



THE UNIVERSITY
of
WISCONSIN
MADISON

Mnemosyne

Lightweight Persistent Memory

Haris Volos

Andres Jaan Tack, Michael M. Swift

University of Wisconsin – Madison



- Storage-Class Memory (SCM) enables memory-like storage
- **Persistent Memory** is an abstraction that enables direct access to SCM
- Durable memory transactions allow consistent in-place updates



- Features

- Memory-like interface (load/store)
- Short access time (1000x faster than flash)
- Non-volatile

- Technologies

- Phase Change Memory (PCM)
- Spin Torque Transfer RAM (STT-RAM)
- Memristors
- Flash-backed DRAM (+ supercapacitor)

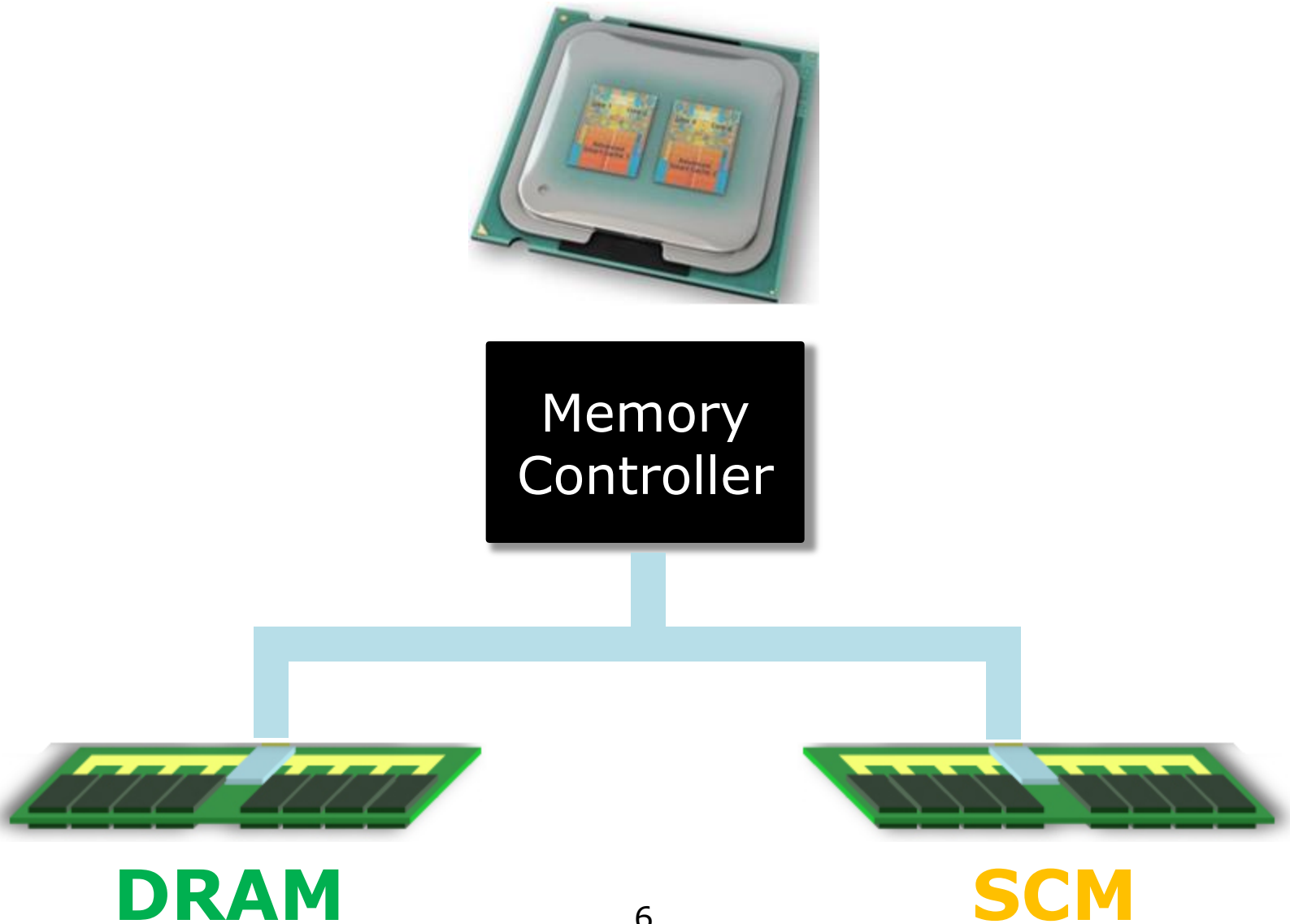


- Web applications
 - Amazon
 - Facebook
- Desktop applications
 - Firefox
- Other
 - Distributed agreement protocols
 - High-frequency trading



- Idea 1: File System
 - Layering overheads
- Idea 2: Persistent Object Stores
 - Single data-storage model
- Our idea: **Persistent Memory**
 - +Flexible and fine-grain
 - +Low-latency

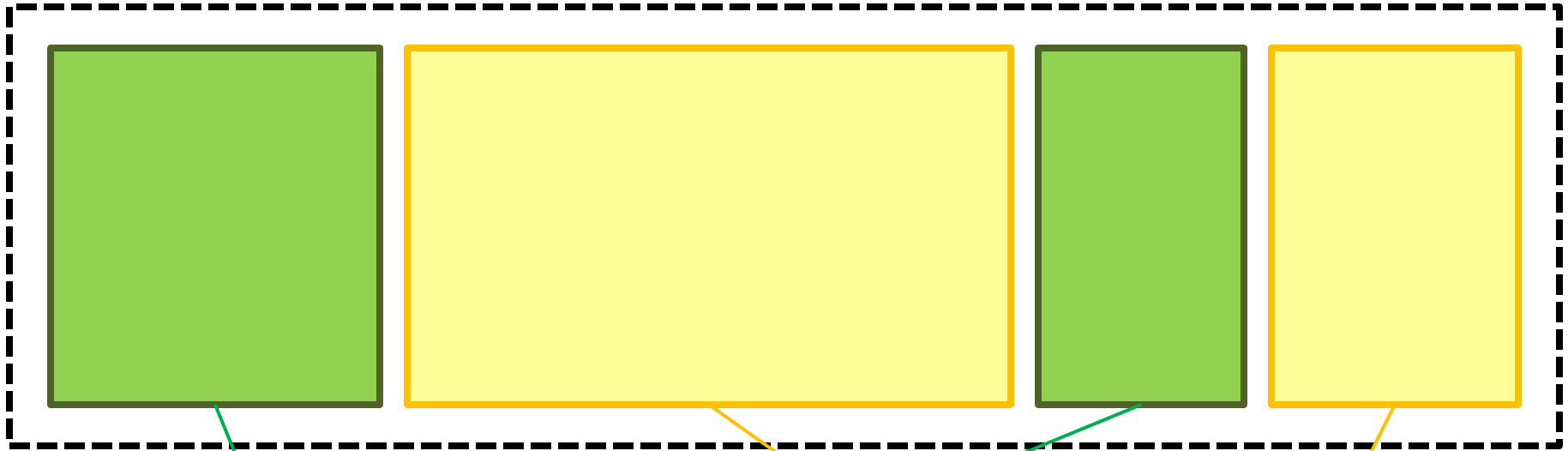
Persistent Memory



Persistent Memory



Address Space



DRAM

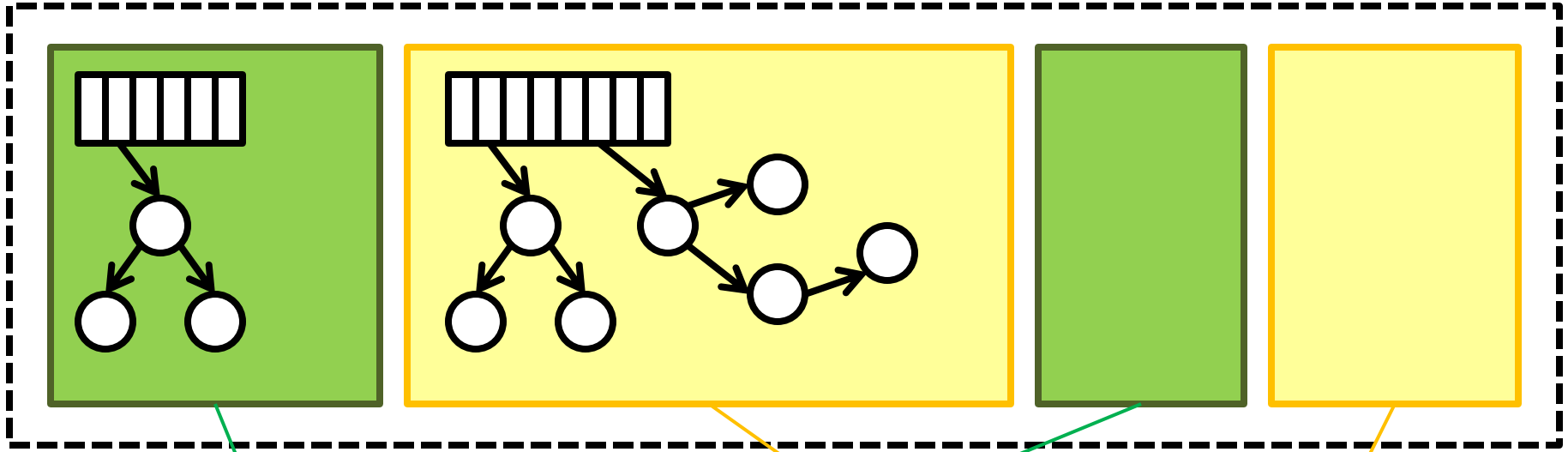


SCM



Persistent Memory

Address Space



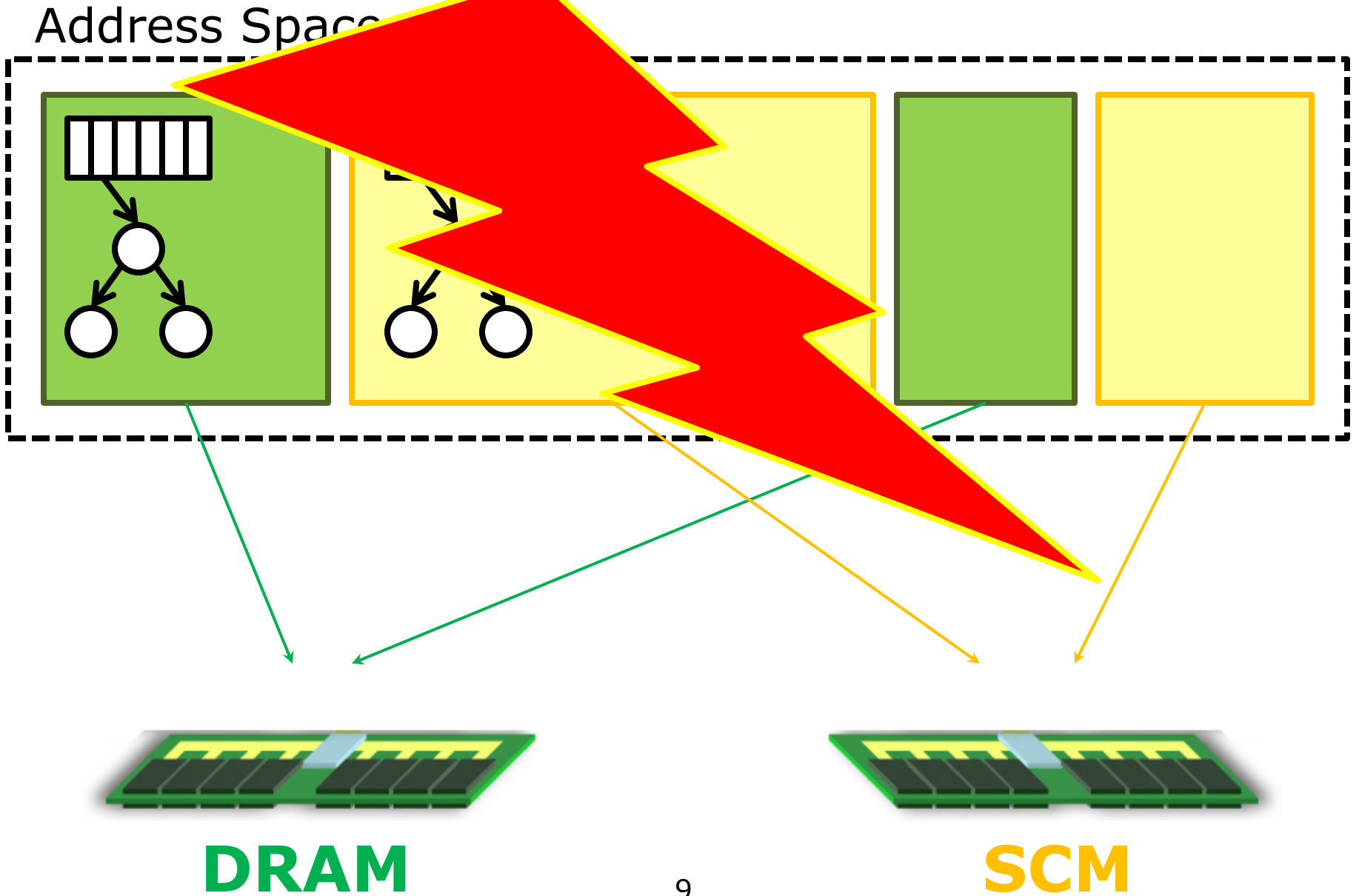
DRAM



SCM



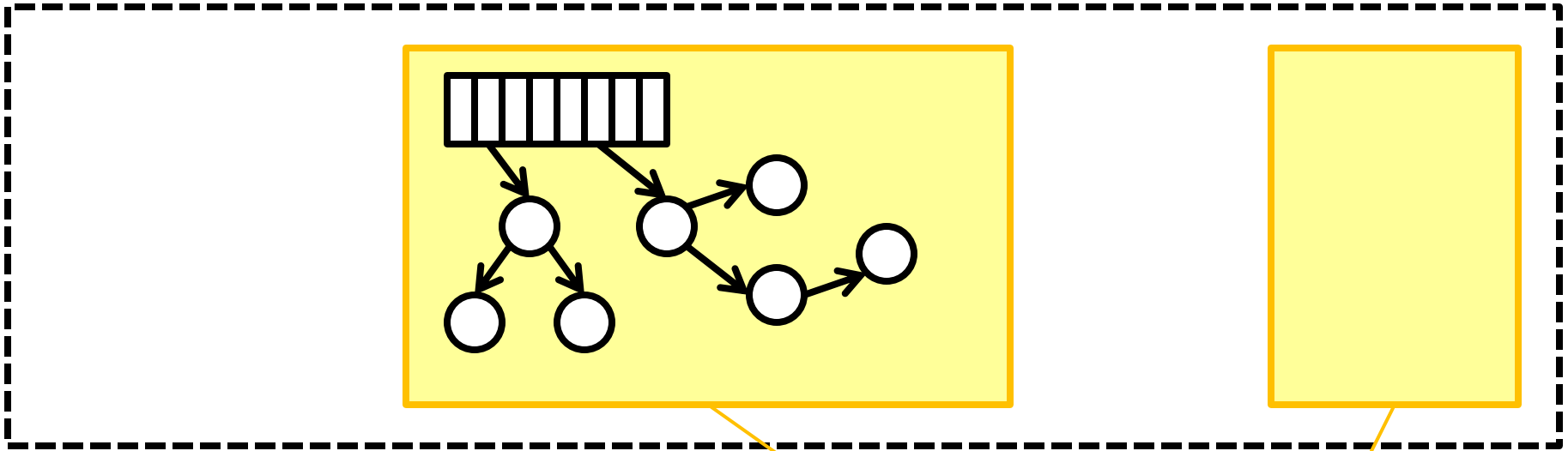
Persistent Memory





Persistent Memory

Address Space



DRAM

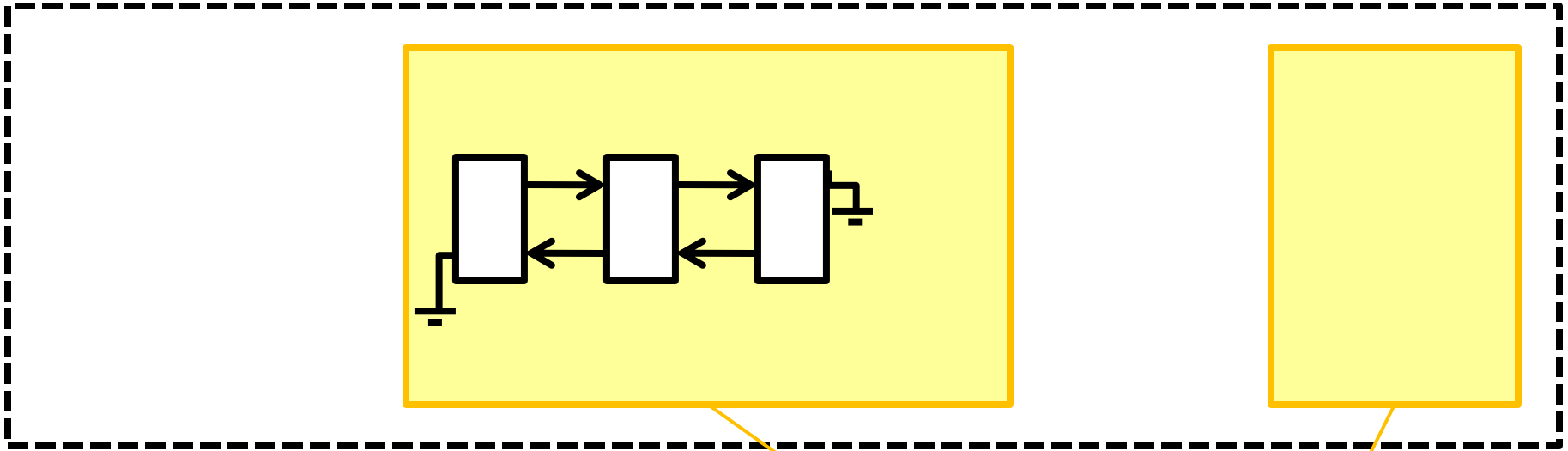


SCM



Persistent Memory

Address Space



DRAM

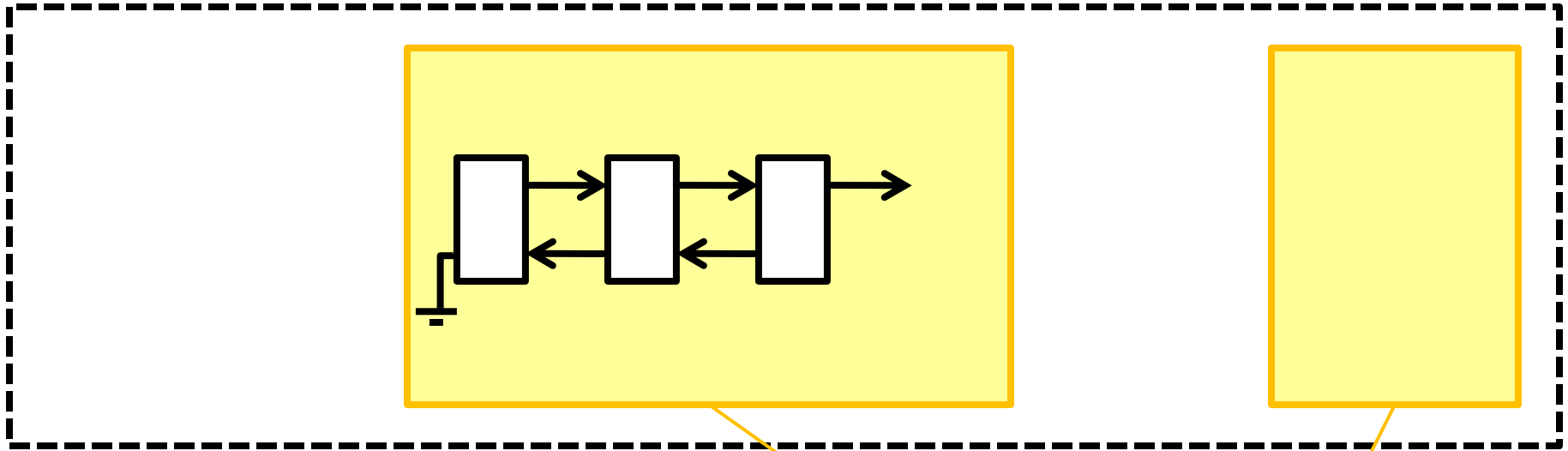


SCM



Persistent Memory

Address Space



DRAM

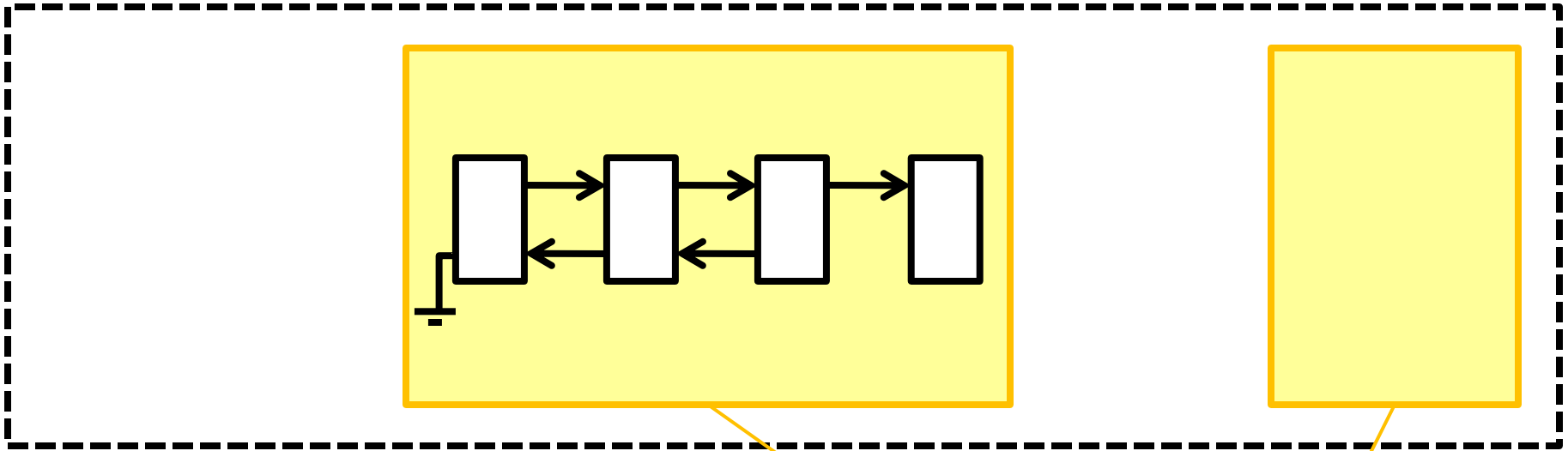


SCM



Persistent Memory

Address Space



DRAM

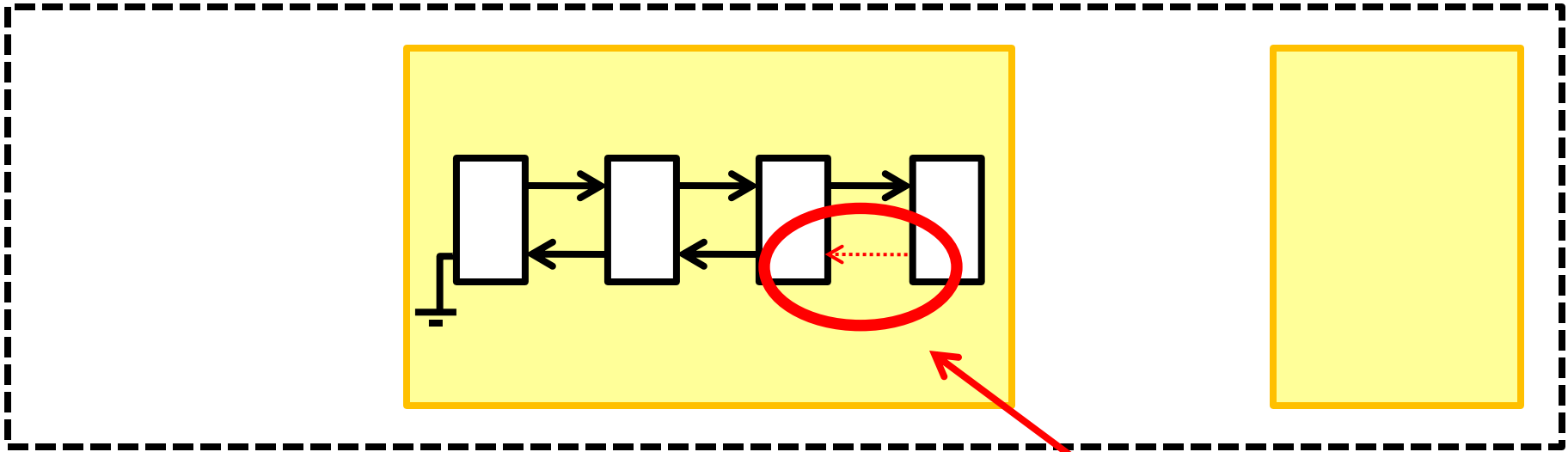


SCM



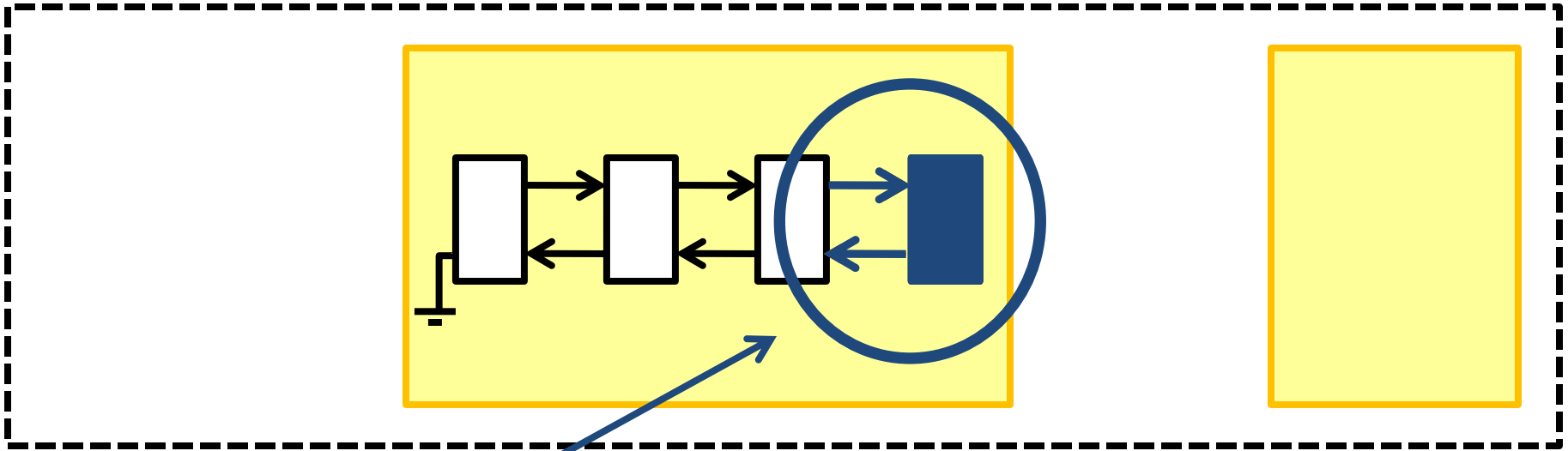
Persistent Memory

Address Space



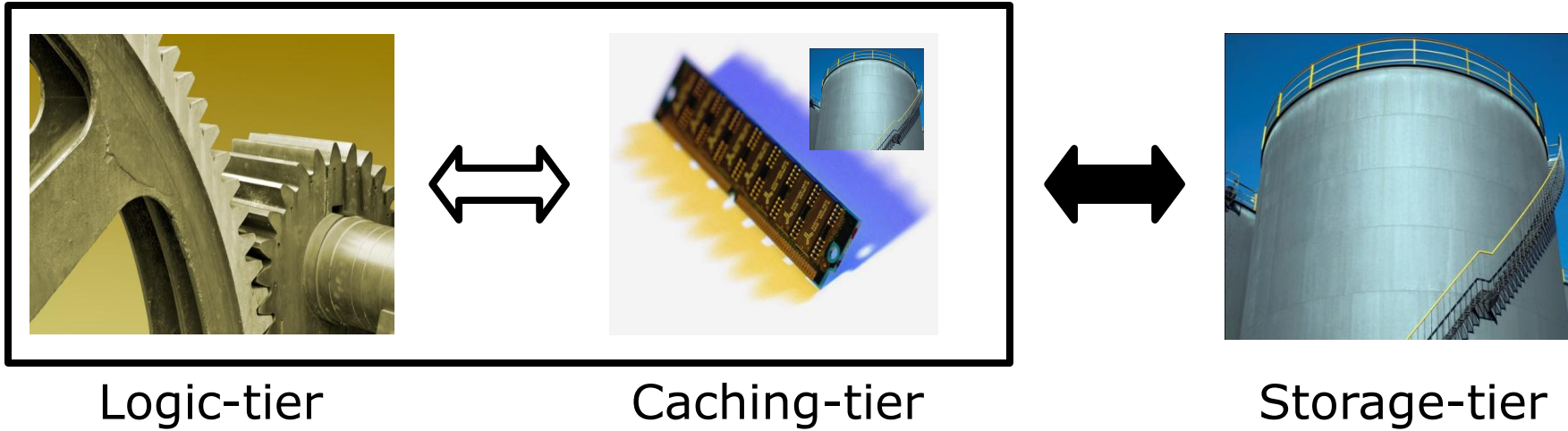
**Inconsistent state
(missing pointer)**

Address Space

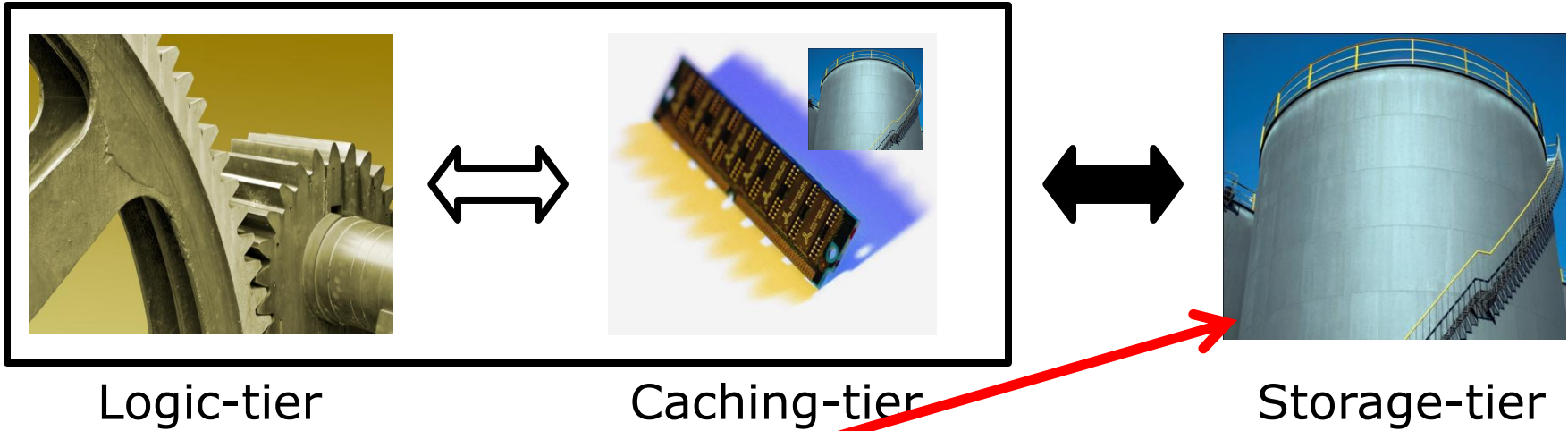


Consistent Updates
prevent inconsistent state

- In-vitro vs. In-vivo persistence

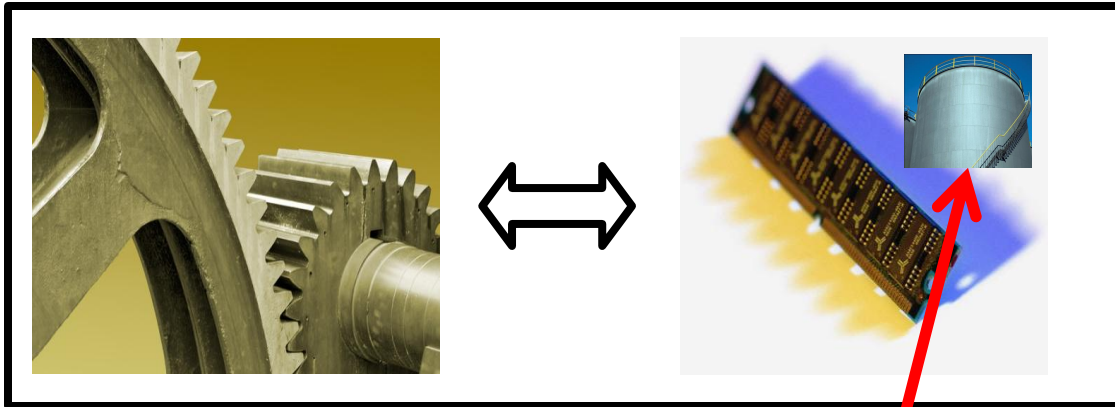


- In-vitro vs. In-vivo persistence



In-vitro persistence

- In-vitro vs. In-vivo persistence



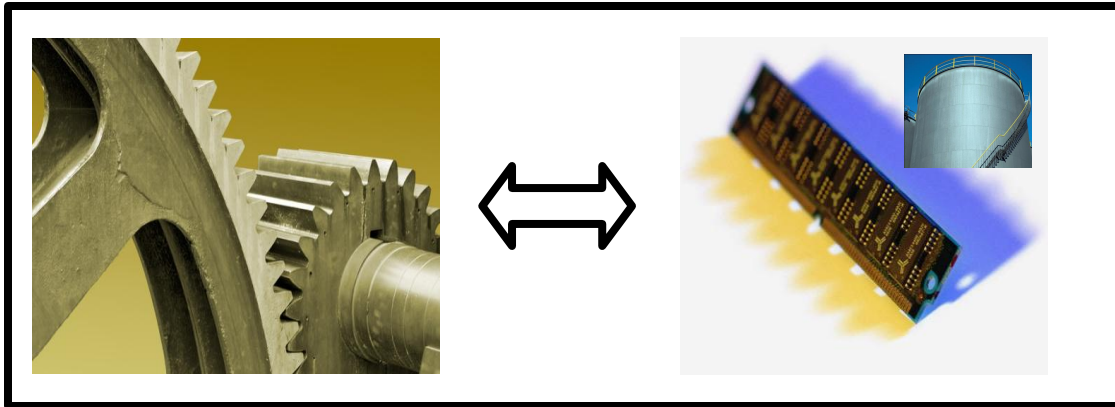
Logic-tier

Caching-tier

In-vivo persistence

In-vivo Persistence

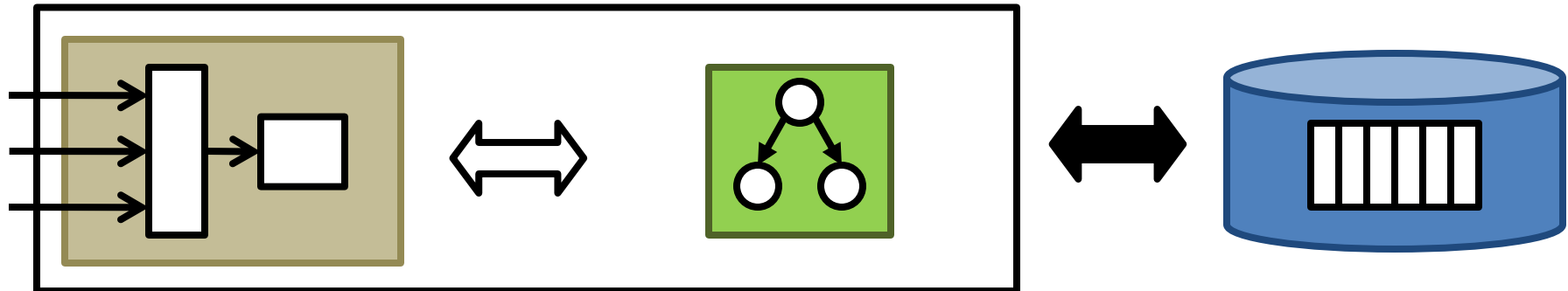
- In-vitro vs. In-vivo persistence



Logic-tier

Caching-tier

- Example: OpenLDAP



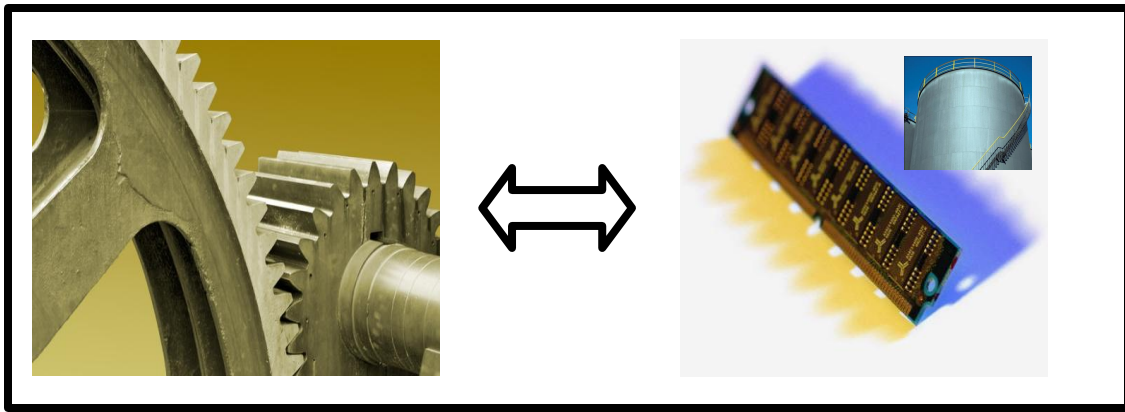
Handle Request

In-memory Directory

Berkeley DB

In-vivo Persistence

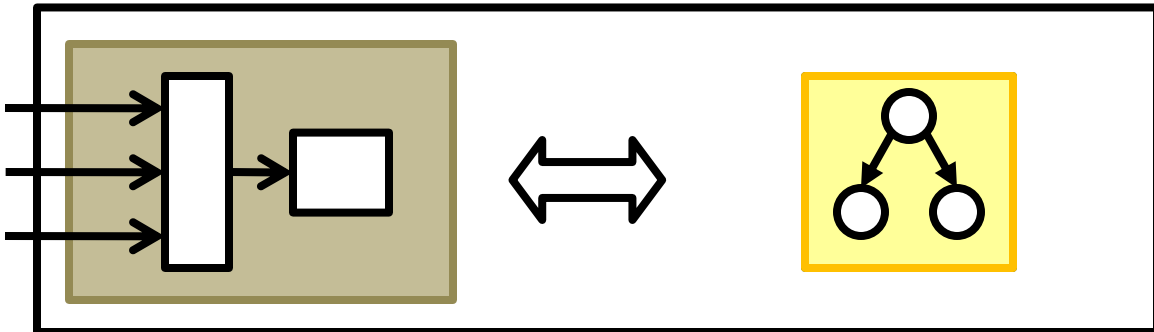
- In-vitro vs. In-vivo persistence



Logic-tier

Caching-tier

- Example: OpenLDAP



Handle Request

In-memory Directory



- Motivation
- **Persistent Memory**
 - **Persistent regions**
 - **Consistent updates**
- Durable Memory Transactions
- Evaluation



Persistent Regions

- Virtual memory segments stored in SCM
 - Enable programming flexibility
 - Persist across restarts

➤ Two region types

Persistent region type	API
Static	<code>pstatic var</code>
Dynamic (persistent heap)	<code>pmalloc</code>

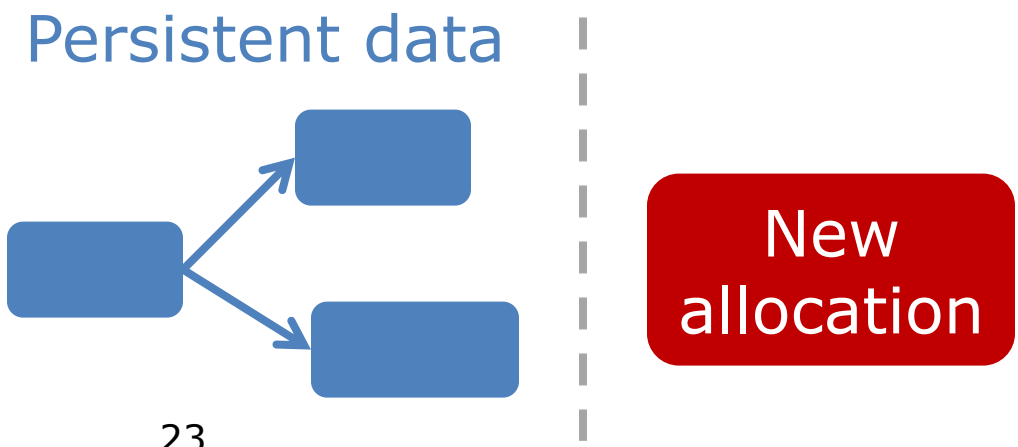


Persistent Regions

- Virtual memory segments stored in SCM
 - Enable programming flexibility
 - Persist across restarts

➤ Two region types

Persistent region type	API
Static	<code>pstatic var</code>
Dynamic (persistent heap)	<code>pmalloc</code>





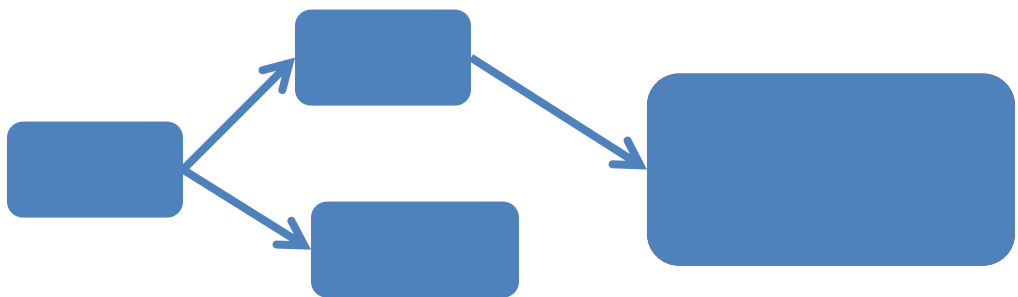
Persistent Regions

- Virtual memory segments stored in SCM
 - Enable programming flexibility
 - Persist across restarts

➤ Two region types

Persistent region type	API
Static	<code>pstatic var</code>
Dynamic (persistent heap)	<code>pmalloc</code>

Persistent data

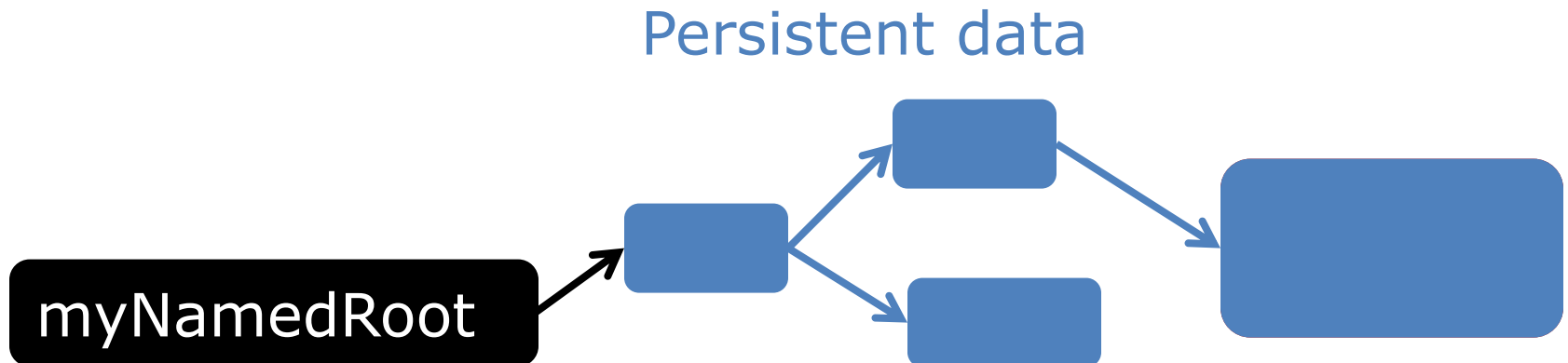




Persistent Regions

- Virtual memory segments stored in SCM
 - Enable programming flexibility
 - Persist across restarts
- Two region types

Persistent region type	API
Static	<code>pstatic var</code>
Dynamic (persistent heap)	<code>pmalloc</code>

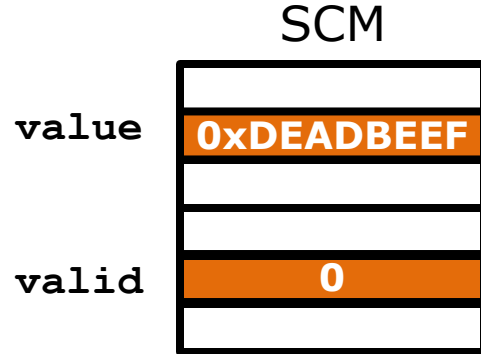
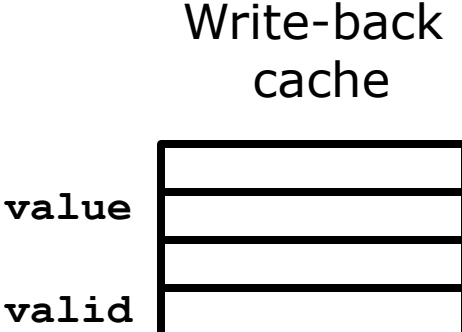




Consistent Updates

- Support updating data without risking correctness after a failure

R.value = 0xC0FFEE
R.valid = 1

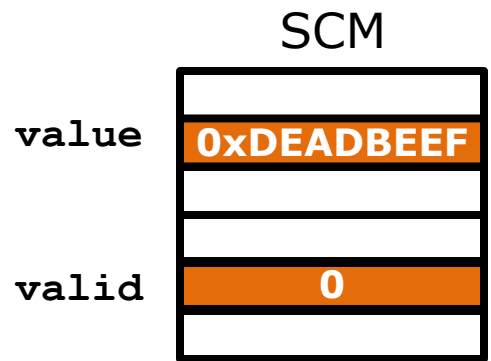
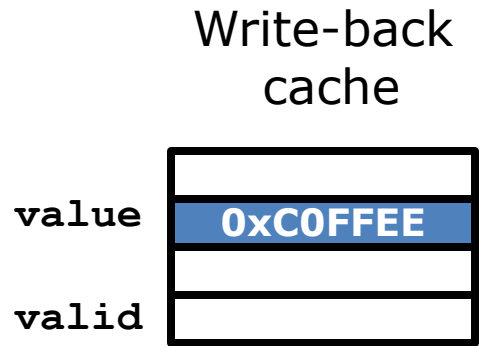




Consistent Updates

- Support updating data without risking correctness after a failure

R.value = 0xC0FFEE
R.valid = 1

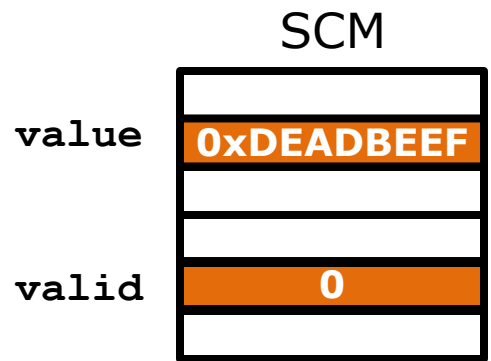
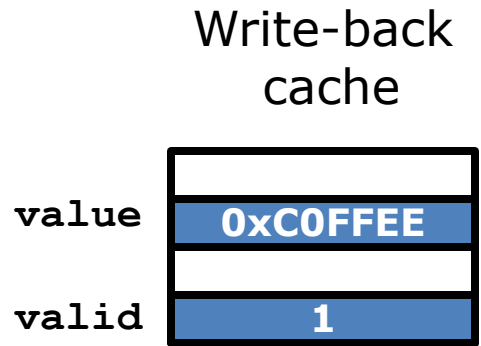




Consistent Updates

- Support updating data without risking correctness after a failure

R.value = 0xC0FFEE
R.valid = 1

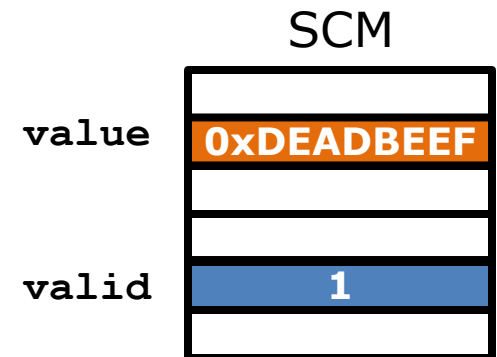
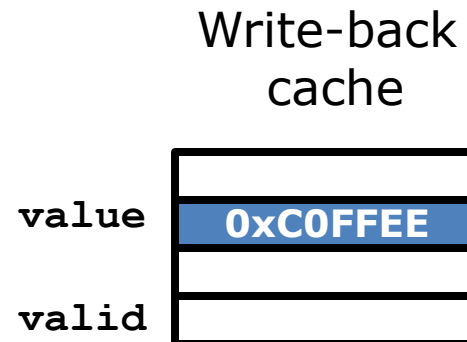


Consistent Updates



- Support updating data without risking correctness after a failure

`R.value = 0xC0FFEE`
`R.valid = 1`

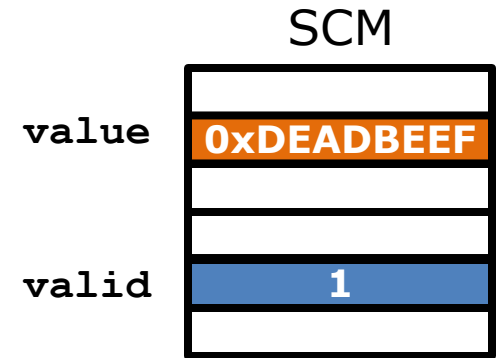
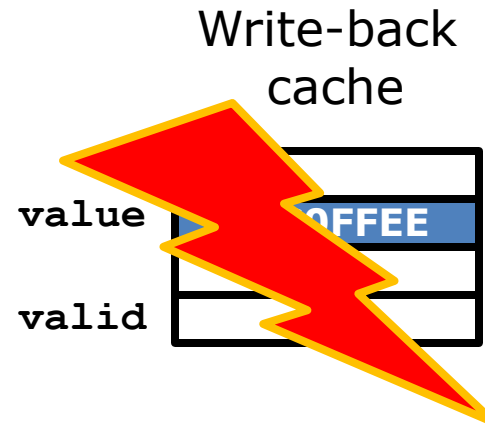


Consistent Updates



- Support updating data without risking correctness after a failure

```
R.value = 0xC0FFEE  
R.valid = 1
```

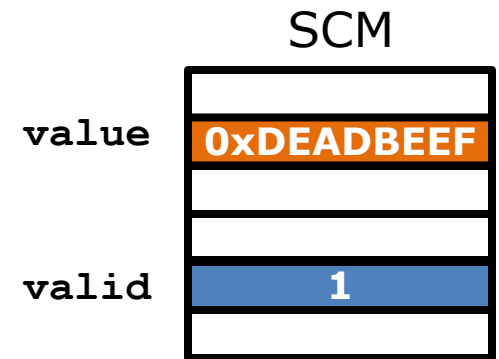
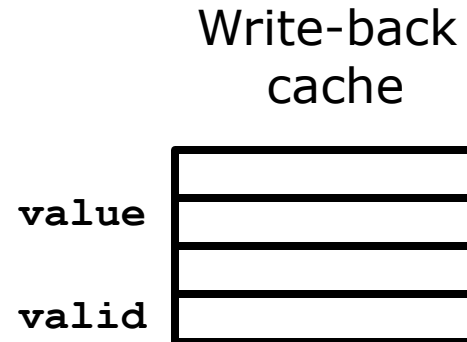


Consistent Updates



- Support updating data without risking correctness after a failure

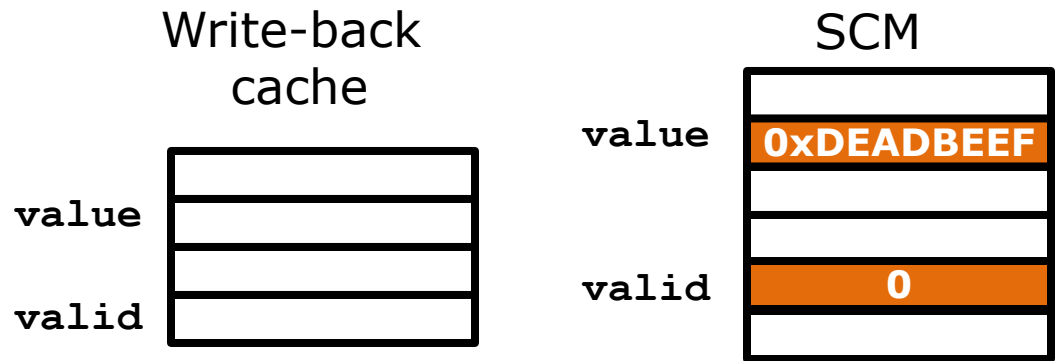
`R.value = 0xC0FFEE`
`R.valid = 1`





Consistent Updates

- Support updating data without risking correctness after a failure
- Rely on conventional CPU primitives for ordering
 - Flush
 - Fence

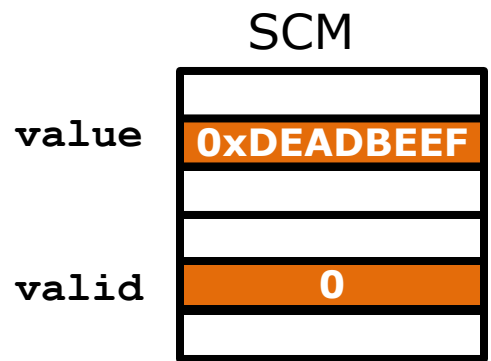
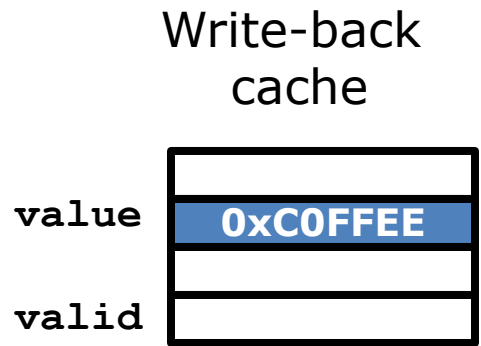




Consistent Updates

- Support updating data without risking correctness after a failure
- Rely on conventional CPU primitives for ordering
 - Flush
 - Fence

R.value ← 0xC0FFEE

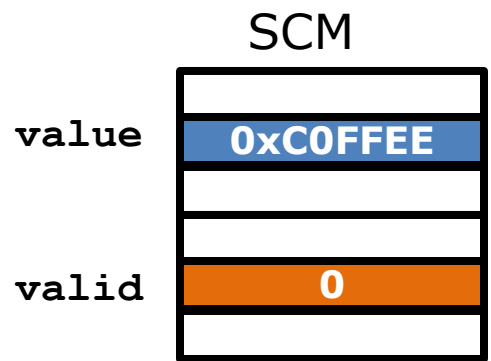
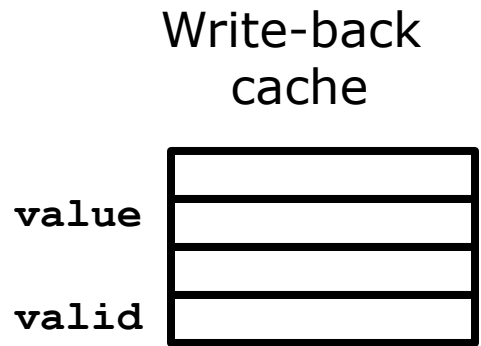




Consistent Updates

- Support updating data without risking correctness after a failure
- Rely on conventional CPU primitives for ordering
 - Flush
 - Fence

```
R.value ← 0xC0FFEE  
FLUSH (&R.value)
```

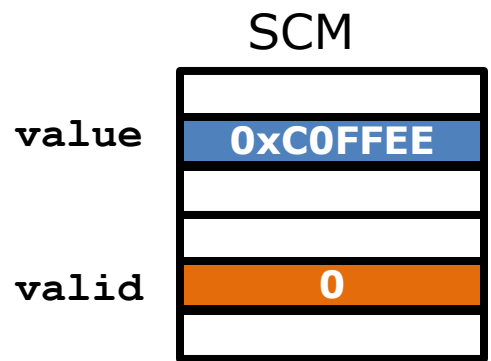
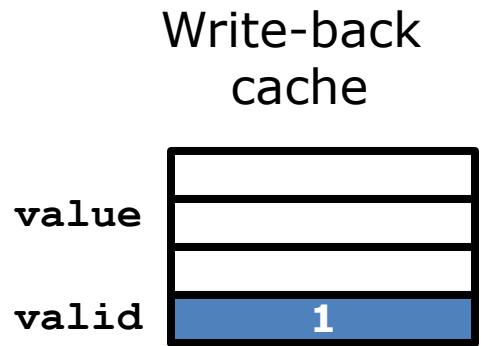




Consistent Updates

- Support updating data without risking correctness after a failure
- Rely on conventional CPU primitives for ordering
 - Flush
 - Fence

```
R.value ← 0xC0FFEE  
FLUSH (&R.value)  
FENCE  
R.valid ← 1
```

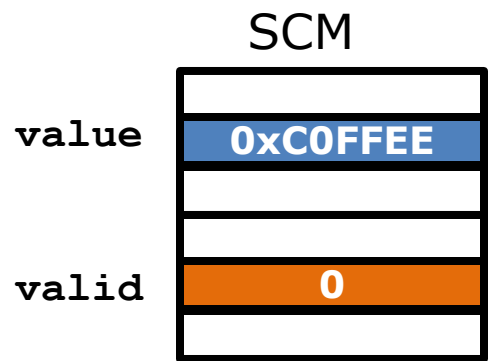
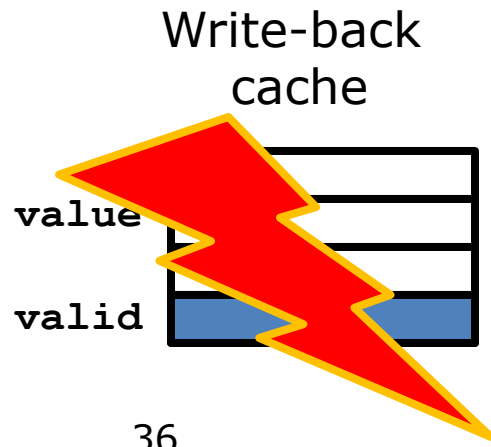




Consistent Updates

- Support updating data without risking correctness after a failure
- Rely on conventional CPU primitives for ordering
 - Flush
 - Fence

```
R.value ← 0xC0FFEE  
FLUSH (&R.value)  
FENCE  
R.valid ← 1
```

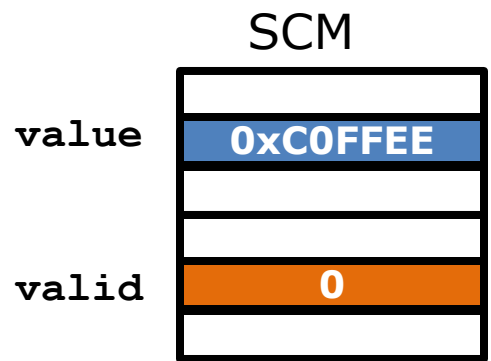
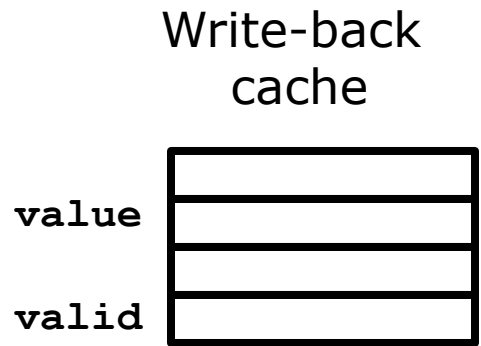




Consistent Updates

- Support updating data without risking correctness after a failure
- Rely on conventional CPU primitives for ordering
 - Flush
 - Fence

```
R.value ← 0xC0FFEE  
FLUSH (&R.value)  
FENCE  
R.valid ← 1
```





- Motivation
- Persistent Memory
 - Persistent regions
 - Consistent updates
- **Durable Memory Transactions**
- Evaluation



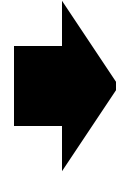
- Compiler instruments **atomic** blocks

```
atomic {  
  R.value = 0xC0F;  
  R.valid = 1;  
}
```



- Compiler instruments **atomic** blocks

```
atomic {  
  R.value = 0xC0F;  
  R.valid = 1;  
}
```



```
begin_transaction();
```

```
commit_transaction();
```




- Compiler instruments **atomic** blocks

```
atomic {  
    R.value = 0xC0F;  
    R.valid = 1;  
}  
    →  
begin_transaction();  
stm_store(&R.value, 0xC0F);  
stm_store(&R.valid, 1);  
commit_transaction();
```



Durable Memory Transactions

- Compiler instruments **atomic** blocks

```
atomic {  
  R.value = 0xC0F;  
  R.valid = 1;  
}  
begin_transaction();  
stm_store(&R.value, 0xC0F);  
stm_store(&R.valid, 1);  
commit_transaction();
```

- Runtime supports ACID transactions
 - Based on TinySTM

Hash Table Example



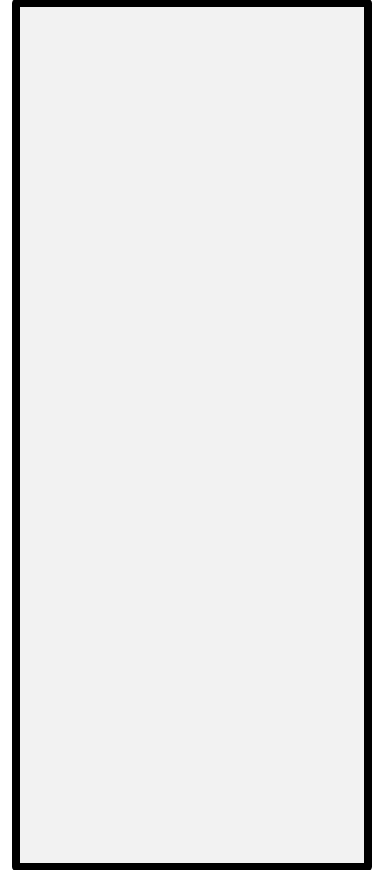
```
pstatic htRoot = NULL;
```

htRoot

0000000000

Static

Persistent Heap

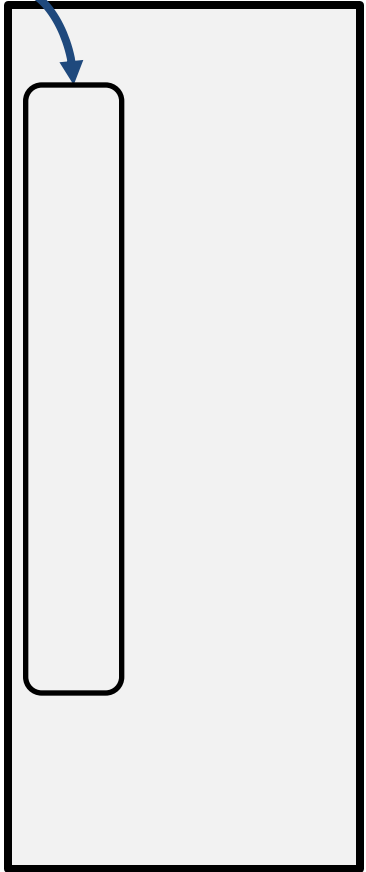




Hash Table Example

```
pstatic htRoot = NULL;
```

```
main() {  
    if (!htRoot)  
        pmalloc(N*sizeof(*bucket), &htRoot);  
}
```



Static

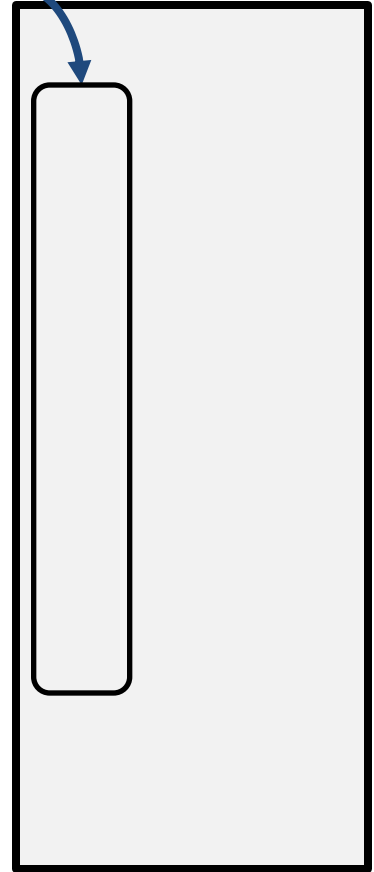
Persistent Heap



Hash Table Example

```
pstatic htRoot = NULL;
```

```
main() {  
    if (!htRoot)  
        pmalloc(N*sizeof(*bucket), &htRoot);  
}
```



Static

Persistent Heap



Static

Persistent Heap

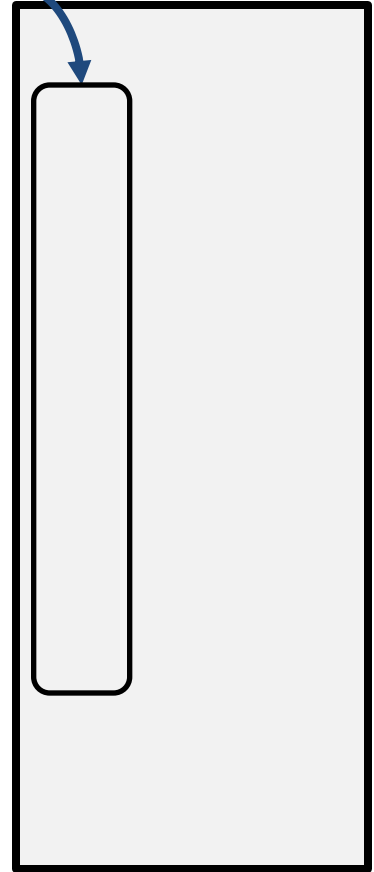
Hash Table Example

```
pstatic htRoot = NULL;
```



```
main() {  
    if (!htRoot)  
        pmalloc(N*sizeof(*bucket), &htRoot);  
}
```

```
update_hash(key, value) {  
    atomic {  
        pmalloc(sizeof(*bucket), &bucket);  
        bucket->key = key;  
        bucket->value = value;  
        insert(hash, bucket);  
    }  
}
```



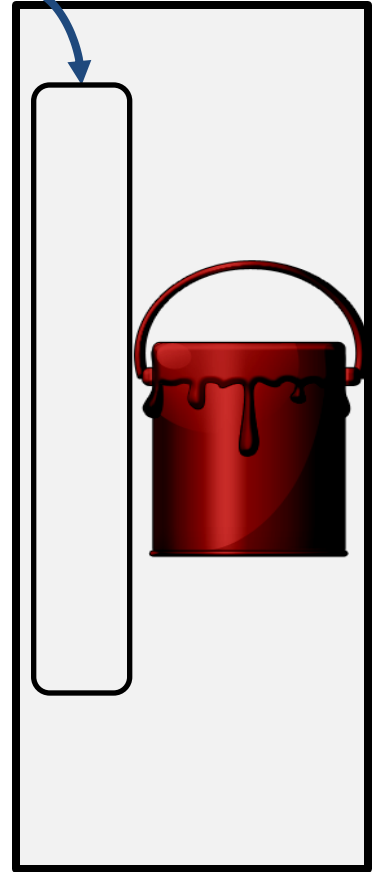


Hash Table Example

```
pstatic htRoot = NULL;
```

```
main() {  
    if (!htRoot)  
        pmalloc(N*sizeof(*bucket), &htRoot);  
}
```

```
update_hash(key, value) {  
    atomic {  
        pmalloc(sizeof(*bucket), &bucket);  
        bucket->key = key;  
        bucket->value = value;  
        insert(hash, bucket);  
    }  
}
```



Static

Persistent Heap

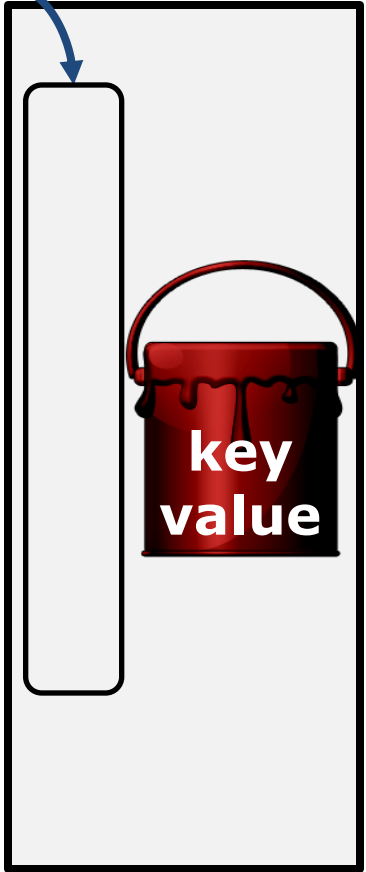


Hash Table Example

```
pstatic htRoot = NULL;
```

```
main() {  
    if (!htRoot)  
        pmalloc(N*sizeof(*bucket), &htRoot);  
}
```

```
update_hash(key, value) {  
    atomic {  
        pmalloc(sizeof(*bucket), &bucket);  
        bucket->key = key;  
        bucket->value = value;  
        insert(hash, bucket);  
    }  
}
```



Static
Persistent Heap

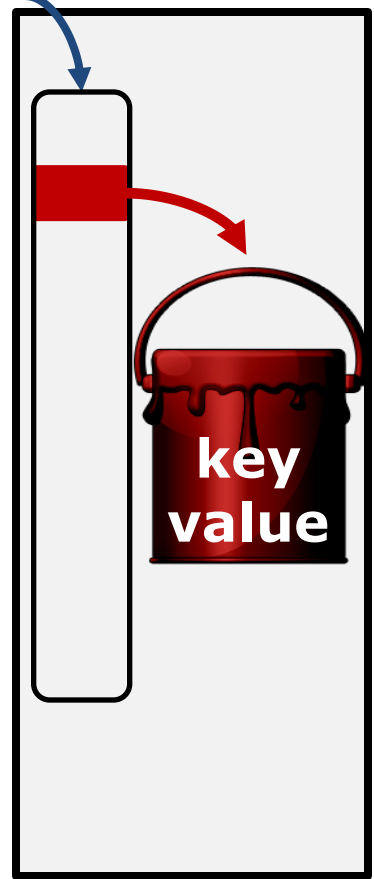


Hash Table Example

```
pstatic htRoot = NULL;
```

```
main() {  
    if (!htRoot)  
        pmalloc(N*sizeof(*bucket), &htRoot);  
}
```

```
update_hash(key, value) {  
    atomic {  
        pmalloc(sizeof(*bucket), &bucket);  
        bucket->key = key;  
        bucket->value = value;  
        insert(hash, bucket);  
    }  
}
```



Static
Persistent Heap



Static

Persistent Heap

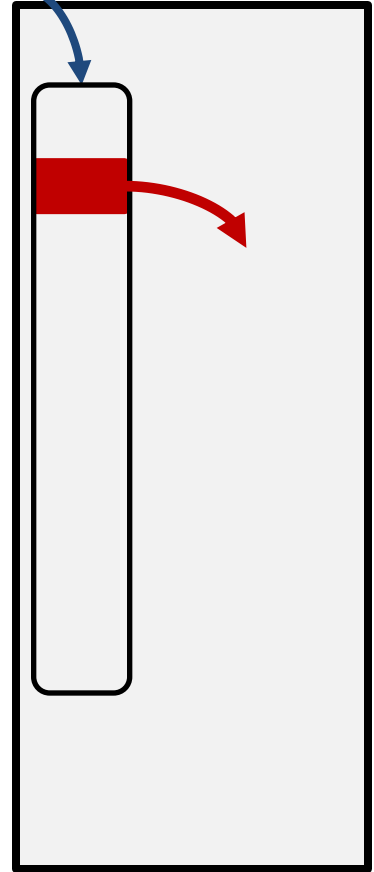
Hash Table Example

```
pstatic htRoot = NULL;
```



```
main() {  
    if (!htRoot)  
        pmalloc(N*sizeof(*bucket), &htRoot);  
}
```

```
update_hash(key, value) {  
    atomic {  
        pmalloc(sizeof(*bucket), &bucket);  
        bucket->key = key;  
        bucket->value = value;  
        insert(hash, bucket);  
    }  
}
```



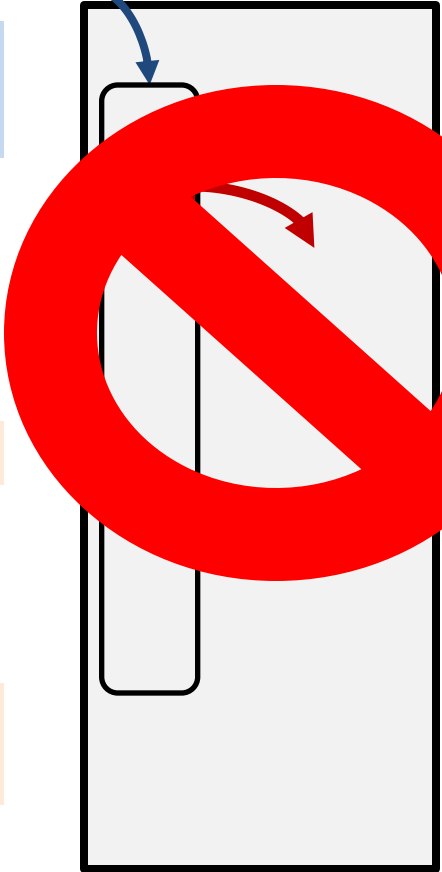


Hash Table Example

```
pstatic htRoot = NULL;
```

```
main() {  
    if (!htRoot)  
        pmalloc(N*sizeof(*bucket), &htRoot);  
}
```

```
update_hash(key, value) {  
    atomic {  
        pmalloc(sizeof(*bucket), &bucket);  
        bucket->key = key;  
        bucket->value = value;  
        insert(hash, bucket);  
    }  
}
```



Static
Heap



- Motivation
- Persistent Memory
 - Persistent regions
 - Consistent updates
- Durable Memory Transactions
- **Evaluation**



- Is it easy to use?
 - Applications

- Is it fast?
 - Microbenchmarks
 - Applications



- **TokyoCabinet**: B-tree based key-value store
 - Original: syncs B-tree to a memory-mapped file
 - Modified: keeps B-tree in persistent memory

- **OpenLDAP**: Lightweight Directory
 - Original: stores dir-entries in Berkeley DB
 - Modified: keeps dir-entries in persistent memory

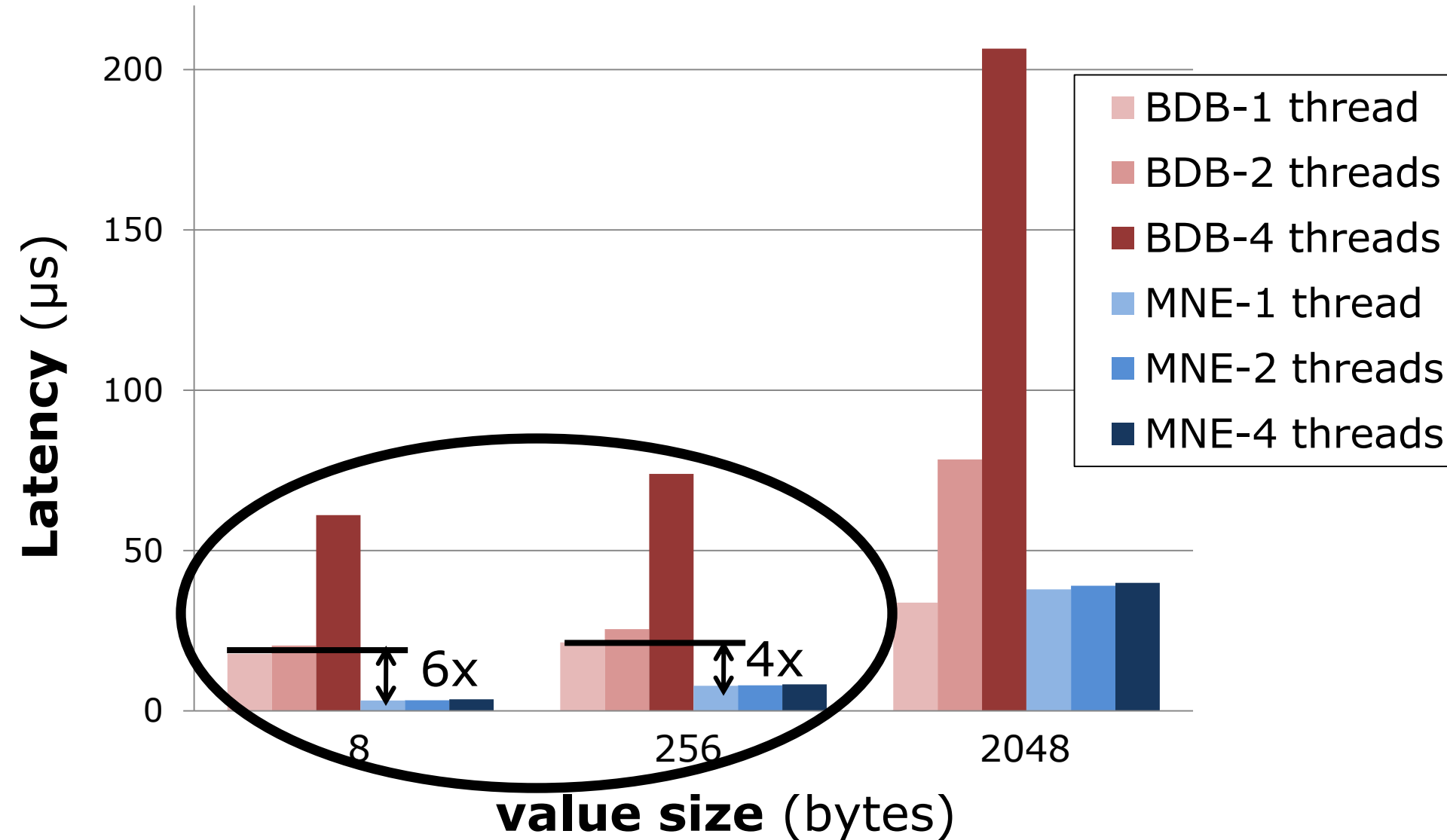


- Platform
 - Intel Core 2 2.5GHz (4 cores)
 - x86-64 Linux 2.6.33
 - Intel STM C/C++ compiler
- Performance Model
 - Writes: DRAM + extra fixed delay (150ns)
 - Reads: DRAM
- Configurations
 - PCM-disk: ext2fs + RAM-disk
 - Mnemosyne: Persistent memory

Microbenchmark – Hash Table



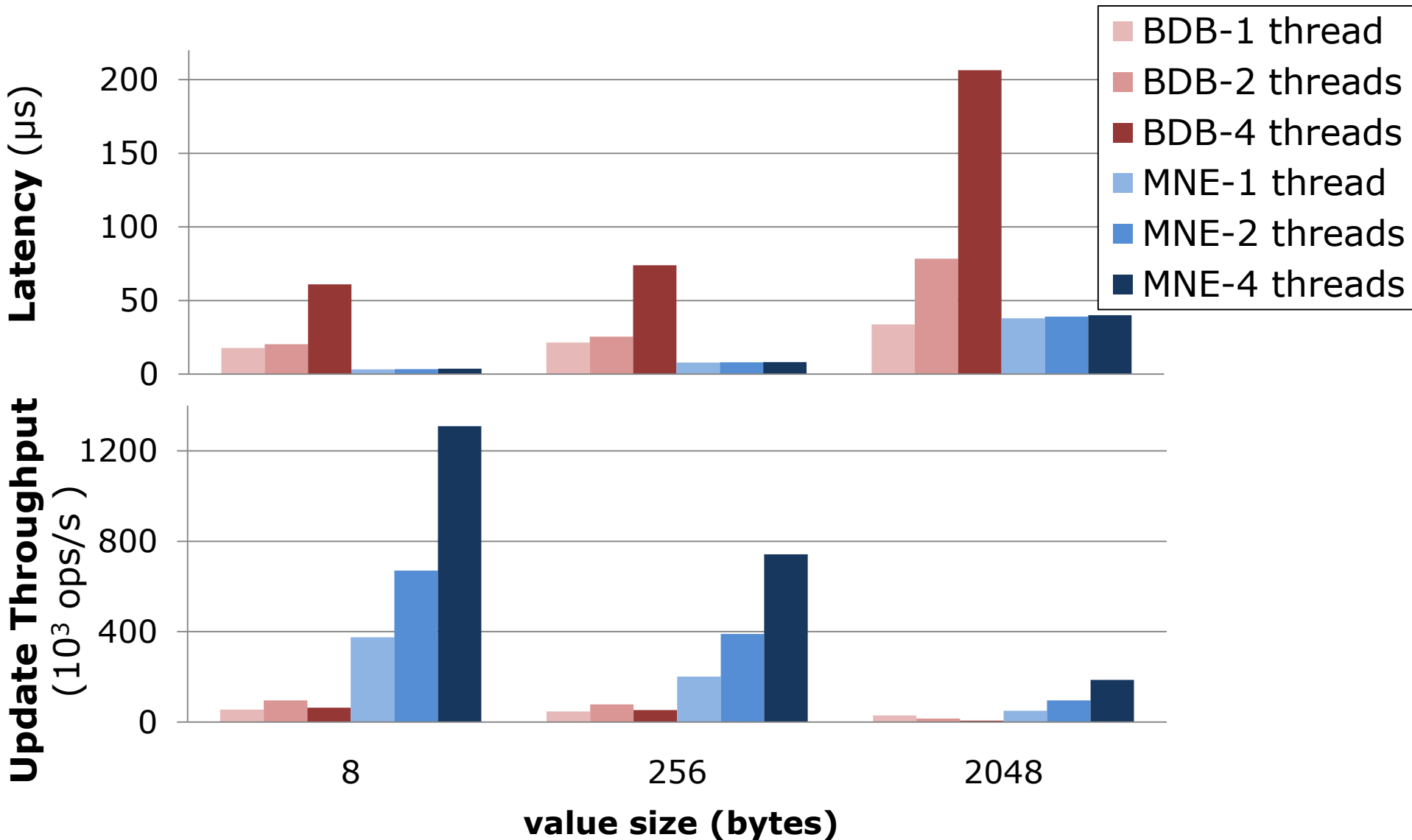
Berkeley DB (PCM-disk) vs Mnemosyne



Microbenchmark – Hash Table



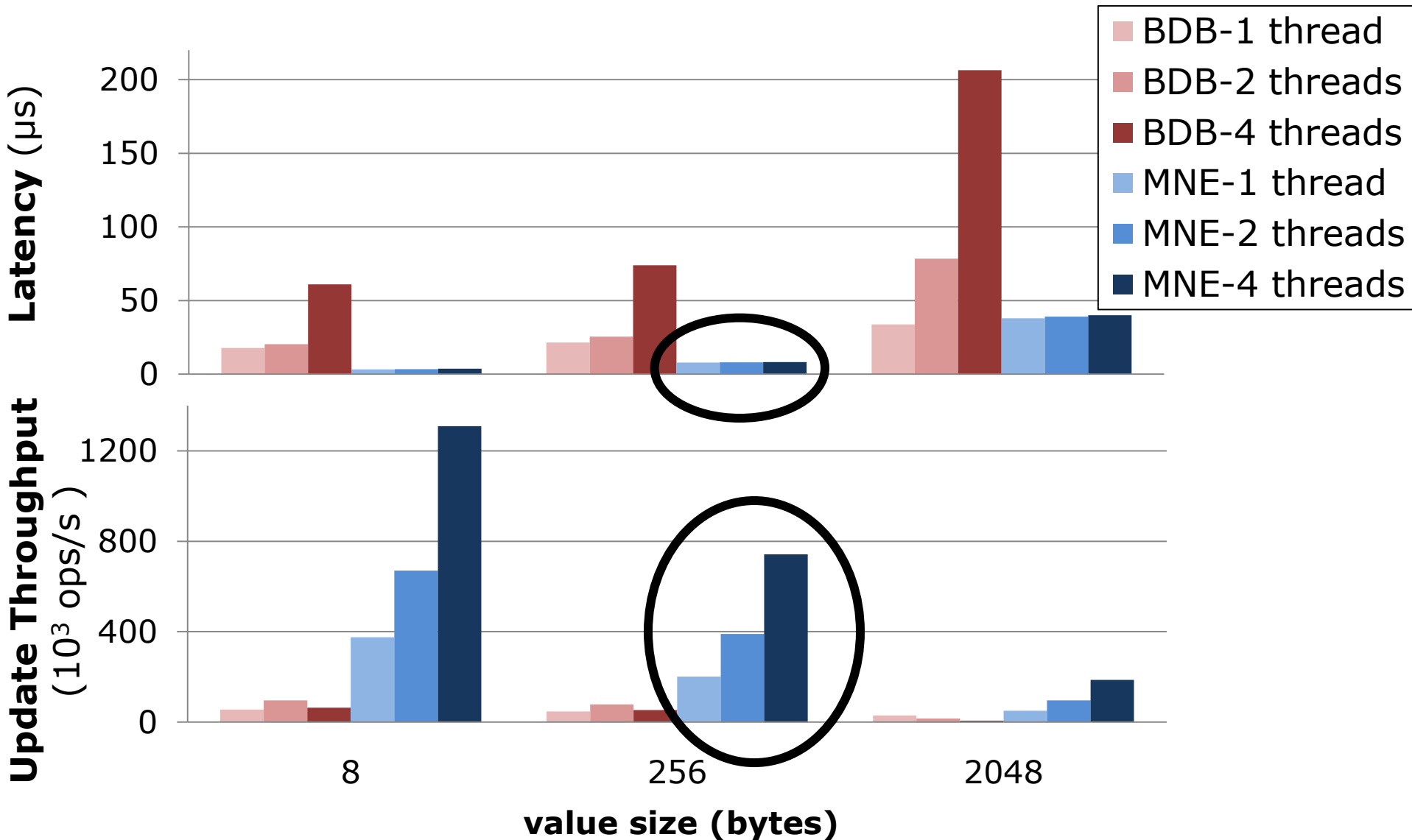
Berkeley DB (PCM-disk) vs Mnemosyne



Microbenchmark – Hash Table



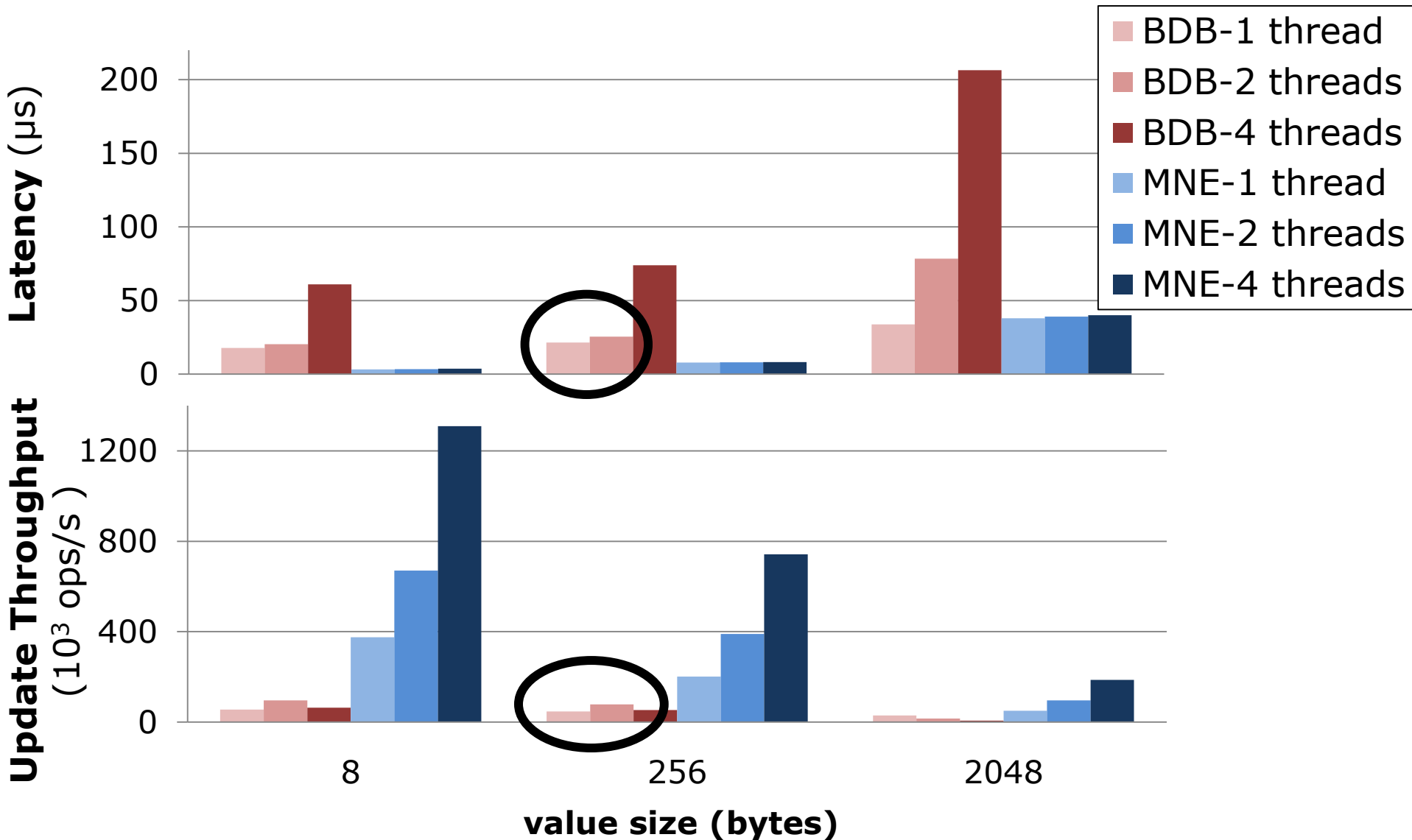
Berkeley DB (PCM-disk) vs Mnemosyne



Microbenchmark – Hash Table

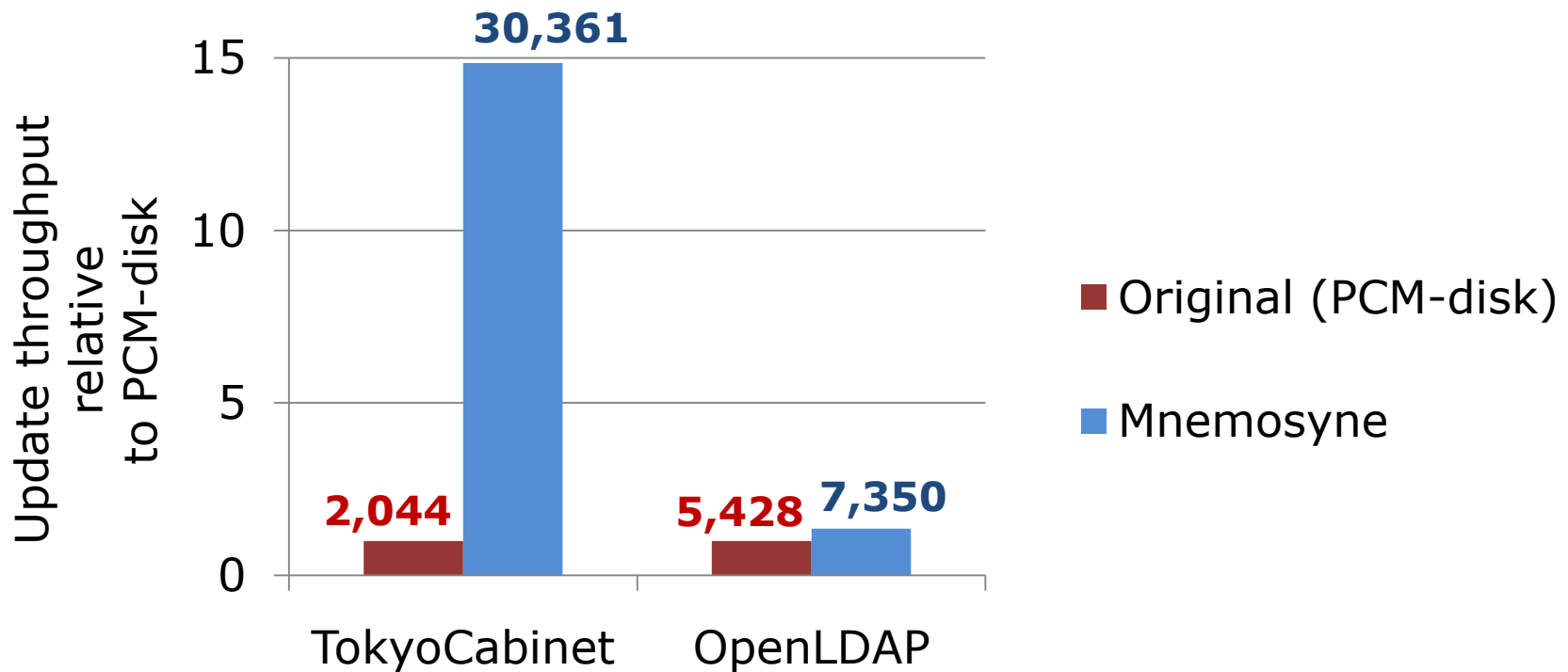


Berkeley DB (PCM-disk) vs Mnemosyne





- Write-mostly workloads
 - TokyoCabinet: 1024-byte ins/del queries
 - OpenLDAP: template-based update queries





- Exposes SCM directly to programmers as persistent memory
- Relies on conventional CPU primitives for ordering
- Enables low-latency, consistent, in-place updates via durable memory transactions