## Beyond Weak Memory Consistency: The Challenges of Memory Persistency

Part I: Low-Level Persistency Models

Azalea Raad

Imperial College London

Viktor Vafeiadis MPI-SWS

azalea@imperial.ac.uk





## Computer Storage





## Computer Storage



✓ fastX volatile





## Computer Storage



✓ fastX volatile





X slow √ persistent

#### What is Non-Volatile Memory (NVM)?





#### What is Non-Volatile Memory (NVM)?



#### NVM: Hybrid Storage + Memory

Best of both worlds:

- ✓ *persistent* (like HDD)
- √ fast, random access (like RAM)



# INTEL® OPTANETM TECHNOLOGY









Non-Volatile memory // x = 0 x := 1 // x = 1 // x = 1 // recovery routine // x = 1



## A: Program Verification





Sequential (1940s)





#### Weak Memory Consistency (WMC)

**No** total execution order (*to*)  $\Rightarrow$ 

weak behaviour absent under SC, caused by:

- instruction *reordering* by compiler
- write propagation across *cache hierarchy*

## WMC: Store Buffering

## WMC: Store Buffering

## WMC: Store Buffering

#### Weak Memory Consistency (WMC)

**No** total execution order (*to*)  $\Rightarrow$ 

weak behaviour absent under SC, caused by:

- instruction *reordering* by compiler
- write propagation across *cache hierarchy*

#### Weak Memory Consistency (WMC)













**!!** Execution continues *ahead of persistence* 

- asynchronous persists



- !! Writes may persist out of order
  - *relaxed* persists



## the *order* in which writes are *made visible* to other threads





#### **Consistency Model**

the *order* in which writes are *made visible* to other threads

#### **Persistency Model**

the *order* in which writes are *persisted* to NVM

// x=

!! E>

 $\mathbf{V}$ 



#### **Consistency Model**

the *order* in which writes are *made visible* to other threads

#### **Persistency Model**

the **order** in which writes are **persisted** to NVM

// x=

!! Ex

#### **NVM Semantics** Consistency + Persistency Model

### Outline

- 1. An *intuitive* account of Intel-x86 persistency: Px86
  - ✤ Warmup: Sequential Px86
  - + Concurrent Px86
- 2. A *formal* account Px86: *operational* semantics
- 3. A *formal* account Px86: *declarative* semantics
- 4. Other Low-level (hardware) persistency models
- 5. Further reading

## 1. An *intuitive* account of Px86

## Warmup: *Sequential* Px86





x:=1 : adds x:=1 to memory



x:=1 : adds x:=1 to memory

a:=x : reads x from memory



x:=1 : adds x:=1 to memory

a:=x : reads x from memory




x:=1 : adds x:=1 to memory

a:=x : reads x from memory









x := 1 : adds x := 1 to p-buffer





- x := 1 : adds x := 1 to p-buffer
- a:=x : if p-buffer contains x, reads latest entry else reads from memory





CPU

(Volatile) Memory

read

write

- x := 1 : adds x := 1 to p-buffer
- a:=x : if p-buffer contains x, reads latest entry else reads from memory







CPU

(Volatile) Memory

read

write

- x := 1 : adds x := 1 to p-buffer
- a:=x : if p-buffer contains x, reads latest entry else reads from memory
  - p-buffer lost; memory retained

unbuffer\* : p-buffer to memory



Unbuffered at *non-deterministic* points in time Buffering & unbuffering orders may disagree

\* at non-deterministic times

### Handling Relaxed Persists



**!! out of order persists** 

explicit persists?

### Explicit Persists: Desiderata



**!! out of order persists** 

explicit persists?

Strength (ordering constraints)

Performance



#### \* clwb and clflushopt: same ordering constraints



- \* clwb and clflushopt: same ordering constraints
- \* clwb does not invalidate cache line
- \* clflushopt invalidates cache line



- \* clwb and clflushopt: same ordering constraints
- \* clwb does not invalidate cache line
- \* clflushopt invalidates cache line
- Clflush: strongest ordering constraints; invalidates cache line

### Strong (Synchronous) Explicit Persists: clflush



Weak (Asynchronous) Explicit Persists: clflushopt & clwb



#### Weak (Asynchronous) Explicit Persists: clflushopt & clwb



### Solution: Persist Sequences



### Solution: Persist Sequences



### Solution: Persist Sequences



Waits until earlier writes on x are persisted
Disallows reordering

✓ synchronous persists
✓ no out of order persists

## 1. An *intuitive* account of Px86

### *Concurrent* Px86





x := 1 : adds x := 1 to buffer



x := 1 : adds x := 1 to buffer

unbuffer\* : buffer to memory



- x := 1 : adds x := 1 to buffer
- unbuffer\* : buffer to memory
  - a:=x : if buffer contains x, reads latest entry else reads from memory

\* at non-deterministic times



- x := 1 : adds x := 1 to buffer
- unbuffer\* : buffer to memory
  - a:=x : if buffer contains x, reads latest entry else reads from memory



buffer and memory lost

# Software WMC: Store Buffering






































### Hardware (Intel x86) WMC: Store Buffering





### Hardware (Intel x86) WMC: Store Buffering





#### Px86: Persistent & Concurrent x86





#### Px86: Persistent & Concurrent x86











buffer/unbuffer order: *consistency* model





# 2. A *formal* account of Px86

## **Operational** Semantics

#### **Basic domains**

- $a \in REG$  Registers
- $v \in VAL$  Values
- $\tau \in TID$  Thread IDs

#### **Basic domains**

- $a \in REG$  Registers
- $v \in VAL$  Values
- $\tau \in TID$  Thread IDs

#### Expressions and sequential commands

 $Exp \ni e ::= v \mid a \mid e+e \mid \cdots$   $PCOM \ni c ::= load(x) \mid store(x, e) \mid CAS(x, e, e') \mid FAA(x, e)$   $\mid mfence \mid sfence \mid flush_{opt} x \mid flush x$   $COM \ni C ::= e \mid c \mid let a := C in C$   $\mid if (C) then C else C \mid repeat C$ 

#### **Basic domains**

- $a \in REG$  Registers
- $v \in VAL$  Values
- $\tau \in TID$  Thread IDs

#### Expressions and sequential commands

 $Exp \ni e ::= v \mid a \mid e+e \mid \cdots$   $PCOM \ni c ::= load(x) \mid store(x, e) \mid CAS(x, e, e') \mid FAA(x, e)$   $\mid mfence \mid sfence \mid flush_{opt} x \mid flush x$   $COM \ni C ::= e \mid c \mid let a := C in C$   $\mid if (C) then C else C \mid repeat C$ 

used for *persistency* 

#### **Basic domains**

- $a \in REG$  Registers
- $v \in VAL$  Values
- $\tau \in TID$  Thread IDs

#### Expressions and sequential commands

 $Exp \ni e ::= v \mid a \mid e+e \mid \cdots$   $PCOM \ni c ::= load(x) \mid store(x, e) \mid CAS(x, e, e') \mid FAA(x, e)$   $\mid mfence \mid sfence \mid flush_{opt} x \mid flush x$   $COM \ni C ::= e \mid c \mid let a := C in C$   $\mid if (C) then C else C \mid repeat C$ 

used for *persistency* 

**Programs**  $P \in PROG \triangleq TID \xrightarrow{fin} COM$ 

# Px86 Operational Semantics = Px86 **Program** Transitions + Px86 **Storage** Transitions

- ➡ First formulated by Raad et al. [2]
- ➡ Later simplified by Khyzha and Lahav [3]

**Thread transitions:** Com  $\xrightarrow{\text{TID:Lab} \cup \{\epsilon\}}$  Com

**Thread transitions:** Com  $\xrightarrow{\text{TID:Lab} \cup \{\epsilon\}}$  Com

$$\boxed{\frac{1}{10ad(x) \xrightarrow{\tau:(\mathsf{R},x,v)} v}} (\text{T-READ})$$

$$\left(\frac{1}{\operatorname{store}(x,v)} \xrightarrow{\tau:(\mathbb{W},x,v)} v\right)$$
 (T-WRITE)

**Thread transitions:** Com  $\xrightarrow{\text{TID:Lab} \cup \{\epsilon\}}$  Com



**Thread transitions:** Com  $\xrightarrow{\text{TID:LAB} \cup \{\epsilon\}}$  Com



**Thread transitions:** Com  $\xrightarrow{\text{TID:Lab} \cup \{\epsilon\}}$  Com



**Thread transitions:** Com  $\xrightarrow{\text{TID:Lab} \cup \{\epsilon\}}$  Com



**Thread transitions:** Com  $\xrightarrow{\text{TID:LAB} \cup \{\epsilon\}}$  Com



let a:=
$$v$$
 in  $C \xrightarrow{\tau:\epsilon} C[v/a]$  (1-LET2

**Thread transitions:** COM  $\xrightarrow{\text{TID:LAB} \cup \{\epsilon\}}$  COM



**Thread transitions:** Com  $\xrightarrow{\text{TID:Lab} \cup \{\epsilon\}}$  Com



**Thread transitions:** Com  $\xrightarrow{\text{TID:LAB} \cup \{\epsilon\}}$  Com

 $LAB \triangleq \{(\mathsf{R}, x, v), (\mathsf{W}, x, v), (\mathsf{U}, x, v, v'), \mathsf{MF}, \mathsf{SF}, (\mathsf{FO}, x), (\mathsf{FL}, x) \mid x \in \mathrm{Loc} \land v, v' \in \mathrm{VAL}\}$ 

**Program transitions:** PROG  $\xrightarrow{\text{TID:LAB} \cup \{\epsilon\}}$  PROG

$$\begin{array}{c} \begin{array}{c} P(\tau) \xrightarrow{\tau:l} C \\ \hline P \xrightarrow{\tau:l} P \xrightarrow{\tau:l} P[\tau \mapsto C] \end{array} \end{array} (PROG) \end{array} \end{array}$$





 $M \in Mem \triangleq Loc \xrightarrow{fin} Val$ 



$$M \in Mem \triangleq Loc \xrightarrow{fin} Val$$
$$B \in BMap \triangleq TID \xrightarrow{fin} Buff$$



buffer/unbuffer order: consistency model → execution reordering

$$M \in Mem \triangleq Loc \xrightarrow{fin} Val$$
$$B \in BMap \triangleq TID \xrightarrow{fin} Buff$$

Later in Program Order

				1	2	3	4	5	6	7
Inread I Inread2	Read Write RMW mfer	mfence	stence	flush <sub>opt</sub>	flush					
$\begin{array}{c c} \bullet \bullet$	rde	Α	Read			<b>√</b>	<ul> <li>Image: A set of the set of the</li></ul>	$\checkmark$		
	О О	B	Write	×	<b>√</b>	1	<ul> <li>Image: A set of the set of the</li></ul>	~	sloc	<ul> <li>Image: A start of the start of</li></ul>
Buffer Buffer	gran	С	RMW	<ul> <li>Image: A set of the set of the</li></ul>	<b>√</b>	<b>√</b>	<b>~</b>		<ul> <li>Image: A set of the set of the</li></ul>	<ul> <li>Image: A set of the set of the</li></ul>
	rog	D	mfence	<ul> <li>Image: A start of the start of</li></ul>	<b>√</b>	<b>√</b>	<b>√</b>		<b>√</b>	<b>√</b>
Persistence Buffer	ier in F	E	sfence	×	<b>√</b>	<b>√</b>	<b>√</b>		<b>√</b>	<b>√</b>
		F	flush <sub>opt</sub>	×	×	1	<b>√</b>		×	sloc
	larl	G	flush	×	1	1	<b>√</b>		sloc	
buffer/unbuffer order: consistency model → execution reordering	Ur	aer p	reserved	Y	: yes	. IIC	) 5100			Joanon
fin										

 $M \in Mem \triangleq Loc \xrightarrow{\min} Val$ 

 $\mathsf{B} \in \mathsf{BMap} \triangleq \mathsf{TId} \xrightarrow{\mathsf{fin}} \mathsf{Buff}$ 



buffer/unbuffer order: consistency model → execution reordering

	Later in Program Order									
			1	2	3	4	5	6	7	
ram Order			Read	Write	RMW	mfence	sfence	<b>flush</b> opt	flush	
	Α	Read	<ul> <li>Image: A set of the set of the</li></ul>	<b>√</b>	<b>√</b>	<b>√</b>	<b>√</b>	<b>√</b>	<ul> <li>Image: A start of the start of</li></ul>	
	B	Write	×	<b>√</b>	<b>√</b>	<b>√</b>	<ul> <li>Image: A set of the set of the</li></ul>	sloc	<ul> <li>Image: A set of the set of the</li></ul>	
	С	RMW	<ul> <li>Image: A set of the set of the</li></ul>	<b>√</b>	<b>√</b>	<b>√</b>	<b>√</b>	<b>√</b>	<ul> <li>Image: A start of the start of</li></ul>	
rog	D	mfence	<ul> <li>Image: A set of the set of the</li></ul>	✓ ✓	<b>√</b>	<b>√</b>	<b>√</b>	✓ ✓	\ \	
in P	E	sfence	×		<b>√</b>	<b>√</b>	<b>√</b>			
ier	F	<b>flush</b> opt	×	×	<b>~</b>			×	sloc	
larl	G	flush	×	<b>√</b>		<b>√</b>	<b>√</b>	sloc	<ul> <li>Image: A set of the set of the</li></ul>	
Order preserved? ✓: yes X: no sloc : iff on the same location ★ Reads reordered before writes/sfence/flush <sub>opt</sub> /flush → delay writes/sfence/flush <sub>opt</sub> /flush execution in thread buffers										

 $M \in Mem \triangleq Loc \xrightarrow{fin} Val$ 

 $B \in BMAP \triangleq TID \xrightarrow{\text{fin}} BUFF$  $b \in BUFF \triangleq Seq \left\langle \left\{ (W, x, v), (FL, x), (F0, x), SF \mid x \in Loc \land v \in VAL \right\} \right.$ 



buffer/unbuffer order: consistency model → *execution reorderin* 

		Later in Program Order									
				1	2	3	4	5	6	7	
	5			Read	Write	RMW	mfence	sfence	<b>flush</b> opt	flush	
	rde	Α	Read	<b>√</b>		<b>√</b>	1	<b>√</b>	<b>√</b>	<b>√</b>	
	0 u	B	Write	×		<b>√</b>	<b>√</b>	<b>√</b>	sloc	<b>√</b>	
	ran	C	RMW	<ul> <li>Image: A set of the set of the</li></ul>	<b>~</b>	<b>√</b>	<b>√</b>	<b>~</b>	<ul> <li>Image: A set of the set of the</li></ul>	<ul> <li>Image: A set of the set of the</li></ul>	
	rog	D	mfence	<b>~</b>		<b>√</b>	<b>√</b>	<b>√</b>	<b>√</b>	<b>√</b>	
	in F	E	sfence	X	<b>√</b>	<b>√</b>	<b>√</b>	<b>√</b>	<b>√</b>	$\checkmark$	
	ier	F	<b>flush</b> opt	×	×	<b>√</b>	<b>√</b>	<b>√</b>	×	sloc	)
	Earl	G	flush	X		-			sloc		
g	Order preserved?       ✓: yes       X: no       sloc: iff on the same location         ◆ Reads reordered before writes/sfence/flush <sub>opt</sub> /flush       →       delay writes/sfence/flush <sub>opt</sub> /flush execution in thread buffers         ◆ flush <sub>opt</sub> reordered w.r.t. writes/flush <sub>opt</sub> /flush on diff. locations       →       their buffer/unbuffer orders (in thread buffers) can disagree         ◆ writes/flush/sfence ordered w.r.t. one another       →       their buffer/unbuffer orders agree (FIFO)										

an in Duamana Andar

 $M \in Mem \triangleq Loc \xrightarrow{fin} Val$ 

 $B \in BMAP \triangleq TID \xrightarrow{\text{fin}} BUFF$  $b \in BUFF \triangleq Seq \left\langle \left\{ (W, x, v), (FL, x), (F0, x), SF \mid x \in Loc \land v \in VAL \right\} \right.$ 



buffer/unbuffer order: persistency model → *persist reordering* 



- Persisting writes may be *delayed*
- Writes on different locations persist in different orders
  - → per-location persist buffers
  - → record & delay writes in pbuff

buffer/unbuffer order: persistency model → *persist reordering* 

 $PB \in PBMAP \triangleq Loc \xrightarrow{fin} PBUFF$  $pb \in PBUFF \triangleq Seq \left\langle \left\{ w(v) \qquad | v \in VAL \right\} \right.$ 



buffer/unbuffer order: persistency model → *persist reordering* 

- Persisting writes may be *delayed*
- Writes on different locations persist in different orders
  - $\rightarrow$  per-location persist buffers
  - → record & delay writes in pbuff
- Is flush executed synchronously
   → no need to delay/record them in pbuff

 $PB \in PBMAP \triangleq Loc \xrightarrow{fin} PBUFF$  $pb \in PBUFF \triangleq Seq \left\langle \left\{ w(v) \qquad \middle| v \in VAL \right\} \right.$ 



buffer/unbuffer order: persistency model → *persist reordering* 

- Persisting writes may be *delayed*
- Writes on different locations persist in different orders
  - → per-location persist buffers
  - → record & delay writes in pbuff
- Is flush executed synchronously
   → no need to delay/record them in pbuff
- Isolation flushopt executed asynchronously
   → record & delay them in pbuff

 $PB \in PBMAP \triangleq Loc \xrightarrow{fin} PBUFF$  $pb \in PBUFF \triangleq Seq \left\langle \left\{ w(v), fo(\tau) \mid v \in VAL \land \tau \in TID \right\} \right.$
# Px86 Storage System



 $PB \in PBMAP \triangleq Loc \xrightarrow{fin} PBUFF$  $pb \in PBUFF \triangleq Seq \left\langle \left\{ w(v), fo(\tau) \mid v \in VAL \land \tau \in TID \right\} \right.$ 

# Px86 Storage System



 $\mathsf{pb} \in \mathsf{PBuff} \triangleq \mathsf{Seq}\left\langle \left\{ \mathsf{w}(v), \mathsf{fo}(\tau) \mid v \in \mathsf{VAl} \land \tau \in \mathsf{TId} \right\} \right.$ 

simplification due to Khyzha & Lahav [3]

# Px86 Storage System



$$\mathsf{pb} \in \mathsf{PBuff} \triangleq \mathsf{Seq}\left\langle \left\{ \mathsf{w}(v), \mathsf{fo}(\tau) \mid v \in \mathsf{VAL} \land \tau \in \mathsf{TID} \right\} \right.$$

simplification due to Khyzha & Lahav [3]

Original model by Raad et al. [2] : one pbuff for all locations





$B(\tau)=b$	(M-WRITE)
$\overline{M,PB,B} \xrightarrow{\tau:(W,x,v)} M,PB,B[\tau \mapsto b.(W,x,v)]$	(IVI- VV RITE)











$$\mathsf{rd}(\mathsf{M},\mathsf{PB},\mathsf{b},x) \triangleq \begin{cases} v & \text{if } \exists S_1, S_2. \ \mathsf{b}=S_1.(\mathsf{W},x,v).S_2 \land \forall v'. \ (\mathsf{W},x,v') \notin S_2 \\ v & \text{else if } \exists S_1, S_2. \ \mathsf{PB}(x)=S_1.\mathsf{w}(v).S_2 \land \forall v'. \ \mathsf{w}(v') \notin S_2 \\ \mathsf{M}(x) & \text{otherwise} \end{cases}$$

(M-FL)





$$\mathsf{rd}(\mathsf{M},\mathsf{PB},\mathsf{b},x) \triangleq \begin{cases} v & \text{if } \exists S_1, S_2. \ \mathsf{b}=S_1.(\mathsf{W},x,v).S_2 \land \forall v'. \ (\mathsf{W},x,v') \notin S_2 \\ v & \text{else if } \exists S_1, S_2. \ \mathsf{PB}(x)=S_1.\mathsf{w}(v).S_2 \land \forall v'. \ \mathsf{w}(v') \notin S_2 \\ \mathsf{M}(x) & \text{otherwise} \end{cases}$$





 $\rightarrow$  ensure no flush<sub>opt</sub> in pbuff





$$\frac{B(\tau) = (W, x, v).b' \quad PB(x) = pb}{M, PB, B \xrightarrow{\tau:\epsilon} M, PB[x \mapsto pb.w(v)], B[\tau \mapsto b']} \quad (M-PROPW)$$

- \* writes/flush/sfence ordered w.r.t. one another  $\rightarrow$  their buffer/unbuffer orders agree (FIFO)
- ✤ Persisting writes may be *delayed* → record & delay writes in pbuff



$$\frac{B(\tau) = (W, x, v).b' \quad PB(x) = pb}{M, PB, B \xrightarrow{\tau:\epsilon} M, PB[x \mapsto pb.w(v)], B[\tau \mapsto b']} \quad (M-PROPW)$$

$$\frac{B(\tau) = (FL, x).b' \quad PB(x) = \epsilon}{M, PB, B \xrightarrow{\tau:\epsilon} M, PB, B[\tau \mapsto b']} \quad (M-PROPFL)$$

- \* writes/flush/sfence ordered w.r.t. one another  $\rightarrow$  their buffer/unbuffer orders agree (FIFO)
- ✤ Persisting writes may be *delayed* → record & delay writes in pbuff
- \* flush executed *synchronously*  $\rightarrow$  no need to delay/record them in pbuff



- \* writes/flush/sfence ordered w.r.t. one another  $\rightarrow$  their buffer/unbuffer orders agree (FIFO)
- ✤ Persisting writes may be *delayed* → record & delay writes in pbuff
- \* flush executed synchronously  $\rightarrow$  no need to delay/record them in pbuff
- Ilushopt + sfence/mfence/RMW = persist sequence → ensure no flushopt in pbuff



- ✤ writes/flush/sfence ordered w.r.t. one another → their buffer/unbuffer orders agree (FIFO)
- ◆ Persisting writes may be *delayed* → record & delay writes in pbuff
- \* flush executed synchronously  $\rightarrow$  no need to delay/record them in pbuff
- Ilushopt + sfence/mfence/RMW = persist sequence → ensure no flushopt in pbuff
- \* flush<sub>opt</sub> executed *asynchronously*  $\rightarrow$  record & delay them in pbuff
- \* flush<sub>opt</sub> reordered w.r.t. writes/flush<sub>opt</sub>/flush on diff. locations  $\rightarrow$  their buffer/unbuffer orders can disagree

## Px86 Storage Transitions: Delayed Persists



$$\frac{\mathsf{PB}(x) = \mathsf{w}(v).\mathsf{pb}}{\mathsf{M}, \mathsf{PB}, \mathsf{B} \xrightarrow{\tau:\epsilon} \mathsf{M}[x \mapsto v], \mathsf{PB}[x \mapsto \mathsf{pb}], \mathsf{B}} \quad (\mathsf{M}\text{-}\mathsf{PersistW})$$

$$\frac{\mathsf{PB}(x) = \mathsf{fo}(\tau).\mathsf{pb}}{\mathsf{M}, \mathsf{PB}, \mathsf{B} \xrightarrow{\tau:\epsilon} \mathsf{M}, \mathsf{PB}[x \mapsto \mathsf{pb}], \mathsf{B}} \quad (\mathsf{M}\text{-}\mathsf{PersistFO})$$

**Program transitions:** PROG  $\xrightarrow{\text{TID:LAB} \cup \{\epsilon\}}$  PROG

**Storage transitions:** Mem × PBMap × BMap  $\xrightarrow{\text{TID:Lab} \cup \{\epsilon\}}$  Mem × PBMap × BMap

**Program transitions:** PROG  $\xrightarrow{\text{TID:LAB} \cup \{\epsilon\}}$  PROG

**Storage transitions:** Mem × PBMap × BMap  $\xrightarrow{\text{TID:Lab} \cup \{\epsilon\}}$  Mem × PBMap × BMap

**Operational semantics:**  $PROG \times MEM \times PBMAP \times BMAP \Rightarrow PROG \times MEM \times PBMAP \times BMAP$ 

**Program transitions:** PROG  $\xrightarrow{\text{TID:LAB} \cup \{\epsilon\}}$  PROG

**Storage transitions:** Mem × PBMAP × BMAP  $\xrightarrow{\text{TId:Lab} \cup \{\epsilon\}}$  Mem × PBMAP × BMAP

**Operational semantics:**  $PROG \times MEM \times PBMAP \times BMAP \Rightarrow PROG \times MEM \times PBMAP \times BMAP$ 

 $\frac{P \xrightarrow{\tau:\epsilon} P'}{P, M, PB, B \Rightarrow P', M, PB, B}$ (SilentP)

**Program transitions:** Prog  $\xrightarrow{\text{TID:LAB} \cup \{\epsilon\}}$  Prog

**Storage transitions:** MEM × PBMAP × BMAP  $\xrightarrow{\text{TID:LAB} \cup \{\epsilon\}}$  MEM × PBMAP × BMAP

**Operational semantics:**  $PROG \times MEM \times PBMAP \times BMAP \Rightarrow PROG \times MEM \times PBMAP \times BMAP$ 

$$\frac{P \xrightarrow{\tau:\epsilon} P'}{P, M, PB, B \Rightarrow P', M, PB, B}$$
(SilentP)

 $\frac{M, PB, B \xrightarrow{\tau:\epsilon} M', PB', B'}{P, M, PB, B \Rightarrow P, M', PB', B'}$ (Silents)

**Program transitions:** Prog  $\xrightarrow{\text{TID:LAB} \cup \{\epsilon\}}$  Prog

**Storage transitions:** MEM × PBMAP × BMAP  $\xrightarrow{\text{TID:LAB} \cup \{\epsilon\}}$  MEM × PBMAP × BMAP

**Operational semantics:**  $PROG \times MEM \times PBMAP \times BMAP \Rightarrow PROG \times MEM \times PBMAP \times BMAP$ 

$$\frac{P \xrightarrow{\tau:\epsilon} P'}{P, M, PB, B \Rightarrow P', M, PB, B}$$
(SilentP)

 $\begin{array}{c} \underbrace{M, PB, B \xrightarrow{\tau:\epsilon} M', PB', B'}_{P, M, PB, B \Rightarrow P, M', PB', B'} (SilentS) \end{array}$ 

$$\begin{pmatrix} \underline{P \xrightarrow{\tau:l} P' & M, PB, B \xrightarrow{\tau:l} M', PB', B'} \\ P, M, PB, B \Rightarrow P', M', PB', B' \end{pmatrix} (STEP)$$

# 3. A *formal* account of Px86

## Declarative Semantics

Represent program behaviours as a set of *consistent executions (graphs)*

- Represent program behaviours as a set of *consistent executions (graphs)*
- An execution graph is a tuple:  $\langle E, po, rf, mo \rangle$
- *E* is the set of *events* (graph nodes), including initialisation writes

- Represent program behaviours as a set of *consistent executions (graphs)*
- An execution graph is a tuple: < E, po, rf, mo >
- ✤ E is the set of events (graph nodes), including initialisation writes
  - each event if of the form (n, τ, l)
    unique event id
    thread id

- Represent program behaviours as a set of *consistent executions (graphs)*
- An execution graph is a tuple: < E, po, rf, mo >
- ✤ E is the set of events (graph nodes), including initialisation writes
  - each event if of the form  $(n, \tau, l)$ unique event id

thread id

l event label: W(x, v), R(x, v), U(x, v, v'), MF, SF, FL(x), FO(x)

- Represent program behaviours as a set of *consistent executions (graphs)*
- An execution graph is a tuple:  $\langle E, po, rf, mo \rangle$
- ✤ E is the set of events (graph nodes), including initialisation writes
  - each event if of the form  $(n, \tau, l)$

unique event id I thread id

l event label: W(x, v), R(x, v), U(x, v, v'), MF, SF, FL(x), FO(x)

- Represent program behaviours as a set of consistent executions (graphs)
- An execution graph is a tuple:  $\langle E, po, rf, mo \rangle$
- ✤ E is the set of events (graph nodes), including initialisation writes
  - each event if of the form  $(n, \tau, l)$

unique event id thread id

l event label: W(x, v), R(x, v), U(x, v, v'), MF, SF, FL(x), FO(x)





- Represent program behaviours as a set of consistent executions (graphs)
- An execution graph is a tuple:  $\langle E, po, rf, mo \rangle$
- ✤ E is the set of events (graph nodes), including initialisation writes
- po, rf and mo are *relations* on events (graph edges)
  - po is the program order: strict total order on events of the same thread





- Represent program behaviours as a set of consistent executions (graphs)
- An execution graph is a tuple:  $\langle E, po, rf, mo \rangle$
- ✤ E is the set of events (graph nodes), including initialisation writes
- po, rf and mo are *relations* on events (graph edges)
  - > po is the program order: strict total order on events of the same thread





- Represent program behaviours as a set of consistent executions (graphs)
- An execution graph is a tuple:  $\langle E, po, rf, mo \rangle$
- ✤ E is the set of events (graph nodes), including initialisation writes
- po, rf and mo are *relations* on events (graph edges)
  - ▶ po is the *program order*: strict total order on events of the same thread
  - rf is the *reads-from* relation: relating each read/update to exactly one write/update on the same location with the same value





- Represent program behaviours as a set of consistent executions (graphs)
- An execution graph is a tuple:  $\langle E, po, rf, mo \rangle$
- ✤ E is the set of events (graph nodes), including initialisation writes
- po, rf and mo are *relations* on events (graph edges)
  - po is the program order: strict total order on events of the same thread
  - rf is the *reads-from* relation: relating each read/update to exactly one write/update on the same location with the same value





- Represent program behaviours as a set of consistent executions (graphs)
- An execution graph is a tuple:  $\langle E, po, rf, mo \rangle$
- ✤ E is the set of events (graph nodes), including initialisation writes
- po, rf and mo are *relations* on events (graph edges)
  - ▶ po is the *program order*: strict total order on events of the same thread
  - rf is the *reads-from* relation: relating each read/update to exactly one write/update on the same location with the same value
  - mo is the *modification order*: strict total order on the writes/updates of the same location



- Represent program behaviours as a set of consistent executions (graphs)
- An execution graph is a tuple:  $\langle E, po, rf, mo \rangle$
- ✤ E is the set of events (graph nodes), including initialisation writes
- po, rf and mo are *relations* on events (graph edges)
  - ▶ po is the *program order*: strict total order on events of the same thread
  - rf is the reads-from relation: relating each read/update to exactly one write/update on the same location with the same value
  - mo is the modification order: strict total order on the writes/updates of the same location
  - Derived relation,  $rb = \frac{rf^{-1}}{r} \frac{mo}{r}$ , is the *reads-before* relation



- Represent program behaviours as a set of consistent executions (graphs)
- An execution graph is a tuple:  $\langle E, po, rf, mo \rangle$
- ✤ E is the set of events (graph nodes), including initialisation writes
- po, rf and mo are *relations* on events (graph edges)
  - ▶ po is the *program order*: strict total order on events of the same thread
  - rf is the reads-from relation: relating each read/update to exactly one write/update on the same location with the same value
  - mo is the *modification order*: strict total order on the writes/updates of the same location
  - Derived relation,  $rb = \frac{rf^{-1}}{r} \frac{mo}{r}$ , is the *reads-before* relation





## Consistent Executions (w/o Persistency)

What is a *consistent* execution?

## Consistent Executions (w/o Persistency)

- What is a *consistent* execution?
  - Depends on the (concurrency) memory model
- What is a *consistent* execution?
  - Depends on the (concurrency) memory model
  - Intel-x86 consistency:

rfi∪moi ∪rbi ⊆ po

(Internal)

Ri : internal (same-thread) subset of R Re : external (diff.-thread) subset of R: R  $\ Ri$ 

- What is a *consistent* execution?
  - Depends on the (concurrency) memory model
  - Intel-x86 consistency:

rfi∪moi ∪rbi ⊆ po

irreflexive(ob)  $ob=(ppo \cup rfe \cup moe \cup rbe)^+$ 

(Internal) (External)

Ri : internal (same-thread) subset of R Re : external (diff.-thread) subset of R: R  $\$  Ri

- What is a *consistent* execution?
  - Depends on the (concurrency) memory model
  - Intel-x86 consistency:

rfi∪moi∪rbi ⊆ po irreflexive(ob)  $ob=(ppo)Urfe∪moe∪rbe)^+$ preserved program order: sloc or ✓ in table



					Late	er in Progra	ım Order		
			1	2	3	4	5	6	7
<u>ل</u>			Read	Write	RMW	mfence	sfence	flush <sub>opt</sub>	flush
rde	Α	Read	<ul> <li>Image: A set of the set of the</li></ul>	1	1	1	1	<b>√</b>	1
Ö	В	Write	×	<b>√</b>	<b>√</b>	<b>√</b>	<b>√</b>	sloc	<b>√</b>
ran	C	RMW	<ul> <li>Image: A set of the set of the</li></ul>	1	1	<b>√</b>	1	<b>√</b>	1
rog	D	mfence	<ul> <li>Image: A start of the start of</li></ul>	1	1	1	1	1	1
ц Ц	E	sfence	×	1	1	1	1	1	1
ler	F	<b>flush</b> opt	×	×	1	1	1	×	sloc
farl	G	flush	×	<b>√</b>	1	<b>√</b>	1	sloc	<b>√</b>

Ri : internal (same-thread) subset of R Re : external (diff.-thread) subset of R: R  $\$  Ri

- What is a *consistent* execution? \*
  - Depends on the (concurrency) memory model
  - Intel-x86 consistency:



Ri : internal (same-thread) subset of R Re : external (diff.-thread) subset of R: R \ Ri

41

(Internal)

(External)

flushopt

1

sloc

1

1

1

X

sloc

W(y, 1)

R(x, 0)

flush

1

1

1

1

1

sloc

1

sfence

1

1

/

1

1

/

3

- What is a *consistent* execution?
  - Depends on the (concurrency) memory model
  - Intel-x86 consistency:



Ri : internal (same-thread) subset of R Re : external (diff.-thread) subset of R: R  $\$  Ri



				Late	er in Progra	ım Order		
		1	2	3	4	5	6	7
		Read	Write	RMW	mfence	sfence	flush <sub>opt</sub>	flush
Α	Read	<ul> <li>✓</li> </ul>	1	1	1	1	<b>√</b>	<ul> <li>Image: A set of the set of the</li></ul>
B	Write	×	1	1	1	1	sloc	1
C	RMW	<ul> <li>Image: A start of the start of</li></ul>	1	<b>√</b>	<b>√</b>	<b>√</b>	<b>√</b>	<ul> <li>Image: A set of the set of the</li></ul>
D	mfence	<ul> <li>Image: A start of the start of</li></ul>	1	1	1	1	<b>√</b>	<ul> <li>Image: A set of the set of the</li></ul>
E	sfence	×	1	1	1	1	<b>√</b>	<ul> <li>Image: A start of the start of</li></ul>
F	<b>flush</b> opt	×	×	1	1	1	×	sloc
G	flush	×	1	1	1	1	sloc	1



x86-consistent execution

Represent program behaviours as a set of consistent & persistent executions

- Represent program behaviours as a set of consistent & persistent executions
- An execution graph is a tuple:  $\langle E, po, rf, mo, P, nvo \rangle$ 
  - Let D ⊆ E be the set of *durable events*: events whose effect *can reach* NVM model-specific for Intel-x86: D = W ∪ U ∪ FL ∪ FO

- Represent program behaviours as a set of consistent & persistent executions
- ✤ An execution graph is a tuple: < E, po, rf, mo, P, nvo >
  - Let D ⊆ E be the set of *durable events*: events whose effect *can reach* NVM model-specific for Intel-x86: D = W ∪ U ∪ FL ∪ FO





- Represent program behaviours as a set of consistent & persistent executions
- An execution graph is a tuple:  $\langle E, po, rf, mo, P, nvo \rangle$ 
  - Let D ⊆ E be the set of *durable events*: events whose effect *can reach* NVM model-specific for Intel-x86: D = W ∪ U ∪ FL ∪ FO
  - $P \subseteq D$  is the set of *persisted events*: events whose effect *have reached* NVM, s.t. init  $\subseteq P$





- Represent program behaviours as a set of consistent & persistent executions
- An execution graph is a tuple:  $\langle E, po, rf, mo, P, nvo \rangle$ 
  - Let D ⊆ E be the set of *durable events*: events whose effect *can reach* NVM model-specific for Intel-x86: D = W ∪ U ∪ FL ∪ FO
  - $P \subseteq D$  is the set of *persisted events*: events whose effect *have reached* NVM, s.t. init  $\subseteq P$







- Represent program behaviours as a set of consistent & persistent executions
- An execution graph is a tuple:  $\langle E, po, rf, mo, P, nvo \rangle$ 
  - Let D ⊆ E be the set of *durable events*: events whose effect *can reach* NVM model-specific for Intel-x86: D = W ∪ U ∪ FL ∪ FO
  - $P \subseteq D$  is the set of *persisted events*: events whose effect *have reached* NVM, s.t. init  $\subseteq P$
  - $nvo \subseteq D \times D$  is the *non-volatile-order*: a strict (partial) order on D that is downward-closed on P: if (e, e')  $\in$  nvo and e'  $\in$  P, then  $e \in P$







Intel-x86 consistency: rfi∪moi∪rbi ⊆ po irreflexive(ob) ob=(ppo∪rfe∪moe∪rbe)+ preserved program order (sloc or ✓ in table)

				Late	er in Progra	ım Order		
		1	2	3	4	5	6	7
		Read	Write	RMW	mfence	sfence	<b>flush</b> opt	flush
Α	Read	<ul> <li>Image: A start of the start of</li></ul>	1	1	1	1	<b>√</b>	<b>√</b>
B	Write	×	1	1	1	<b>√</b>	sloc	1
С	RMW	<ul> <li>Image: A set of the set of the</li></ul>	1	1	1	1	1	1
D	mfence	<ul> <li>Image: A set of the set of the</li></ul>	1	1	1	1	<b>√</b>	<b>√</b>
E	sfence	×	1	1	1	<b>√</b>	<b>√</b>	1
F	<b>flush</b> opt	×	×	1	1	<b>√</b>	×	sloc
G	flush	X	1	1	1	1	sloc	1

- Intel-x86 consistency: rfi∪moi∪rbi ⊆ po irreflexive(ob) ob=(ppoUrfe∪moe∪rbe)+ preserved program order (sloc or ✓ in table)
- Intel-x86 persistency (simplified):

					Late	er in Progra	ım Order		
			1	2	3	4	5	6	7
L			Read	Write	RMW	mfence	sfence	<b>flush</b> opt	flush
rdeı	Α	Read	1	1	<b>√</b>	1	<ul> <li>Image: A start of the start of</li></ul>	1	<ul> <li>Image: A set of the set of the</li></ul>
Ö	B	Write	×	1	1	1	1	sloc	<ul> <li>Image: A set of the set of the</li></ul>
ran	С	RMW	<ul> <li>Image: A set of the set of the</li></ul>	1	1	1	1	1	<ul> <li>Image: A start of the start of</li></ul>
rog	D	mfence	<ul> <li>Image: A start of the start of</li></ul>	1	1	1	1	1	<ul> <li>Image: A set of the set of the</li></ul>
in P	E	sfence	×	1	1	1	1	<b>√</b>	<ul> <li>Image: A set of the set of the</li></ul>
ier	F	<b>flush</b> opt	×	×	1	1	1	×	sloc
Earl	G	flush	×	<b>√</b>	<b>√</b>	<b>√</b>	<b>√</b>	sloc	<b>√</b>





- Intel-x86 consistency: rfi∪moi∪rbi ⊆ po irreflexive(ob) ob=(ppo∪rfe∪moe∪rbe)+ preserved program order (sloc or ✓ in table)
- Intel-x86 persistency (simplified):

 $\mathsf{FL} \cup \mathsf{dom}(\mathsf{FO}_{\mathsf{PO}} \mathsf{SF} \cup \mathsf{MF} \cup \mathsf{U}) \subseteq \mathsf{P}$ 

					Late	er in Progra	ım Order		
			1	2	3	4	5	6	7
<u>ب</u>			Read	Write	RMW	mfence	sfence	<b>flush</b> opt	flush
rdeı	Α	Read	<ul> <li>Image: A set of the set of the</li></ul>	1	1	1	1	1	<ul> <li>Image: A second s</li></ul>
Ö	В	Write	×	<b>√</b>	<b>√</b>	<b>√</b>	<b>√</b>	sloc	<ul> <li>Image: A start of the start of</li></ul>
ran	С	RMW	<ul> <li>Image: A second s</li></ul>	1	1	1	1	1	<ul> <li>Image: A start of the start of</li></ul>
rog	D	mfence	<ul> <li>Image: A start of the start of</li></ul>	<b>√</b>	1	1	1	<b>√</b>	<ul> <li>Image: A start of the start of</li></ul>
in P	E	sfence	×	<b>√</b>	1	1	1	<b>√</b>	<ul> <li>Image: A set of the set of the</li></ul>
ier	F	<b>flush</b> opt	×	×	1	1	1	×	sloc
Earl	G	flush	×	<b>~</b>	<b>√</b>	<b>√</b>	<b>√</b>	sloc	<ul> <li>Image: A start of the start of</li></ul>





Intel-x86 consistency: rfi∪moi∪rbi ⊆ po irreflexive(ob) ob=(ppo)Urfe∪moe ∪rbe)+ preserved program order (sloc or ✓ in table)

					Late	er in Progra	ım Order		
			1	2	3	4	5	6	7
5			Read	Write	RMW	mfence	sfence	<b>flush</b> opt	flush
rdeı	Α	Read	1	<b>√</b>	<b>√</b>	<b>√</b>	<b>√</b>	<b>√</b>	<b>√</b>
Ö	В	Write	×	<b>√</b>	<b>√</b>	<b>√</b>	<b>√</b>	sloc	<b>√</b>
ran	C	RMW	<b>√</b>	1	1	<b>√</b>	<b>√</b>	1	<b>√</b>
rog	D	mfence	1	1	1	1	1	1	<b>√</b>
in P	E	sfence	×	1	1	1	1	1	1
ier	F	<b>flush</b> opt	×	×	1	1	1	×	sloc
arl	G	flush	X	1	1	1	1	sloc	1

Intel-x86 persistency (simplified):

```
(FL) \cup dom(FO \xrightarrow{po} SF \cup MF \cup U) \subseteq P
strong persists
```





Intel-x86 consistency: rfi∪moi∪rbi ⊆ po irreflexive(ob) ob=(ppo)Urfe∪moe Urbe)+ preserved program order (sloc or ✓ in table)

					Late	er in Progra	im Order		
			1	2	3	4	5	6	7
<u>ب</u>			Read	Write	RMW	mfence	sfence	<b>flush</b> opt	flush
rde	Α	Read	<ul> <li>Image: A set of the set of the</li></ul>	<b>√</b>	<b>√</b>	<b>√</b>	<b>√</b>	<b>√</b>	<ul> <li>Image: A set of the set of the</li></ul>
Ö	B	Write	×	<b>√</b>	<b>√</b>	<b>√</b>	<b>√</b>	sloc	<ul> <li>Image: A set of the set of the</li></ul>
ran	C	RMW	<ul> <li>Image: A start of the start of</li></ul>	<b>√</b>	1	<b>√</b>	<b>√</b>	<b>√</b>	<ul> <li>Image: A set of the set of the</li></ul>
rog	D	mfence	<ul> <li>Image: A start of the start of</li></ul>	<b>√</b>	<b>√</b>	<b>√</b>	<b>√</b>	<b>√</b>	<b>~</b>
in F	E	sfence	×		<b>√</b>	<b>√</b>	<b>√</b>	<b>√</b>	<ul> <li>Image: A set of the set of the</li></ul>
ier	F	<b>flush</b> opt	×	×	<b>√</b>	<b>√</b>	<b>√</b>	×	sloc
Earl	G	flush	×		<b>√</b>	<b>√</b>	<b>√</b>	sloc	<b>~</b>

Intel-x86 persistency (simplified):

 $FL \cup (dom(FO \xrightarrow{po} SF \cup MF \cup U)) \subseteq P$ strong persists persist sequences

![](_page_159_Figure_5.jpeg)

![](_page_159_Figure_6.jpeg)

![](_page_159_Picture_7.jpeg)

- Intel-x86 consistency: rfi∪moi∪rbi ⊆ po irreflexive(ob) ob=(ppo∪rfe∪moe∪rbe)+ preserved program order (sloc or ✓ in table)
- Intel-x86 persistency (simplified):

```
FL \cup dom(FO \xrightarrow{po} SF \cup MF \cup U) \subseteq PD \xrightarrow{ob} \cap sloc} D \subseteq nvo
```

```
Later in Program Order
                                         3
                                                                         6
                                 2
                                                                     flushopt
                      Read
                               Write
                                       RMW
                                                mfence
                                                           sfence
                                                                                 flush
3arlier in Program Order
                        1
                                                              /
             Read
                                 1
                                         1
                                                   1
      A
                                                                                   \checkmark
      B
                                         /
                        X
                                 1
                                                   1
                                                              1
                                                                       sloc
            Write
                                                                                   1
      C
            RMW
                        1
                                1
                                         1
                                                   1
                                                              1
                                                                         1
                                                                                   1
      D
                        1
                                1
                                         1
                                                   1
                                                              1
                                                                                   1
           mfence
                                                                         \checkmark
                        X
                                                              /
      E
                                         1
                                                   1
                                                                         1
                                                                                   1
            sfence
                                1
                        X
                                X
                                         1
                                                   1
                                                              1
                                                                         X
                                                                                  sloc
           flush<sub>opt</sub>
                                                              1
       G
            flush
                        X
                                1
                                                   1
                                                                       sloc
                                                                                   1
```

1 x:=1;
2 clflushopt x;
3 sfence;
4 y:=1

![](_page_160_Figure_6.jpeg)

![](_page_160_Picture_7.jpeg)

- Intel-x86 consistency: rfi∪moi∪rbi ⊆ po irreflexive(ob) ob=(ppo∪rfe∪moe∪rbe)+ preserved program order (sloc or ✓ in table)
- Intel-x86 persistency (simplified):

 $\mathsf{FL} \cup \mathsf{dom}(\mathsf{FO} \xrightarrow{\mathsf{PO}} \mathsf{SF} \cup \mathsf{MF} \cup \mathsf{U}) \subseteq \mathsf{P}$ 

 $D \xrightarrow{\text{ob}} \cap \text{sloc} D \subseteq \text{nvo}$ 

![](_page_161_Figure_5.jpeg)

![](_page_161_Figure_6.jpeg)

![](_page_161_Figure_7.jpeg)

- Intel-x86 consistency: rfi∪moi∪rbi ⊆ po irreflexive(ob) ob=(ppo)Urfe∪moe ∪rbe)+ preserved program order (sloc or ✓ in table)
- Intel-x86 persistency (simplified):

```
FL U dom(FO \xrightarrow{po} SFUMFUU) \subseteq P

D \xrightarrow{ob} \cap sloc} D \subseteq nvo

(FO U FL) \xrightarrow{ob} D \subseteq nvo
```

![](_page_162_Picture_4.jpeg)

![](_page_162_Figure_5.jpeg)

![](_page_162_Picture_6.jpeg)

					Late	er in Progra	ım Order		
			1	2	3	4	5	6	7
5			Read	Write	RMW	mfence	sfence	<b>flush</b> opt	flush
rde	Α	Read	<ul> <li>Image: A start of the start of</li></ul>	<b>√</b>	1	1	1	<b>√</b>	<ul> <li>Image: A start of the start of</li></ul>
Ö	B	Write	×	1	1	1	1	sloc	<ul> <li>Image: A start of the start of</li></ul>
ran	С	RMW	<ul> <li>Image: A second s</li></ul>	1	1	1	1	1	<ul> <li>Image: A start of the start of</li></ul>
rog	D	mfence	<ul> <li>Image: A set of the set of the</li></ul>	1	1	1	1	1	<ul> <li>Image: A start of the start of</li></ul>
in P	Ε	sfence	×	1	1	1	1	1	<ul> <li>Image: A start of the start of</li></ul>
ier	F	<b>flush</b> opt	×	×	1	1	1	×	sloc
Earl	G	flush	×	<b>√</b>	1	<ul> <li>Image: A set of the set of the</li></ul>	<b>√</b>	sloc	<ul> <li>Image: A start of the start of</li></ul>

- Intel-x86 consistency: rfi∪moi∪rbi ⊆ po irreflexive(ob) ob=(ppo)Urfe∪moe ∪rbe)+ preserved program order (sloc or ✓ in table)
- Intel-x86 persistency (simplified):

```
FL U dom(FO \xrightarrow{po} SFUMFUU) \subseteq P

D \xrightarrow{ob} \cap sloc} D \subseteq nvo

(FO U FL) \xrightarrow{ob} D \subseteq nvo
```

![](_page_163_Picture_4.jpeg)

![](_page_163_Figure_5.jpeg)

					Late	er in Progra	ım Order		
			1	2	3	4	5	6	7
L			Read	Write	RMW	mfence	sfence	flush <sub>opt</sub>	flush
rdei	Α	Read	1	1	1	1	1	1	<ul> <li>Image: A set of the set of the</li></ul>
Ő	В	Write	×	1	1	1	1	sloc	<ul> <li>Image: A start of the start of</li></ul>
ran	С	RMW	<ul> <li>Image: A start of the start of</li></ul>	1	1	1	1	1	<ul> <li>Image: A start of the start of</li></ul>
rog	D	mfence	<ul> <li>Image: A start of the start of</li></ul>	1	1	1	1	1	<ul> <li>Image: A start of the start of</li></ul>
ц Ч	Ε	sfence	×	1	1	1	1	1	<ul> <li>Image: A start of the start of</li></ul>
ler	F	<b>flush</b> opt	X	×	1	1		×	sloc
arb	G	flush	X	1	1	1	1	sloc	<ul> <li>Image: A set of the set of the</li></ul>

- Intel-x86 consistency: rfi∪moi∪rbi ⊆ po irreflexive(ob) ob=(ppo)Urfe∪moe ∪rbe)+ preserved program order (sloc or ✓ in table)
- Intel-x86 persistency (simplified):

```
FL U dom(FO \xrightarrow{po} SFUMFUU) \subseteq P

D \xrightarrow{ob} \cap sloc} D \subseteq nvo

(FO U FL) \xrightarrow{ob} D \subseteq nvo
```

![](_page_164_Picture_4.jpeg)

![](_page_164_Figure_5.jpeg)

	Later in Program Order									
			1	2	3	4	5	6	7	
L			Read	Write	RMW	mfence	sfence	<b>flush</b> opt	flush	
rae	Α	Read	<ul> <li>Image: A start of the start of</li></ul>	1	1	1	<ul> <li>Image: A start of the start of</li></ul>	<b>√</b>	1	
Ď	B	Write	×	1	1	1	<ul> <li>Image: A start of the start of</li></ul>	sloc	1	
Tan	C	RMW	<ul> <li>✓</li> </ul>	1	1	1	<ul> <li>Image: A start of the start of</li></ul>	<b>√</b>	1	
	D	mfence	<ul> <li>✓</li> </ul>	1	1	1	<ul> <li>Image: A set of the set of the</li></ul>	<b>√</b>	1	
ц Ц	E	sfence	×	$\checkmark$	1	1	1	<b>√</b>	1	
ler	F	<b>flush</b> opt	×	×	1	1	<ul> <li>Image: A set of the set of the</li></ul>	×	sloc	
Carn	G	flush	×	1	1	<b>√</b>	<b>√</b>	sloc	<b>√</b>	

- Intel-x86 consistency: rfi∪moi∪rbi ⊆ po irreflexive(ob) ob=(ppo)Urfe∪moe ∪rbe)+ preserved program order (sloc or ✓ in table)
- Intel-x86 persistency (simplified):

```
FL U dom(FO \xrightarrow{po} SFUMFUU) \subseteq P

D \xrightarrow{ob} \cap sloc} D \subseteq nvo

(FO U FL) \xrightarrow{ob} D \subseteq nvo
```

![](_page_165_Picture_4.jpeg)

![](_page_165_Figure_5.jpeg)

![](_page_165_Picture_6.jpeg)

					Late	er in Progra	ım Order		
			1	2	3	4	5	6	7
<u>ب</u>			Read	Write	RMW	mfence	sfence	<b>flush</b> opt	flush
rdeı	Α	Read	<ul> <li>Image: A start of the start of</li></ul>	<b>√</b>	1	1	1	<b>√</b>	<ul> <li>Image: A set of the set of the</li></ul>
Ö	B	Write	×	1	1	1	1	sloc	<ul> <li>Image: A set of the set of the</li></ul>
ran	С	RMW	<ul> <li>Image: A second s</li></ul>	1	1	1	1	1	<ul> <li>Image: A set of the set of the</li></ul>
rog	D	mfence	<ul> <li>Image: A set of the set of the</li></ul>	1	1	1	1	1	<ul> <li>Image: A set of the set of the</li></ul>
in P	Ε	sfence	×	1	1	1	1	1	<ul> <li>Image: A set of the set of the</li></ul>
ier	F	<b>flush</b> opt	×	×	1	1	1	×	sloc
Earl	G	flush	×	<b>√</b>	<b>√</b>	<b>√</b>	<b>√</b>	sloc	<ul> <li>Image: A start of the start of</li></ul>

- Intel-x86 consistency: rfi∪moi∪rbi ⊆ po irreflexive(ob) ob=(ppo)Urfe∪moe ∪rbe)+ preserved program order (sloc or ✓ in table)
- Intel-x86 persistency (simplified):

```
FL U dom(FO \xrightarrow{po} SFUMFUU) \subseteq P

D \xrightarrow{ob} \cap sloc} D \subseteq nvo

(FO U FL) \xrightarrow{ob} D \subseteq nvo
```

![](_page_166_Picture_4.jpeg)

![](_page_166_Figure_5.jpeg)

Later in Program Order									
			1	2	3	4	5	6	7
L			Read	Write	RMW	mfence	sfence	<b>flush</b> opt	flush
rdei	Α	Read	<ul> <li>Image: A second s</li></ul>	1	1	1	1	1	<ul> <li>Image: A start of the start of</li></ul>
5 C	В	Write	×	1	1	1	1	sloc	<ul> <li>Image: A start of the start of</li></ul>
ran	С	RMW	<ul> <li>Image: A start of the start of</li></ul>	<b>√</b>	1	1	1	1	<ul> <li>Image: A start of the start of</li></ul>
rog	D	mfence	<ul> <li>Image: A start of the start of</li></ul>	<b>√</b>	1	1	1	1	<ul> <li>Image: A start of the start of</li></ul>
ц Ц	Ε	sfence	×	<b>√</b>	1	1	1	1	<ul> <li>Image: A start of the start of</li></ul>
ler	F	<b>flush</b> opt	×	×	1	1	1	×	sloc
tarl	G	flush	×	<b>√</b>	1	<ul> <li>Image: A set of the set of the</li></ul>	<b>√</b>	sloc	<ul> <li>Image: A start of the start of</li></ul>

- Intel-x86 consistency: rfi∪moi∪rbi ⊆ po irreflexive(ob) ob=(ppo)∪rfe∪moe ∪rbe)+ preserved program order (sloc or ✓ in table)
- Intel-x86 persistency (simplified):

```
FL U dom(FO \xrightarrow{po} SFUMFUU) \subseteq P

D \xrightarrow{ob} \cap sloc} D \subseteq nvo

(FO U FL) \xrightarrow{ob} D \subseteq nvo
```

1	x:=1;
2	clflushopt x;
3	sfence;
4	y:=1

persisted event (in P)

![](_page_167_Figure_5.jpeg)

			Later in Program Order						
			1	2	3	4	5	6	7
L			Read	Write	RMW	mfence	sfence	flush <sub>opt</sub>	flush
rde	Α	Read	<ul> <li>✓</li> </ul>	1	<ul> <li>✓</li> </ul>	1	1	<ul> <li>Image: A set of the set of the</li></ul>	<ul> <li>Image: A start of the start of</li></ul>
Ö	B	Write	×	1	1	1	1	sloc	<ul> <li>✓</li> </ul>
ran	С	RMW	1	1	1	1	1	1	<ul> <li>✓</li> </ul>
rog	D	mfence	<ul> <li>✓</li> </ul>	1	1	1	1	1	<ul> <li>✓</li> </ul>
пЪ	Ε	sfence	×	1	1	1	1	<ul> <li>✓</li> </ul>	<ul> <li>✓</li> </ul>
ier	F	<b>flush</b> opt	×	×	1	1	1	×	sloc
L	G	fluch	X				1	sloc	

# 4. Other hardware persistency models

#### ARMv8 Consistency

• A complex model; much weaker than Intel-x86 consistency

#### ARMv8 Consistency

- ► A complex model; much weaker than Intel-x86 consistency
- Operational model: promising semantics by Kang et al.

#### ARMv8 Consistency

- A complex model; much weaker than Intel-x86 consistency
- Operational model: promising semantics by Kang et al.
- ▶ Declarative model: by Pulte et al. defines the ob relation

#### ARMv8 Consistency

- A complex model; much weaker than Intel-x86 consistency
- Operational model: promising semantics by Kang et al.
- ▶ Declarative model: by Pulte et al. defines the ob relation

#### ARMv8 Persistency

▶ Weak explicit persists: DC CVAP x — analogous to clflushopt x

#### ARMv8 Consistency

- A complex model; much weaker than Intel-x86 consistency
- Operational model: promising semantics by Kang et al.
- ▶ Declarative model: by Pulte et al. defines the ob relation

- ▶ Weak explicit persists: DC CVAP x analogous to clflushopt x
- Persist sequences: DC CVAP x; DSB SY DSB SY is a strong fence (analogous to mfence)

#### ARMv8 Consistency

- A complex model; much weaker than Intel-x86 consistency
- Operational model: promising semantics by Kang et al.
- ▶ Declarative model: by Pulte et al. defines the ob relation

- Weak explicit persists: DC CVAP x analogous to clflushopt x
- Persist sequences: DC CVAP x; DSB SY DSB SY is a strong fence (analogous to mfence)
- Operational model: an extension of promising semantics by Cho et al. [5]

#### ARMv8 Consistency

- A complex model; much weaker than Intel-x86 consistency
- Operational model: promising semantics by Kang et al.
- ▶ Declarative model: by Pulte et al. defines the ob relation

- Weak explicit persists: DC CVAP x analogous to clflushopt x
- Persist sequences: DC CVAP x; DSB SY DSB SY is a strong fence (analogous to mfence)
- Operational model: an extension of promising semantics by Cho et al. [5]
- ▶ Declarative model: by Raad et al. [4], and Cho et al. [5] (simplified)

#### ARMv8 Consistency

- A complex model; much weaker than Intel-x86 consistency
- Operational model: promising semantics by Kang et al.
- ▶ Declarative model: by Pulte et al. defines the ob relation

- Weak explicit persists: DC CVAP x analogous to clflushopt x
- ▶ Persist sequences: DC CVAP x; DSB SY DSB SY is a strong fence (analogous to mfence)
- Operational model: an extension of promising semantics by Cho et al. [5]
- ▶ Declarative model: by Raad et al. [4], and Cho et al. [5] (simplified)

ARMv8 Persistency	Intel-x86 Persistency

#### ARMv8 Consistency

- A complex model; much weaker than Intel-x86 consistency
- Operational model: promising semantics by Kang et al.
- ▶ Declarative model: by Pulte et al. defines the ob relation

- ▶ Weak explicit persists: DC CVAP x analogous to clflushopt x
- ▶ Persist sequences: DC CVAP x; DSB SY DSB SY is a strong fence (analogous to mfence)
- Operational model: an extension of promising semantics by Cho et al. [5]
- ▶ Declarative model: by Raad et al. [4], and Cho et al. [5] (simplified)

ARMv8 Persistency
$dom(DC_CVAP_{\text{DSB}}^{\text{DO}} DSB_SY) \subseteq P$

Intel-x86 Persistency
$FL \cup dom(FO \xrightarrow{PO} SFUMFUU) \subseteq P$

#### ARMv8 Consistency

- ► A complex model; much weaker than Intel-x86 consistency
- Operational model: promising semantics by Kang et al.
- ▶ Declarative model: by Pulte et al. defines the ob relation

- ▶ Weak explicit persists: DC CVAP x analogous to clflushopt x
- Persist sequences: DC CVAP x; DSB SY DSB SY is a strong fence (analogous to mfence)
- Operational model: an extension of promising semantics by Cho et al. [5]
- ▶ Declarative model: by Raad et al. [4], and Cho et al. [5] (simplified)

**ARMv8 Persistency**  $dom(DC_CVAP_{\rightarrow}^{po}DSB_SY) \subseteq P$  $\cap$  sloc  $D \subseteq$  nvo

Intel-x86 Persistency
$FL \cup dom(FO \xrightarrow{PO} SF \cup MF \cup U) \subseteq P$
$D \xrightarrow{\text{ob}} D \subseteq \text{nvo}$

#### ARMv8 Consistency

- A complex model; much weaker than Intel-x86 consistency
- Operational model: promising semantics by Kang et al.
- ▶ Declarative model: by Pulte et al. defines the ob relation

- ▶ Weak explicit persists: DC CVAP x analogous to clflushopt x
- ▶ Persist sequences: DC CVAP x; DSB SY DSB SY is a strong fence (analogous to mfence)
- Operational model: an extension of promising semantics by Cho et al. [5]
- ▶ Declarative model: by Raad et al. [4], and Cho et al. [5] (simplified)

ARMv8 Persistency  $dom(DC_CVAP_{OO}^{OO} DSB_SY) \subseteq P$  $D \xrightarrow{ob} \cap sloc} D \subseteq nvo$  $(DC_CVAP) \xrightarrow{ob} D \subseteq nvo$ 

Intel-x86 Persistency
$FL \cup dom(FO \xrightarrow{PO} SF \cup MF \cup U) \subseteq P$
$D \xrightarrow{\text{ob}} \cap \text{sloc} D \subseteq \text{nvo}$
(FL U FO) $\xrightarrow{ob} D \subseteq nvo$
## 5. Further Reading

## A. Intel-x86 / TSO Persistency Models

- [1] Persistence Semantics for Weak Memory: Integrating <u>Epoch Persistency</u> with the <u>TSO</u> Memory Model Azalea Raad, Viktor Vafeiadis
- [2] Persistency Semantics of the <u>Intel-x86</u> Architecture Azalea Raad, John Wickerson, Gil Neiger, Viktor Vafeiadis
- [3] *Taming <u>x86-TSO Persistency</u>* Artem Khyzha, Ori Lahav

## B. ARMv8 Persistency Models

- [4] Weak Persistency Semantics from the Ground Up: Formalising the Persistency Semantics of <u>ARMv8</u> and <u>Transactional Models</u> Azalea Raad, John Wickerson, Viktor Vafeiadis
- [5] Revamping Hardware Persistency Models: View-Based and Axiomatic Persistency Models for <u>Intel-x86</u> and <u>Armv8</u> Kyeongmin Cho, Sung-Hwan Lee, Azalea Raad, Jeehoon Kang

## C. Persistent Verification

- [6] <u>Linearizability</u> of Persistent Memory Objects Under a Full-System-Crash Failure Model Joseph Izraelevitz, Hammurabi Mendes, Michael L. Scott
- [7] <u>Persistent</u> Owicki-Gries Reasoning: A <u>Program Logic</u> for Reasoning about Persistent Programs on Intel-x86 Azalea Raad, Ori Lahav, Viktor Vafeiadis
- [8] Defining and Verifying Durable Opacity: Correctness for <u>Persistent Software Transactional Memory</u> Eleni Bila, Simon Doherty, Brijesh Dongol, John Derrick, Gerhard Schellhorn, Heike Wehrheim
- [9] PerSeVerE: Persistency Semantics for Verification under Ext4 Michalis Kokologiannakis, Ilya Kaisin, Azalea Raad, Viktor Vafeiadis
- [10] Deciding reachability under <u>persistent x86-TSO</u> Parosh Aziz Abdulla, Mohamed Faouzi Atig, Ahmed Bouajjani, K. Narayan Kumar, Prakash Saivasan