

Script generated by TTT

Title: Petter: Programmiersprachen (23.10.2019)

Date: Wed Oct 23 12:30:52 CEST 2019

Duration: 76:13 min

Pages: 27

MESI Example



Consider how the following code might execute:

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>

- 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
 - ▶ Mx: modified, with value x
 - ▶ Ex: exclusive, with value x
 - ▶ Sx: shared, with value x
 - ▶ I: invalid

MESI Example

Consider how the following code might execute:

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>

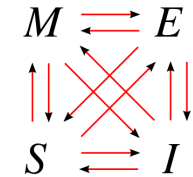
- in all examples, the initial values of variables are assumed to be 0
- suppose that `a` and `b` reside in different cache lines
- assume that a cache line is larger than the variable itself
- we write the content of a cache line as
 - ▶ Mx: modified, with value x
 - ▶ Ex: exclusive, with value x
 -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

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	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>

- 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
 - ▶ Mx: modified, with value x
 - ▶ Ex: exclusive, with value x
 - ▶ Sx: shared, with value x
 - ▶ I: invalid

MESI Example

Consider how the following code might execute:

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>

- 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
 - ▶ Mx: modified, with value x
 - ▶ Ex: exclusive, with value x
 - ▶ Sx: shared, with value x
 - ▶ I: invalid



MESI Example (I)

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>

statement	CPU A		CPU B		RAM		message
	a	b	a	b	a	b	
A.1	I	I	I	I	0	0	read invalidate of a from CPU A invalidate ack. of a from CPU B read response of a=0 from RAM
	I	I	I	I	0	0	
	I	I	I	I	0	0	
B.1	M1	I	I	I	0	0	read of b from CPU B read response with b=0 from RAM
	M1	I	I	I	0	0	
B.1	M1	I	I	E0	0	0	read invalidate of b from CPU A read response of b=0 from CPU B invalidate ack. of b from CPU B
A.2	M1	I	I	E0	0	0	
	M1	I	I	E0	0	0	
	M1	S0	I	S0	0	0	
	M1	M1	I	I	0	0	



MESI Example (I)

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>

statement	CPU A		CPU B		RAM		message
	a	b	a	b	a	b	
A.1	I	I	I	I	0	0	read invalidate of a from CPU A invalidate ack. of a from CPU B read response of a=0 from RAM
	I	I	I	I	0	0	
	I	I	I	I	0	0	
B.1	M1	I	I	I	0	0	read of b from CPU B read response with b=0 from RAM
	M1	I	I	I	0	0	
B.1	M1	I	I	E0	0	0	read invalidate of b from CPU A read response of b=0 from CPU B invalidate ack. of b from CPU B
A.2	M1	I	I	E0	0	0	
	M1	I	I	E0	0	0	
	M1	S0	I	S0	0	0	
	M1	M1	I	I	0	0	



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
```

statement	CPU A		CPU B		RAM		message
	a	b	a	b	a	b	
A.1	I	I	I	I	0	0	} read invalidate of a from CPU A } invalidate ack. of a from CPU B } read response of a=0 from RAM
	I	I	I	I	0	0	
	I	I	I	I	0	0	
B.1	M1	I	I	I	0	0	} read of b from CPU B } read response with b=0 from RAM
	M1	I	I	I	0	0	
B.1 A.2	M1	I	I	E0	0	0	} read invalidate of b from CPU A } read response of b=0 from CPU B } invalidate ack. of b from CPU B
	M1	I	I	E0	0	0	
	M1	I	I	E0	0	0	
	M1	S0	I	S0	0	0	
	M1	M1	I	I	0	0	

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
```

statement	CPU A		CPU B		RAM		message
	a	b	a	b	a	b	
A.1	I	I	I	I	0	0	} read invalidate of a from CPU A } invalidate ack. of a from CPU B } read response of a=0 from RAM
	I	I	I	I	0	0	
	I	I	I	I	0	0	
B.1	M1	I	I	I	0	0	} read of b from CPU B } read response with b=0 from RAM
	M1	I	I	I	0	0	
B.1 A.2	M1	I	I	E0	0	0	} read invalidate of b from CPU A } read response of b=0 from CPU B } invalidate ack. of b from CPU B
	M1	I	I	E0	0	0	
	M1	I	I	E0	0	0	
	M1	S0	I	S0	0	0	
	M1	M1	I	I	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
```

statement	CPU A		CPU B		RAM		message
	a	b	a	b	a	b	
B.1	M1	M1	I	I	0	0	} read of b from CPU B } write back of b=1 from CPU A
	M1	M1	I	I	0	0	
B.2	M1	S1	I	S1	0	1	} read of a from CPU B } write back of a=1 from CPU A
	M1	S1	I	S1	0	1	
...
A.1	S1	S1	S1	S1	1	1	} invalidate of a from CPU A } invalidate ack. of a from CPU B
	S1	S1	I	S1	1	1	
	M1	S1	I	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
```

statement	CPU A		CPU B		RAM		message
	a	b	a	b	a	b	
B.1	M1	M1	I	I	0	0	} read of b from CPU B } write back of b=1 from CPU A
	M1	M1	I	I	0	0	
B.2	M1	S1	I	S1	0	1	} read of a from CPU B } write back of a=1 from CPU A
	M1	S1	I	S1	0	1	
...
A.1	S1	S1	S1	S1	1	1	} invalidate of a from CPU A } invalidate ack. of a from CPU B
	S1	S1	I	S1	1	1	
	M1	S1	I	S1	1	1	

MESI Example (II)

```

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

```

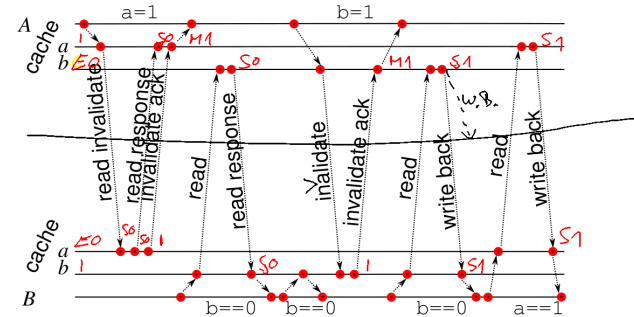
Thread B
while (b == 0) {}; // B.1
assert(a == 1); // B.2
    
```

statement	CPU A		CPU B		RAM		message
	a	b	a	b	a	b	
B.1	M1	M1	I	I	0	0	} read of b from CPU B
	M1	M1	I	I	0	0	
B.2	M1	S1	I	S1	0	1	} write back of b=1 from CPU A
	M1	S1	I	S1	0	1	
A.1	S1	S1	S1	S1	1	1	} read of a from CPU B
	S1	S1	S1	S1	1	1	
A.2	S1	S1	I	S1	1	1	} write back of a=1 from CPU A
	M1	S1	I	S1	1	1	
A.1	S1	S1	S1	S1	1	1	} invalidate of a from CPU A
	S1	S1	I	S1	1	1	
A.2	S1	S1	I	S1	1	1	} invalidate ack. of a from CPU B
	M1	S1	I	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

Introducing Store Buffers: Out-Of-Order Stores

Out-of-Order Execution

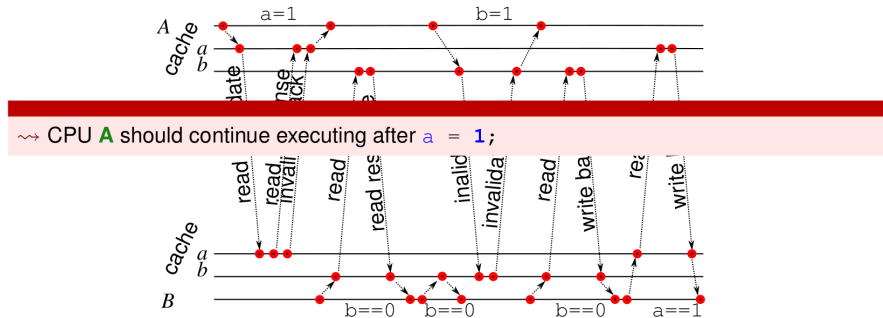
⚠ performance problem: writes always stall

```

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

```

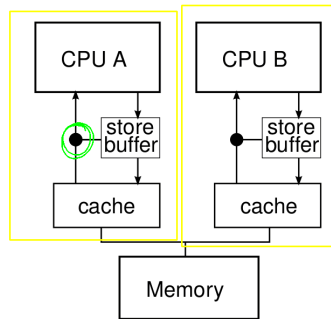
Thread B
while (b == 0) {}; // B.1
assert(a == 1); // B.2
    
```



Store Buffers



⚠ *Abstract Machine Model*: defines semantics of memory accesses

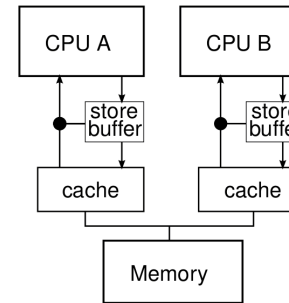


- put *each* store into a *store buffer* and continue execution
- Store buffers apply stores in various orders:
 - ▶ FIFO (Sparc/x86-*TSO*)
 - ▶ unordered (Sparc *PSO*)
- ⚠ program order still needs to be observed locally
 - ▶ store buffer snoops read channel and
 - ▶ on matching address, returns the youngest value in buffer

Store Buffers



⚠ *Abstract Machine Model*: defines semantics of memory accesses

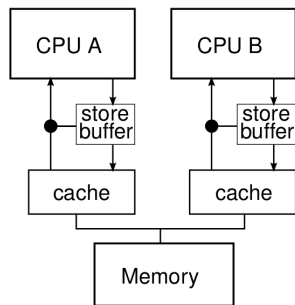


- put *each* store into a *store buffer* and continue execution
- Store buffers apply stores in various orders:
 - ▶ FIFO (Sparc/x86-*TSO*)
 - ▶ unordered (Sparc *PSO*)
- ⚠ program order still needs to be observed locally
 - ▶ store buffer snoops read channel and
 - ▶ on matching address, returns the youngest value in buffer

Store Buffers



⚠ *Abstract Machine Model*: defines semantics of memory accesses



- put *each* store into a *store buffer* and continue execution
- Store buffers apply stores in various orders:
 - ▶ FIFO (Sparc/x86-*TSO*)
 - ▶ unordered (Sparc *PSO*)
- ⚠ program order still needs to be observed locally
 - ▶ store buffer snoops read channel and
 - ▶ on matching address, returns the youngest value in buffer

TSO Model: Formal Spec [SI92]



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)

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

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

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

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

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

Particularly, one ordering property is not guaranteed:

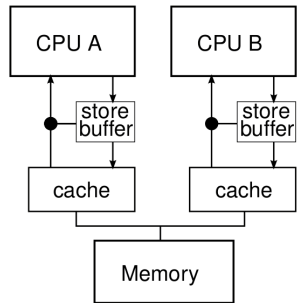
$$St_i[a] \leq Ld_i[b] \not\Rightarrow St_i[a] \sqsubseteq Ld_i[b]$$

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

Store Buffers



⚠ *Abstract Machine Model*: defines semantics of memory accesses

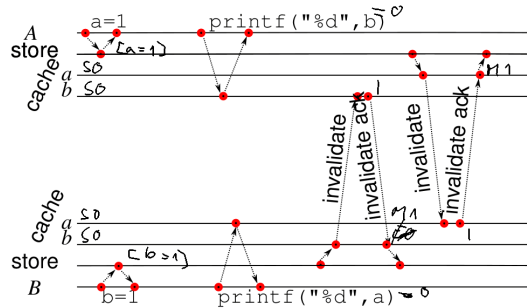


- put *each* store into a *store buffer* and continue execution
- Store buffers apply stores in various orders:
 - ▶ FIFO (Sparc/x86-*TSO*)
 - ▶ unordered (Sparc *PSO*)
- ⚠ program order still needs to be observed locally
 - ▶ store buffer snoops read channel and
 - ▶ on matching address, returns the youngest value in buffer

Happened-Before Model for TSO



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



TSO Model: Formal Spec [SI92]



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)

$$st_i[a] \leq st_i[b] \Rightarrow st_i[a] \sqsubseteq st_i[b]$$

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

$$ld_i[a] \leq op_i[b] \Rightarrow ld_i[a] \sqsubseteq op_i[b]$$

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

$$val(ld_i[a]) = val(st_j[a] \mid st_j[a] \Rightarrow \max_{\sqsubseteq} (\{st_k[a] \mid st_k[a] \sqsubseteq ld_i[a]\} \cup \{st_i[a] \mid st_i[a] \leq ld_i[a]\}))$$

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!

TSO in the Wild: x86



The x86 CPU, powering desktops and servers around the world is a common representative of a TSO Memory Model based CPU.

- FIFO store buffers keep quite strong consistency properties
- The major obstacle to Sequential Consistency is

$$st_i[a] \leq ld_i[b] \not\Rightarrow st_i[a] \sqsubseteq ld_i[b]$$

- ▶ modern x86 CPUs provide the `mfence` instruction
- ▶ `mfence` orders all memory instructions:

$$op_i \leq mfence() \leq op_i' \Rightarrow op_i \sqsubseteq op_i'$$

- a fence between write and loads gives sequentially consistent CPU behavior (and is as slow as a CPU without store buffer)

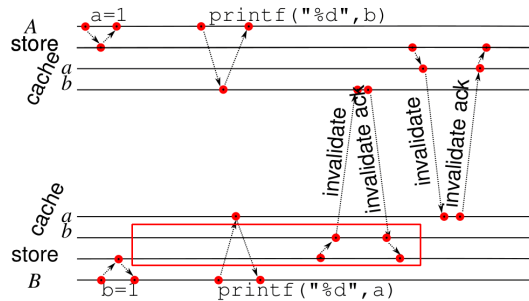
~ use fences only when necessary

Happened-Before Model for TSO



Thread A	Thread B
<pre>a = 1; printf("%d", b);</pre>	<pre>b = 1; printf("%d", a);</pre>

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



TSO in the Wild: x86



The x86 CPU, powering desktops and servers around the world is a common representative of a TSO Memory Model based CPU.

- FIFO store buffers keep quite strong consistency properties
- The major obstacle to Sequential Consistency is

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

- ▶ modern x86 CPUs provide the `m fence` instruction
- ▶ `m fence` orders all memory instructions:

$$Op_i \leq mfence() \leq Op_i' \Rightarrow Op_i \sqsubseteq Op_i'$$

- a fence between write and loads gives sequentially consistent CPU behavior (and is as slow as a CPU without store buffer)
- ↪ use fences only when necessary

PSO Model: Formal Spec [S192]



Definition (Partial 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 Fenced stores in program order (\leq) are embedded into the memory order (\sqsubseteq)

$$St_i[a] \leq sfence() \leq St_i[b] \Rightarrow St_i[a] \sqsubseteq St_i[b]$$
- 3 Stores to the same address in program order (\leq) are embedded into the memory order (\sqsubseteq)

$$St_i[a] \leq St_i[a'] \Rightarrow St_i[a] \sqsubseteq St_i[a']$$
- 4 Loads preceding another operation (wrt. program order \leq) are embedded into the memory order (\sqsubseteq)

$$Ld_i[a] \leq Op_i[b] \Rightarrow Ld_i[a] \sqsubseteq Op_i[b]$$
- 5 A load's value is determined by the latest write as observed by the local CPU

$$val(Ld_i[a]) = val(St_j[a] \mid St_j[a] \sqsubseteq_{\text{CPU } i} \{St_k[a] \mid St_k[a] \sqsubseteq Ld_i[a]\} \cup \{St_i[a] \mid St_i[a] \leq Ld_i[a]\})$$

⚠ Now also stores are not guaranteed to be in order any more:

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

↪ What about sequential consistency for the whole system?

Explicit Synchronization: Write Barrier



Overtaking of messages *may be desirable* and does not need to be prohibited in general.

- generalized store buffers render programs incorrect that assume sequential consistency between *different* CPUs
- whenever a store in front of another operation in one CPU must be observable in this order *by 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