

# Relaxed Persist Ordering Using Strand Persistency

Vaibhav Gogte, William Wang<sup>§</sup>, Stephan Diestelhorst<sup>§</sup>,  
Peter M. Chen, Satish Narayanasamy, Thomas F. Wenisch



ISCA 2020



# Promise of persistent memory (PM)

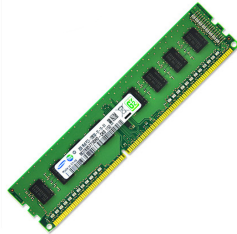
**Performance**



**Density**



**Non-volatility**



# Promise of persistent memory (PM)

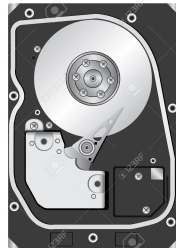
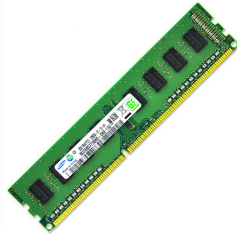
Performance



Density



Non-volatility



## Intel Announces New Optane DC Persistent Memory \*

By Joel Hruska on May 31, 2018 at 8:15 am | [1 Comment](#)

*“Optane DC Persistent Memory will be offered in packages of up to 512GB per stick.”*

*“... expanding memory per CPU socket to as much as 3TB.”*

\* Source: [www.extremetech.com](http://www.extremetech.com)

# Promise of persistent memory (PM)

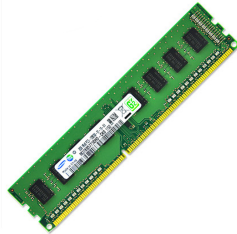
Performance



Density



Non-volatility



## Intel Announces New Optane DC Persistent Memory \*

By Joel Hruska on May 31, 2018 at 8:15 am | [1 Comment](#)

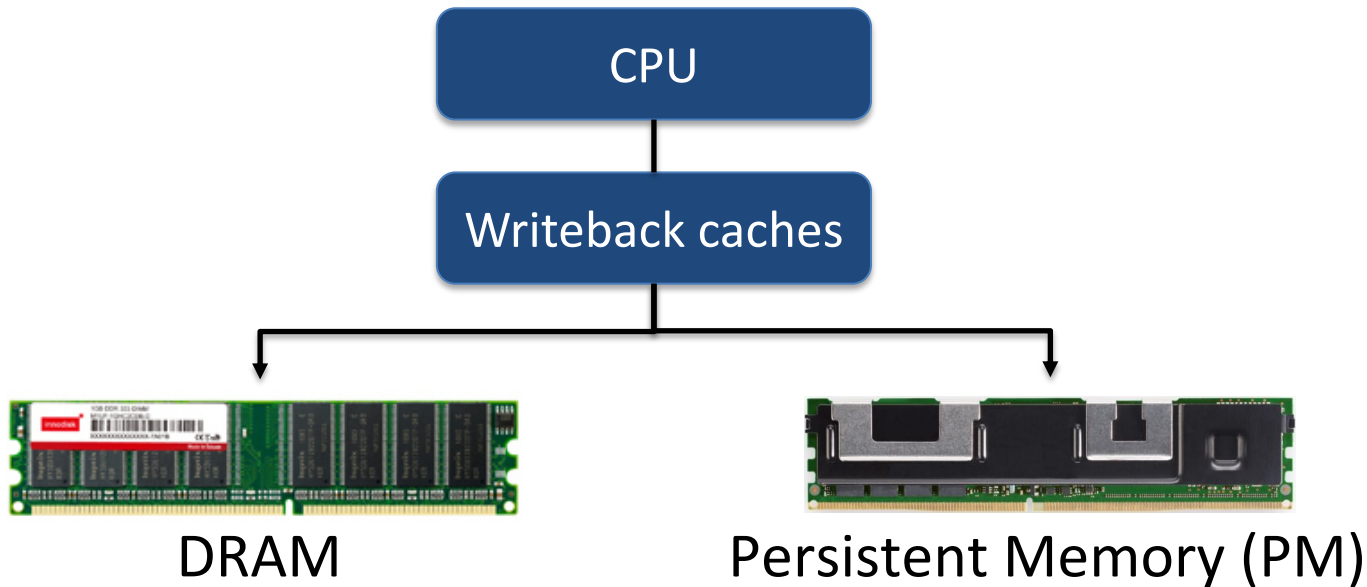
*“Optane DC Persistent Memory will be offered in packages of up to 512GB per stick.”*

*“... expanding memory per CPU socket to as much as 3TB.”*

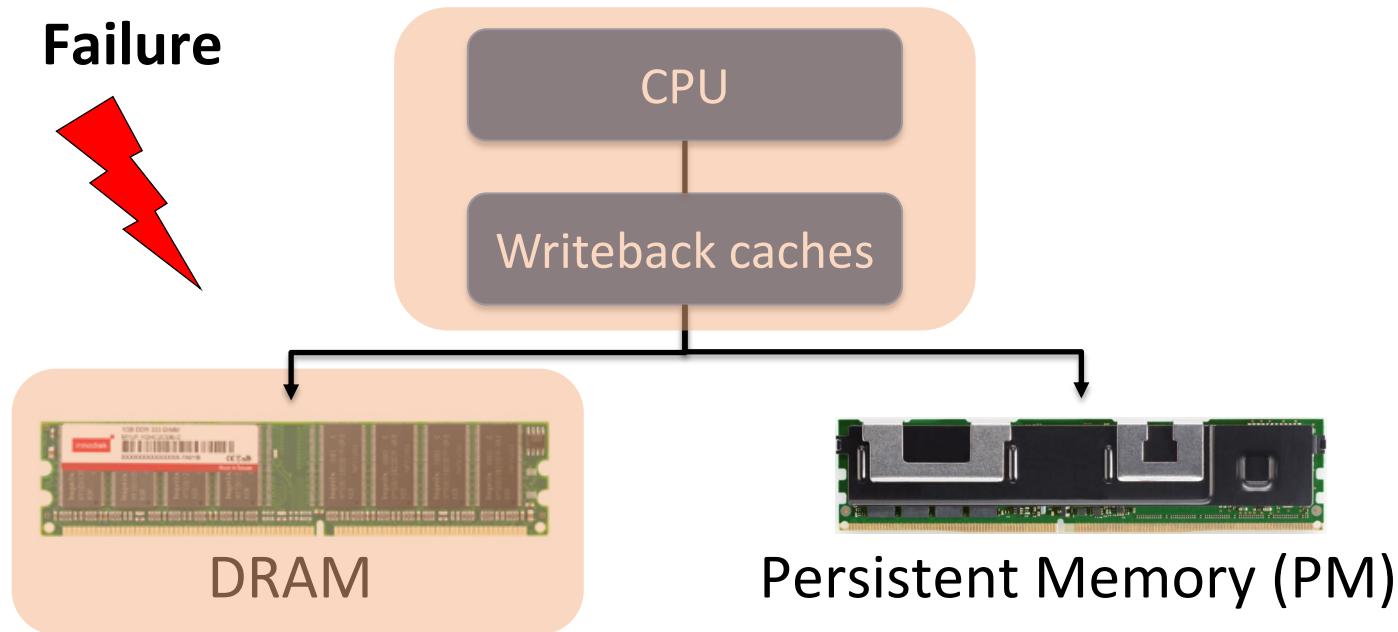
\* Source: [www.extremetech.com](http://www.extremetech.com)

Byte-addressable, load-store interface to durable storage

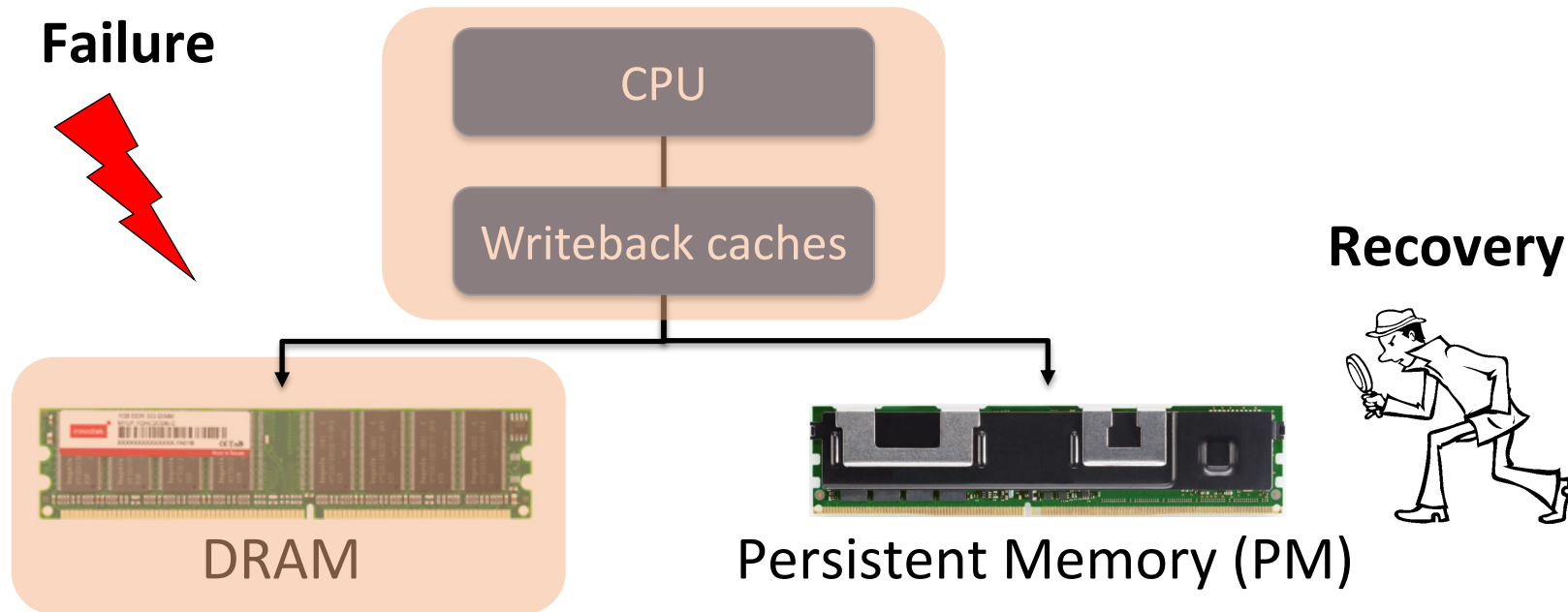
# Persistent memory system



# Persistent memory system



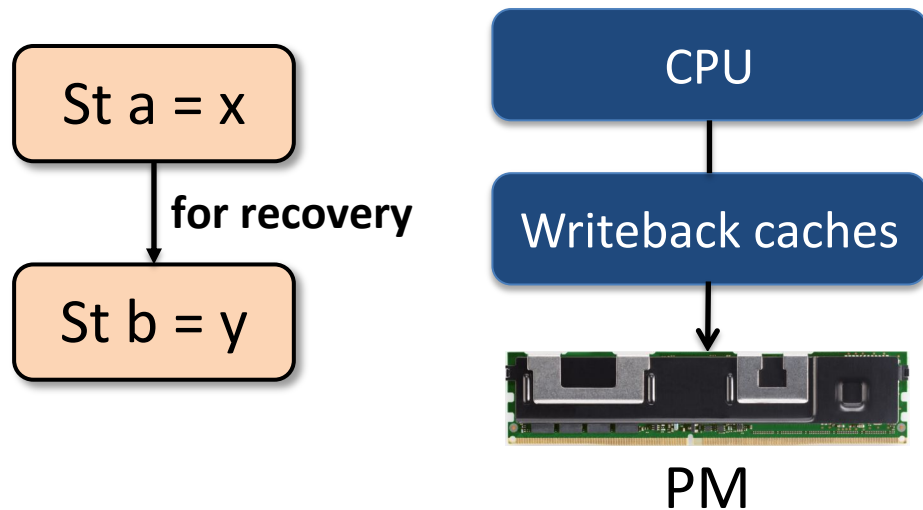
# Persistent memory system



Recovery can inspect PM data-structures to restore system to a consistent state

# Recovery requires PM access ordering

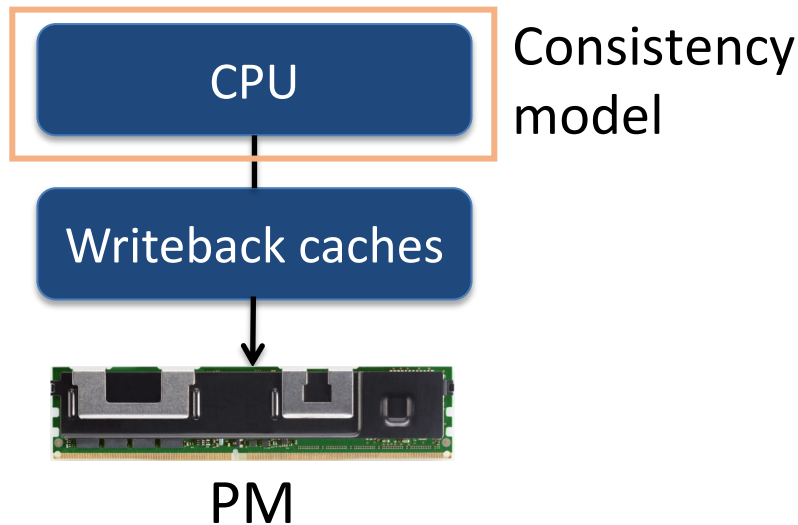
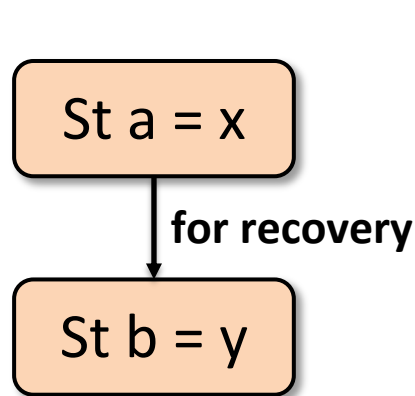
Intel x86 primitives







# Recovery requires PM access ordering

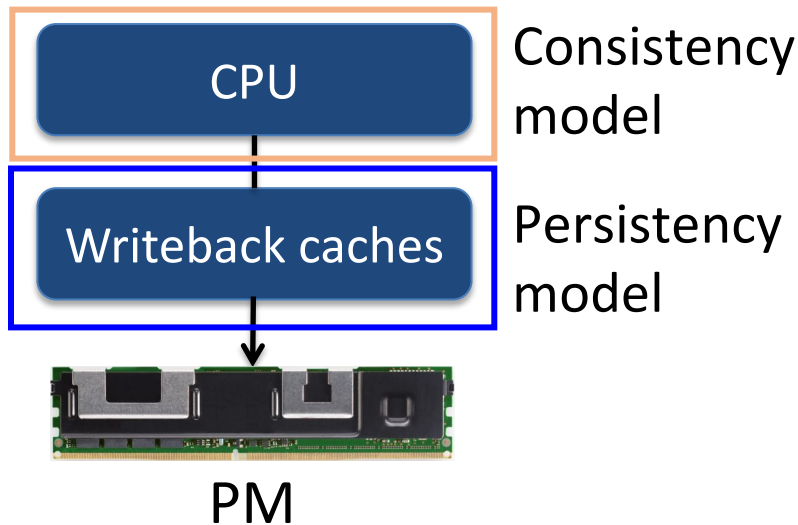
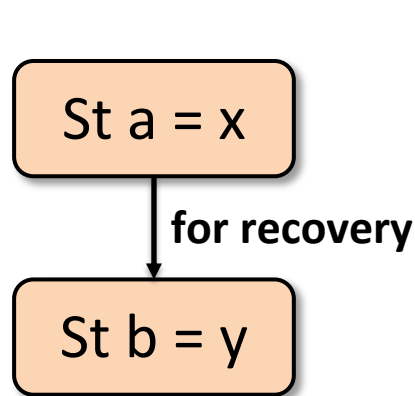


## Intel x86 primitives

St a = x

St b = y

# Recovery requires PM access ordering

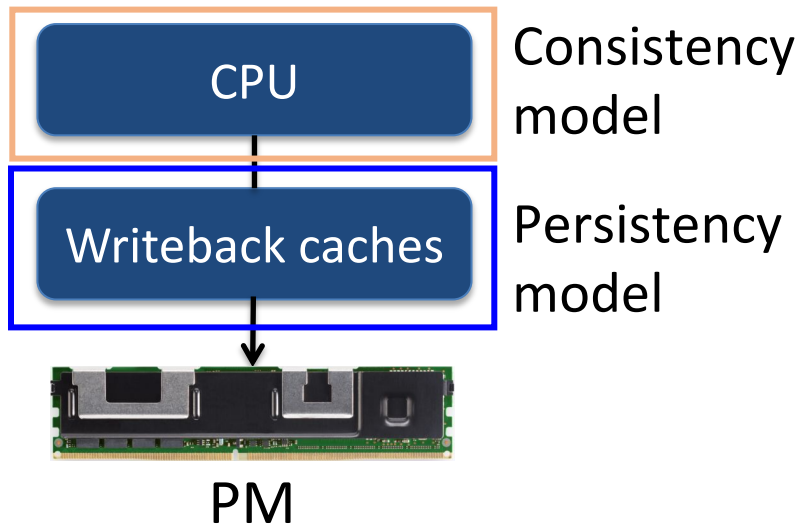
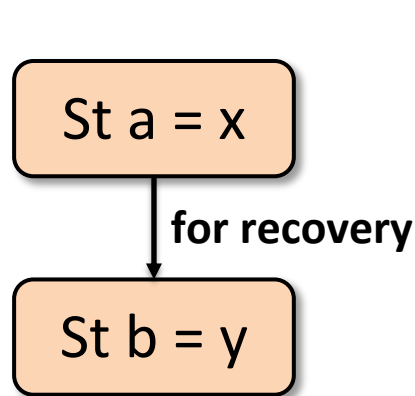


## Intel x86 primitives

St a = x

St b = y

# Recovery requires PM access ordering



## Intel x86 primitives

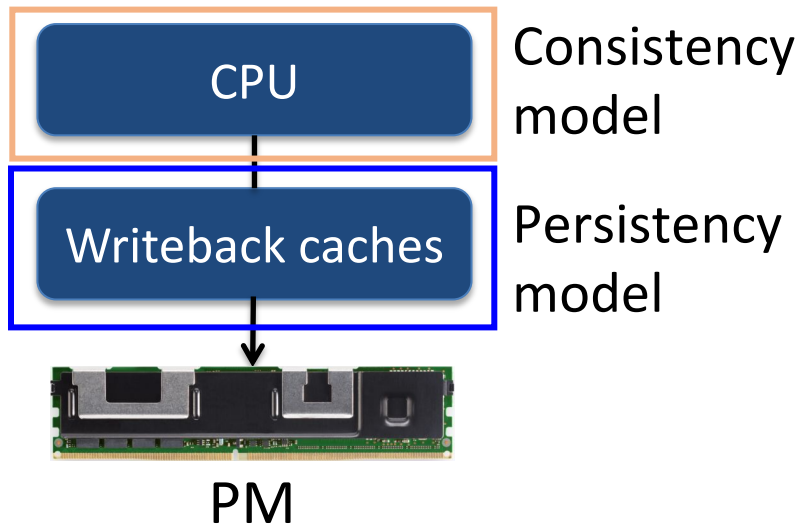
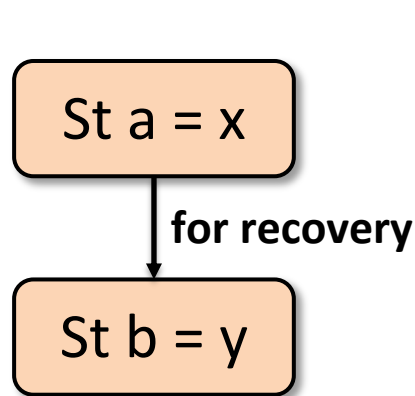
St a = x

CLWB(a)

St b = y

CLWB(b)

# Recovery requires PM access ordering

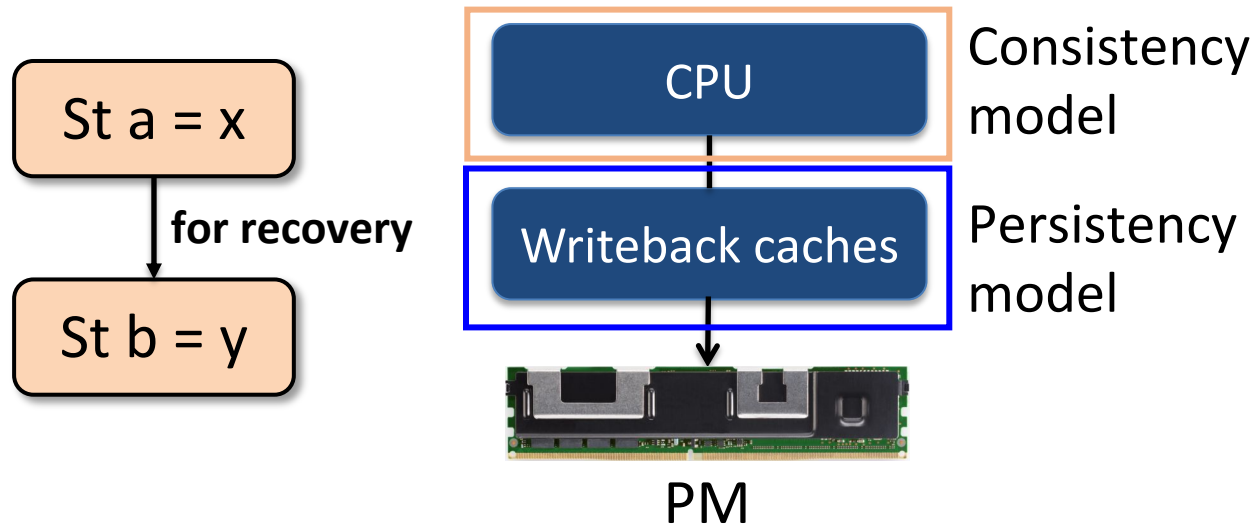


## Intel x86 primitives

- St a = x
- CLWB(a)
- SFENCE
- St b = y
- CLWB(b)



# Recovery requires PM access ordering



## Intel x86 primitives

St a = x

CLWB(a)

SFENCE

St b = y

CLWB(b)

Hardware systems provide primitives to express *persist* order to PM

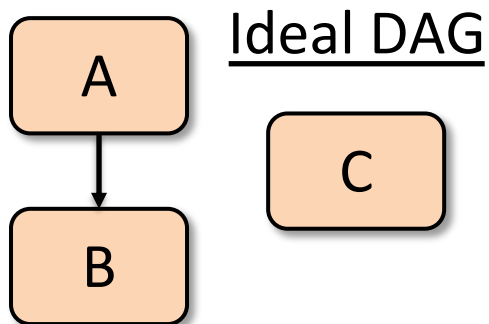


# Hardware imposes overly strict constraints

St A = 1; CLWB (A)

St B = 2; CLWB (B)

St C = 3; CLWB (C)





# Hardware imposes overly strict constraints

St A = 1; CLWB (A)

St B = 2; CLWB (B)

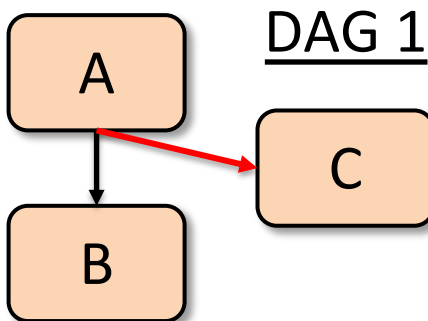
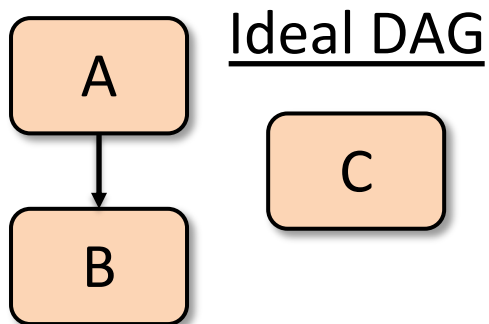
St C = 3; CLWB (C)

St A = 1; CLWB (A)

**SFENCE**

St B = 2; CLWB (B)

St C = 3; CLWB (C)





# Hardware imposes overly strict constraints

St A = 1; CLWB (A)

St B = 2; CLWB (B)

St C = 3; CLWB (C)

St A = 1; CLWB (A)

**SFENCE**

St B = 2; CLWB (B)

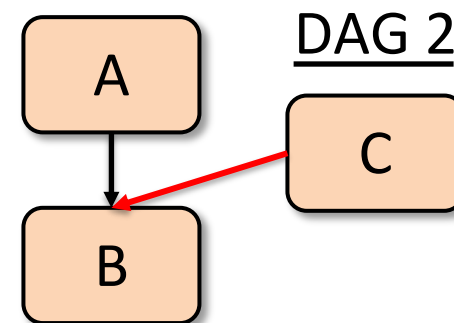
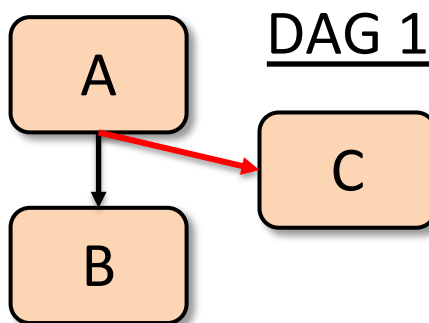
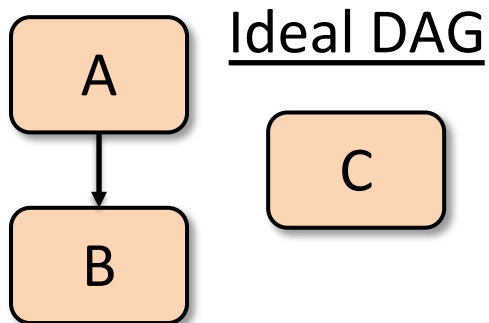
St C = 3; CLWB (C)

St A = 1 ; CLWB (A)

St C = 3; CLWB (C)

**SFENCE**

St B = 2; CLWB (B)





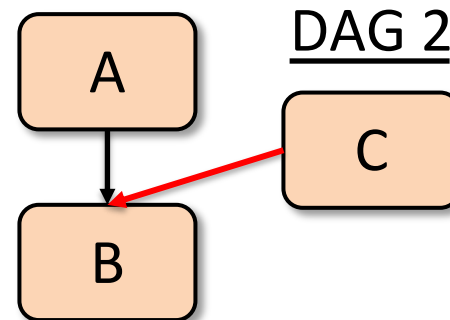
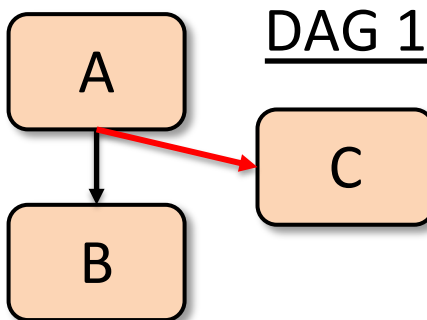
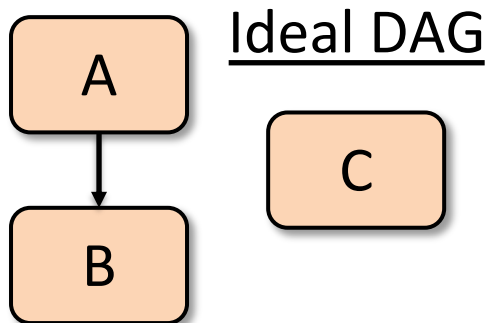


# Hardware imposes overly strict constraints

St A = 1; CLWB (A)  
 St B = 2; CLWB (B)  
 St C = 3; CLWB (C)

St A = 1; CLWB (A)  
**SFENCE**  
 St B = 2; CLWB (B)  
 St C = 3; CLWB (C)

St A = 1 ; CLWB (A)  
 St C = 3; CLWB (C)  
**SFENCE**  
 St B = 2; CLWB (B)



Primitives in existing hardware systems overconstrain PM accesses

# Contributions

- Our proposal: StrandWeaver
  - Builds strand persistency model in hardware
  - Specifies precise persist ordering constraints
- Comprises primitives: **PersistBarrier**, **NewStrand**, and **JoinStrand**
  - Can encode an arbitrary DAG
- Map language-level persistency models to ISA level primitives
  - Leverage hw primitives to build persistency models efficiently



# Contributions

- Our proposal: StrandWeaver
  - Builds strand persistency model in hardware
  - Specifies precise persist ordering constraints
- Comprises primitives: **PersistBarrier**, **NewStrand**, and **JoinStrand**
  - Can encode an arbitrary DAG
- Map language-level persistency models to ISA level primitives
  - Leverage hw primitives to build persistency models efficiently

StrandWeaver results in 1.45x (avg.) speedup over Intel x86



# Outline

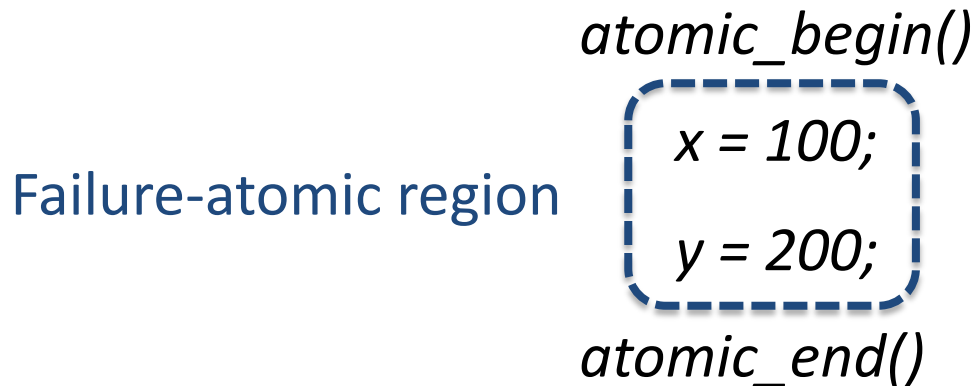
- Contributions
- Example: Failure atomicity
- Existing hardware vs. strand persistency model
- Our proposal: StrandWeaver
- Evaluation



# Example: Failure atomicity

## Failure atomicity:

Which group of stores persist atomically?





# Example: Failure atomicity

## Failure atomicity:

Which group of stores persist atomically?



Failure atomicity limits state that recovery can observe after failure



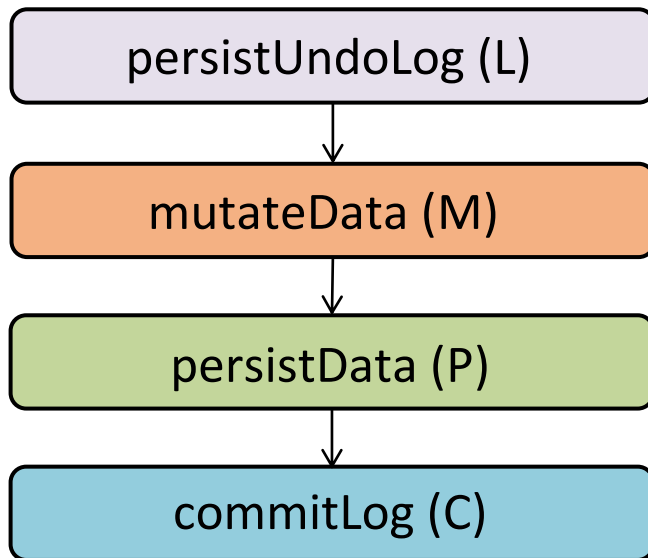
# Undo logging for failure atomicity

*Init: x = 0; y = 0*

atomic\_begin()

```
x = 1;  
y = 2;
```

atomic\_end()





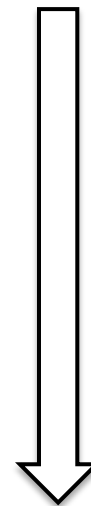
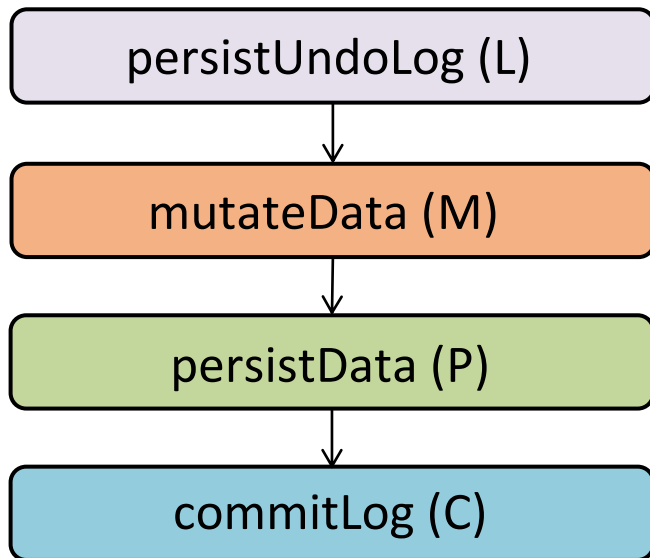
# Undo logging for failure atomicity

*Init: x = 0; y = 0*

atomic\_begin()

```
x = 1;  
y = 2;
```

atomic\_end()



**Failure-atomic**

Undo logging steps ordered to ensure failure atomicity





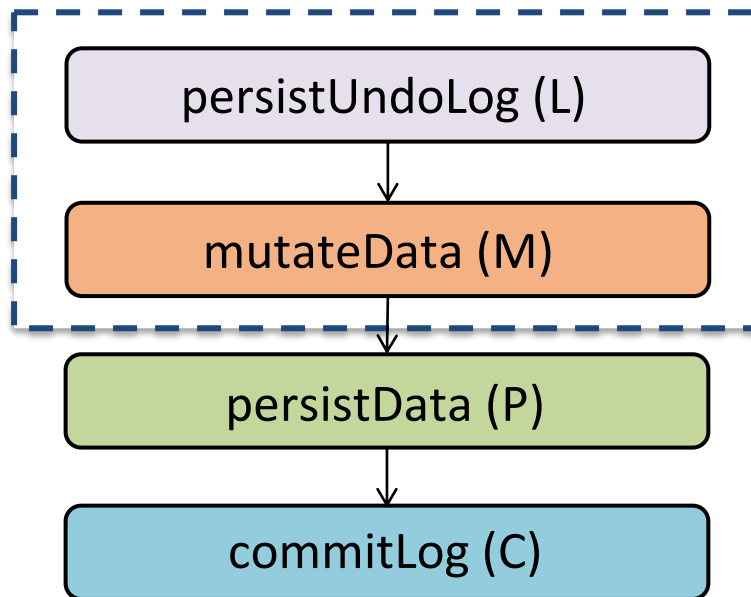
# Undo logging for failure atomicity

*Init: x = 0; y = 0*

atomic\_begin()

```
x = 1;  
y = 2;
```

atomic\_end()



**Failure-atomic**

Undo logging steps ordered to ensure failure atomicity

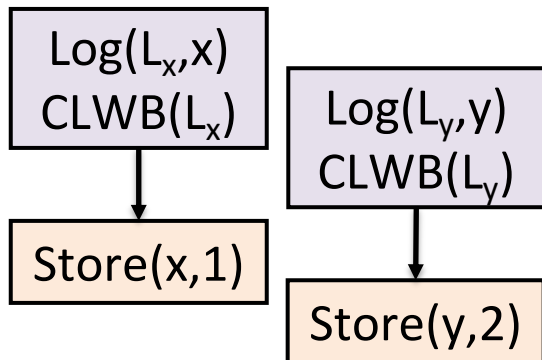
# Hardware imposes stricter constraints

## Ideal ordering

## SFENCE ordering

```

atomic_begin()
  x = 1;
  y = 2;
atomic_end()
  
```



Log(L<sub>x</sub>,x)  
CLWB(L<sub>x</sub>)

**SFENCE**

Store(x,1)

Log(L<sub>y</sub>,y)  
CLWB(L<sub>y</sub>)

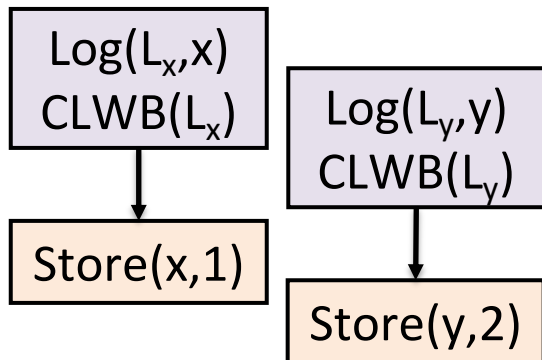
**SFENCE**

Store(y,2)

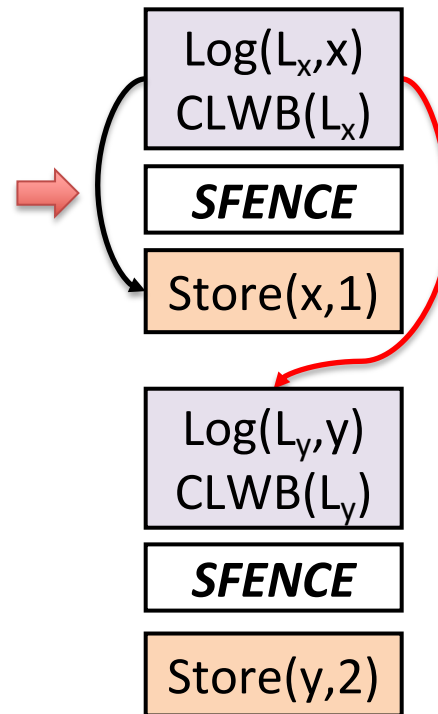
# Hardware imposes stricter constraints

## Ideal ordering

```
atomic_begin()
  x = 1;
  y = 2;
atomic_end()
```



## SFENCE ordering





# Hardware imposes stricter constraints

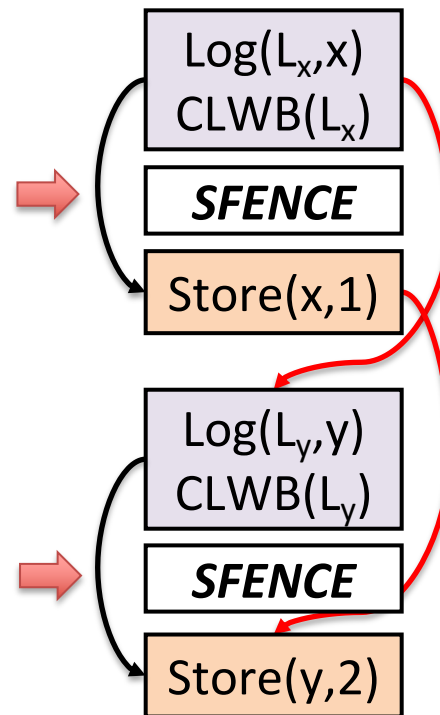
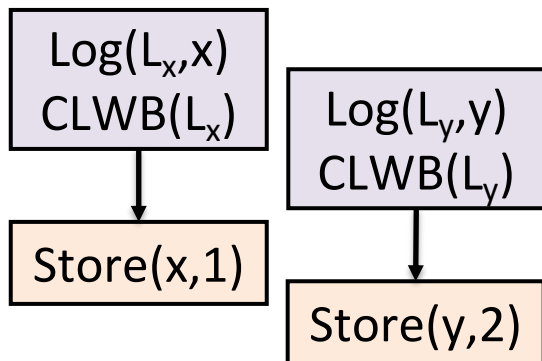
## Ideal ordering

## SFENCE ordering

```

atomic_begin()
  x = 1;
  y = 2;
atomic_end()

```





# StrandWeaver: Hardware Strand Persistency Model

High-level languages

Failure atomicity for language-level persistency models

Compiler

Logging impl. that map to hardware primitives

Hardware ISA

ISA primitives: PersistBarrier, NewStrand, JoinStrand



# StrandWeaver: Hardware Strand Persistency Model

High-level languages

Failure atomicity for language-level persistency models

Compiler

Logging impl. that map to hardware primitives

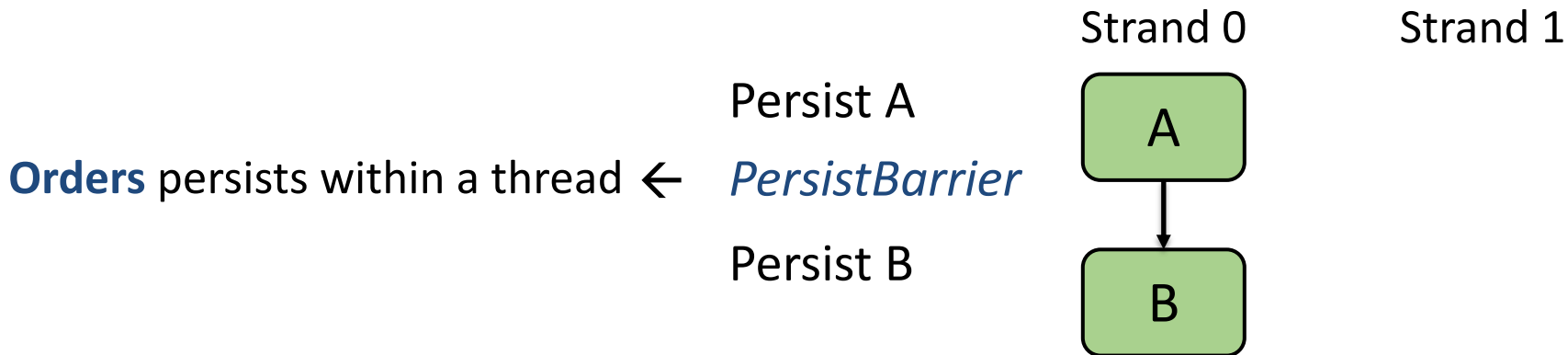
Hardware ISA

ISA primitives: PersistBarrier, NewStrand, JoinStrand



# StrandWeaver enables persist concurrency

- Provides primitives to express precise persist order





# StrandWeaver enables persist concurrency

- Provides primitives to express precise persist order

**Orders** persists within a thread ←

Persist A

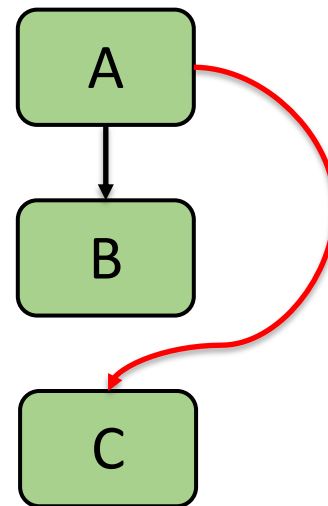
Persist B

Persist C

*PersistBarrier*

Strand 0

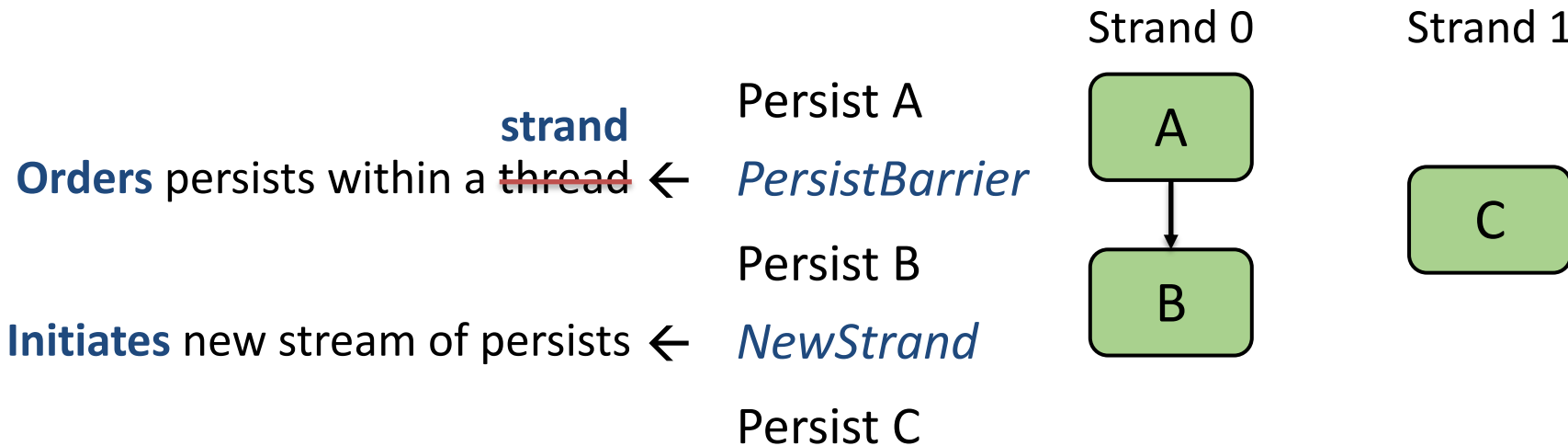
Strand 1





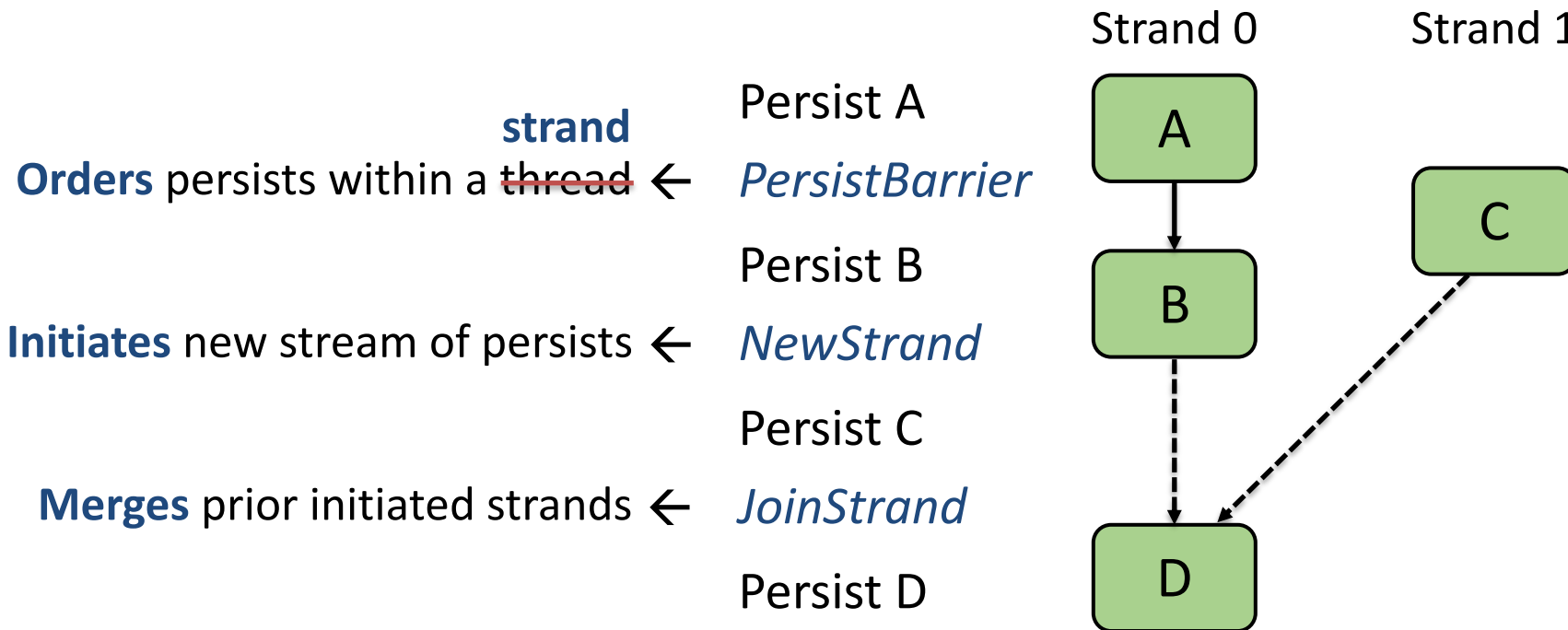
# StrandWeaver enables persist concurrency

- Provides primitives to express precise persist order

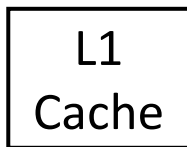
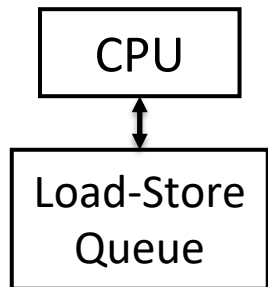


# StrandWeaver enables persist concurrency

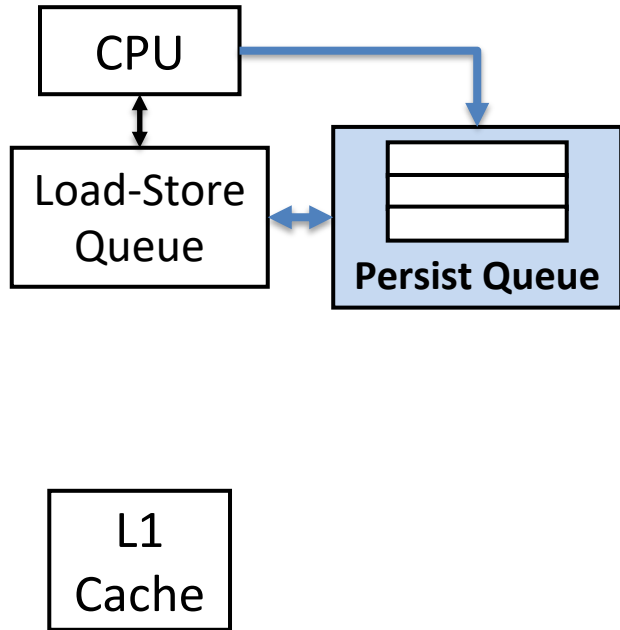
- Provides primitives to express precise persist order



# StrandWeaver architecture



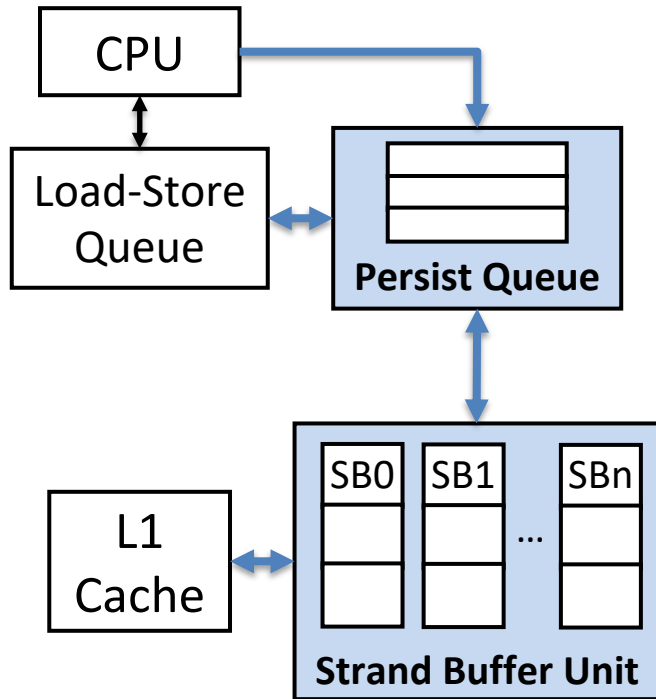
# StrandWeaver architecture



## Persist queue

- Manages ongoing StrandWeaver primitives
- Orders CLWBs separated by *JoinStrand*

# StrandWeaver architecture



## Persist queue

- Manages ongoing StrandWeaver primitives
- Orders CLWBs separated by *JoinStrand*

## Strand Buffer Unit

- Issues CLWBs and flushes dirty cache lines
- Ensures CLWBs on diff. strands are concurrent
- Monitors coherence reqs. for inter-thread order

# Running example

## Example code

CLWB(A)

*NewStrand*

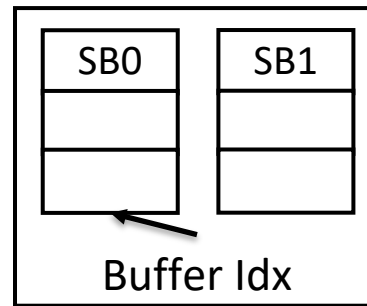
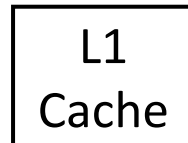
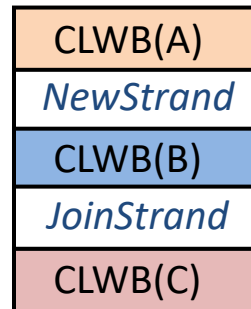
CLWB(B)

*JoinStrand*

CLWB(C)



## Persist Queue



## Strand Buffer Unit

# Running example

## Example code

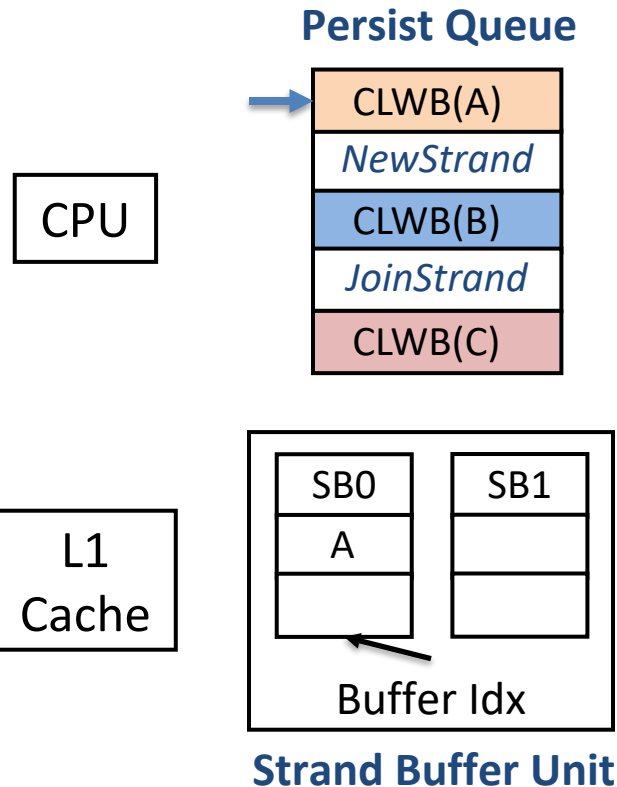
CLWB(A)

*NewStrand*

CLWB(B)

*JoinStrand*

CLWB(C)



# Running example

## Example code

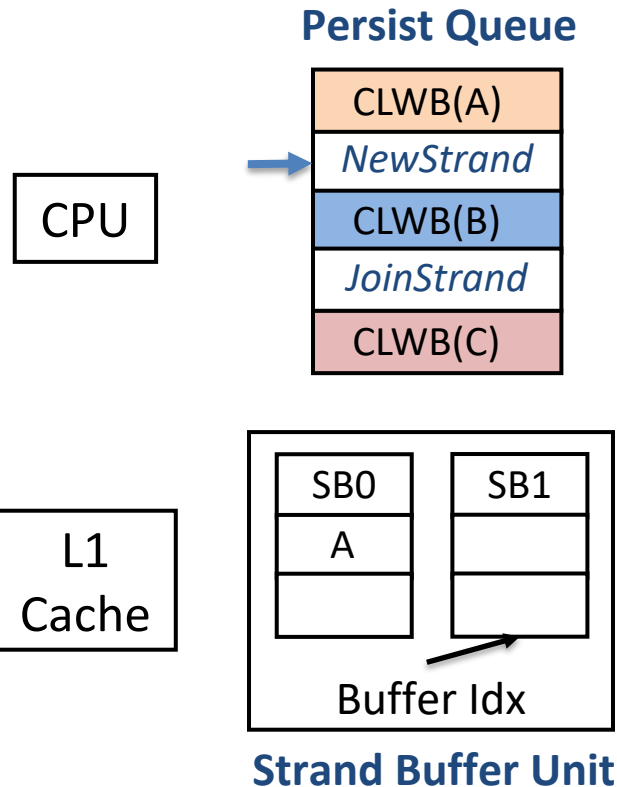
CLWB(A)

*NewStrand*

CLWB(B)

*JoinStrand*

CLWB(C)





# Running example

## Example code

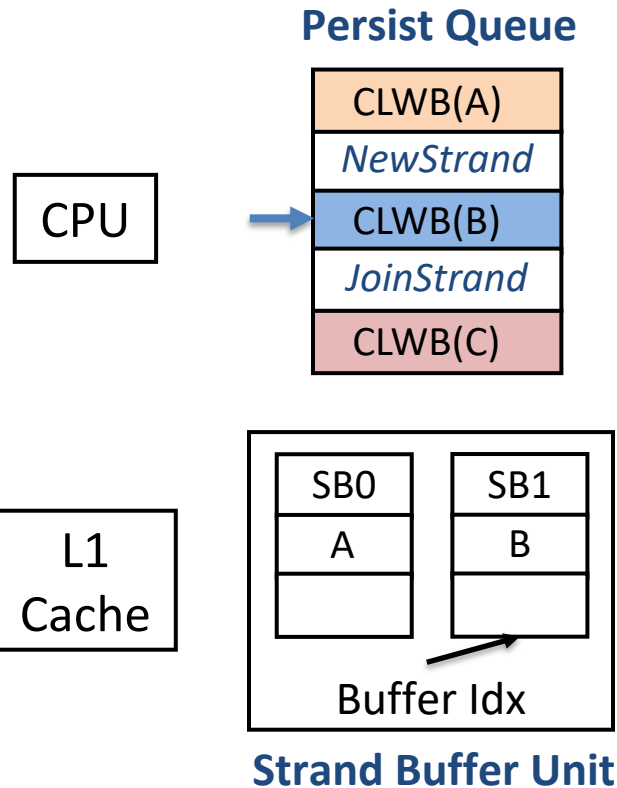
CLWB(A)

*NewStrand*

CLWB(B)

*JoinStrand*

CLWB(C)



# Running example

## Example code

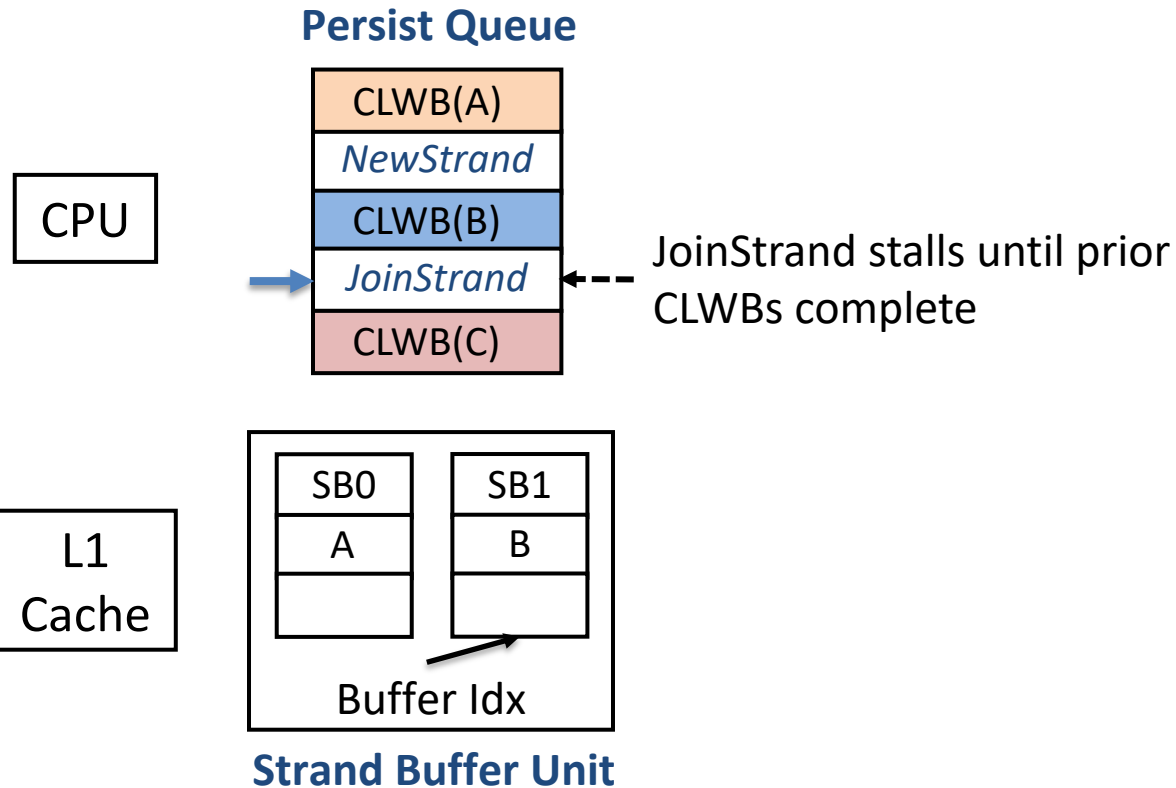
CLWB(A)

*NewStrand*

CLWB(B)

*JoinStrand*

CLWB(C)



# Running example

## Example code

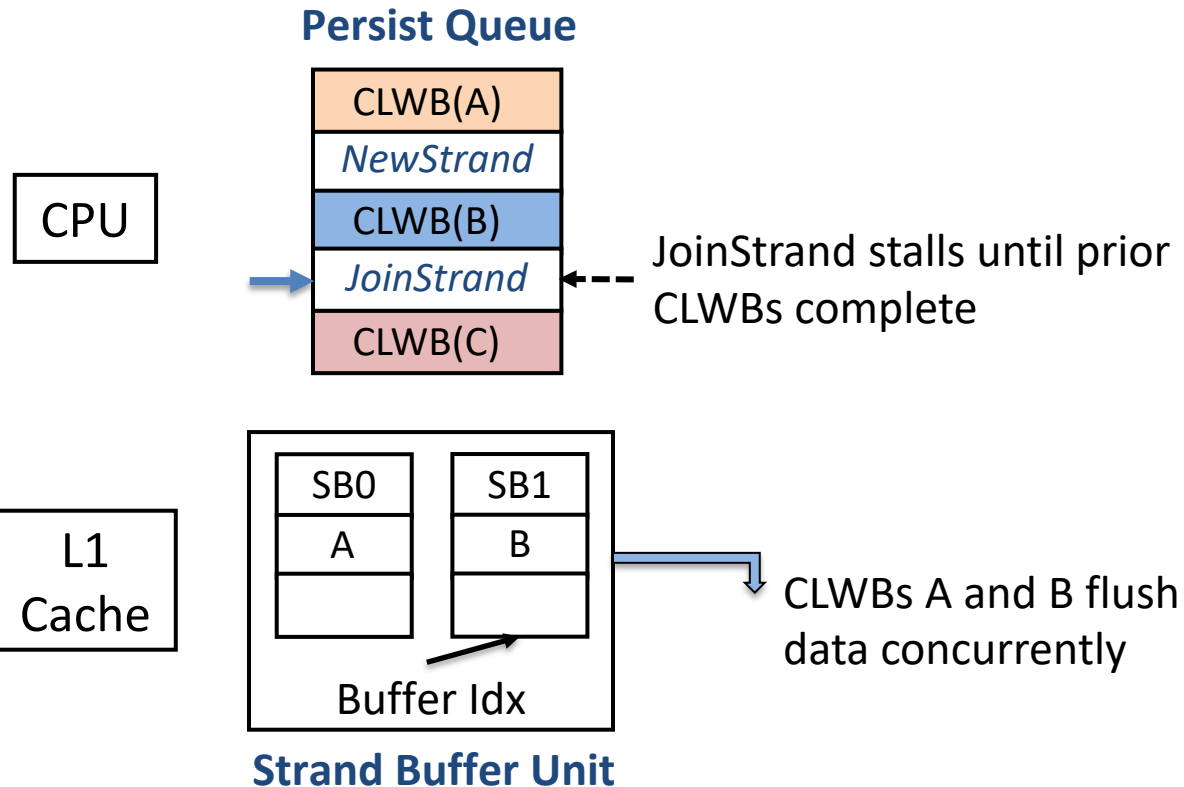
CLWB(A)

*NewStrand*

CLWB(B)

*JoinStrand*

CLWB(C)



# Running example

## Example code

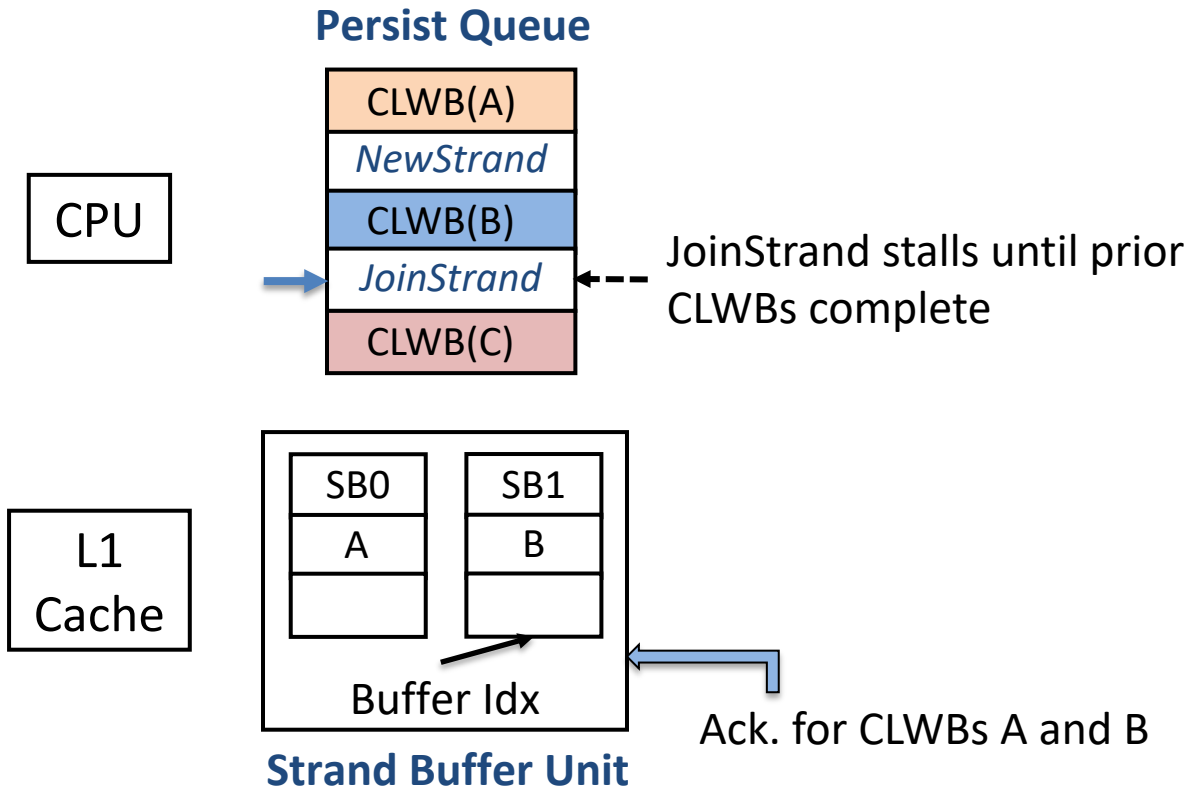
CLWB(A)

*NewStrand*

CLWB(B)

*JoinStrand*

CLWB(C)



# Running example

## Example code

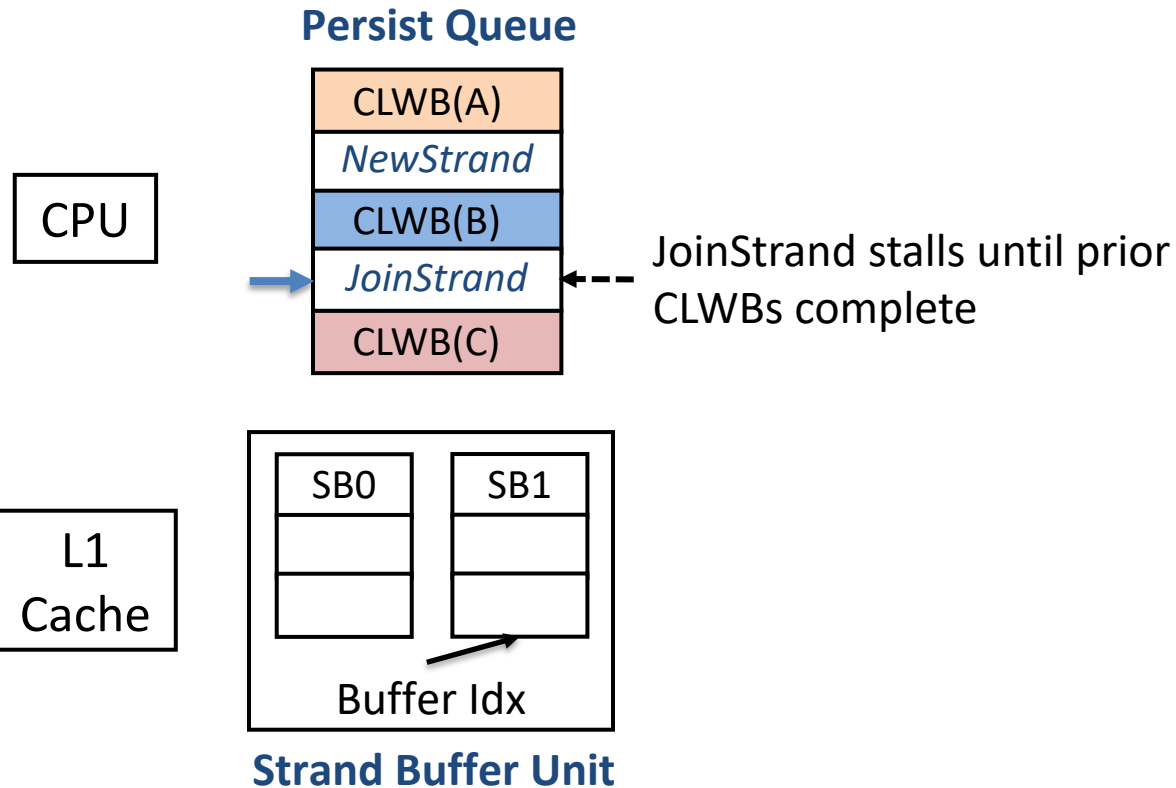
CLWB(A)

*NewStrand*

CLWB(B)

*JoinStrand*

CLWB(C)



# Running example

## Example code

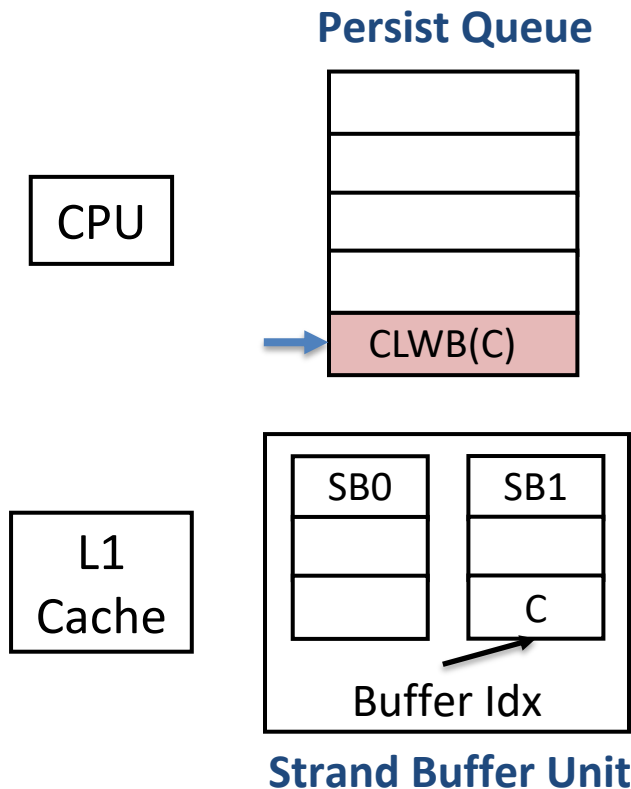
CLWB(A)

*NewStrand*

CLWB(B)

*JoinStrand*

CLWB(C)





# StrandWeaver: From ISA to high-level language

High-level languages

Failure atomicity for language-level persistency models

Compiler

Logging impl. that map to hardware primitives

Hardware ISA

ISA primitives: PersistBarrier, NewStrand, JoinStrand



# Logging using StrandWeaver primitives

```
atomic_begin()
```

```
    x = 1;
```

```
    y = 2;
```

```
atomic_end()
```

```
Log(Lx,x)
```

```
CLWB(Lx)
```

```
PersistBarrier
```

```
Store(x,1)
```

```
NewStrand
```

```
Log(Ly,y)
```

```
CLWB(Ly)
```

```
PersistBarrier
```

```
Store(y,2)
```

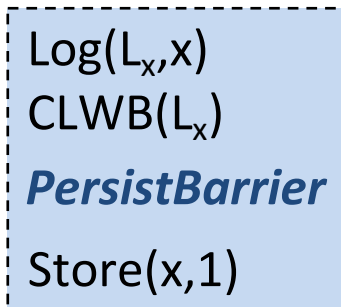
```
JoinStrand
```



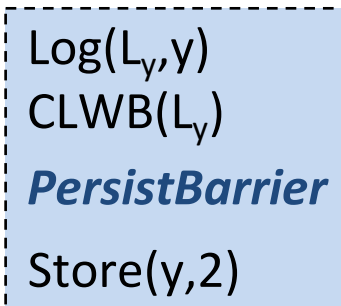


# Logging using StrandWeaver primitives

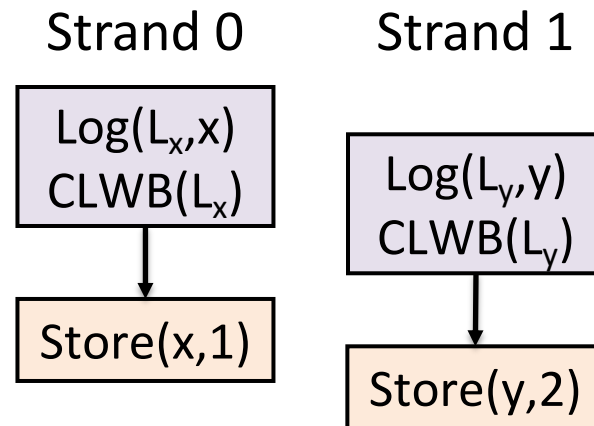
```
atomic_begin()  
  x = 1;  
  y = 2;  
atomic_end()
```



*NewStrand*



*JoinStrand*



# Logging using StrandWeaver primitives

```
atomic_begin()  
  x = 1;  
  y = 2;  
atomic_end()
```

Log( $L_x, x$ )

CLWB( $L_x$ )

**PersistBarrier**

Store( $x, 1$ )

**NewStrand**

Log( $L_y, y$ )

CLWB( $L_y$ )

**PersistBarrier**

Store( $y, 2$ )

**JoinStrand**

Strand 0

Log( $L_x, x$ )  
CLWB( $L_x$ )

Store( $x, 1$ )

Strand 1

Log( $L_y, y$ )  
CLWB( $L_y$ )

Store( $y, 2$ )



# StrandWeaver: From ISA to high-level language

High-level languages

Failure atomicity for language-level persistency models

Compiler

Logging impl. that map to hardware primitives

Hardware ISA

ISA primitives: PersistBarrier, NewStrand, JoinStrand



# High-level language implementations

```
L1.lock();  
    x -= 100;  
    y += 100;  
    L2.lock();  
        a -= 100;  
        b += 100;  
    L2.unlock();  
L1.unlock();
```

## **ATLAS** [Chakrabarti14]

- Failure-atomic outermost critical sections



# High-level language implementations

```
L1.lock();
```

```
    x -= 100;
```

```
    y += 100;
```

```
L2.lock();
```

```
    a -= 100;
```

```
    b += 100;
```

```
L2.unlock();
```

```
L1.unlock();
```

**ATLAS** [Chakrabarti14]

- Failure-atomic outermost critical sections

**Coupled-SFR** [Gogte18]

- Failure-atomic synchronization-free regions

**Decoupled-SFR** [Gogte18]

- Failure-atomic synchronization-free regions



# High-level language implementations

```
L1.lock();
```

```
x -= 100;
```

```
y += 100;
```

```
L2.lock();
```

```
a -= 100;
```

```
b += 100;
```

```
L2.unlock();
```

```
L1.unlock();
```

**ATLAS** [Chakrabarti14]

- Failure-atomic outermost critical sections

**Coupled-SFR** [Gogte18]

- Failure-atomic synchronization-free regions

**Decoupled-SFR** [Gogte18]

- Failure-atomic synchronization-free regions

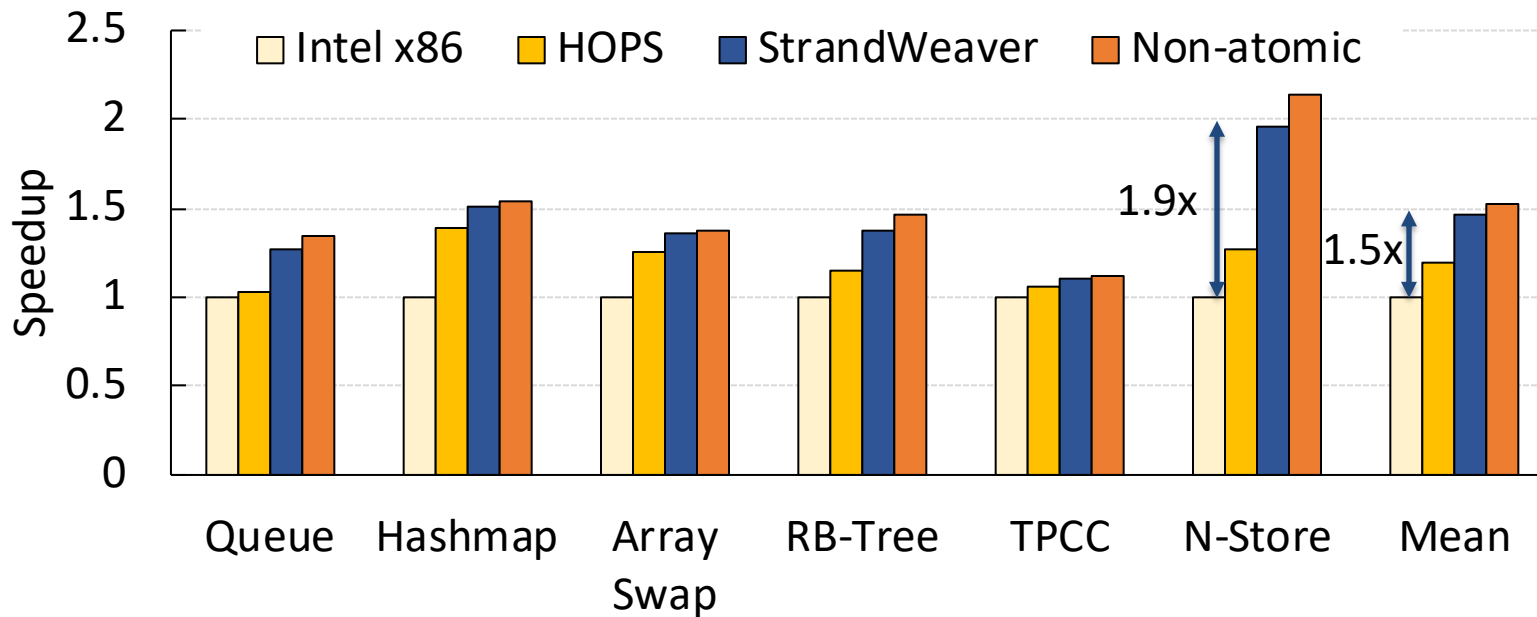


# Methodology

- Gem5 simulator
- Micro-benchmarks:
  - **Queue**: insert/delete entries in a queue
  - **Hashmap**: update values in persistent hash table
  - **Array swaps**: random swaps of array elements
  - **RBTree**: insert/delete entries in red-black tree
  - **TPCC**: new order transaction from TPCC
- Benchmarks:
  - **N-Store** [Arulraj15]: persistent KV-Store benchmark



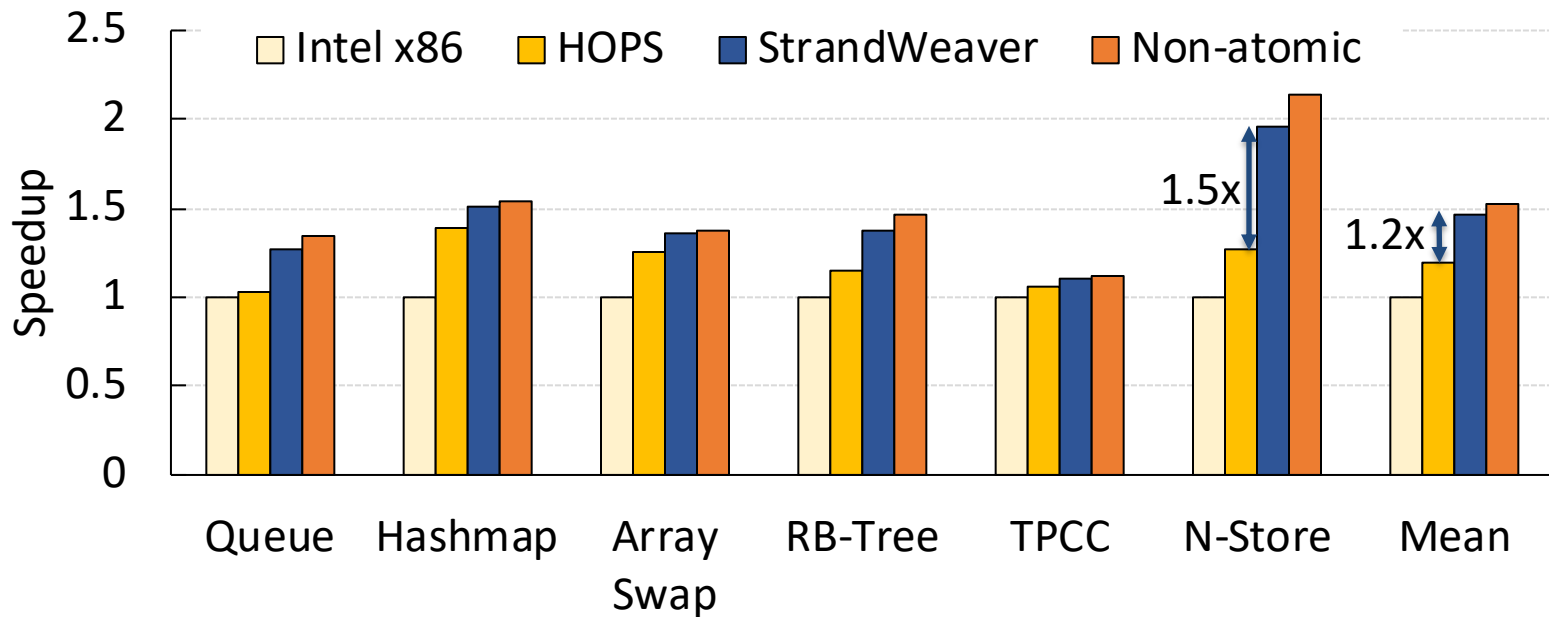
# Performance comparison with Intel x86



StrandWeaver achieves avg. speedup of 1.5x compared to the baseline

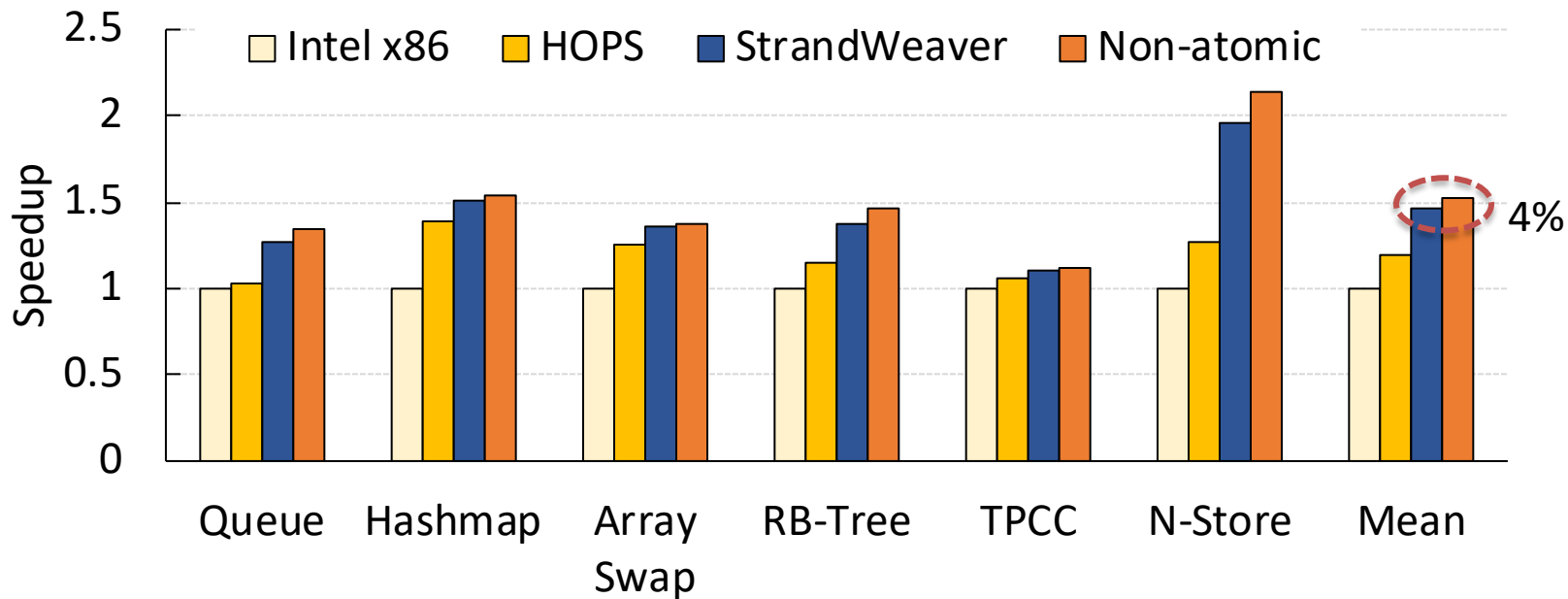


# Performance comparison with Intel x86



StrandWeaver achieves avg. speedup of 1.2x over HOPS

# Performance comparison with Intel x86



StrandWeaver performance is within 4% of non-atomic design

# Conclusion

- Strand persistency to precisely order persists
- Three primitives: **PersistBarrier**, **NewStrand** and **JoinStrand**
  - Work together to relax ordering constraints in undo logging
- Evaluation using language-level persistency models
- Performance improvement of 1.45x average over Intel x86

# Relaxed Persist Ordering Using Strand Persistency

Vaibhav Gogte, William Wang<sup>§</sup>, Stephan Diestelhorst<sup>§</sup>,  
Peter M. Chen, Satish Narayanasamy, Thomas F. Wenisch



ISCA 2020

