# Persistence Semantics for Weak Memory

Integrating Epoch Persistency with the TSO Memory Model

**Azalea Raad**       Viktor Vafeiadis

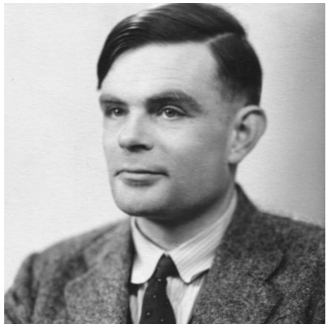Max Planck Institute for Software Systems (MPI-SWS)

Thursday 8 November

OOPSLA 2018

Boston, USA

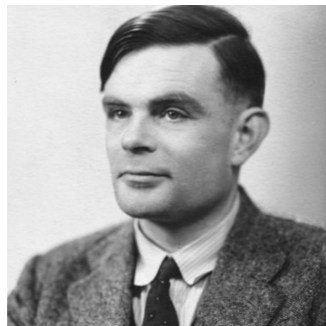✉ azalea@mpi-sws.org          🔗 SoundAndComplete.org          🐦 @azalearaad

# History

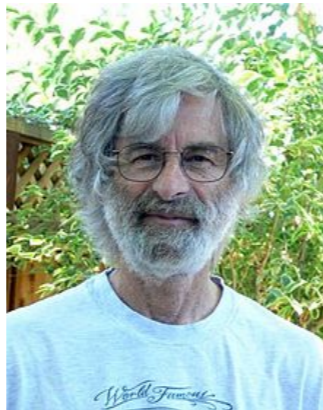Difficulty



😊

Sequential

time

# History

# History

Difficulty



time

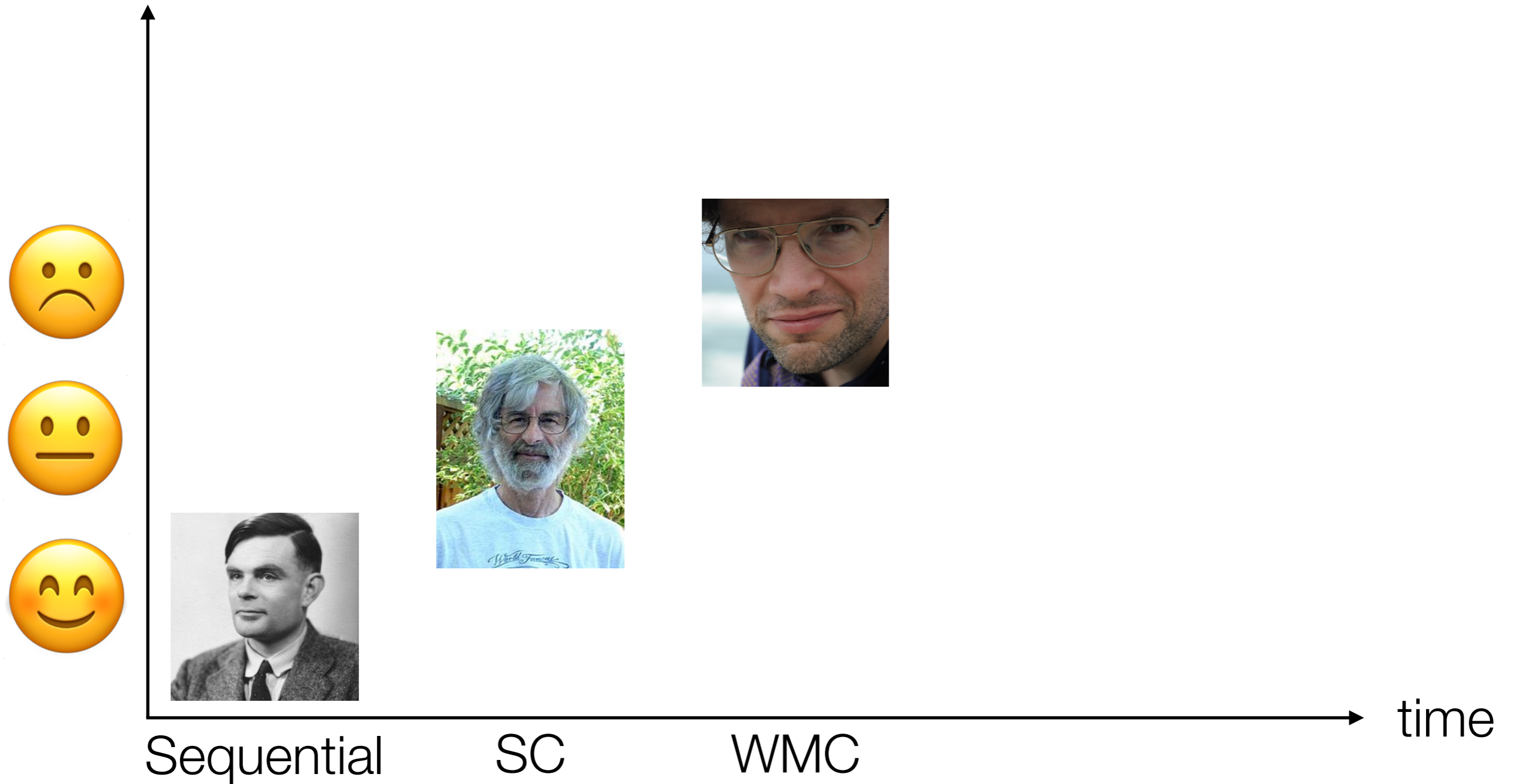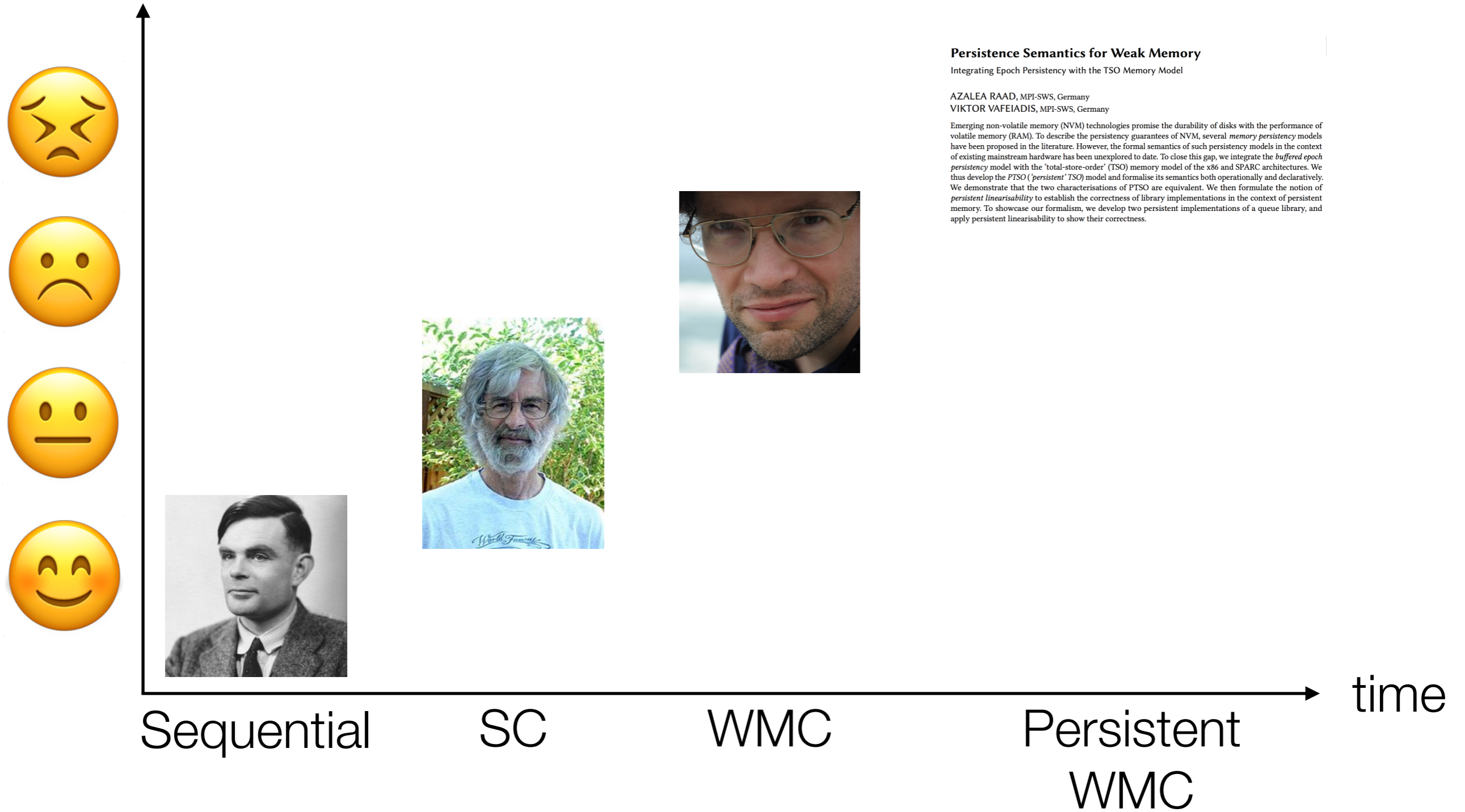Sequential        SC        WMC

# History



Difficulty

Persistence Semantics for Weak Memory
Integrating Epoch Persistency with the TSO Memory Model

AZALEA RAAD, MPI-SWS, Germany
VIKTOR VAFEIADIS, MPI-SWS, Germany

Emerging non-volatile memory (NVM) technologies promise the durability of disks with the performance of volatile memory (RAM). To describe the persistency guarantees of NVM, several *memory persistency* models have been proposed in the literature. However, the formal semantics of such persistency models in the context of existing mainstream hardware has been unexplored to date. To close this gap, we integrate the *buffered epoch persistency* model with the 'total-store-order' (TSO) memory model of the x86 and SPARC architectures. We thus develop the *PTSO* ('*persistent*' *TSO*) model and formalise its semantics both operationally and declaratively. We demonstrate that the two characterisations of PTSO are equivalent. We then formulate the notion of *persistent linearisability* to establish the correctness of library implementations in the context of persistent memory. To showcase our formalism, we develop two persistent implementations of a queue library, and apply persistent linearisability to show their correctness.

Sequential      SC      WMC      Persistent WMC

time

# What is Persistent Memory?

**_Volatile_** memory
// x = 0
x := 1
// x = 1

// x = v : reading x yields v

# What is Persistent Memory?

**_Volatile_** memory
// x = 0
x := 1
// x = 1

// no recovery
// x = 0

// x = v : reading x yields v

# What is Persistent Memory?

**Volatile** memory
// x = 0
x := 1
// x = 1

// no recovery
// x = 0

**Persistent** memory
// x = 0
x := 1
// x = 1

// recovery routine
// x = 0 OR x = 1

// x = v : reading x yields v

# What is Persistent Memory?

**Volatile** memory
```
// x = 0
x := 1
// x = 1
```
// no recovery
// x = 0

**Persistent** memory
```
// x = 0
x := 1
// x = 1
```
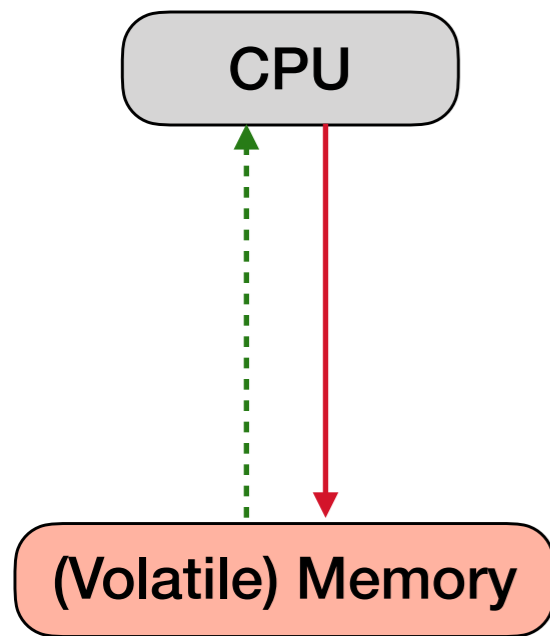// recovery routine
// x = 0 OR x = 1

persists are **asynchronous** (buffered): may not persist immediately
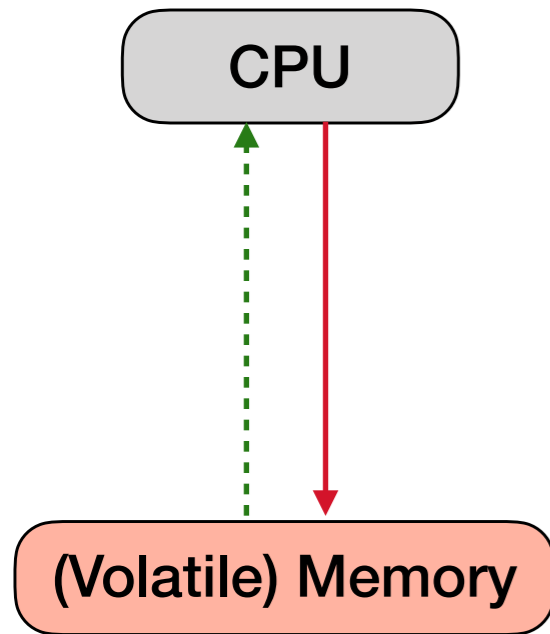
// x = v : reading x yields v

# (Sequential) Hardware

# (Sequential) Hardware
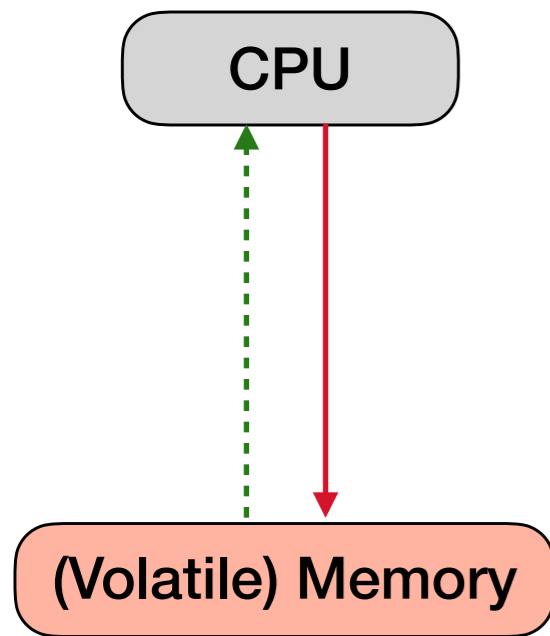
**CPU**

**(Volatile) Memory**

# (Sequential) Hardware



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

# (Sequential) Hardware



$x:=1$ : adds $x:=1$ to memory
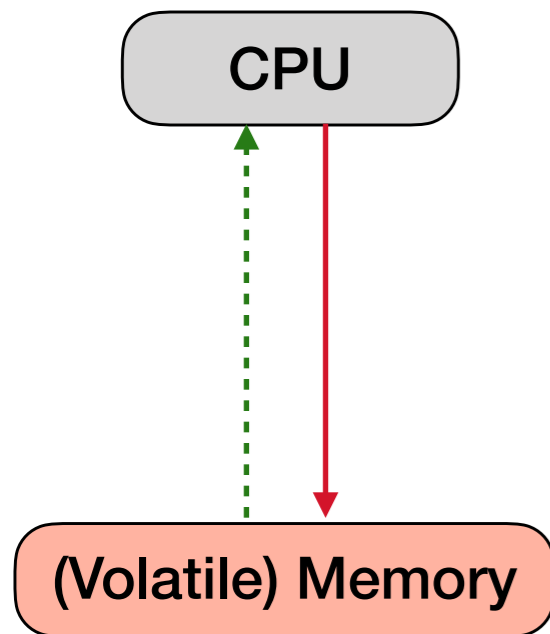
$a:=x$ : reads $x$ from memory

# (Sequential) Hardware

CPU

(Volatile) Memory

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

$a:=x$ : reads $x$ from memory

memory lost

# (Sequential) Hardware

CPU

(Volatile) Memory

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

$a:=x$ : reads $x$ from memory

memory lost

CPU

Persistence Buffer

(Persistent) Memory

# (Sequential) Hardware

CPU

(Volatile) Memory

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

$a:=x$ : reads $x$ from memory

memory lost

CPU

Persistence Buffer

(Persistent) Memory

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

# (Sequential) Hardware

CPU

(Volatile) Memory

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

$a:=x$ : reads $x$ from memory

memory lost

CPU

Persistence Buffer

(Persistent) Memory

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

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

# (Sequential) Hardware

**CPU**

**(Volatile) Memory**

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

$a:=x$ : reads $x$ from memory
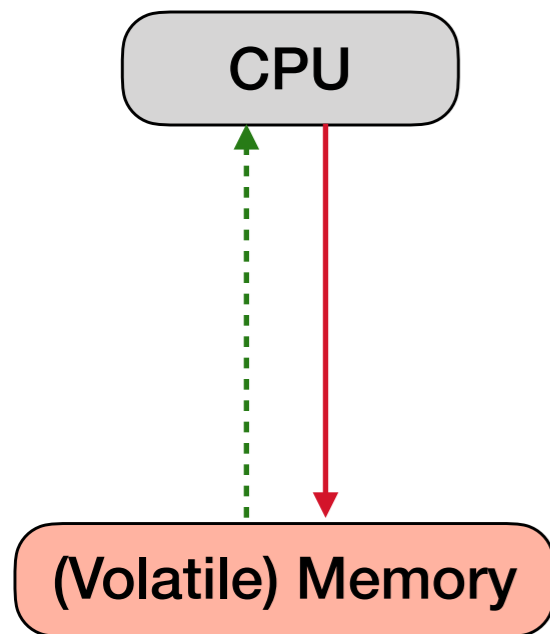
⚡ memory lost

---

**CPU**

**Persistence Buffer**

**(Persistent) Memory**

$x:=1$ : 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*

# (Sequential) Hardware



**CPU**

**(Volatile) Memory**

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

$a:=x$  :  reads $x$ from memory

⚡  memory lost

**CPU**

**Persistence Buffer**

**(Persistent) Memory**

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

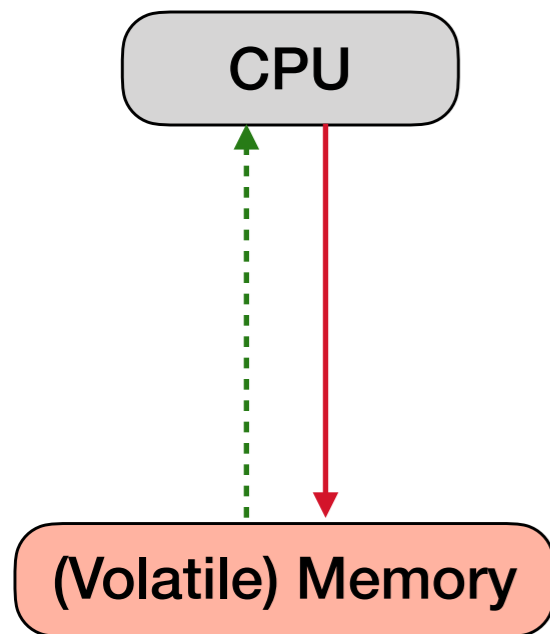$a:=x$  :  if p-buffer contains $x$, reads latest entry
else reads from memory

⚡  p-buffer lost; memory ***retained***

unbuffer* : p-buffer to memory

\* at non-deterministic times

# What is Memory Persistency Model?

- Memory **consistency** model describes:

    the order writes are made visible to other threads
    e.g. SC, TSO, …

# What is Memory Persistency Model?

- Memory ***consistency*** model describes:

  the order writes are made visible to other threads
  e.g. SC, TSO, ...

- Memory ***persistency*** model describes:

  the order writes are persisted to memory
  e.g. Epoch Persistency

# What is Memory Persistency Model?

- Memory **consistency** model describes:

  **Problem**

  **Formal**
  Epoch Persistency Model
  for
  **Mainstream Hardware** (**Weak** Memory Models)

  the order writes are persisted to memory
  e.g. Epoch Persistency

# What Can Go Wrong?

```
            // x=0;y=0

             x := 1;

             y := 1;
```



```
            // recovery routine

// x=0;y=0 OR x=1;y=1 OR x=1;y=0 OR x=0;y=1
```

# What Can Go Wrong?

// x=0;y=0

x := 1;

y := 1;

// recovery routine

// x=0;y=0 OR x=1;y=1 OR x=1;y=0 OR x=0;y=1

‼️ Writes may persist out of order

# What Can Go Wrong?

```
// x=0; y=0

x := 1;

y := 1;
```

// recovery routine

// x=0; y=0 OR x=1; y=1 OR x=1; y=0 OR x=0; y=1

‼️ Writes may persist out of order

☛ *persistent fence* `pfence`

# Persistent Fence

```
            // x=0;y=0
              x := 1;
    ☞        pfence;
              y := 1;
```



// recovery routine

// x=0;y=0 OR x=1;y=1 OR x=1;y=0 OR ~~x=0;y=1~~

# Persistent Fence

```
 a     x := 1;

 b     y := 2;

 c     x := 3;

       pfence;

 d     z := 4;
```

# Persistent Fence

- writes on **same locations** persist in <u>execution order</u>

a x := 1;

b y := 2;

c x := 3;

pfence;

d z := 4;

a persists before c

# Persistent Fence

- writes on **same locations** persist in <u>execution order</u>
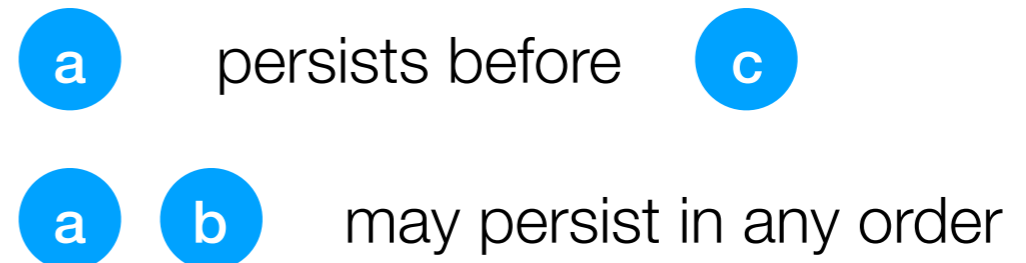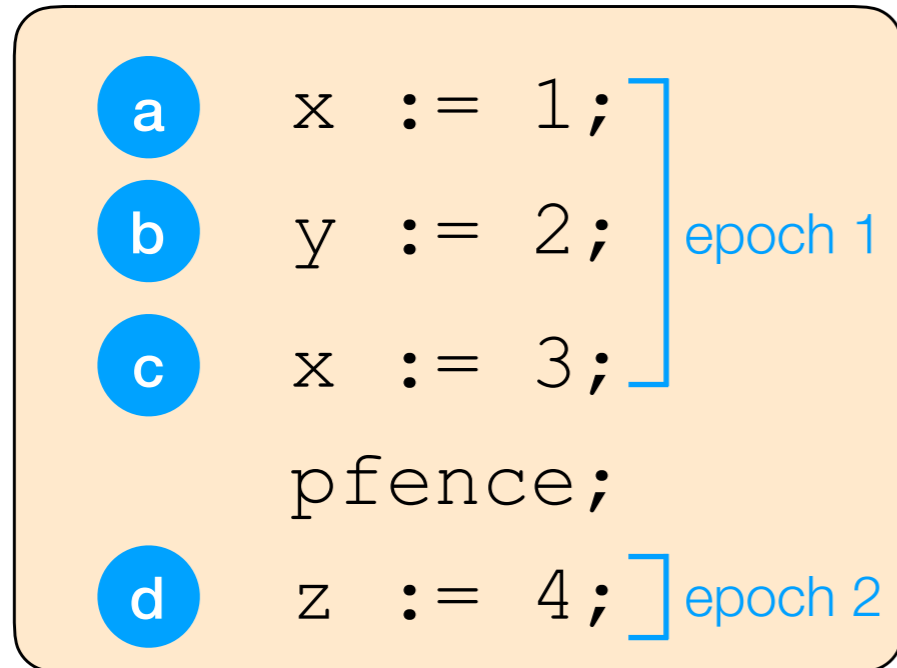- writes on **different locations** are <u>unordered</u>

```
a    x := 1;
b    y := 2;
c    x := 3;
     pfence;
d    z := 4;
```

(a) persists before (c)
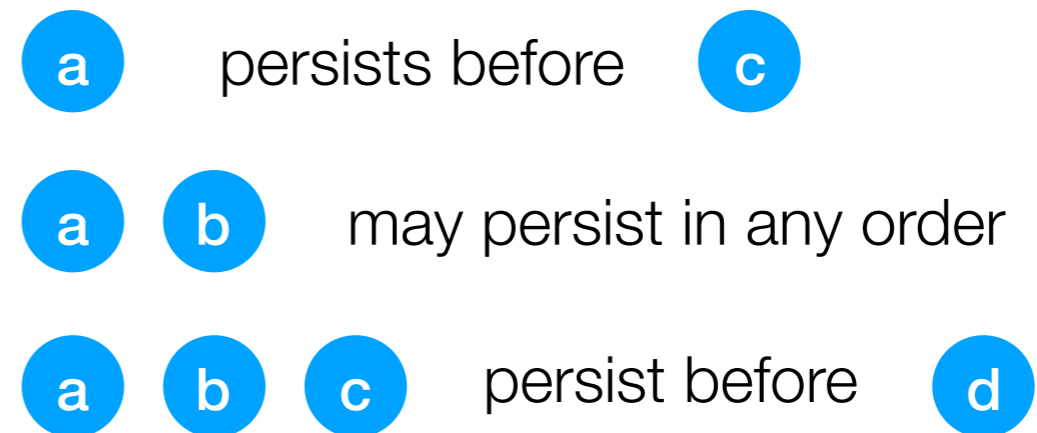
(a) (b) may persist in any order

# Persistent Fence

- writes on **same locations** persist in <u>execution order</u>
- writes on **different locations** are <u>unordered</u>
- `pfence` adds a new **epoch**

```
a   x := 1; ⎤
b   y := 2; ⎥ epoch 1
c   x := 3; ⎦
    pfence;
d   z := 4; ⎤ epoch 2
```

a persists before c

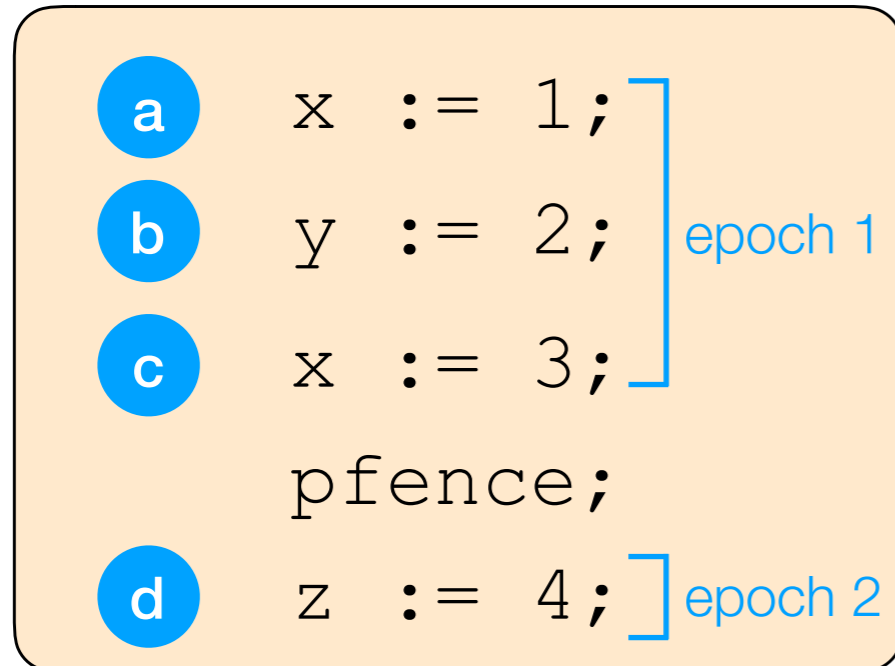a b may persist in any order

# Persistent Fence

- writes on **same locations** persist in <u>execution order</u>
- writes on **different locations** are <u>unordered</u>
- `pfence` adds a new **epoch**
- writes persist in <u>epoch order</u>

```
a   x := 1;
b   y := 2;    epoch 1
c   x := 3;
    pfence;
d   z := 4;    epoch 2
```

a persists before c

a b may persist in any order

a b c persist before d

# What Can Go Wrong (Continued)?

```
            // x=0; y=0

              x := 1;
asynchronous
(buffered)  →  pfence;

              y := 1;
```

```
        // recovery routine
```

```
// x=0; y=0  OR  x=1; y=1  OR  x=1; y=0
```
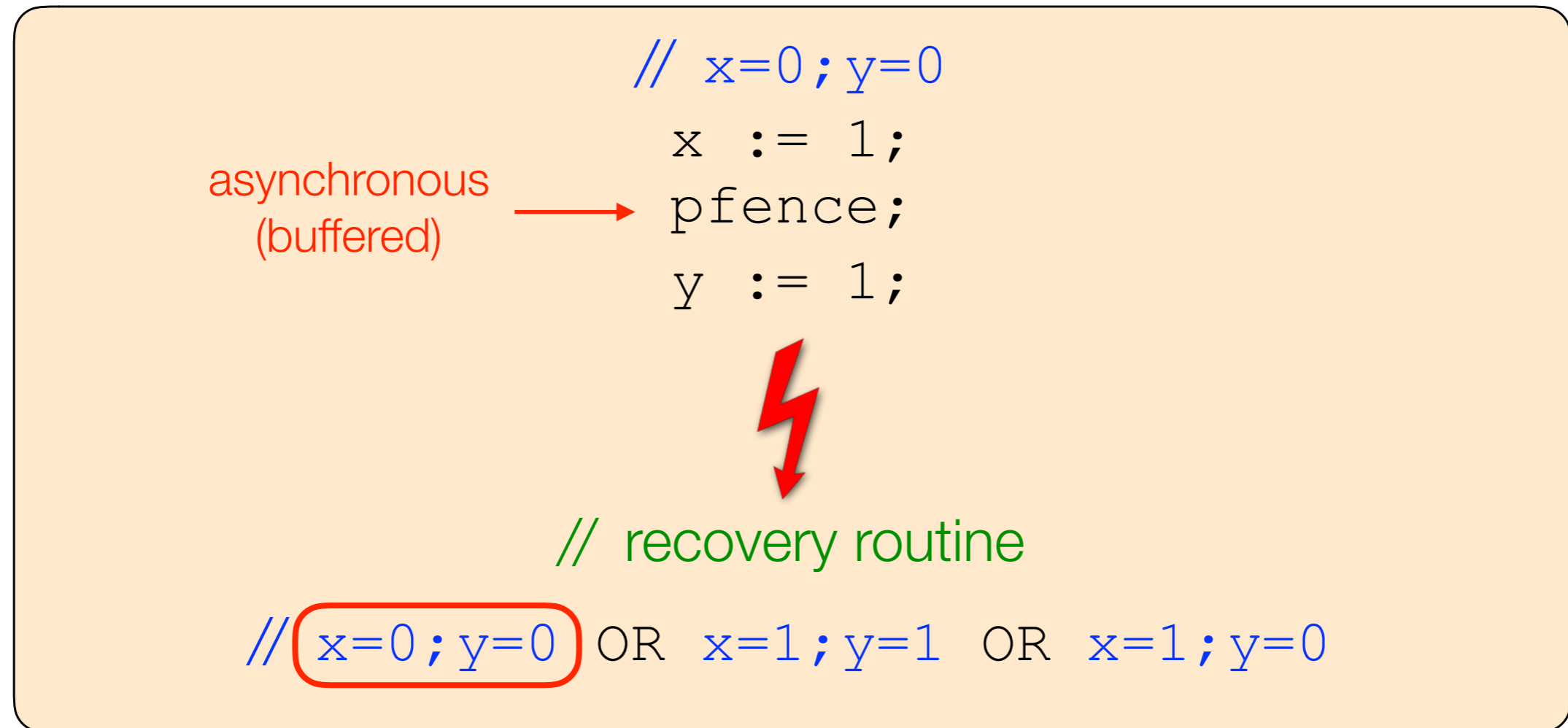
**‼** Execution continues ahead of persistence

# What Can Go Wrong (Continued)?

// x=0; y=0

x := 1;

asynchronous (buffered) →  pfence;

y := 1;

// recovery routine

// (x=0; y=0) OR x=1; y=1 OR x=1; y=0

**‼** Execution continues ahead of persistence

☞ ***persistent sync*** `psync`

# What Can Go Wrong (Continued)?

// x=0; y=0

x := 1;

asynchronous
(buffered) →    pfence;

y := 1;
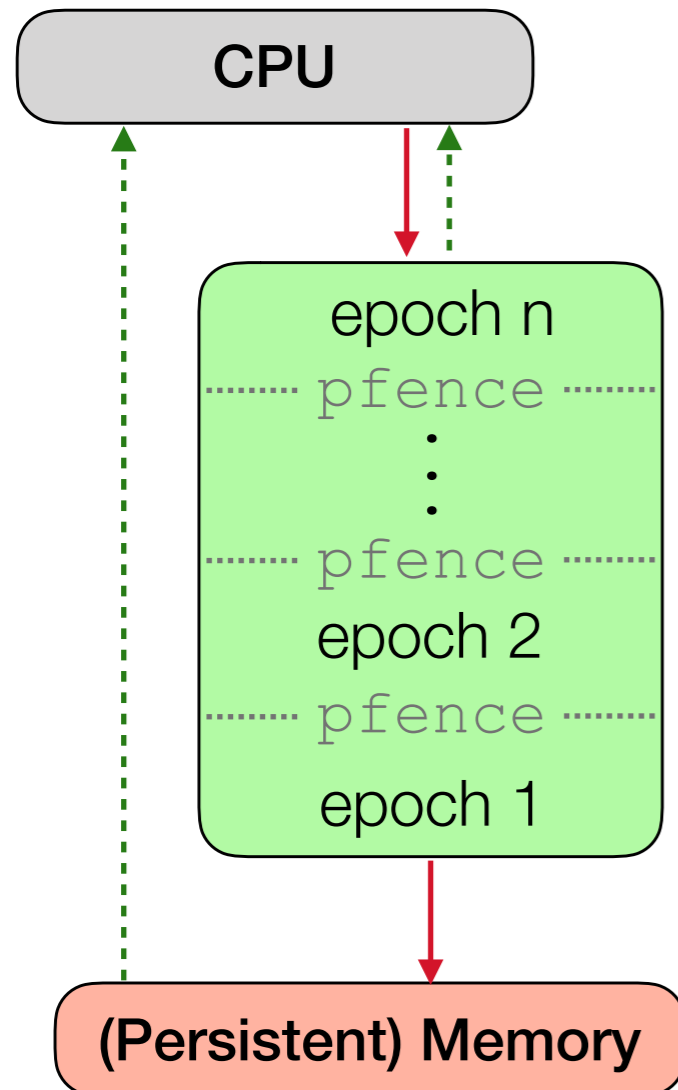
// recovery routine

// x=0; y=0 OR x=1; y=1 OR x=1; y=0

‼️ Execution continues ahead of persistence

☞ **persistent sync** `psync`

`C1; psync; C2`

- same persist-ordering as `pfence`
- `C2` executed only when **all** `C1` writes have persisted

# Persistent Sync

```
// x=0; y=0
      x := 1;
☞     psync;
      y := 1;
```



```
// recovery routine

// x̶=̶0̶;̶y̶=̶0̶ OR  x=1; y=1 OR x=1; y=0
```

**‼️** Execution continues ahead of persistence

☞ ***persistent sync*** `psync`

`C1; psync; C2`

- same persist-ordering as `pfence`
- `C2` executed only when ***all*** `C1` writes have persisted
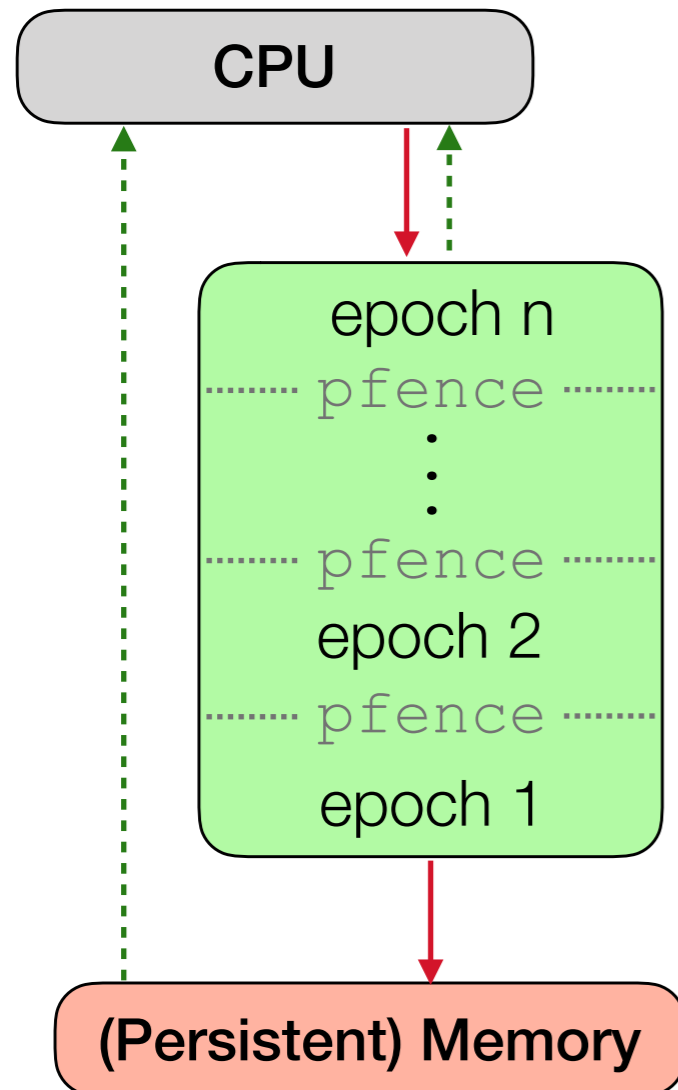
# (Sequential) Hardware



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

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

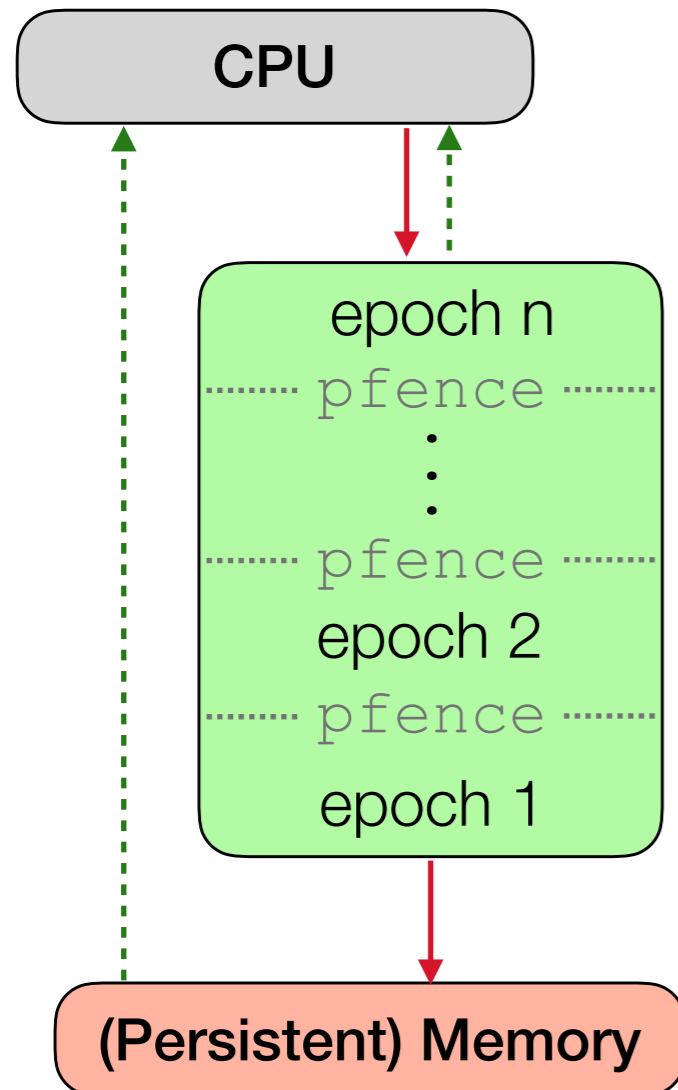p-buffer lost; memory retained

# (Sequential) Hardware



CPU

epoch n
········ pfence ········
⋮
········ pfence ········
epoch 2
········ pfence ········
epoch 1

(Persistent) Memory

`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 (in epoch order)

* at non-deterministic times

11

# (Sequential) Hardware



CPU

epoch n
········ pfence ········
·
·
·
········ pfence ········
epoch 2
pfence
epoch 1

(Persistent) Memory

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

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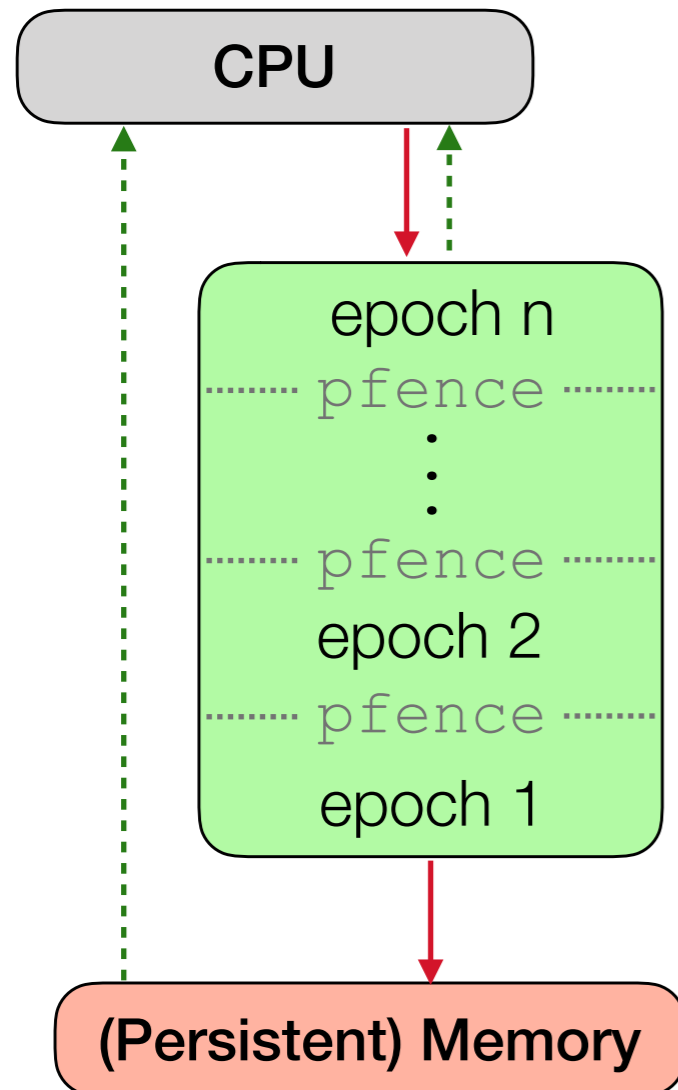else reads from memory

p-buffer lost; memory retained

unbuffer* : p-buffer to memory (in epoch order)

pfence : introduces a new epoch in p-buffer

* at non-deterministic times

# (Sequential) Hardware



**CPU**

epoch n
········· `pfence` ·········
⋮
········· `pfence` ·········
epoch 2
········· `pfence` ·········
epoch 1

**(Persistent) Memory**

$x:=1$ : 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 (in epoch order)

`pfence` : introduces a new epoch in p-buffer

`psync` : flushes the entire p-buffer to memory

\* at non-deterministic times

# What about Concurrency?

TSO

POWER

ARMv8

…

# What about Concurrency?

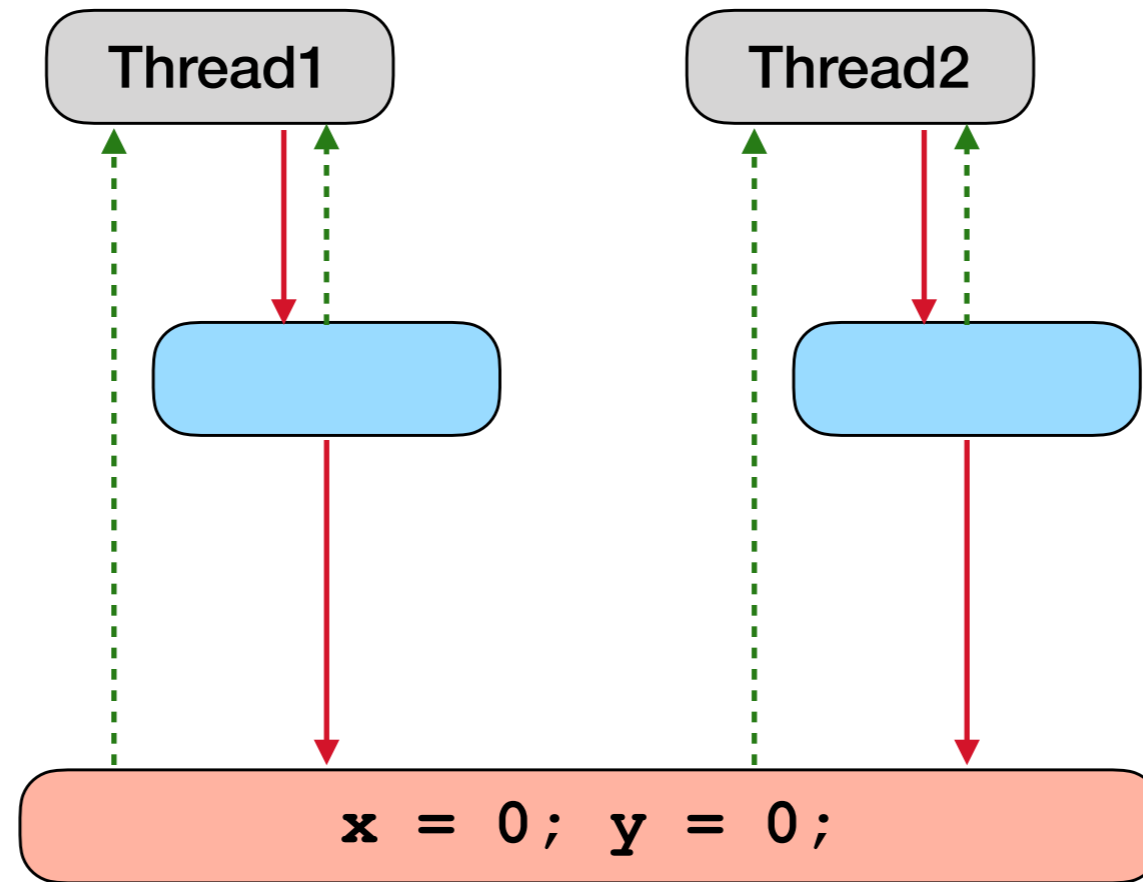TSO · POWER · ARMv8 · ...

# Contributions

# Contributions

- PTSO: <u>First</u> formal epoch persistency semantics under mainstream hardware

  ‣ **Operational** model

  ‣ **Declarative** model

  ‣ **Equivalence** of the two models

# Contributions

- PTSO: <u>First</u> formal epoch persistency semantics under mainstream hardware
  - ‣ ***Operational*** model
  - ‣ ***Declarative*** model
  - ‣ ***Equivalence*** of the two models

# Contributions

- PTSO: <u>First</u> formal epoch persistency semantics under mainstream hardware

  ‣ **Operational** model

  ‣ **Declarative** model

  ‣ **Equivalence** of the two models

- Verifying programs under PTSO

  ‣ PTSO programming **pattern**

  ‣ Correctness condition: **persistent linearisability**

  ‣ Verified several **examples** under PTSO

# Total Store Ordering (TSO)

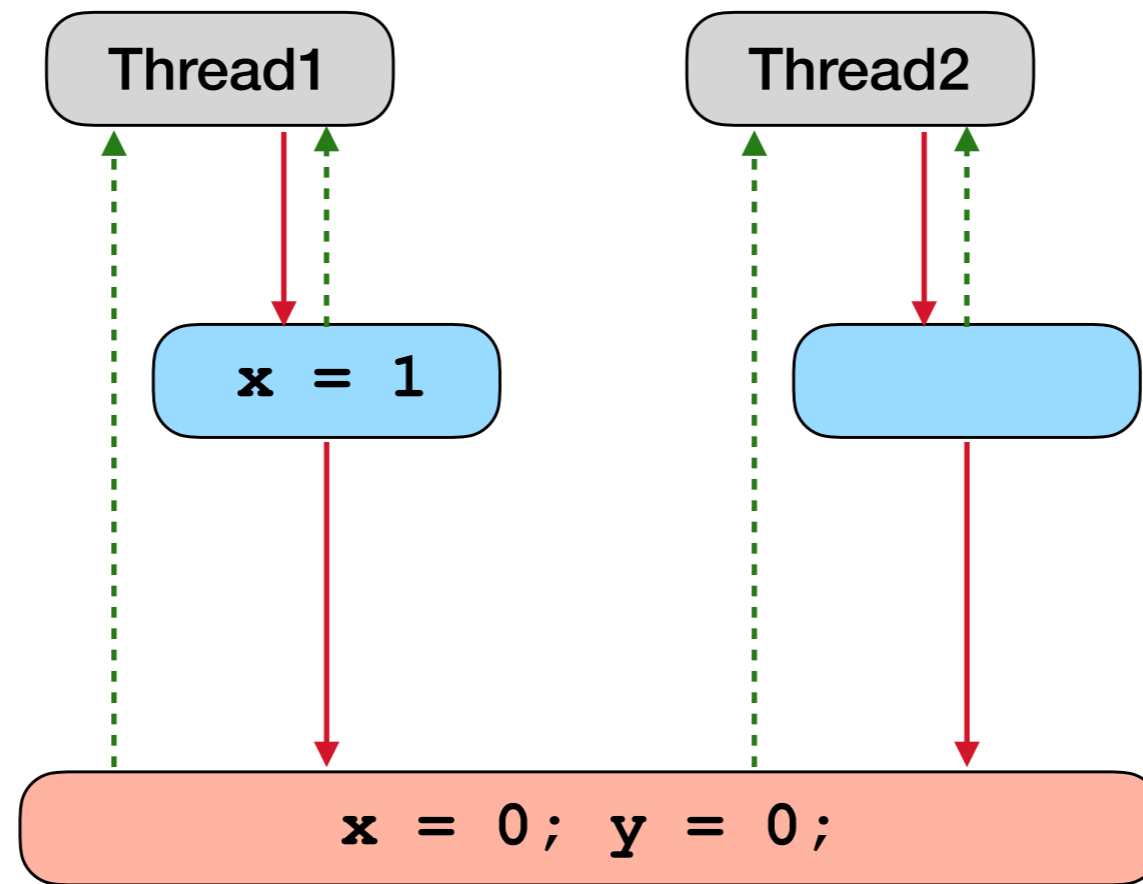# Total Store Ordering (TSO)

# Total Store Ordering (TSO)



Store Buffering (SB)

# Total Store Ordering (TSO)



Store Buffering (SB)

Thread1
```
  x := 1;
  a := y;
```

Thread2
```
  y := 1;
  c := x;
```

# Total Store Ordering (TSO)



Store Buffering (SB)

| Thread1 | Thread2 |
|---------|---------|
| x := 1; | y := 1; |
| a := y; | c := x; |

# Total Store Ordering (TSO)



Thread1    Thread2

**x = 1**    **y = 1**

**x = 0; y = 0;**

Thread1    Thread2

```
x := 1;          y := 1;
a := y;          c := x;
```

Store Buffering (SB)

# Total Store Ordering (TSO)



Store Buffering (SB)

# Total Store Ordering (TSO)



Store Buffering (SB)

| Thread1 | Thread2 |
|---------|---------|
| x := 1; | y := 1; |
| a := y; // 0 | c := x; |

# Total Store Ordering (TSO)



Store Buffering (SB)

```
Thread1                    Thread2

x := 1;                    y := 1;
a := y;  // 0        ☞     c := x;
☞
```

# Total Store Ordering (TSO)



Store Buffering (SB)

# Total Store Ordering (TSO)



Store Buffering (SB)

# Total Store Ordering (TSO)

Thread1    Thread2

`y = 1`

`x = 1; y = 0;`

Thread1

Thread2

```
x := 1;              y := 1;
a := y;  // 0        c := x;  // 0
```

Store Buffering (SB)

# Total Store Ordering (TSO)



Thread1    Thread2
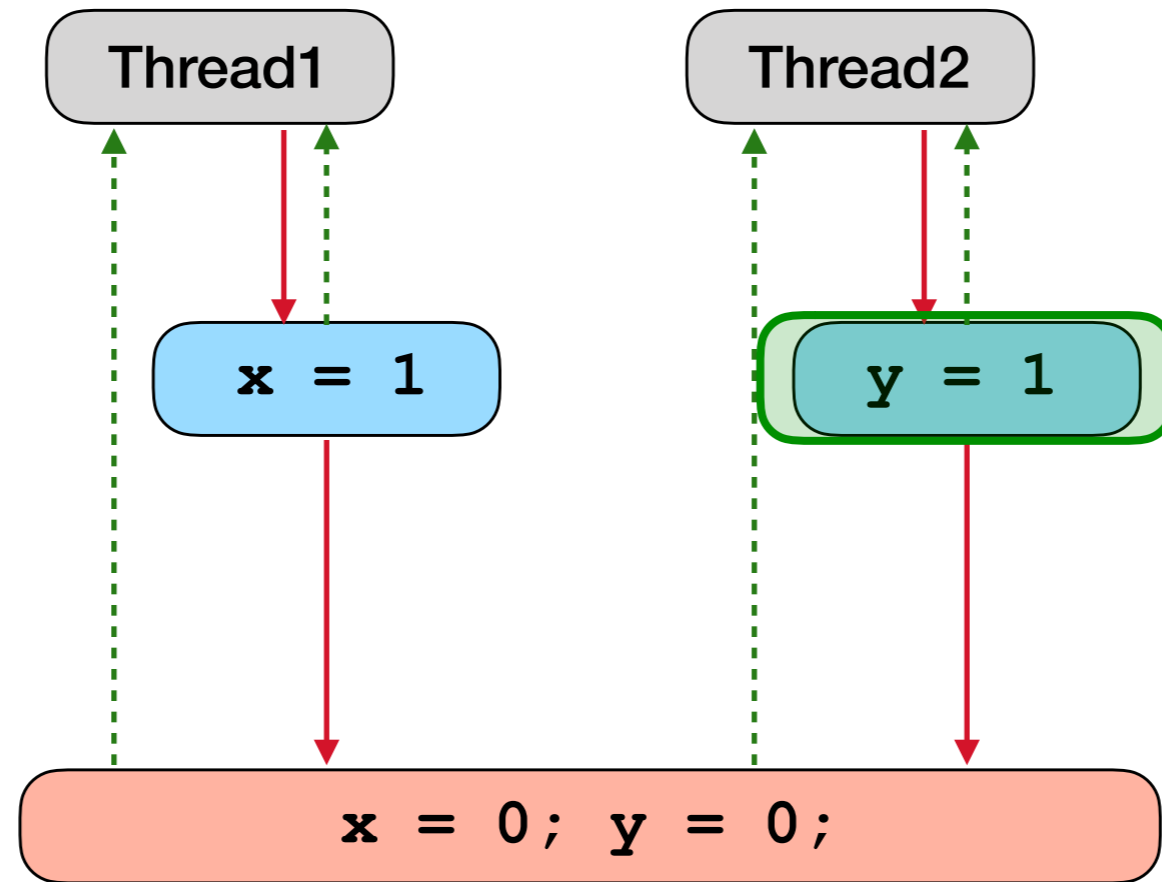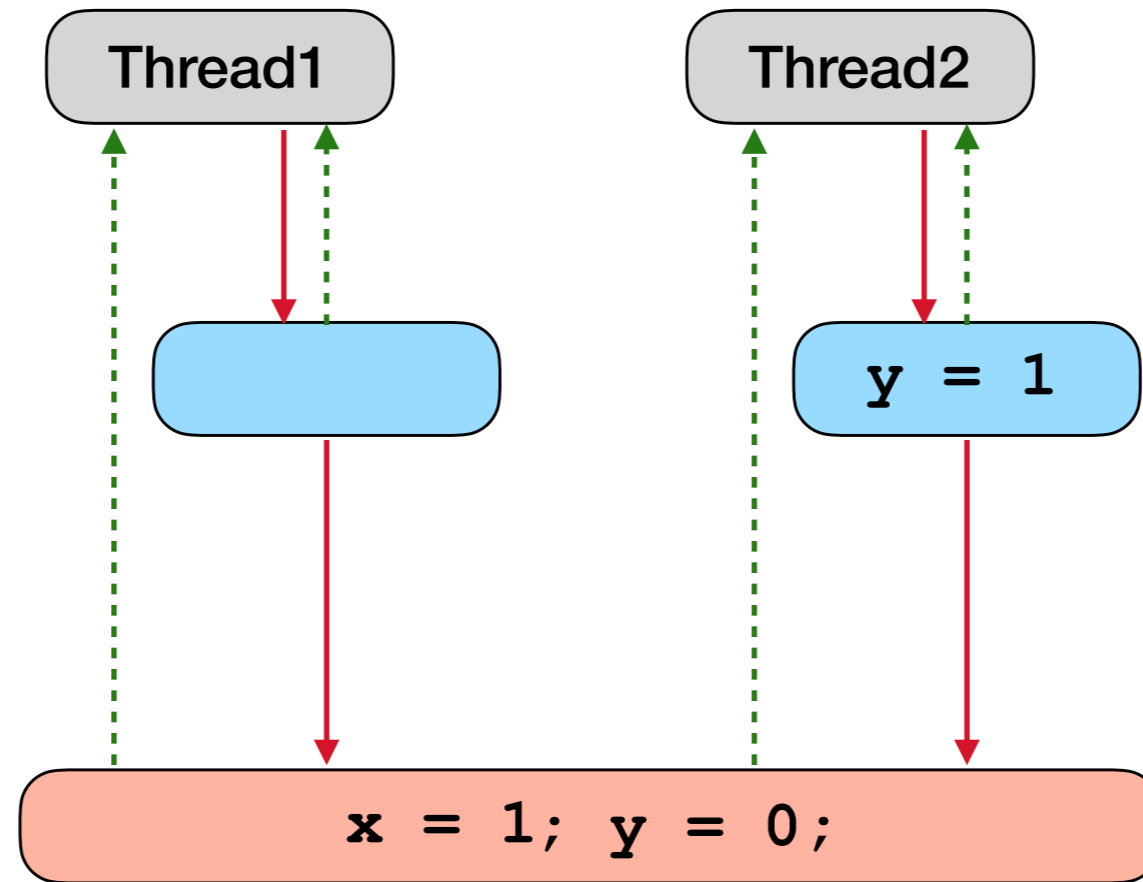
x = 1; y = 1;

Thread1

x := 1;
a := y;  // 0

Thread2

y := 1;
c := x;  // 0

Store Buffering (SB)

# Persistent TSO (PTSO)

# Persistent TSO (PTSO)

# Persistent TSO (PTSO)

# Contributions

- PTSO: <u>First</u> formal epoch persistency semantics under mainstream hardware
  - ▸ ***Operational*** model
  - ▸ ***Declarative*** model
  - ▸ ***Equivalence*** of the two models

# Contributions

- PTSO: <u>First</u> formal epoch persistency semantics under mainstream hardware
  - ‣ *Operational* model
  - ‣ *Declarative* model
  - ‣ *Equivalence* of the two models

- Verifying programs under PTSO

  - ‣ PTSO programming *pattern*
  - ‣ Correctness condition: *persistent linearisability*
  - ‣ Verified several *examples* under PTSO

# Contributions

- PTSO: <u>First</u> formal epoch persistency semantics under mainstream hardware
  - ‣ *Operational* model
  - ‣ *Declarative* model
  - ‣ *Equivalence* of the two models

- Verifying programs under PTSO

  - ‣ PTSO programming *pattern*
  - ‣ Correctness condition: *persistent linearisability*
  - ‣ Verified several *examples* under PTSO

# Verifying programs under PTSO

```
1. q.enq(v) ≜
2.   pc:=getPC(); t:=getTC();
3.   n:=newNode(v,t,pc);
4.   map[t][pc]:=n; pfence;
5.   lock(q); h:=q.head;
6.   while (q.data[h] != null)
7.     h:=h+1;
8.   q.data[h]:=n;
9.   pfence; unlock(q);

10. q.deq() ≜
11.  pc:=getPC(); t:=getTC();
12.  lock(q);h:=q.head;n:=q.data[h];
13.  map[t][pc]:=n;
14.  if (n != null) {
15.    t':=n.t; pc':=n.pc;
16.    map[t'][pc']:=⊤ }
17.  pfence;
18.  if (n != null) {
19.    q.head:=h+1; pfence; }
20.  unlock(q); return n;

21. lock(q) ≜
22.  while (!CAS(q.lock,0,1)) skip;

23. unlock(q) ≜ q.lock:=0;

24. isIn(q,n) ≜
25.  h:=q.head; c:=q.data[h];
26.  while (c != null) {
27.    if (n==c) return true;
28.    else { h:=h+1; c:=q.data[h]; }
29.  } return false;

30. getProgress(t) ≜
31.  pc:=-1; n:=⊥;
32.  while (map[t][pc+1] !=⊥){
33.    pc++; n:=map[t][pc]; }
34.  return (pc,n);
```

```
35. start() ≜
36.  lq:=newQueue();
37.  s:=P.size; lmap:=newMap(s);
38.  for (t in P)
39.    lmap[t]:=newArray(P[t].size,⊥);
40.  pfence;
41.  q:=lq; map:=lmap; run(P);

42. recover() ≜
43.  if (q==null || map==null)
44.    goto start();
45.  for(t in P) enq[t]:=-1;
46.  unlock(q);
47.  for(t in P) { // deq recovery
48.    (pc,n):=getProgress(t);
49.    if (pc>=0 && isDeq(P[t][pc])) {
50.      if (n==null)
51.        P'[t]:=sub(P[t],pc+1);
52.      else {
53.        if (inIn(q,n))
54.          P'[t]:=sub(P[t],pc);
55.        else
56.          P'[t]:=sub(P[t],pc+1);
57.        t':=n.t; pc':=n.pc;
58.        enq[t']:=max(enq[t'],pc'+1);}
59.    }
60.    else if (pc<0) P'[t]:=P[t];
61.  }
62.  for(t in P) { // enq recovery
63.    (pc,n):=getProgress(t);
64.    if (pc>=0 && isEnq(P[t][pc])) {
65.      if (pc < enq[t])
66.        P'[t]:=sub(P[t],enq[t]);
67.      else if (n==⊤ || isIn(q,n))
68.        P'[t]:=sub(P[t],pc+1);
69.      else
70.        P'[t]:=sub(P[t],pc); }
71.  } run(P');
```

The **persistent** variant of the Michael-Scott queue and its **recovery** mechanism

```
1. q.enq(v) ≜
2.   pc:=getPC(); t:=getTC();
3.   n:=newNode(v,t,pc);
4.   map[t][pc]:=n; pfence;
5.   lock(q); h:=q.head;
6.   while (q.data[h] != null)
7.     h:=h+1;
8.   q.data[h]:=n;
9.   pfence; unlock(q);

10. q.deq() ≜
11.   pc:=getPC(); t:=getTC();
12.   lock(q); h:=q.head; n:=q.data[h];
13.   map[t][pc]:=n;
```

```
35. start() ≜
36.   lq:=newQueue();
37.   s:=P.size; lmap:=newMap(s);
38.   for(t in P)
39.     lmap[t]:=newArray(P[t].size,⊥);
40.   pfence;
41.   q:=lq; map:=lmap; run(P);

42. recover() ≜
43.   if (q==null || map==null)
44.     goto start();
45.   for(t in P) enq[t]:=-1;
46.   unlock(q);
47.   for(t in P) { // deq recovery
```

## What constitutes a *correct persistent* implementation?
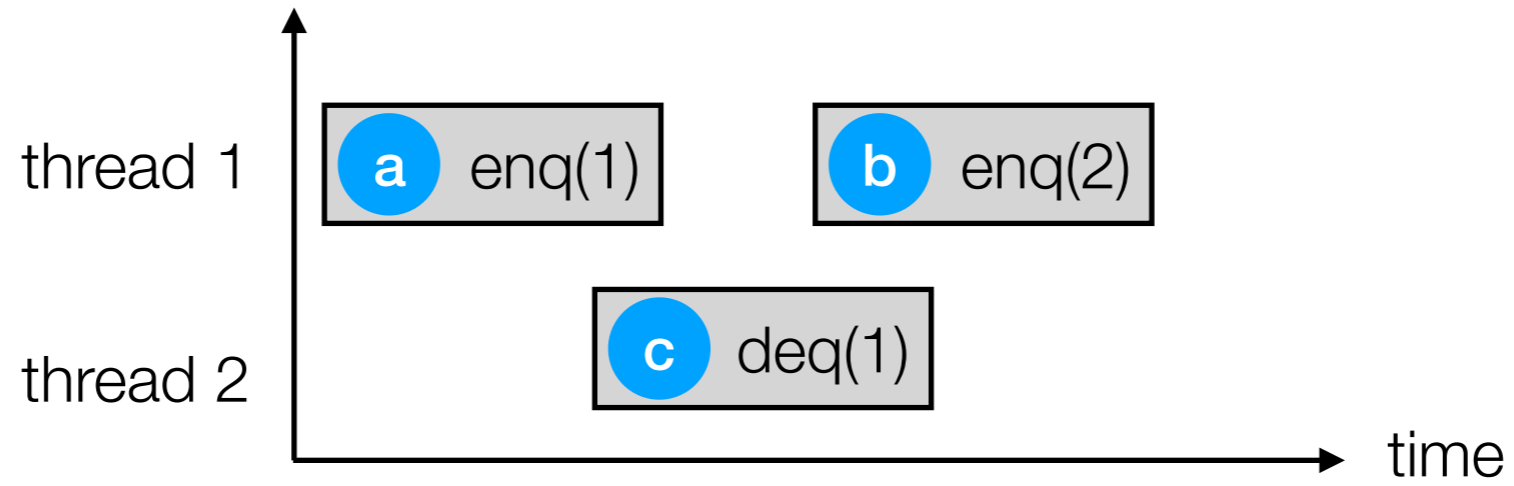
```
23. unlock(q) ≜ q.lock:=0;

24. isIn(q,n) ≜
25.   h:=q.head; c:=q.data[h];
26.   while (c != null) {
27.     if (n==c) return true;
28.     else { h:=h+1; c:=q.data[h]; }
29.   } return false;

30. getProgress(t) ≜
31.   pc:=-1; n:=⊥;
32.   while (map[t][pc+1] != ⊥) {
33.     pc++; n:=map[t][pc]; }
34.   return (pc,n);
```

```
58.       enq[t']:=max(enq[t'],pc'+1);}
59.     }
60.     else if (pc<0) P'[t]:=P[t];
61.   }
62.   for(t in P) { // enq recovery
63.     (pc,n):=getProgress(t);
64.     if (pc>=0 && isEnq(P[t][pc])) {
65.       if (pc < enq[t])
66.         P'[t]:=sub(P[t],enq[t]);
67.       else if (n==⊤ || isIn(q,n))
68.         P'[t]:=sub(P[t],pc+1);
69.       else
70.         P'[t]:=sub(P[t],pc); }
71.   } run(P');
```

The *persistent* variant of the Michael-Scott queue and its *recovery* mechanism

# Linearisability



thread 1    **a** enq(1)      **b** enq(2)

thread 2    **c** deq(1)

time

# Linearisability



- Define happens-before relation *hb*
  - ▸ $(e_1, e_2) \in hb \iff e_1.end <_{time} e_2.begin$

# Linearisability



- Define happens-before relation *hb*
  - ▸ $(e_1, e_2) \in hb \iff e_1.end <_{\textbf{time}} e_2.begin$
  
    -- e.g. $(a, b) \in hb$      $(a, c) \notin hb$

# Linearisability



- Define happens-before relation *hb*

  ▸ $(e_1, e_2) \in$ *hb* $\iff e_1$.end $<_{\text{time}}$ $e_2$.begin

  -- e.g. ( **a** , **b** ) $\in$ *hb*      ( **a** , **c** ) $\notin$ *hb*

- ***Linearisable*** $\iff \exists$ *H. H* totally orders events

  ▸ *H* respects *hb*

  ▸ *H* is a ***legal*** sequence (library-specific)

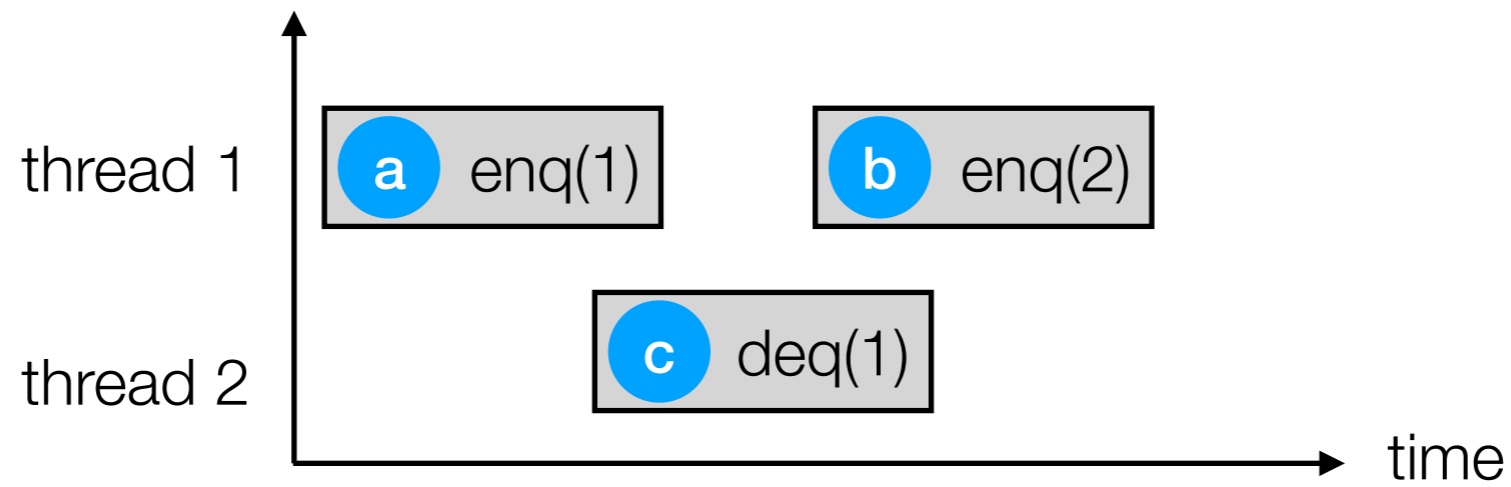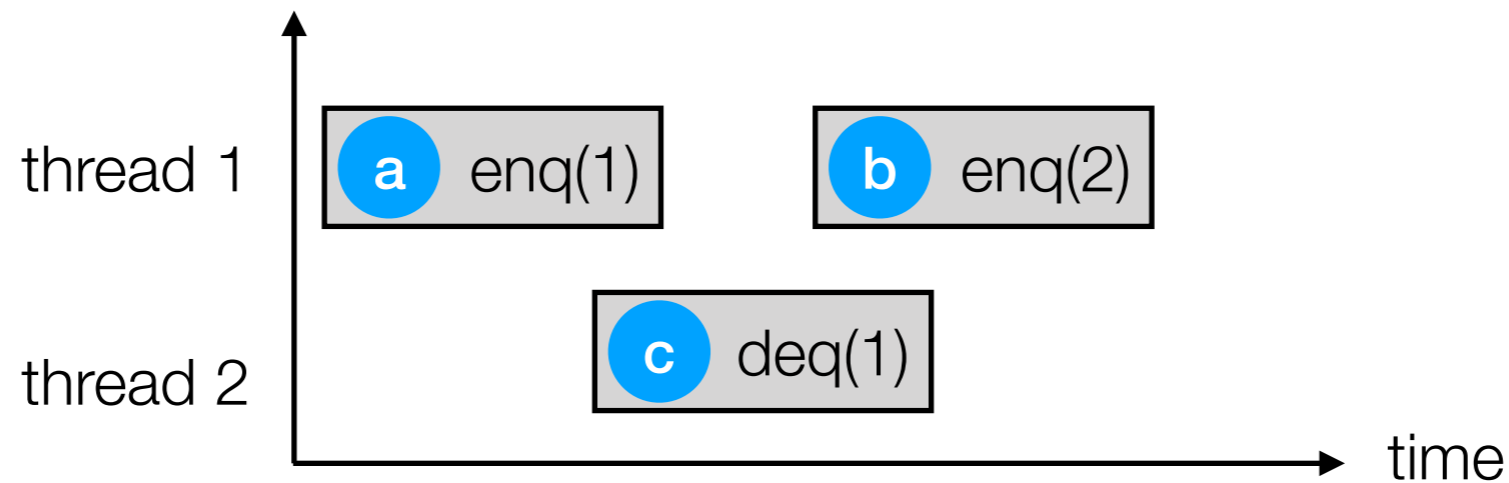# Linearisability



- Define happens-before relation *hb*
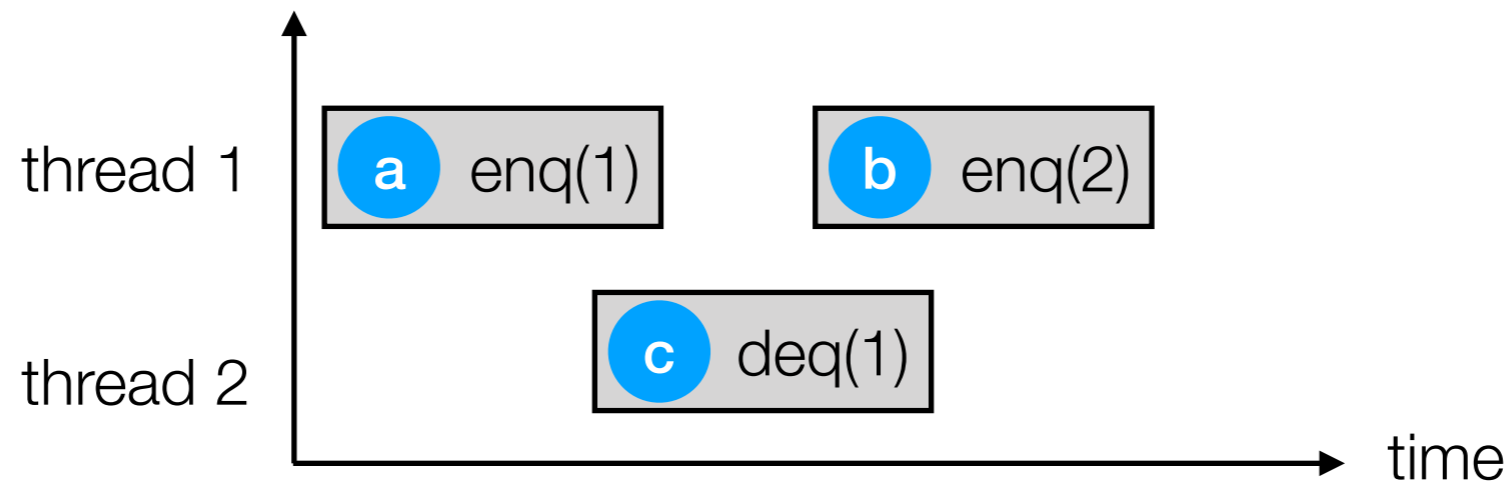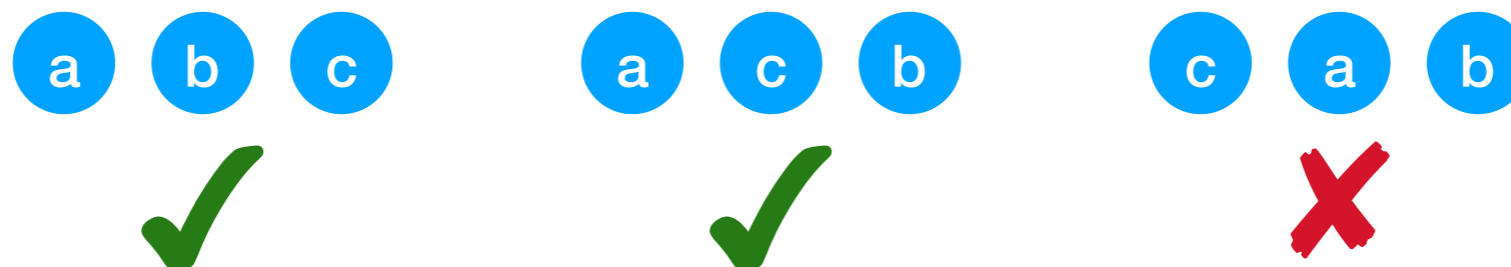  - ‣ $(e_1, e_2) \in hb \iff e_1.\text{end} <_{\text{time}} e_2.\text{begin}$

    -- e.g. $(\text{a}, \text{b}) \in hb$ $\qquad (\text{a}, \text{c}) \notin hb$

- ***Linearisable*** $\iff \exists H.\ H$ totally orders events
  - ‣ *H* respects *hb*
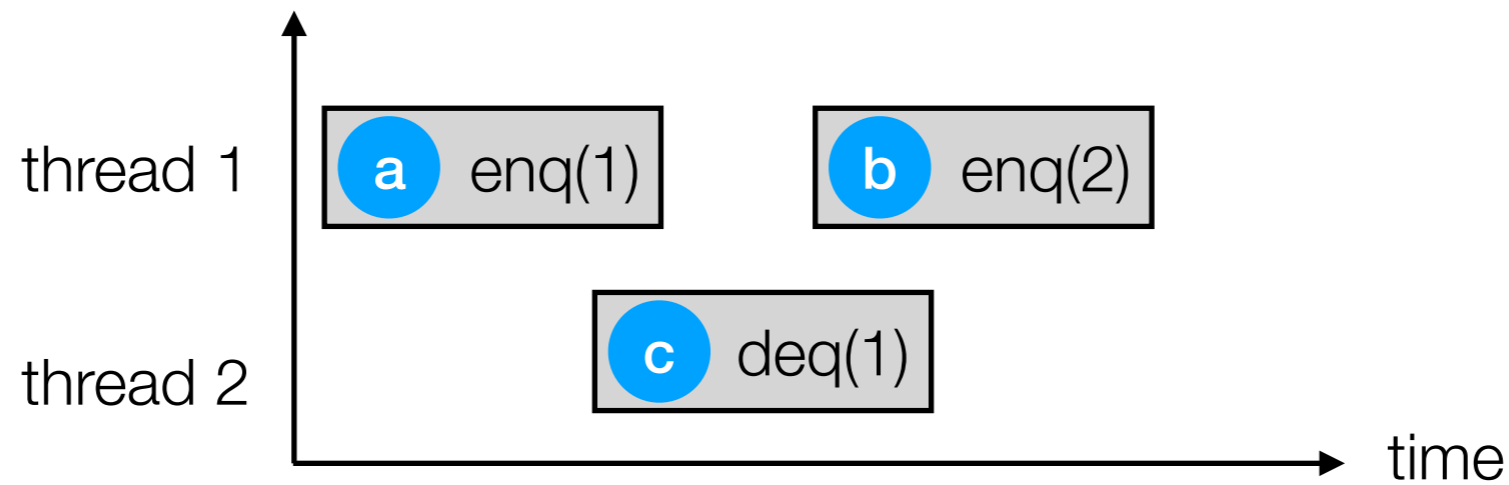  - ‣ *H* is a ***legal*** sequence (library-specific)

    -- e.g. FIFO sequences for queue

# Linearisability



- Define happens-before relation *hb*

  ‣ $(e_1, e_2) \in$ *hb* $\iff e_1$.end $<_{\textbf{time}}$ $e_2$.begin

  -- e.g. ( **a** , **b** ) $\in$ *hb*        ( **a** , **c** ) $\notin$ *hb*

- ***Linearisable*** $\iff \exists$ *H. H* totally orders events

  ‣ *H* respects *hb*

  ‣ *H* is a ***legal*** sequence (library-specific)
    -- e.g. FIFO sequences for queue

linearisable

**a** **b** **c**    **a** **c** **b**    **c** **a** **b**

✔        ✔        ✘

# Linearisability



- Define happens-before relation *hb*

  ‣ $(e_1, e_2) \in hb \iff e_1.end <_{time} e_2.begin$

    -- e.g. ( a , b ) ∈ *hb*    ( a , c ) ∉ *hb*

- **Linearisable** $\iff \exists H.$ *H* totally orders events

  ‣ *H* respects *hb*

  ‣ *H* is a **legal** sequence (library-specific)
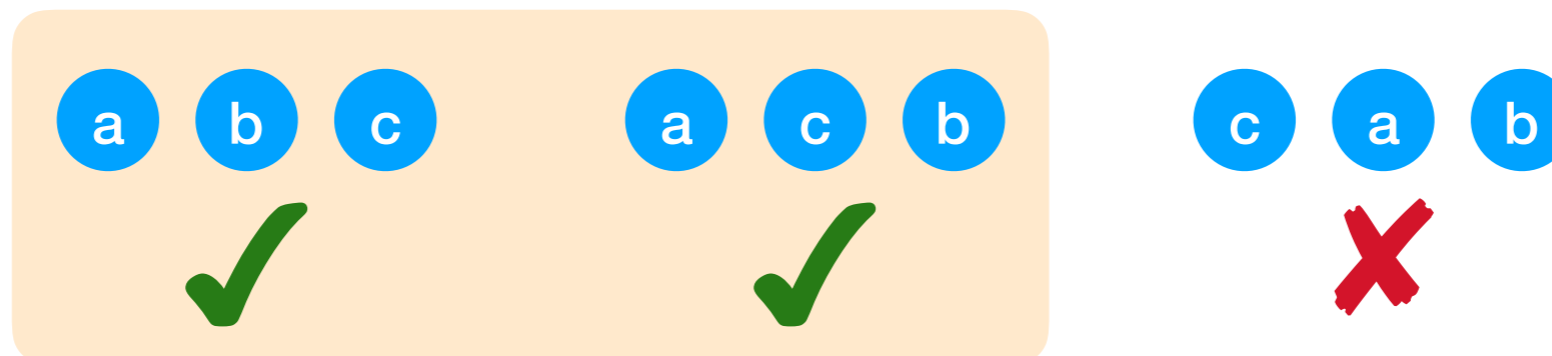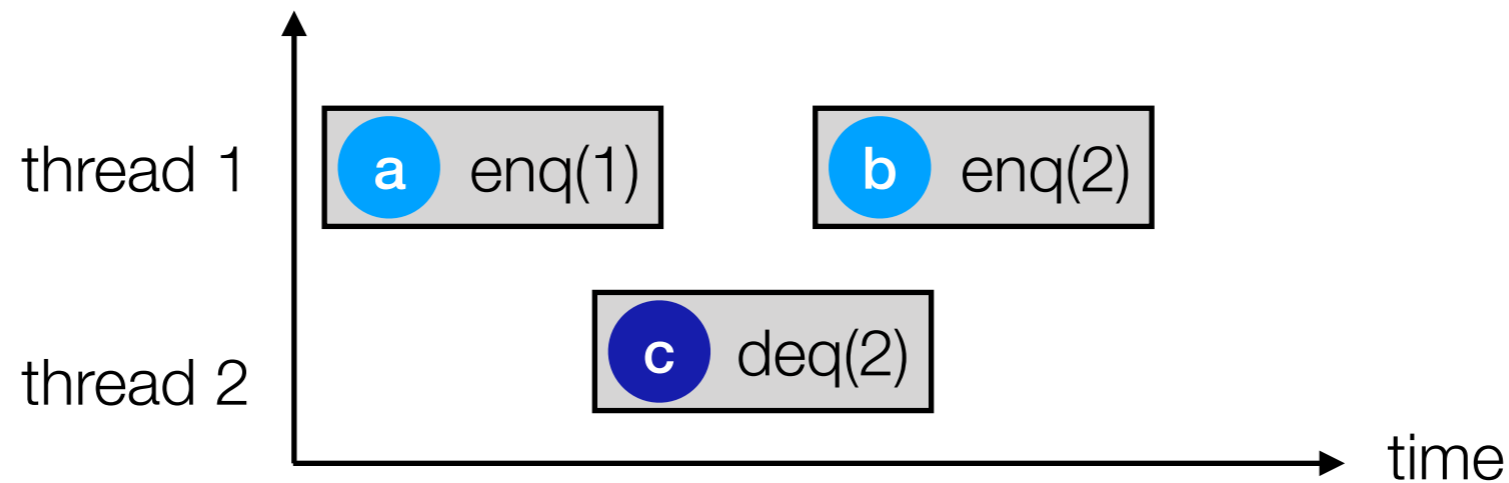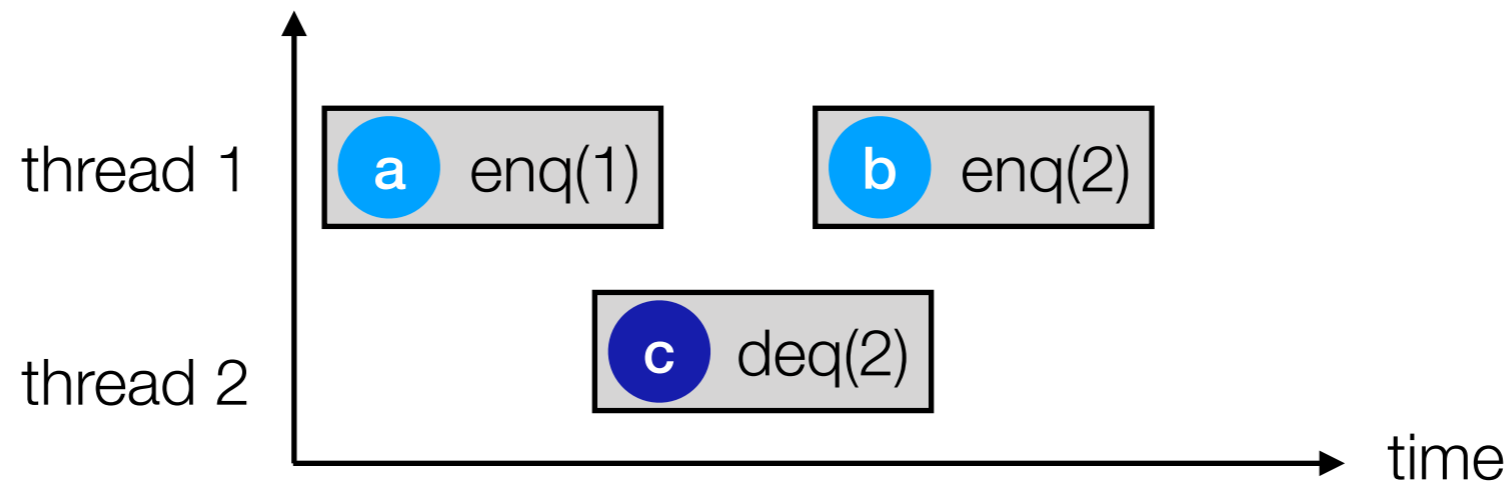    -- e.g. FIFO sequences for queue

# Linearisability



- Define happens-before relation *hb*
  - $(e_1, e_2) \in hb \iff e_1.\text{end} <_{\text{time}} e_2.\text{begin}$

    -- e.g. $(a, b) \in hb \qquad (a, c) \notin hb$

- **Linearisable** $\iff \exists H.\; H$ totally orders events
  - *H* respects *hb*
  - *H* is a **legal** sequence (library-specific)
    -- e.g. FIFO sequences for queue

# Linearisability



thread 1 — a enq(1) — b enq(2)

thread 2 — c deq(2)

time

- Define happens-before relation *hb*
  - ▸ $(e_1, e_2) \in hb \iff e_1.\text{end} <_{\text{time}} e_2.\text{begin}$
    
    -- e.g. ( a , b ) $\in hb$     ( a , c ) $\notin hb$

- ***Linearisable*** $\iff \exists H. H$ totally orders events
  - ▸ *H* respects *hb*
  - ▸ *H* is a ***legal*** sequence (library-specific)
    
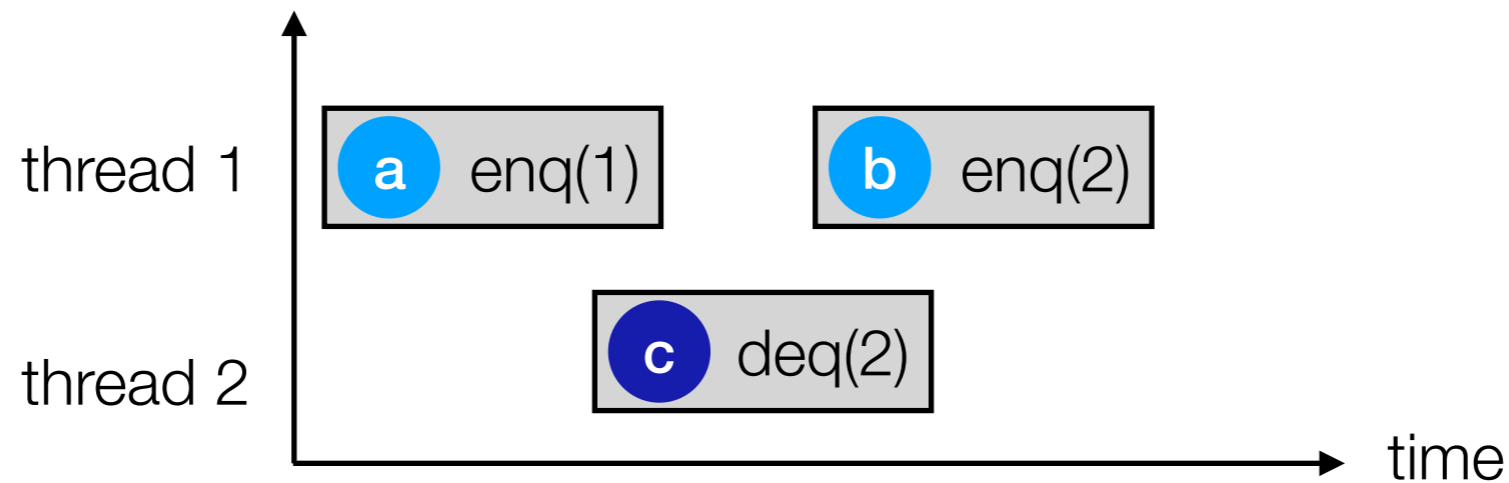    -- e.g. FIFO sequences for queue

non-linearisable
(not legal)
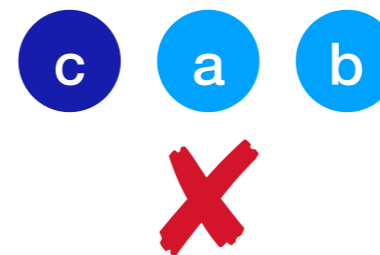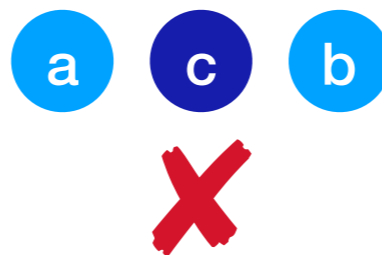
a b c ✗     a c b ✗     c a b ✗

# Persistent Linearisability



- Define happens-before relation *hb*
  - ▸ $(e_1, e_2) \in hb \iff e_1.end <_{\textbf{time}} e_2.begin$
    -- e.g. $(\text{a}, \text{b}) \in hb$      $(\text{a}, \text{c}) \notin hb$
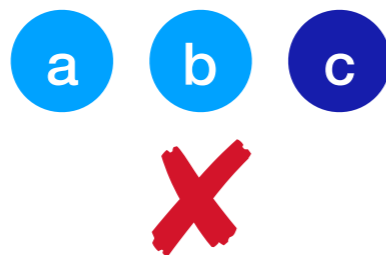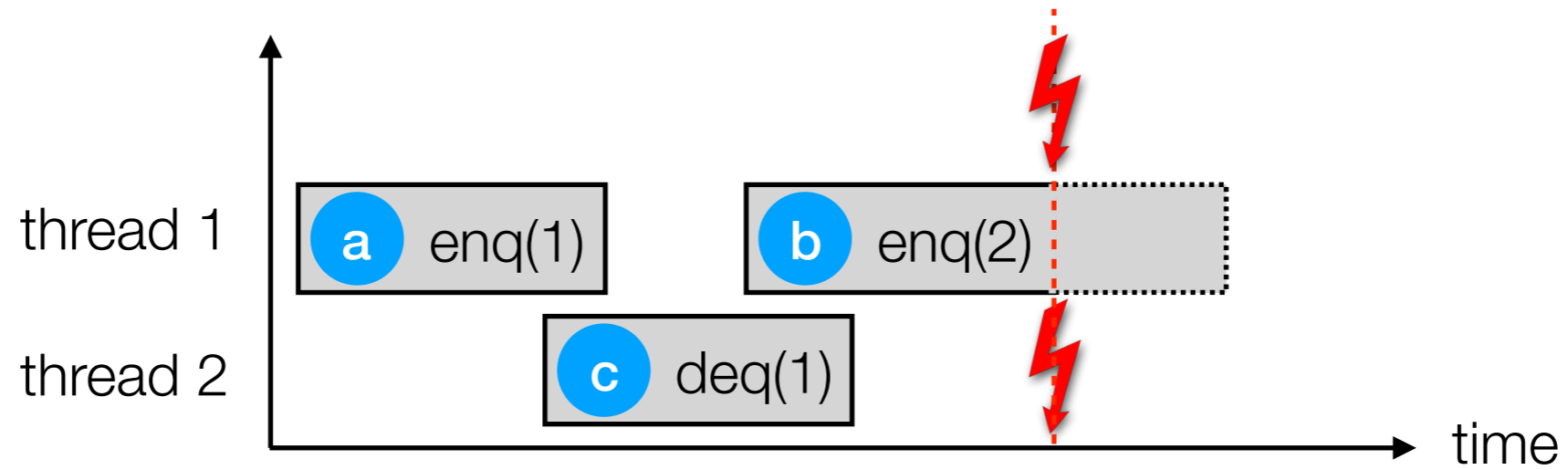
# Persistent Linearisability



- Define happens-before relation *hb*

  ‣ $(e_1, e_2) \in hb \iff e_1.\text{end} <_{\text{time}} e_2.\text{begin}$

  -- e.g. ( a , b ) $\in hb$    ( a , c ) $\notin hb$

- ***Persistently linearisable*** $\iff \exists$ *H. H* totally orders a ***subset*** *S* of events

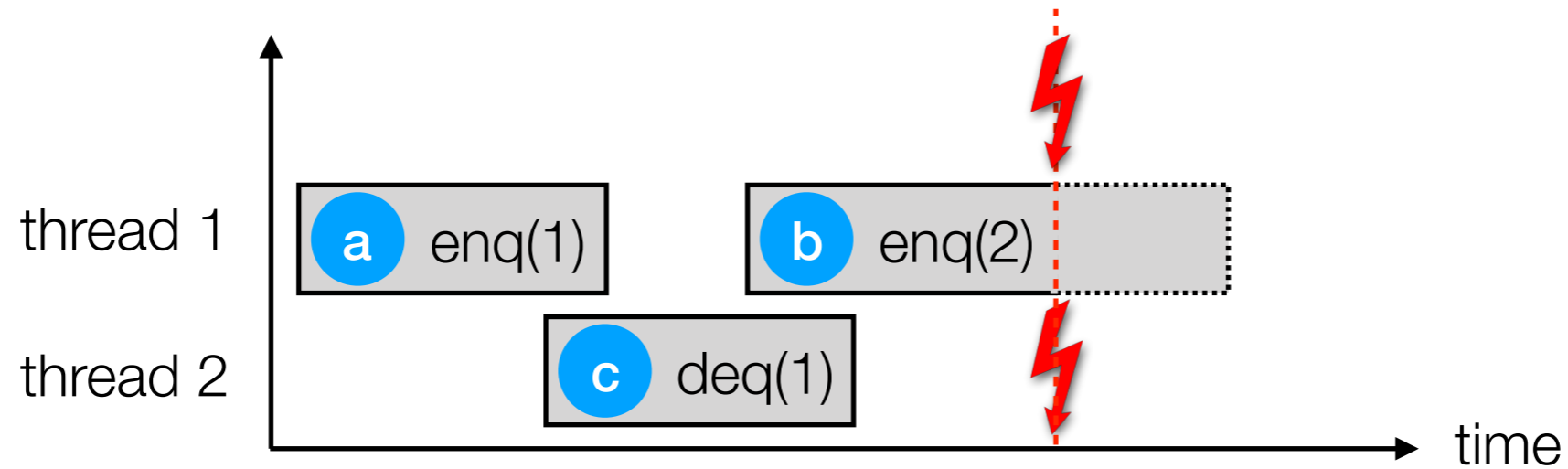  ‣ *H* respects *hb*

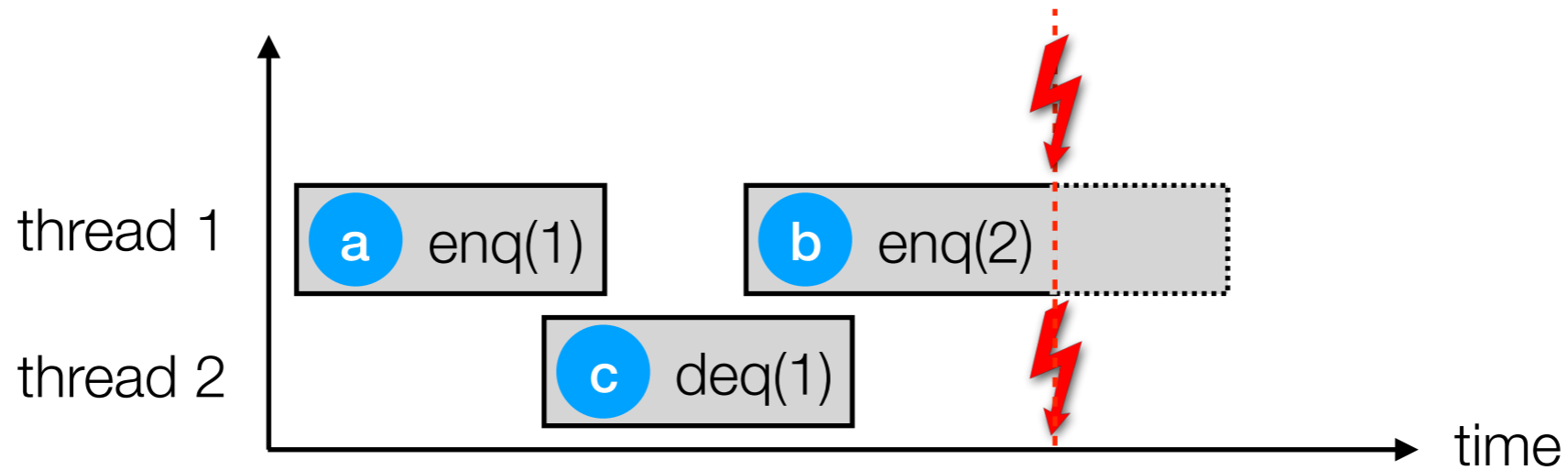  ‣ *H* is a legal sequence

# Persistent Linearisability



- Define happens-before relation *hb*

  ▸ $(e_1, e_2) \in hb \iff e_1.\text{end} <_{\text{time}} e_2.\text{begin}$

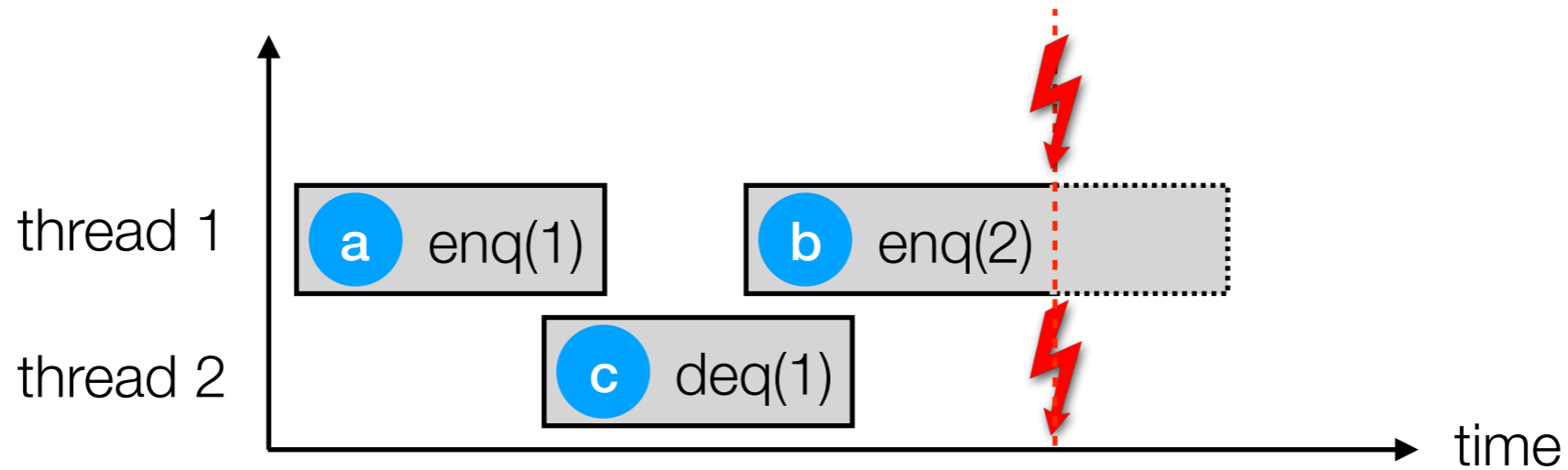    -- e.g. $(a, b) \in hb$       $(a, c) \notin hb$

- **Persistently linearisable** $\iff \exists H. H$ totally orders a **subset** *S* of events

  ▸ *H* respects *hb*

  ▸ *H* is a legal sequence

  ▸ *S* is *hb*-**prefix-closed** : $(a, b) \in hb$   and   $b \in S \implies a \in S$

    -- persists are **asynchronous**: only a **prefix** may persist after a crash

20

# Persistent Linearisability



- Define happens-before relation *hb*
  - ▸ $(e_1, e_2) \in hb \iff e_1.\text{end} <_{\text{time}} e_2.\text{begin}$

    -- e.g. ( a , b ) $\in hb$    ( a , c ) $\notin hb$

- ***Persistently linearisable*** $\iff \exists$ *H. H* totally orders a ***subset*** *S* of events
  - ▸ *H* respects *hb*
  - ▸ *H* is a legal sequence
  - ▸ *S* is *hb*-***prefix-closed*** : $(a, b) \in hb$  and  $b \in S \implies a \in S$

    -- persists are ***asynchronous***: only a ***prefix*** may persist <u>after a crash</u>
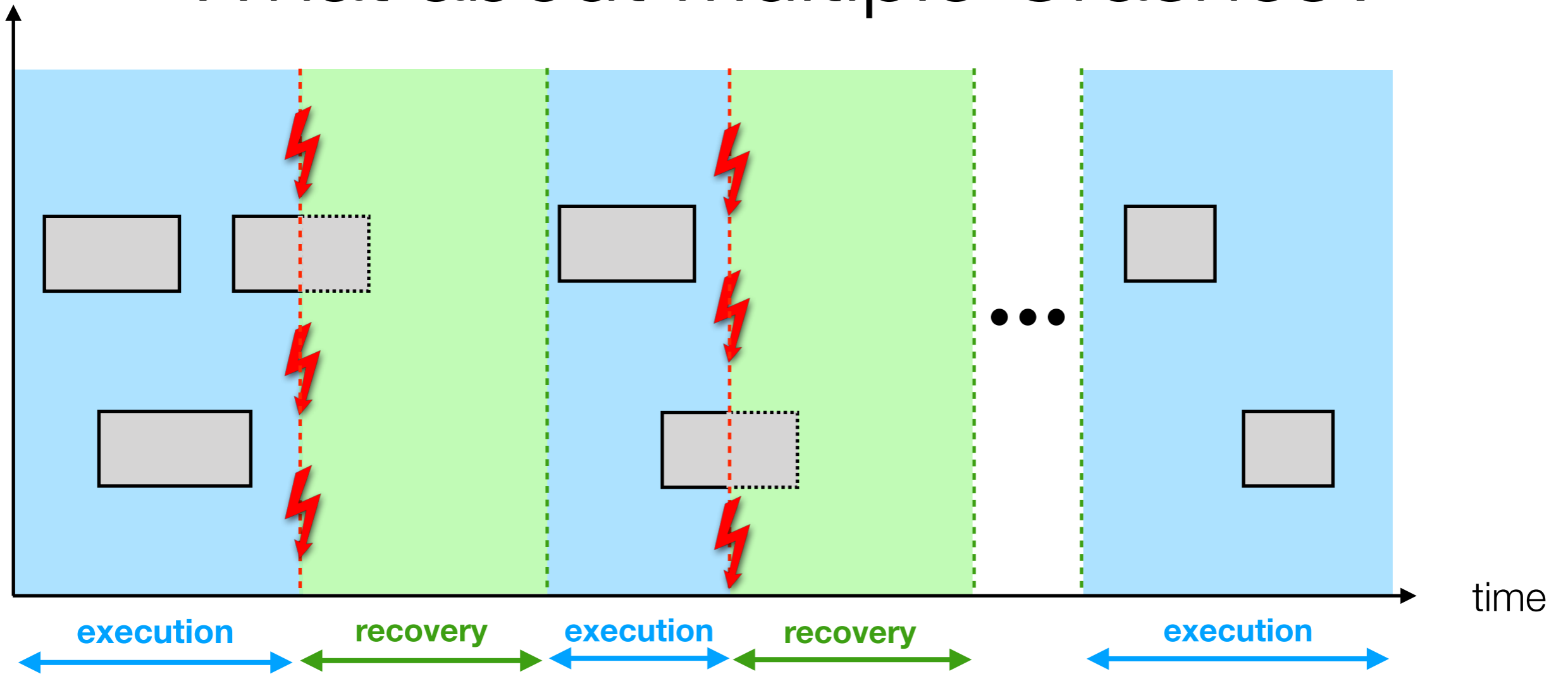
a c

Persistently linearisable ✔

20

# What about Multiple Crashes?



execution    recovery    execution    recovery    execution

time

# What about Multiple Crashes?



no crashes

execution    recovery    execution    recovery    execution

time

21

# What about Multiple Crashes?



no crashes

execution    recovery    execution    recovery    execution

$G_1$    $G_2$    $G_n$

time

# What about Multiple Crashes?



- A **chain** $G_1 \cdots G_n$ is **persistently linearisable** $\iff \exists\, H_1 \cdots H_n$.
  - ‣ $H_i$ persistently linearises $G_i$ — as before
  - ‣ $H_1 ++ \cdots ++ H_n$ is a legal sequence

# Conclusions

- PTSO: <u>First</u> formal epoch persistency semantics under mainstream hardware
  - ‣ **_Operational_** model
  - ‣ **_Declarative_** model
  - ‣ **_Equivalence_** of the two models

- Verifying programs under PTSO

  - ‣ PTSO programming **_pattern_**
  - ‣ Correctness condition: **_persistent linearisability_**
  - ‣ Verified several **_examples_** under PTSO

# Conclusions

- PTSO: <u>First</u> formal epoch persistency semantics under mainstream hardware
  - ‣ ***Operational*** model
  - ‣ ***Declarative*** model
  - ‣ ***Equivalence*** of the two models

- Verifying programs under PTSO

  - ‣ PTSO programming ***pattern***
  - ‣ Correctness condition: ***persistent linearisability***
  - ‣ Verified several ***examples*** under PTSO

## Thank you for listening!

# Programming Pattern

```
1. // log progress
2. pfence
3. // do the work
4. pfence
```

# Programming Pattern

```
1. // log progress
2. pfence
3. // do the work
4. pfence
```

Log *at most one step* ahead of work

# Programming Pattern

```
1. // log progress
2. pfence
3. // do the work
4. pfence
```

Log *at most one step* ahead of work

```
q.enq(v) ≜
1.   pc:=getPC(); t:=getTC();
     n:=newNode(v,t,pc,null);
     map[t][pc]:=n;
2.   pfence;
3.   h:=q.head;
     find: while (q.data[h] != null)
       h:=h+1;
     if (!CAS(q.data[h],null,n))
       goto find;
4.   pfence;
```