

Specifying & Verifying Non-Volatile Memory

Azalea Raad

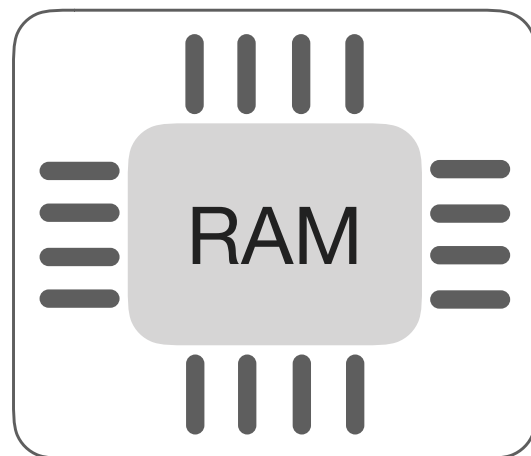
Max Planck Institute for Software Systems (MPI-SWS)

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

Computer Storage

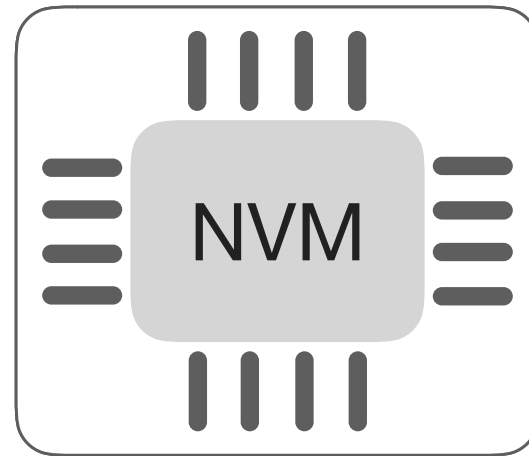


✓ *fast*
✗ *volatile*



✗ *slow*
✓ *persistent*

What is Non-Volatile Memory (NVM)?



NVM: Hybrid Storage + Memory

Best of both worlds:

✓ ***persistent*** (like HDD)

✓ ***fast, random access*** (like RAM)

Why NVM **Now**?

✓ **Persistent**

✓ **Higher capacity** (than RAM)

✓ **Green**

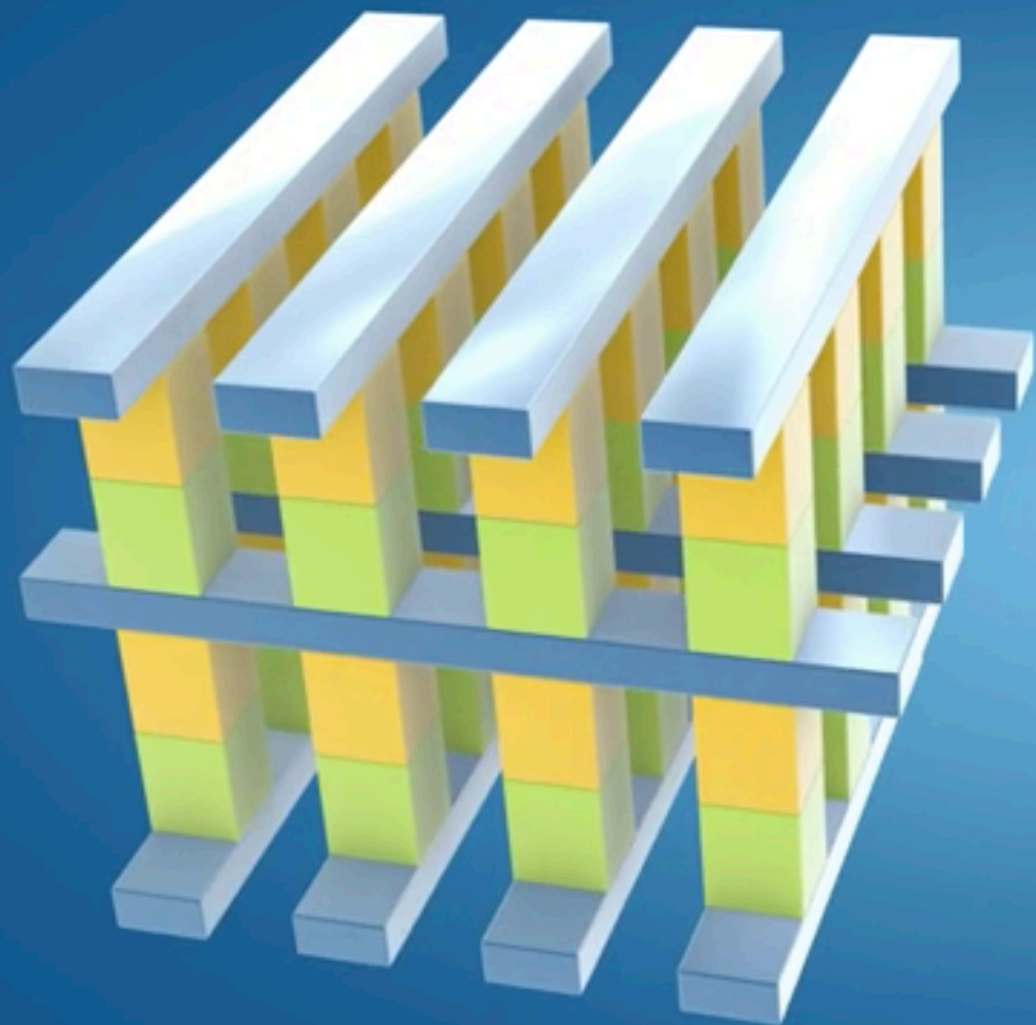
- ▶ 32x capacity for 3x power consumption

✓ **Low latency – no intermediaries**

- ▶ no OS/FS intermediaries
- ▶ no paging, no context switching, no interrupts, no kernel code
- ▶ short instruction path

✓ Many **applications**

- ▶ fast in-memory databases
- ▶ file systems
- ▶ ...



INTEL® OPTANE™ TECHNOLOGY



FAST

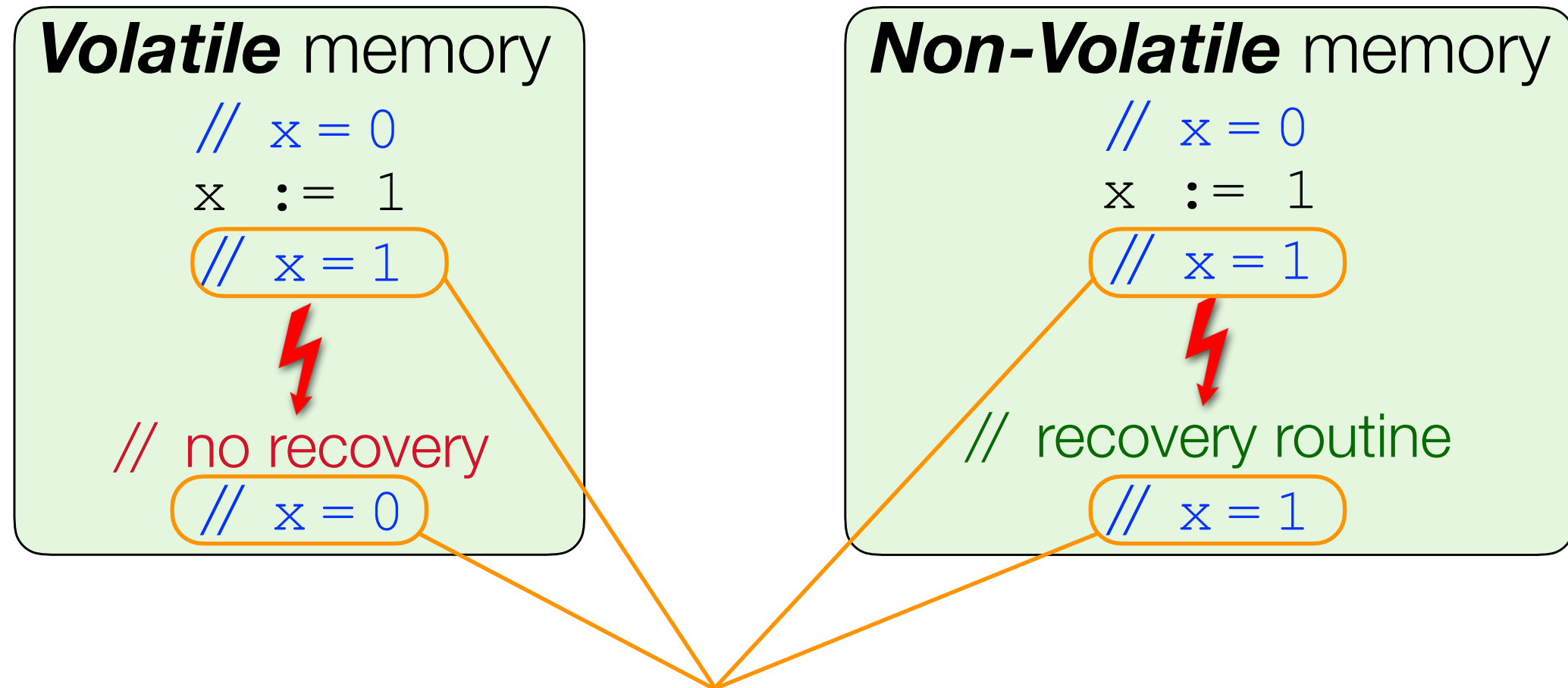


DENSE



NON-VOLATILE

Q: Why **Formal** NVM Semantics?



A: Program **Verification**

Q: Why ***Formal*** NVM Semantics?

What about ***Concurrency***?

// $x = y = \dots = 0$

$C_1 \parallel C_2 \parallel \dots \parallel C_n$

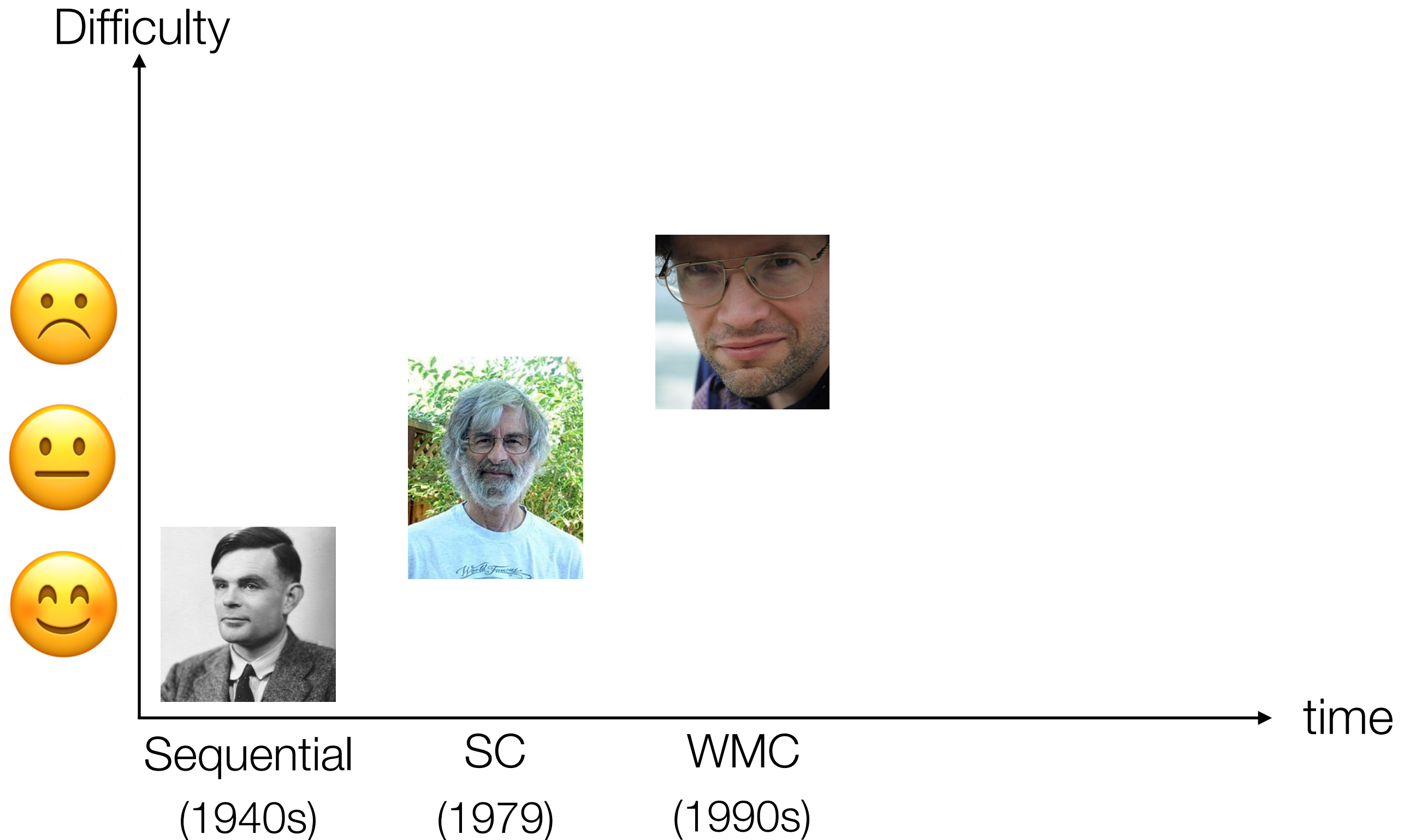
// ???



// recovery routine

// ???

Formal Semantic Models



Weak Memory Consistency (WMC)

No total execution order (*to*) \Rightarrow

weak behaviour absent under SC, caused by:

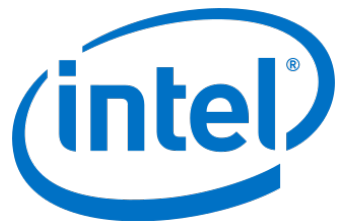
- **software**

instruction ***reordering*** by compiler



- **hardware**

write propagation across ***cache hierarchy***



Weak Memory Consistency (WMC)

No total execution order (*to*) \Rightarrow

weak be

- soft

Consistency Model

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

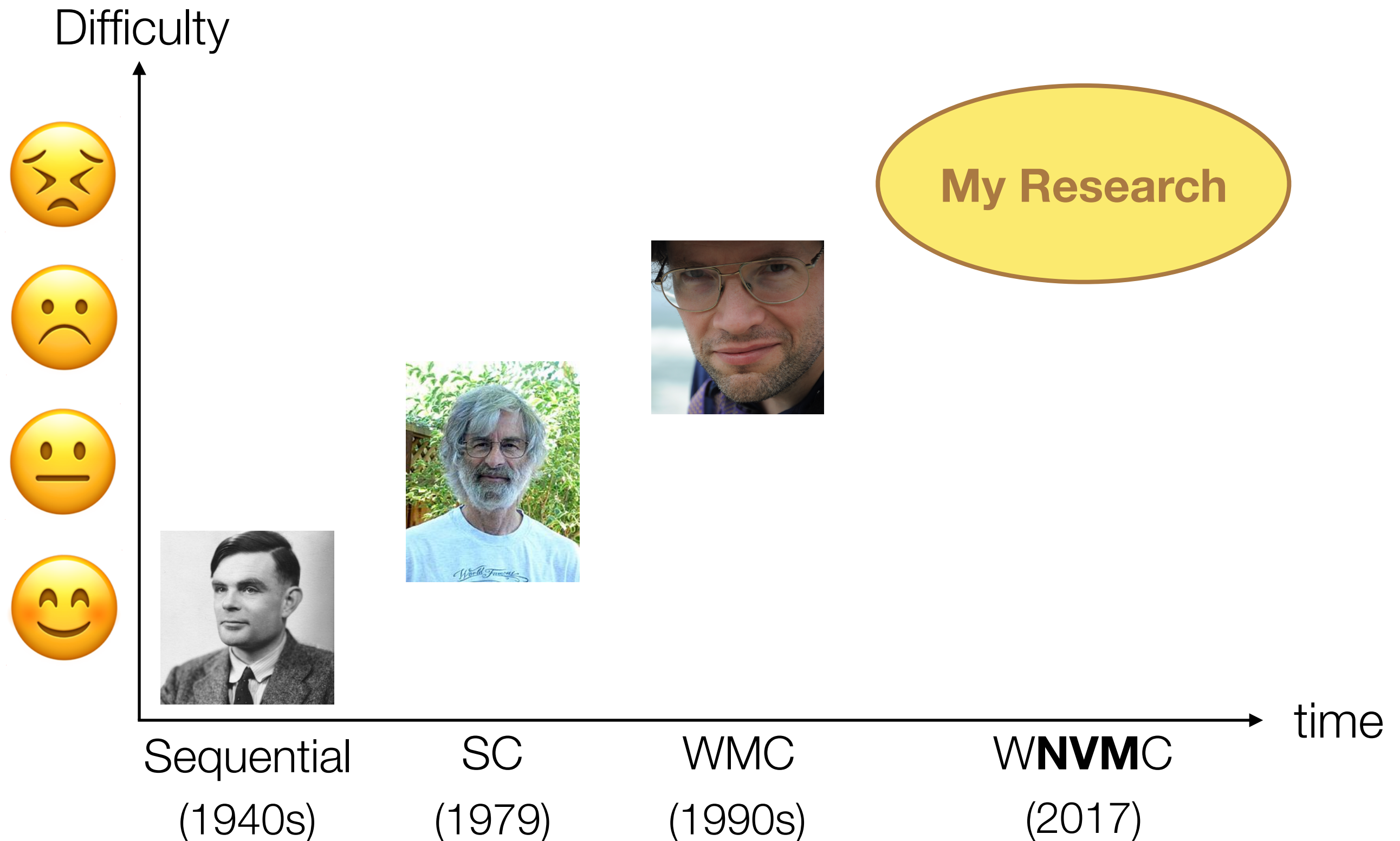
- hard

e.g. x86-TSO, ARM, POWER, C11, Java

intel arm



Formal Semantic Models



What Can Go Wrong?

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

```
x := 1 ;
```

```
y := 1 ;
```



```
// recovery routine
```

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

!! Execution continues ***ahead of persistence***
— ***asynchronous*** persists

What Can Go Wrong?

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

```
x  := 1 ;
```

```
y  := 1 ;
```



```
// recovery routine
```

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

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

— ***asynchronous*** persists

!! Writes may persist ***out of order***

— ***relaxed*** persists

What Can Go Wrong?

Consistency Model

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

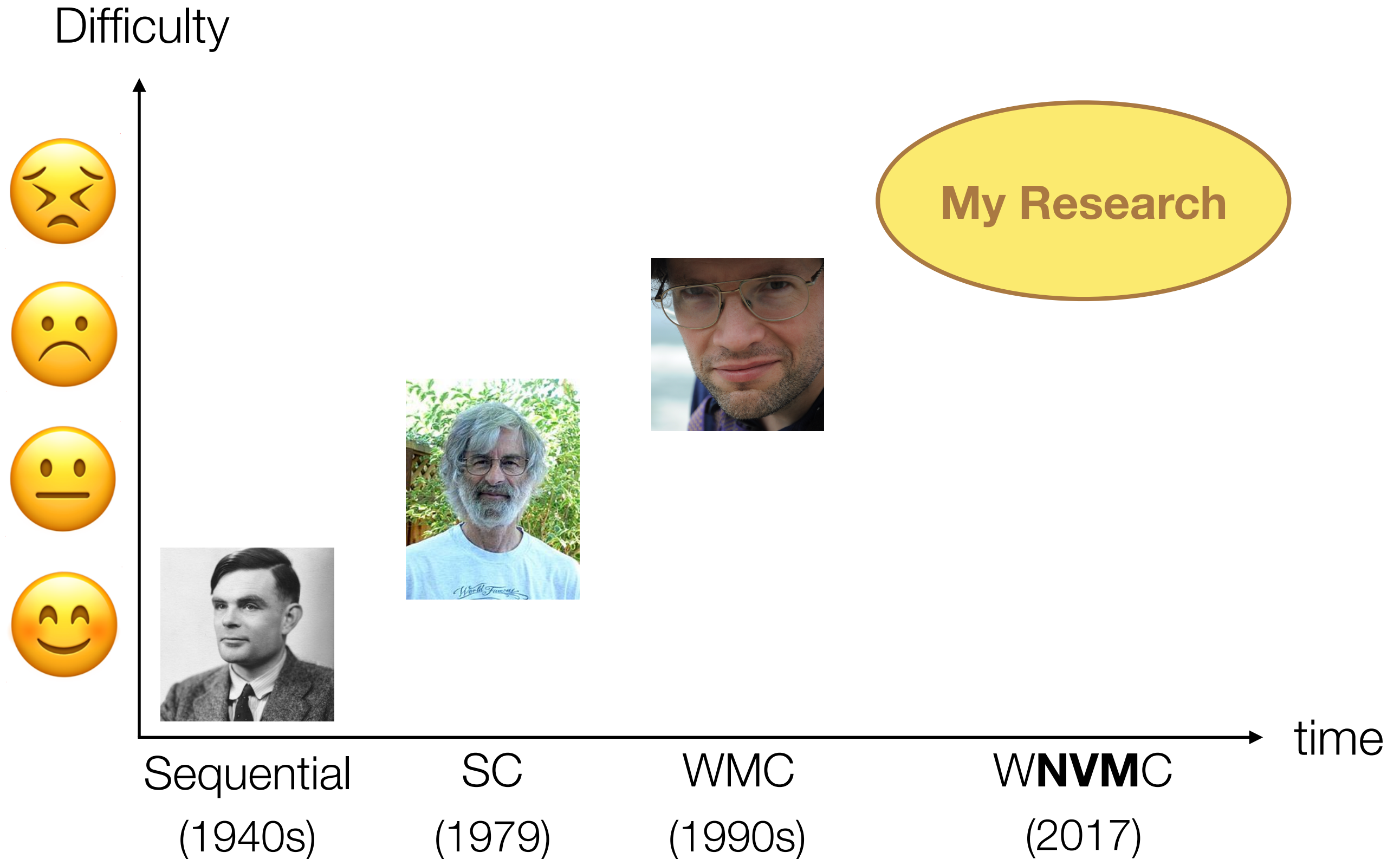
Persistency Model

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

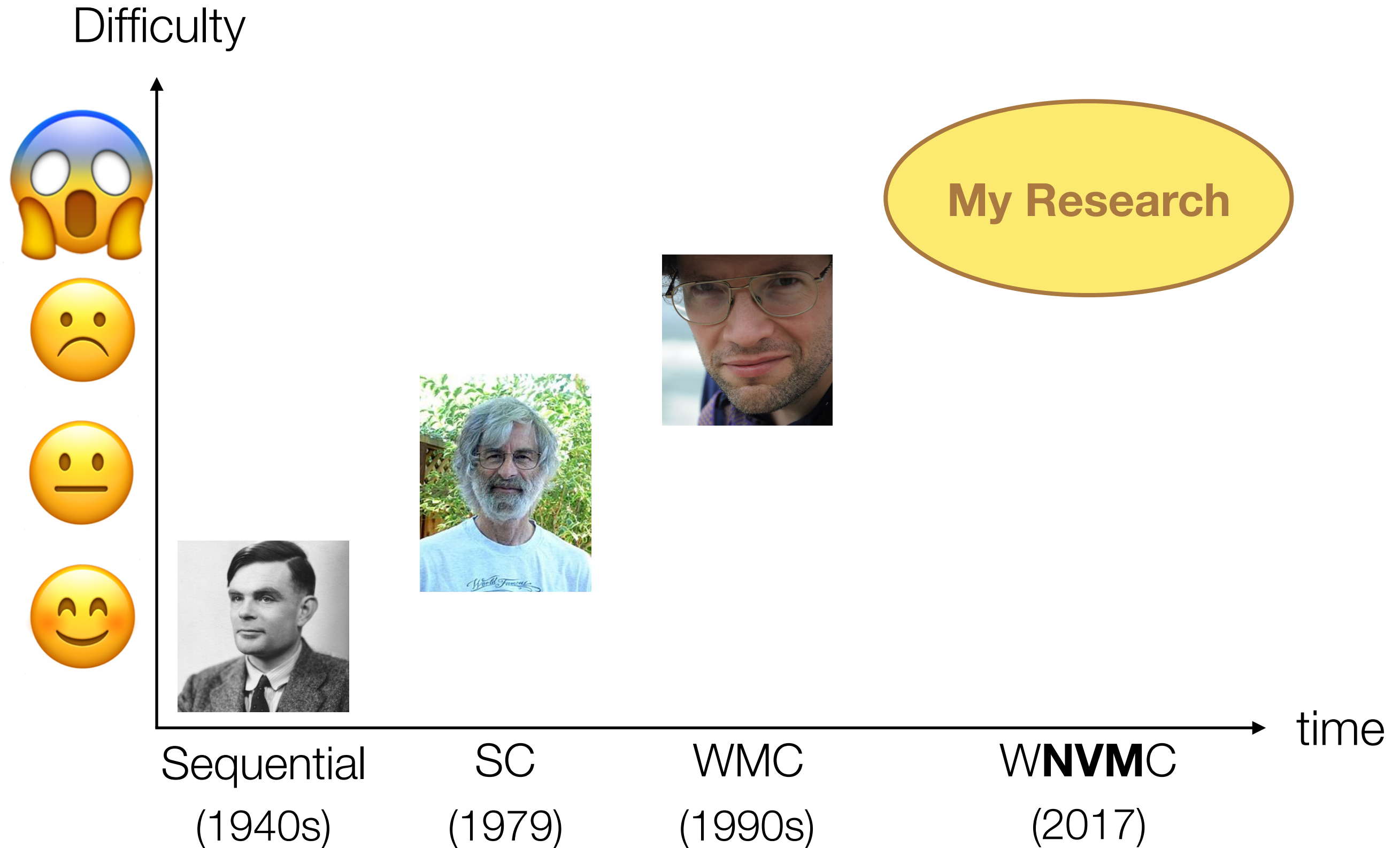
NVM Semantics

Consistency + Persistency Model

Formal Concurrency Models



Formal Concurrency Models



Challenge #1: ***Relaxed*** Persists

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

```
x := 1;
```

```
y := 1;
```


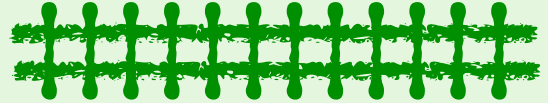



```
// recovery routine
```

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

!! out of order persists

Persist Barriers

```
// x=0; y=0  
  
x := 1;  
   
y := 1;  
  
  
  
// recovery routine  
  
// x=1; y=1 OR x=0; y=0 OR x=1; y=0 OR x=0; y=1
```

!! out of order persists

 ***persist barriers***

e.g. **SFENCE** on Intel

DSB_{full} on ARM

Challenge #2: ***Asynchronous*** Persists

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

```
x := 1 ;
```



```
y := 1 ;
```



```
// recovery routine
```

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

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

Explicit Persists

```
// x=0; y=0  
x := 1;  
➡ persist x;  
|||||  
y := 1;  
⚡  
// recovery routine
```

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

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

➡ **persist** instructions

e.g. **CLWB** / **CLFLUSHOPT** / **CLFLUSH** on Intel (per cache line)

DC-CVAP on ARM (per cache line)

psync under epoch persistency (global)

Here's Some Maths!

ARM Persistence Semantics

- together w. ARM UK
- declarative specification
- discovered *ambiguities in manual*

Intel Persistence Semantics

- Together w. Intel US
- declarative specification
- operational specification
- equivalence theorem
- discovered *ambiguities in manual*

$$\begin{array}{c}
 \text{Thread transitions: } \text{COM} \xrightarrow{\text{TID:LAB}_{\text{Px86}} \cup \{\epsilon\}} \text{COM} \quad \text{Program transitions: } \text{PROG} \xrightarrow{\text{TID:LAB}_{\text{Px86}} \cup \{\epsilon\}} \text{PROG} \\
 \\
 \frac{C_1 \xrightarrow{\tau:l} C'_1}{\text{let } a := C_1 \text{ in } C_2 \xrightarrow{\tau:l} \text{let } a := C'_1 \text{ in } C_2} \text{ (T-LET1)} \quad \frac{}{\text{let } a := v \text{ in } C \xrightarrow{\tau:\epsilon} C[v/a]} \text{ (T-LET2)} \\
 \\
 \frac{C \xrightarrow{\tau:l} C'}{\text{if } (C) \text{ then } C_1 \text{ else } C_2 \xrightarrow{\tau:l} \text{if } (C') \text{ then } C_1 \text{ else } C_2} \text{ (T-IF1)} \quad \frac{v \neq 0 \Rightarrow C = C_1 \quad v = 0 \Rightarrow C = C_2}{\text{if } (v) \text{ then } C_1 \text{ else } C_2 \xrightarrow{\tau:\epsilon} C} \text{ (T-IF2)} \\
 \\
 \frac{}{\text{repeat } C \xrightarrow{\tau:\epsilon} \text{if } (C) \text{ then } (\text{repeat } C) \text{ else } 0} \text{ (T-REPEAT)} \quad \frac{}{\text{store}(x, v) \xrightarrow{\tau:(W, x, v)} v} \text{ (T-WRITE)} \\
 \\
 \frac{}{\text{load}(x) \xrightarrow{\tau:(R, x, v)} v} \text{ (T-READ)} \quad \frac{}{\text{FAA}(x, v) \xrightarrow{\tau:(U, x, v_0, v_0+v)} v_0} \text{ (T-FAA)} \quad \frac{}{\text{mfence} \xrightarrow{\tau:MF} 1} \text{ (T-MF)} \\
 \\
 \frac{v \neq v_1}{\text{CAS}(x, v_1, v_2) \xrightarrow{\tau:(R, x, v)} 0} \text{ (T-CAS0)} \quad \frac{}{\text{CAS}(x, v_1, v_2) \xrightarrow{\tau:(U, x, v_1, v_2)} 1} \text{ (T-CAS1)} \quad \frac{}{\text{sfence} \xrightarrow{\tau:SF} 1} \text{ (T-SF)} \\
 \\
 \frac{}{\text{flush}_{\text{opt}} x \xrightarrow{\tau:(FQ, x)} 1} \text{ (T-FO)} \quad \frac{}{\text{flush } x \xrightarrow{\tau:(FL, x)} 1} \text{ (T-FL)} \quad \frac{P(\tau) \xrightarrow{\tau:l} C}{P \xrightarrow{\tau:l} P[\tau \mapsto C]} \text{ (PROG)} \\
 \hline
 \text{Storage transitions: } \text{MEM} \times \text{PBUF} \times \text{BMAP} \xrightarrow{\text{TID:LAB}_{\text{Px86}} \cup \{\epsilon\}} \text{MEM} \times \text{PBUF} \times \text{BMAP} \\
 \\
 \frac{B(\tau)=b \quad o \in \{\text{pfo}, \text{pfl}\} \quad x \in X \quad \langle \text{sf} \rangle \notin b \quad o=\text{pfo} \Rightarrow \forall v. \forall x' \in X. \langle x', v \rangle, \langle \text{fo}, x' \rangle, \langle \text{fl}, x' \rangle \notin b \quad o=\text{pfl} \Rightarrow \forall y, v. \forall x' \in X. \langle y, v \rangle, \langle \text{fo}, x' \rangle, \langle \text{fl}, y \rangle \notin b}{M, PB, B \xrightarrow{\tau:\epsilon} M, PB, B[\tau \mapsto b.(o, x)]} \text{ (M-BFETCH)} \\
 \\
 \frac{B(\tau)=b_1.(o, x).b_2 \quad o \in \{\text{pfo}, \text{pfl}\}}{M, PB, B \xrightarrow{\tau:\epsilon} M, PB, B[\tau \mapsto b_1.b_2]} \text{ (M-BDROP)} \quad \frac{PB=PB_1.(x, v).PB_2 \quad \forall y, v'. \langle x, v' \rangle, \langle \text{per}, y \rangle \notin PB_1}{M, PB, B \xrightarrow{\tau:\epsilon} M[x \mapsto v], PB_1.PB_2, B} \text{ (M-PROP W)} \\
 \\
 \frac{B(\tau)=b \quad x \in X \quad \forall y. \forall x' \in X. \langle \text{pfl}, y \rangle, \langle \text{pfo}, x' \rangle \notin b}{M, PB, B \xrightarrow{\tau:(W, x, v)} M, PB, B[\tau \mapsto b.(x, v)]} \text{ (M-WRITE)} \quad \frac{B(\tau)=b \quad \text{rd}(M, PB, b, x)=v}{M, PB, B \xrightarrow{\tau:(R, x, v)} M, PB, B} \text{ (M-READ)} \\
 \\
 \frac{B(\tau)=\epsilon \quad \text{rd}(M, PB, \epsilon, x)=v_r}{M, PB, B \xrightarrow{\tau:(U, x, v_r, v_w)} M, PB[x \mapsto PB(x).v_w], B} \text{ (M-RMW)} \quad \frac{B(\tau)=\epsilon}{M, PB, B \xrightarrow{\tau:MF} M, PB, B} \text{ (M-MF)} \\
 \\
 \frac{B(\tau)=b \quad \forall x. \forall o \in \{\text{pfo}, \text{pfl}\}. \langle o, x \rangle \notin b}{M, PB, B \xrightarrow{\tau:SF} M, PB, B[\tau \mapsto b.(sf)]} \text{ (M-SF)} \quad \frac{B(\tau)=\langle \text{sf} \rangle.b}{M, PB, B \xrightarrow{\tau:\epsilon} M, PB, B[\tau \mapsto b]} \text{ (M-BPROP SF)} \\
 \\
 \frac{B(\tau)=b \quad x \in X \quad \forall x' \in X. \langle \text{pfl}, x' \rangle, \langle \text{pfo}, x' \rangle \notin b}{M, PB, B \xrightarrow{\tau:(FQ, x)} M, PB, B[\tau \mapsto b.(fo, x)]} \text{ (M-FO1)} \quad \frac{B(\tau)=b_1.(pfo, x).b_2 \quad x \in X \quad \forall x' \in X. \langle \text{pfo}, x' \rangle, \langle \text{pfl}, x' \rangle \notin b_1}{M, PB, B \xrightarrow{\tau:(FQ, x)} M, PB, B[\tau \mapsto b_1.b_2]} \\
 \\
 \frac{B(\tau)=b \quad x \in X \quad \forall y. \forall x' \in X. \langle \text{pfl}, y \rangle, \langle \text{pfo}, x' \rangle \notin b}{M, PB, B \xrightarrow{\tau:(FL, x)} M, PB, B[\tau \mapsto b.(fl, x)]} \text{ (M-FL1)} \quad \frac{B(\tau)=b_1.(pfl, x).b_2 \quad x \in X \quad \forall x' \in X. \langle \text{pfo}, x' \rangle \notin b_1 \quad \forall y.}{M, PB, B \xrightarrow{\tau:(FL, x)} M, PB, B[\tau \mapsto b_1.b_2]} \\
 \\
 \frac{B(\tau)=b_1.(x, v).b_2 \quad x \in X \quad \forall y, v. \forall x' \in X. \langle \text{sf} \rangle, \langle y, v \rangle, \langle \text{fl}, y \rangle, \langle \text{fo}, x' \rangle \notin b_1}{M, PB, B \xrightarrow{\tau:\epsilon} M, PB.(x, v), B[\tau \mapsto b_1.b_2]} \text{ (M-BPROP W)} \quad \frac{PB=PB_1.(per, x).PB_2 \quad \forall y, v. \langle x, v \rangle, \langle \text{per}, y \rangle \notin PB_1}{M, PB, B \xrightarrow{\tau:\epsilon} M, PB_1.PB_2, B} \text{ (M-BPROP P)} \\
 \\
 \frac{B(\tau)=b_1.(o, x).b_2 \quad o \in \{\text{fo}, \text{fl}\} \quad x \in X \quad \langle \text{sf} \rangle \notin b_1 \quad o=\text{fo} \Rightarrow \forall v. \forall x' \in X. \langle x', v \rangle, \langle \text{fo}, x' \rangle, \langle \text{fl}, x' \rangle \notin b_1 \quad o=\text{fl} \Rightarrow \forall y, v. \forall x' \in X. \langle y, v \rangle, \langle \text{fo}, x' \rangle, \langle \text{fl}, y \rangle \notin b_1}{M, PB, B \xrightarrow{\tau:\epsilon} M, PB.(per, x), B[\tau \mapsto b_1.b_2]} \text{ (M-BPROP P)}
 \end{array}$$

Here's Some Maths!

Problem

counter-intuitive semantics
low-level hardware details

Solution

high-level, hardware-agnostic
NVM libraries:

Persistent Transactions

ARM Persistent

- together w
- declarative
- discovered

Intel Persistent

- Together w
- declarative
- operational
- equivalence
- discovered

$$\begin{array}{c} \text{tid:LAB}_{\text{Px86}} \cup \{\epsilon\} \\ \text{prog: PROG} \xrightarrow{\quad} \text{PROG} \\ \\ \frac{}{\tau:\epsilon \xrightarrow{\quad} C[v/a]} \text{ (T-LET2)} \\ \frac{C=C_1 \quad v=0 \Rightarrow C=C_2}{\text{if } C_1 \text{ else } C_2 \xrightarrow{\tau:\epsilon} C} \text{ (T-IF2)} \\ \frac{}{\tau:(W, x, v) \xrightarrow{\quad} v} \text{ (T-WRITE)} \\ \frac{}{\text{mfence} \xrightarrow{\tau:\text{MF}} 1} \text{ (T-MF)} \\ \frac{}{\text{sfence} \xrightarrow{\tau:\text{SF}} 1} \text{ (T-SF)} \\ \frac{\tau:l \xrightarrow{\quad} C}{\text{P}[\tau \mapsto C]} \text{ (PROG)} \\ \text{-----} \\ \times \text{ BMAP} \\ \frac{\langle \text{fo}, x' \rangle, \langle \text{fl}, y \rangle \notin b}{\quad} \text{ (M-BFETCH)} \\ \frac{\langle x, v' \rangle, \langle \text{per}, y \rangle \notin \text{PB}_1}{\quad} \text{ (M-PROP W)} \\ \frac{}{\text{B}(M, \text{PB}, b, x) = v} \text{ (M-READ)} \\ \frac{}{\tau:\epsilon \xrightarrow{\quad} M, \text{PB}, B} \text{ (M-MF)} \\ \frac{}{\tau:\text{MF} \xrightarrow{\quad} M, \text{PB}, B} \text{ (M-BPROP SF)} \\ \frac{}{\text{B}(\tau) = \epsilon} \text{ (M-BPROP SF)} \\ \frac{}{\text{B}, B[\tau \mapsto b]} \text{ (M-BPROP SF)} \\ \frac{x \in X \quad \forall x' \in X. \langle \text{pfo}, x' \rangle, \langle \text{pfl}, x' \rangle \notin b_1}{\tau:(\text{FO}, x) \xrightarrow{\quad} M, \text{PB}, B[\tau \mapsto b_1.b_2]} \text{ (M-BPROP W)} \\ \frac{x \in X \quad \forall x' \in X. \langle \text{pfo}, x' \rangle, \langle \text{pfl}, x' \rangle \notin b_1 \quad \forall y. \langle x \rangle.b_2}{\tau:(\text{FL}, x) \xrightarrow{\quad} M, \text{PB}, B[\tau \mapsto b_1.b_2]} \text{ (M-BPROP W)} \\ \frac{\text{B}(\tau) = b_1. \langle x, v \rangle. b_2 \quad x \in X \quad \forall y, v. \forall x' \in X. \langle \text{sf} \rangle, \langle y, v \rangle, \langle \text{fl}, y \rangle, \langle \text{fo}, x' \rangle \notin b_1}{M, \text{PB}, B \xrightarrow{\tau:\epsilon} M, \text{PB}, \langle x, v \rangle, B[\tau \mapsto b_1.b_2]} \text{ (M-BPROP W)} \\ \frac{\text{B}(\tau) = b_1. \langle o, x \rangle. b_2 \quad o \in \{\text{fo}, \text{fl}\} \quad x \in X \quad \langle \text{sf} \rangle \notin b_1}{o = \text{fo} \Rightarrow \forall v. \forall x' \in X. \langle x', v \rangle, \langle \text{fo}, x' \rangle, \langle \text{fl}, x' \rangle \notin b_1 \quad o = \text{fl} \Rightarrow \forall y, v. \forall x' \in X. \langle y, v \rangle, \langle \text{fo}, x' \rangle, \langle \text{fl}, y \rangle \notin b_1} \text{ (M-BPROP P)} \\ \frac{}{M, \text{PB}, B \xrightarrow{\tau:\epsilon} M, \text{PB}, \langle \text{per}, x \rangle, B[\tau \mapsto b_1.b_2]} \text{ (M-BPROP P)} \end{array}$$

Transactions: NVM vs. TM/Databases

- ▶ **TM** (Transactional Memory): *run on **volatile** hardware*
 - ➔ **No** persistency guarantees
- ▶ **Database** transactions: *run on **persistent** hardware*
 - ➔ **Specific** persistency guarantees
 - ✦ **only strict** persistency
 - ✦ **only synchronous** persistency
- ▶ **NVM** transactions: *run on **persistent** hardware*
 - ➔ **Range** of persistency guarantees
 - ✦ **strict or relaxed** persistency
 - ✦ **synchronous or asynchronous** persistency

What is a Transaction?

Concurrency control mechanism:

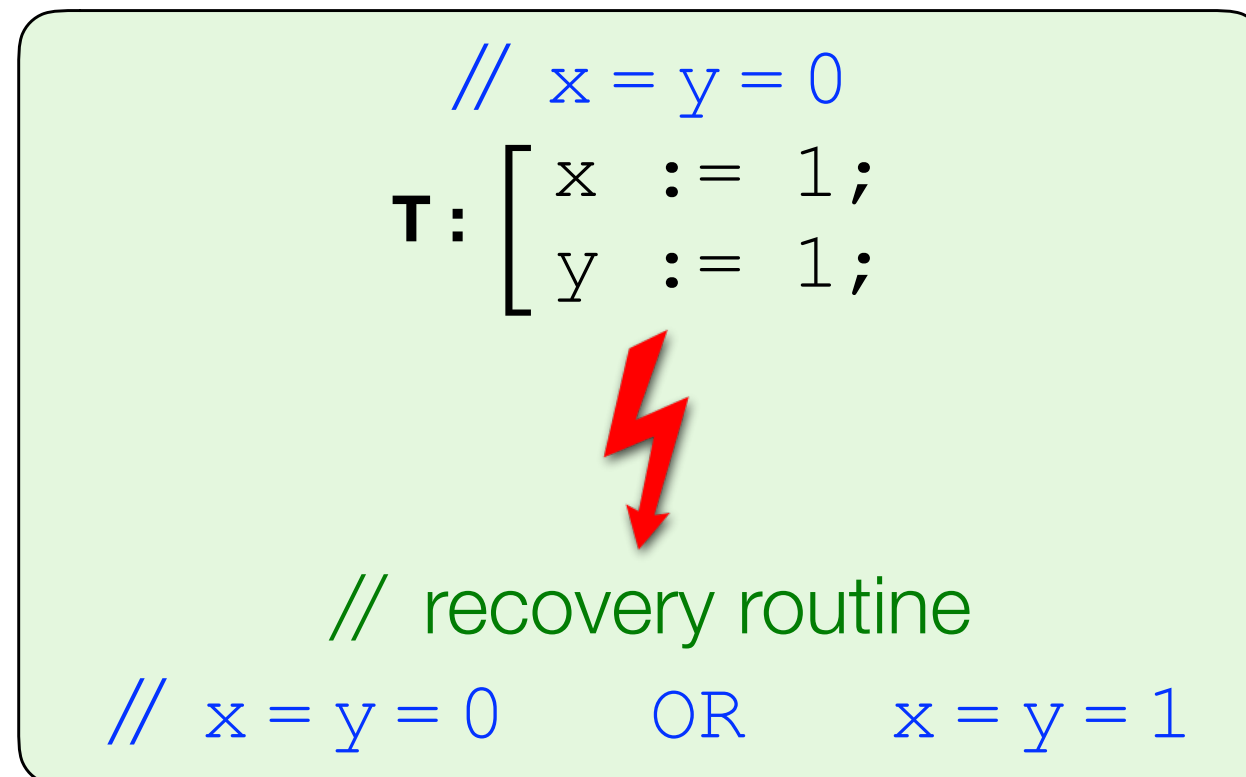
- ▶ **atomic** work unit:
 - ➔ all-or-nothing writes
- ▶ **consistent** (e.g. serialisable)

```
// x = y = 0  
  
T: [ x := 1;  
    y := 1;  
    ]  
  
// x = y = 0    OR    x = y = 1
```

What is a **Persistent** Transaction?

Concurrency & **persistency** control mechanism:

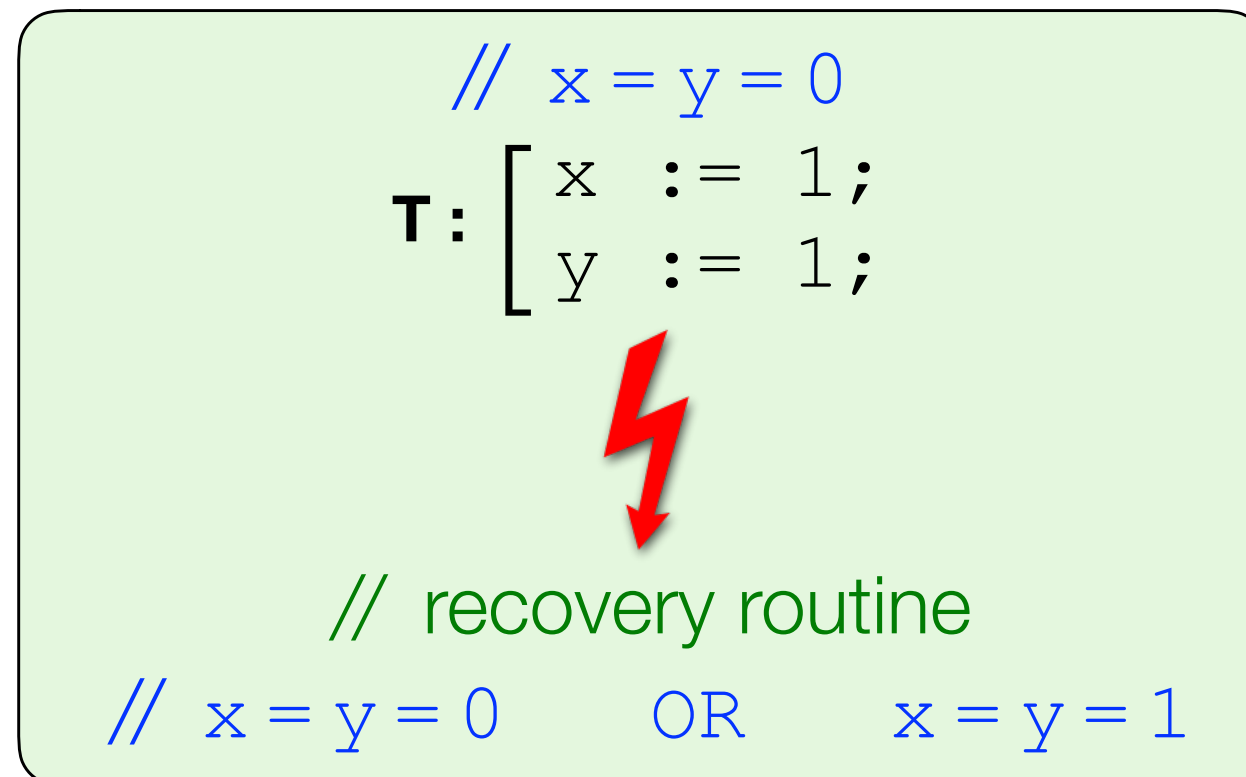
- ▶ **atomic** work unit:
 - ➔ all-or-nothing writes
 - ➔ all-or-nothing **persists**
- ▶ **consistent** (e.g. serialisable)



What is a **Persistent** Transaction?

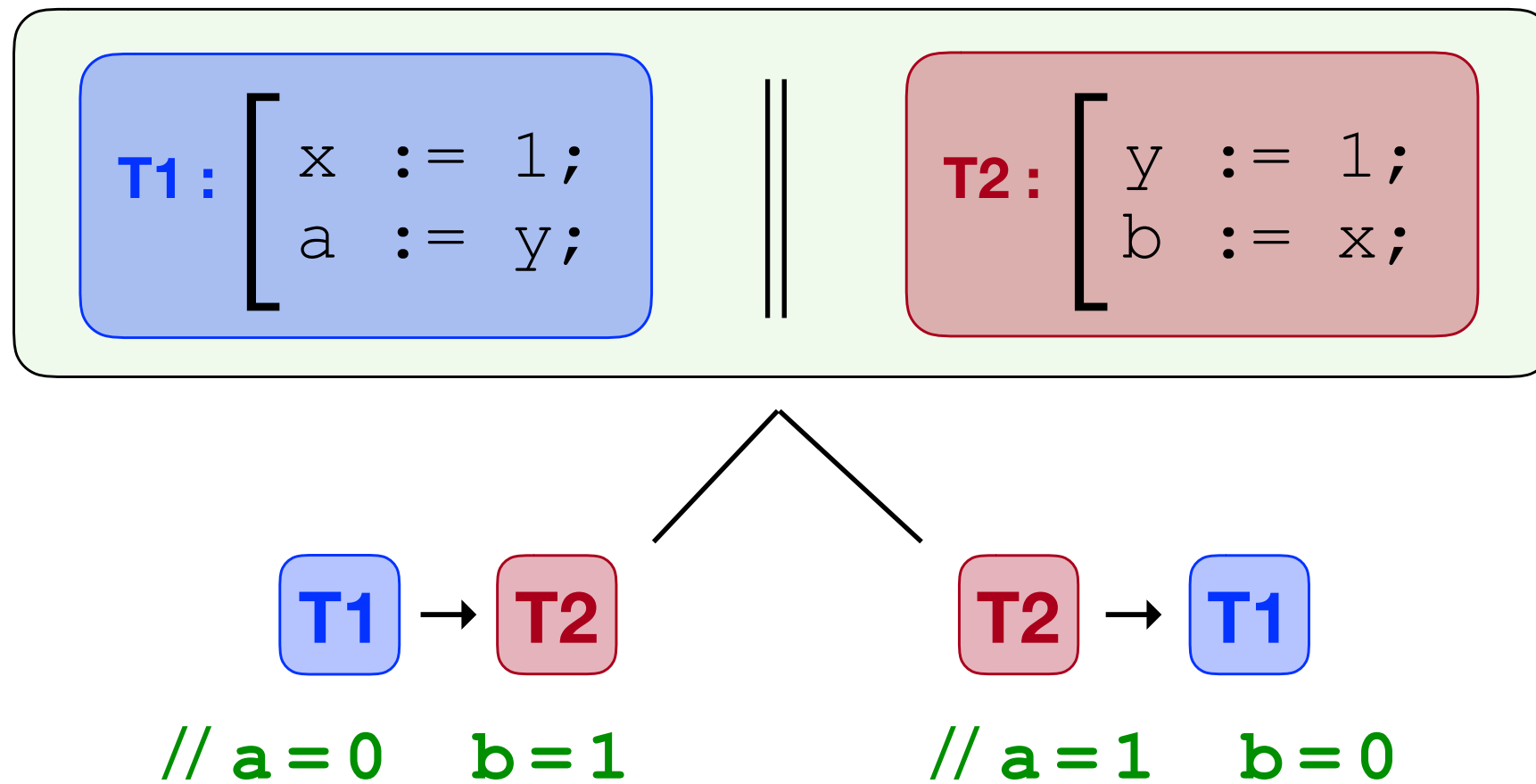
Concurrency & **persistency** control mechanism:

- ▶ **atomic** work unit:
 - ➔ all-or-nothing writes
 - ➔ all-or-nothing **persists**
- ▶ **consistent** (e.g. serialisable)
- ▶ **persistent** (e.g. **persistently serialisable**)



Serialisability (SER)

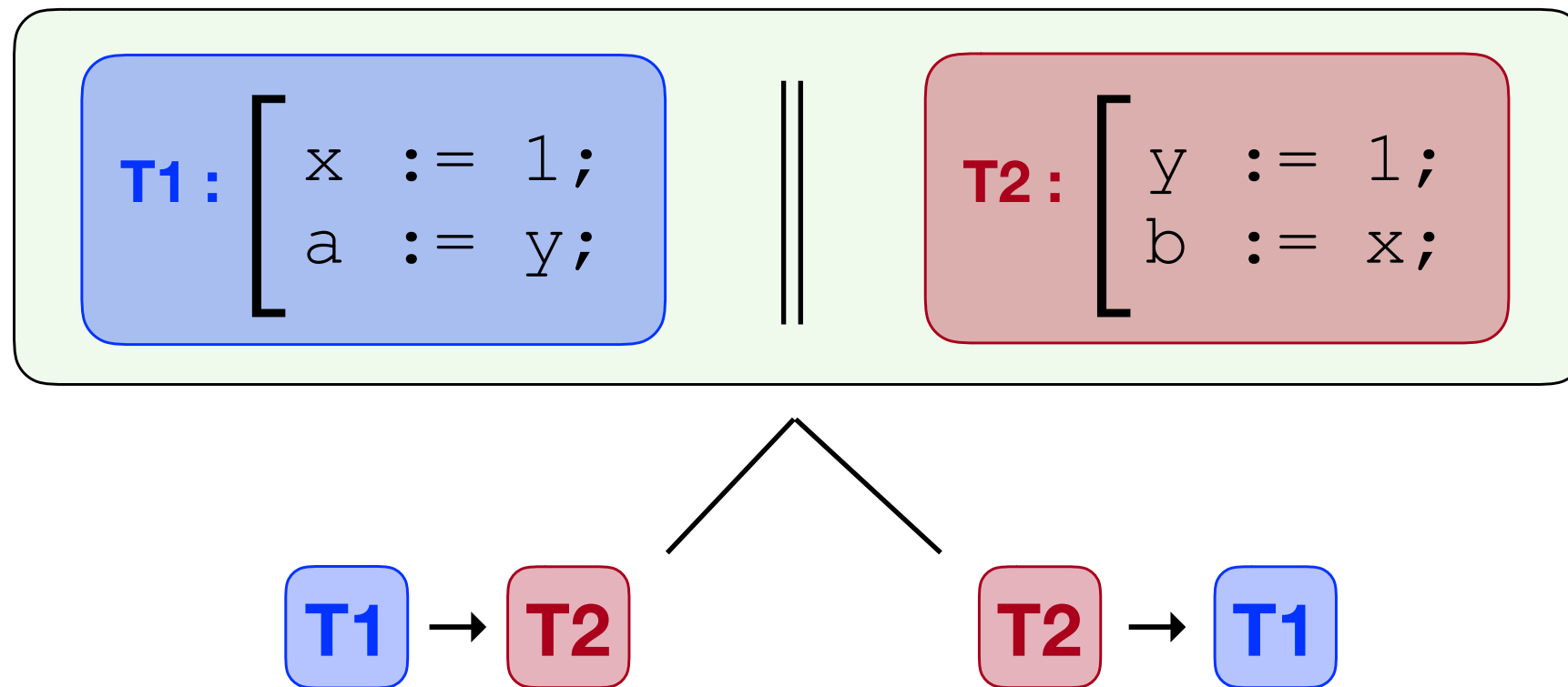
All transactions appear to **execute** in a sequential order



Persistent Serialisability (PSER)

All transactions appear to **execute** in a sequential order

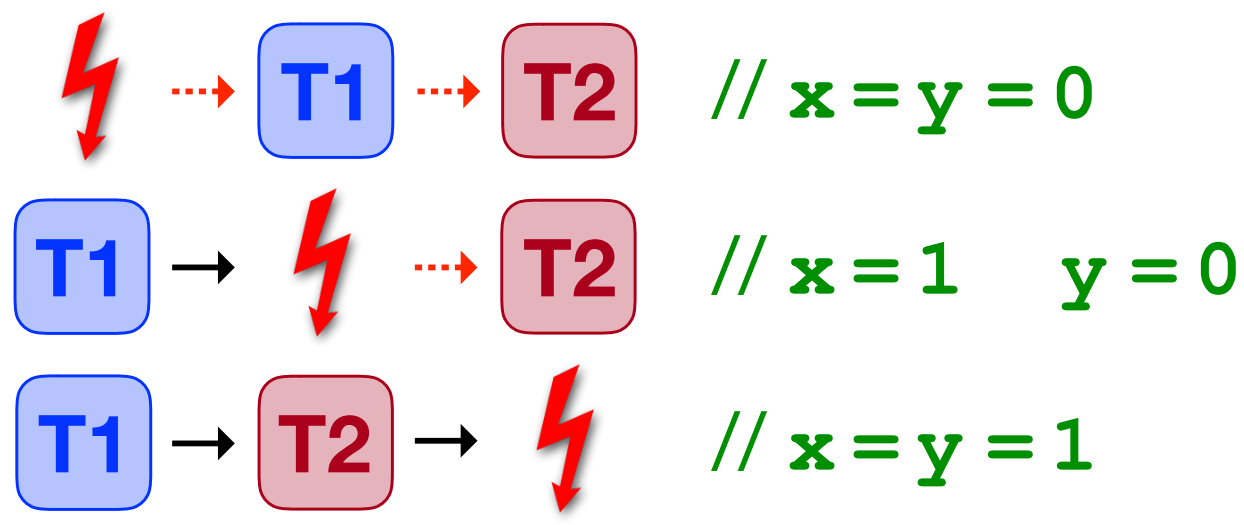
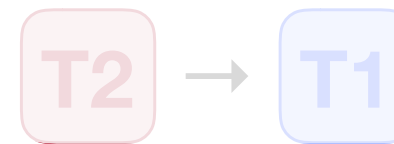
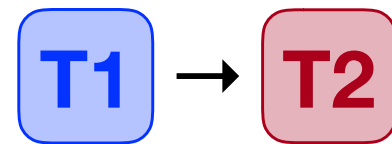
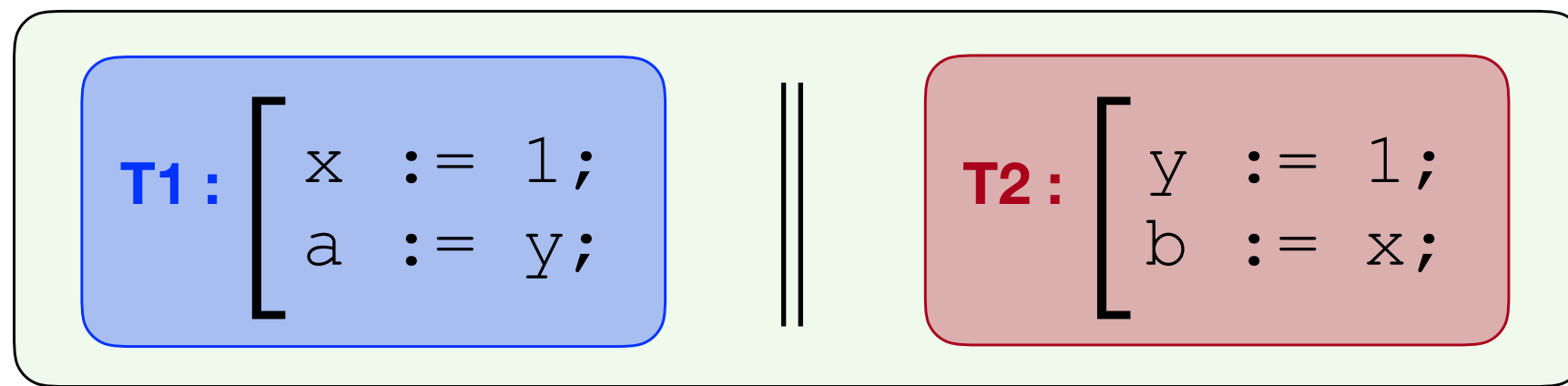
A prefix of transactions appears to **persist** in the **same sequential order**



Persistent Serialisability (PSER)

All transactions appear to **execute** in a sequential order

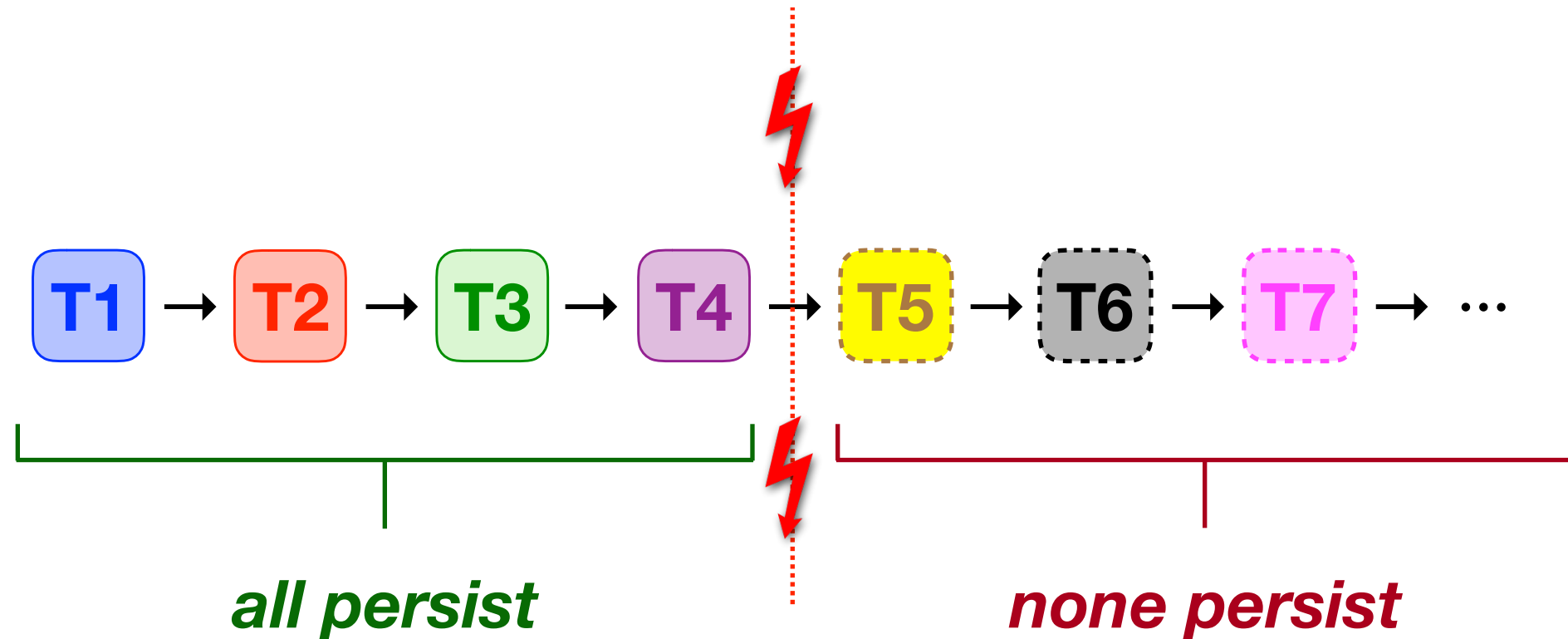
A prefix of transactions appears to **persist** in the **same sequential order**



Persistent Serialisability (PSER)

All transactions appear to **execute** in a sequential order

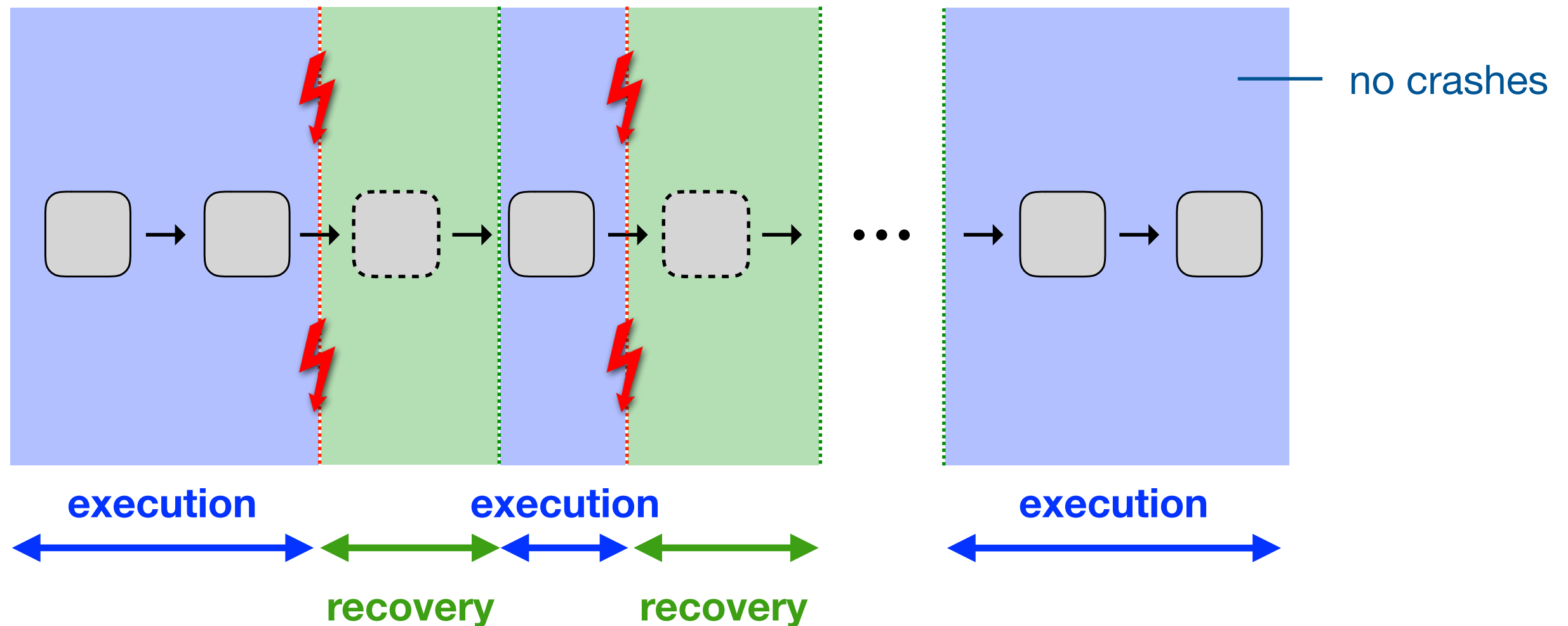
A prefix of transactions appears to **persist** in the **same sequential order**



Persistent Serialisability (PSER)

All transactions appear to **execute** in a sequential order

A prefix of transactions appears to **persist** in the *same sequential order* in **each era**



Persistent Serialisability (PSER)

All transactions appear to ***execute*** in a sequential order

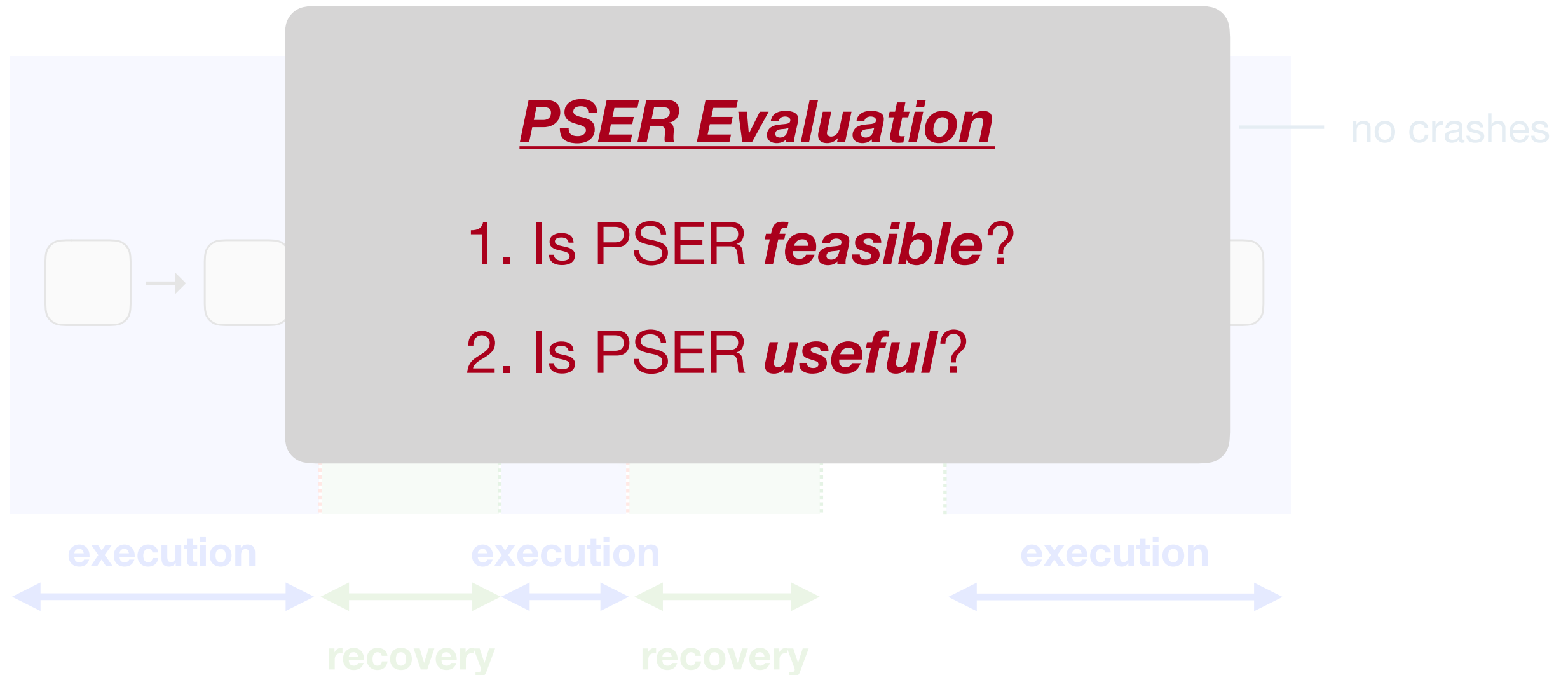
A prefix of transactions appears to ***persist*** in the ***same sequential order*** in *each era*



Persistent Serialisability (PSER)

All transactions appear to ***execute*** in a sequential order

A prefix of transactions appears to ***persist*** in the ***same sequential order*** in *each era*



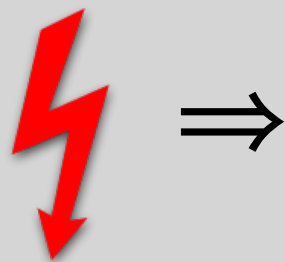
Is PSER *Feasible*?

✓ PSER *implementation* in **ARM**

✓ PSER *implementation* in **Intel**

Take **SER** Implementation — e.g. 2-PL

- ✦ add code for **persistence** — e.g. **psync**
- ✦ add code to **log metadata** for **recovery**
- ✦ add **recovery mechanism**



recovery mechanism

check log for **incomplete** transactions:
either **complete**
or **rollback**

Is PSER *Feasible*?

✓ PSER *implementation* in *ARM*

✓ PSER *implementation* in *Intel*

Take *SER* Implementation — e.g. 2-PL



Yes!

Correct Implementation* in *ARM* & *Intel



check log for *incomplete* transactions:

either ***complete***

or ***rollback***

Is PSER *Useful*?

Given library L (e.g. queue library):

1. Take **any** correct **sequential** implementation of L
2. wrap each operation in a PSER transaction

```
enq(q, v) =  
  < enq_body >
```

```
deq(q) =  
  < deq_body >
```

sequential queue imp.

```
enq(q, v) =  
  psер{  
    < enq_body > }  
  }
```

```
deq(q) =  
  psер{  
    < deq_body > }  
  }
```

Is PSER *Useful*?

Given library L (e.g. queue library):

1. Take **any** correct **sequential** implementation of L
2. wrap each operation in a PSER transaction

\Rightarrow **correct, concurrent & persistent** implementation of L

```
enq(q, v) =  
  < enq_body >
```

```
deq(q) =  
  < deq_body >
```

sequential queue imp.

```
enq(q, v) =  
  pser{  
    < enq_body > }  
deq(q) =  
  pser{  
    < deq_body > }
```

correct
concurrent & persistent
queue imp.

Is PSER *Useful*?

Given library L (e.g. queue library):

1. T

2. w



Yes!

***any* correct *sequential* implementation**



***correct, concurrent & persistent*
implementation**

er

de

< deq_body >

pool

< deq_body > }

correct

***concurrent & persistent*
queue imp.**

sequential queue imp.

Summary

✓ Formalised ***architecture-level*** NVM semantics:

✦ ***ARM***

✦ ***Intel***

✓ Formalised ***language-level*** NVM semantics:

✦ ***PSER***

✦ ***Feasibility***: implemented PSER on ARM and Intel

✦ ***Utility***: PSER for concurrent & persistent library implementation

? Future Work:

✦ other transactional models

✦ model checking algorithms

✦ program logics

Thank You for Listening!