0	0	<i>write back</i> of b=1 from CPU A
B.2	M1	S1		S1	0	1	<i>read</i> of a from CPU B
	M1	S1		S1	0	1	<i>write back</i> of a=1 from CPU A
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
A.1	S1	S1	S1	S1	1	1	<i>invalidate</i> of a from CPU A
	S1	S1		S1	1	1	<i>invalidate ack.</i> of a from CPU B
	M1	S1		S1	1	1	

MESI Example: Happened Before Model



Idea: each cache line one process, A caches b=0 as E, B caches a=0 as E



Observations:

- each memory access must complete before executing next instruction
- ↪ add edge

Summary: MESI cache coherence protocol



Sequential consistency:

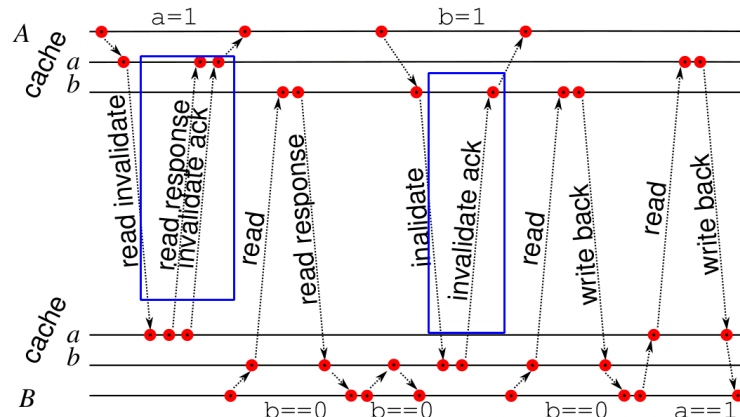
- a characterization of well-behaved programs
- a model for differing speed of execution
- for fixed paths through the threads *and* a total order between accesses to the same variable: executions can be illustrated by happened-before diagram with one process per variable
- MESI cache coherence protocol ensures SC for processors with caches

Out-of-Order Execution



⚠ performance problem: writes always stall

Thread A	Thread B
<code>a = 1; // A.1</code>	<code>while (b == 0) {}; // B.1</code>
<code>b = 1; // A.2</code>	<code>assert(a == 1); // B.2</code>

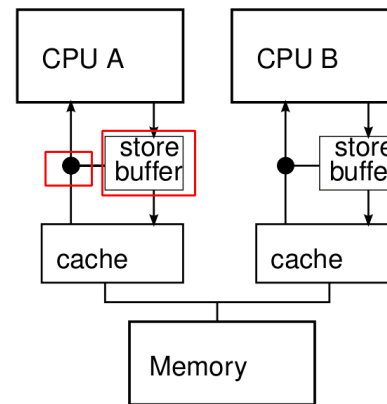


Introducing Store Buffers: Out-Of-Order Stores

Store Buffers and Total Store Ordering [?]



Goal: continue execution after *cache-miss* write operation



- put each write into a *store buffer* and trigger fetching of cache line
- once a cache line has arrived, apply relevant writes
 - ▶ today, a store buffer is always a *queue* [OSS09]
 - ▶ two writes to the same location are not merged

⚠ sequential consistency *per CPU* is violated unless

- ▶ each read checks store buffer before cache
 - ▶ on hit, return the youngest value that is waiting to be written
- ↔ TSO

Excursion : non-FIFO store buffers (↔ Sparc/PSO)

TSO Model: Formal Spec [?]



Definition (Total Store Order)

1 The store order wrt. memory (\sqsubseteq) is total

$$\forall a, b \in \text{addr } i, j \in \text{CPU} \quad (St_i[a] \sqsubseteq St_j[b]) \vee (St_j[b] \sqsubseteq St_i[a])$$

2 Stores in program order (\leq) are embedded into the memory order (\sqsubseteq)

$$St_i[a] \leq St_j[b] \Rightarrow St_i[a] \sqsubseteq St_j[b]$$

3 Loads preceding an other operation (wrt. program order \leq) are embedded into the memory order (\sqsubseteq)

$$Ld_i[a] \leq Op_j[b] \Rightarrow Ld_i[a] \sqsubseteq Op_j[b]$$

4 A load's value is determined by the latest write as observed by the local CPU

$$val(Ld_i[a]) = val(St_j[a] \mid St_j[a] \sqsubseteq Ld_i[a] \mid \max_{\sqsubseteq} \{St_k[a] \mid St_k[a] \sqsubseteq Ld_i[a]\} \cup \{St_i[a] \mid St_i[a] \leq Ld_i[a]\})$$

Particularly, one ordering property is not guaranteed:

$$St_i[a] \leq Ld_i[b] \not\Rightarrow St_i[a] \sqsubseteq Ld_i[b]$$

⚠ Local stores may be observed earlier by local loads than from somewhere else!

↪ What about sequential consistency for the whole system?

Explicit Synchronization: Write Barrier



Overtaking of messages *is desirable* and should not be prohibited in general.

- store buffers render programs incorrect that assume sequential consistency between *different* CPUs
- whenever two stores in one CPU must appear *in sequence at a different CPU*, an explicit *write barrier* has to be inserted
 - a write barrier marks all current store operations in the store buffer
 - the next store operation is only executed when all marked stores in the buffer have completed
- x86 CPUs provide the `sfence` instruction
- a write barrier after each write gives sequentially consistent CPU behavior (and is as slow as a CPU without store buffer)

↪ use (write) barriers only when necessary

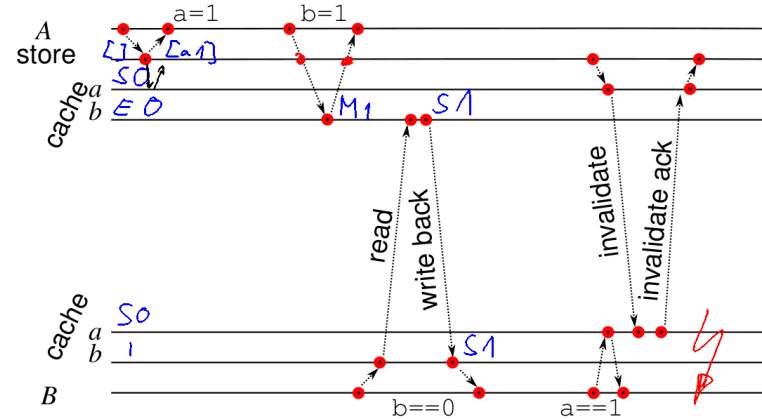
Happened-Before Model for Store Buffers



```
Thread A
a = 1;
b = 1;
```

```
Thread B
while (b == 0) {};
assert(a == 1);
```

Assume cache A contains: a: S0, b: E0, cache B contains: a: S0, b: I



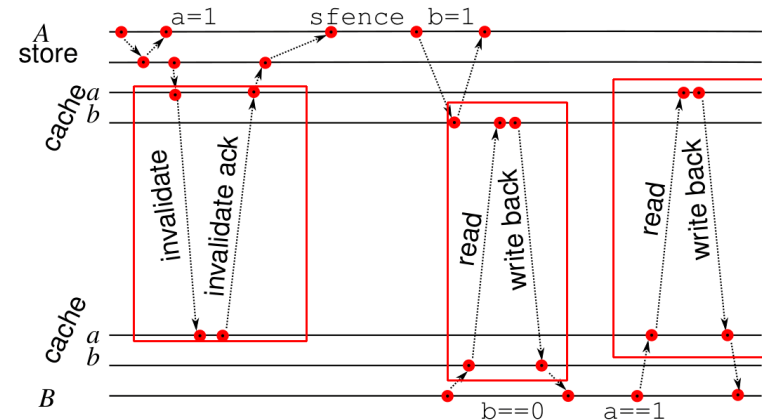
Happened-Before Model for Write Barriers



```
Thread A
a = 1;
sfence();
b = 1;
```

```
Thread B
while (b == 0) {};
assert(a == 1);
```

Assume cache A contains: a: S0, b: E0, cache B contains: a: S0, b: I



Introducing Invalidate Queues: O-o-O Reads