

STM in the Small: Trading Generality for Performance in Software Transactional Memory

Aleksandar Dragojević

Our motivation

- Concurrent data structures
- Hash-tables, skip-lists
 - In-memory database indices²
 - Key-value stores

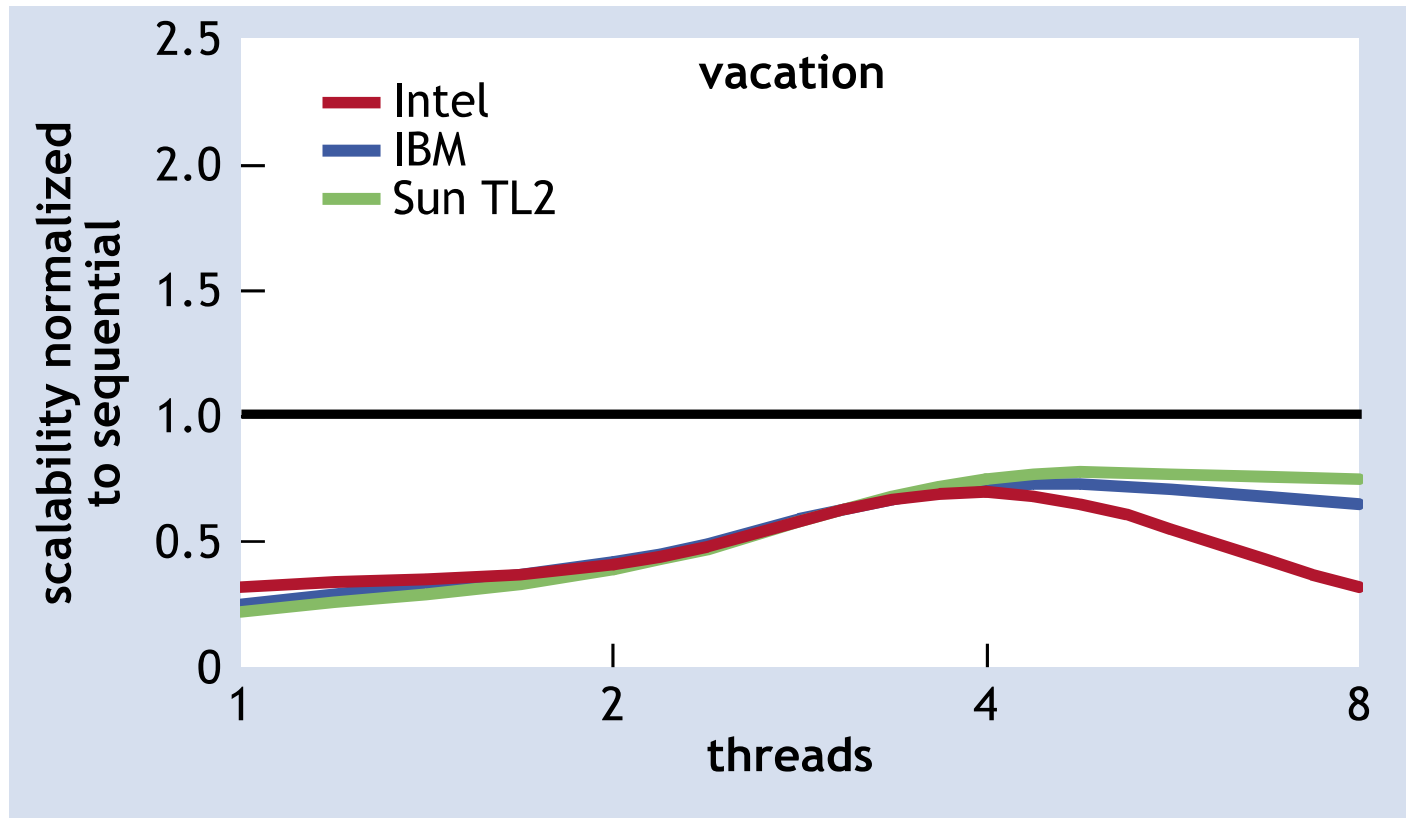
[2] *Larson et al.* “High-performance concurrency control mechanisms for main-memory databases” VLDB Dec. 2011

Software Transactional Memory

- Simplifies data structures
 - Ensures correct implementations
- Enables use of complex data structures
 - More complex than when using CAS directly

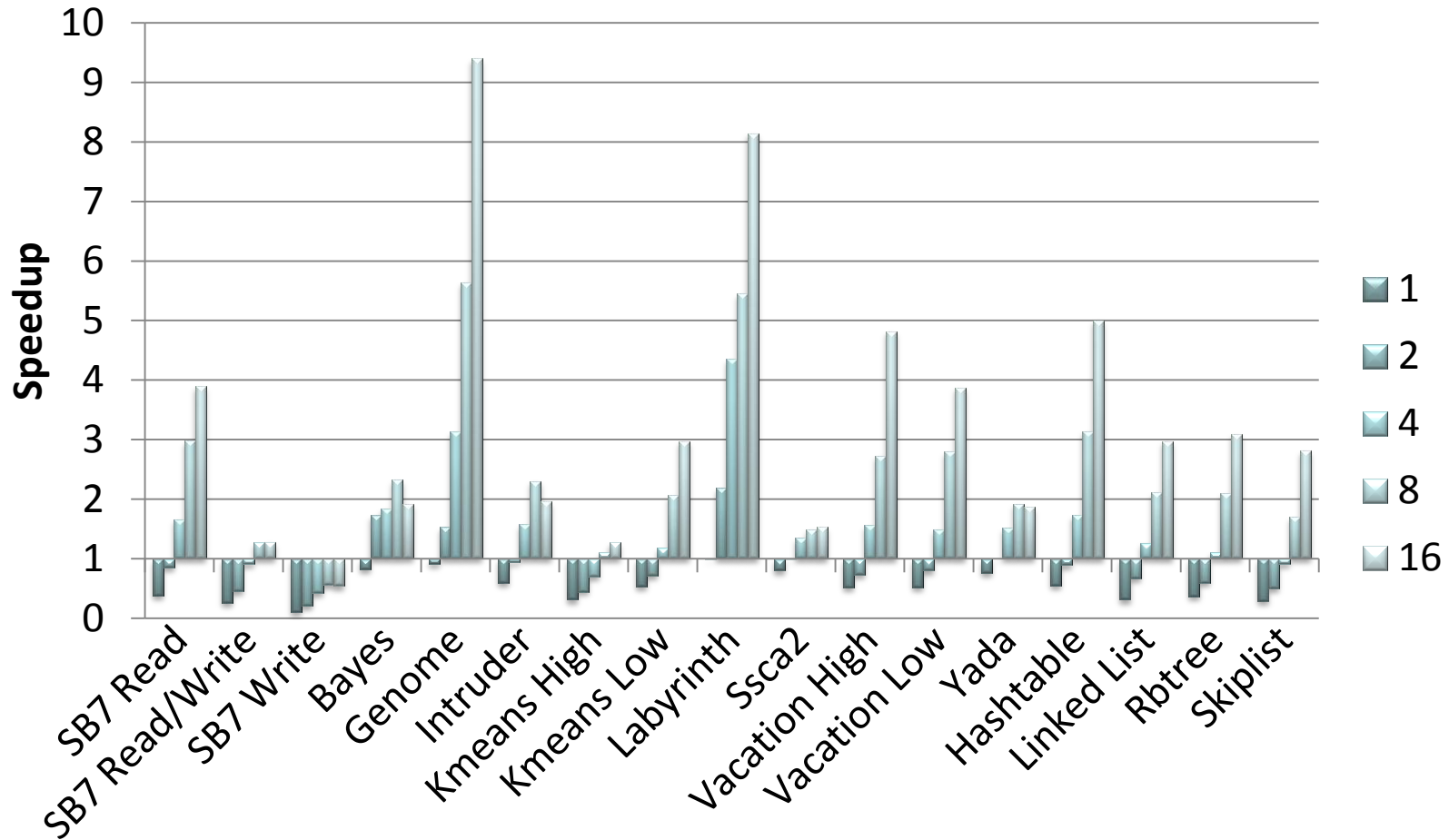
What is STM performance like?

What is STM performance like?

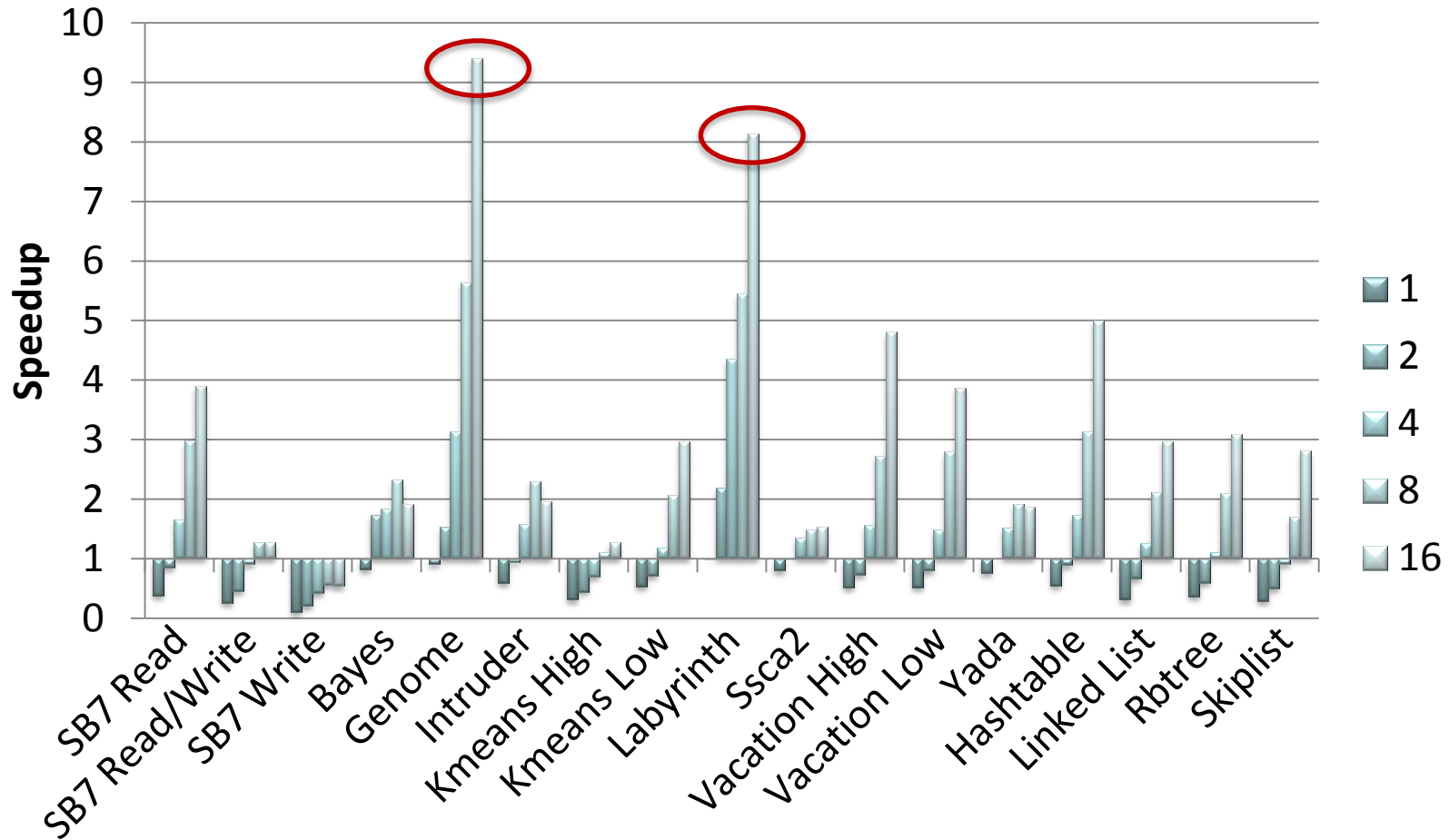


Graph from: Cascaval et al. "Software transactional memory: why is it only a research toy" ACMQ September 2008

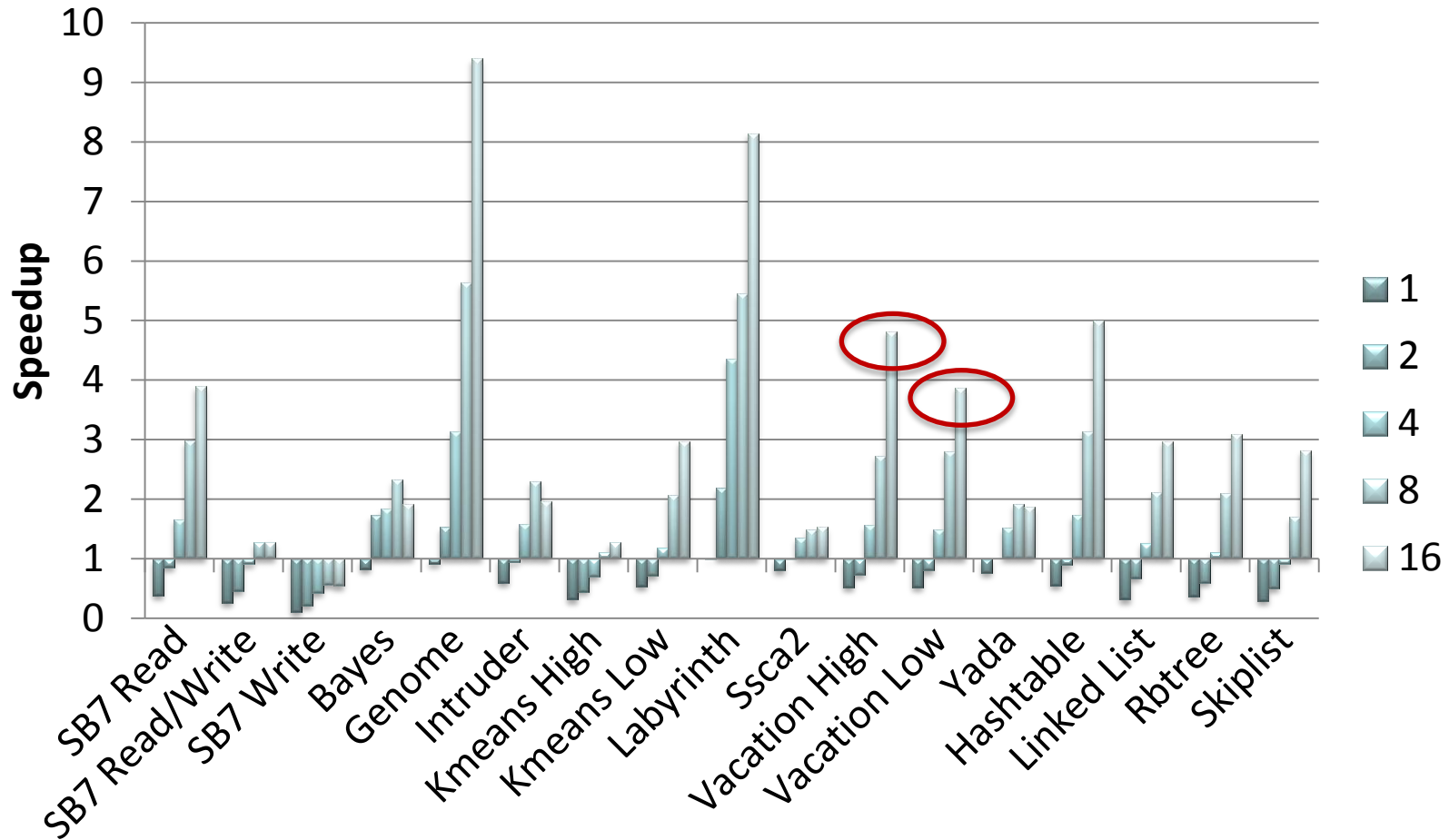
What is STM performance like?



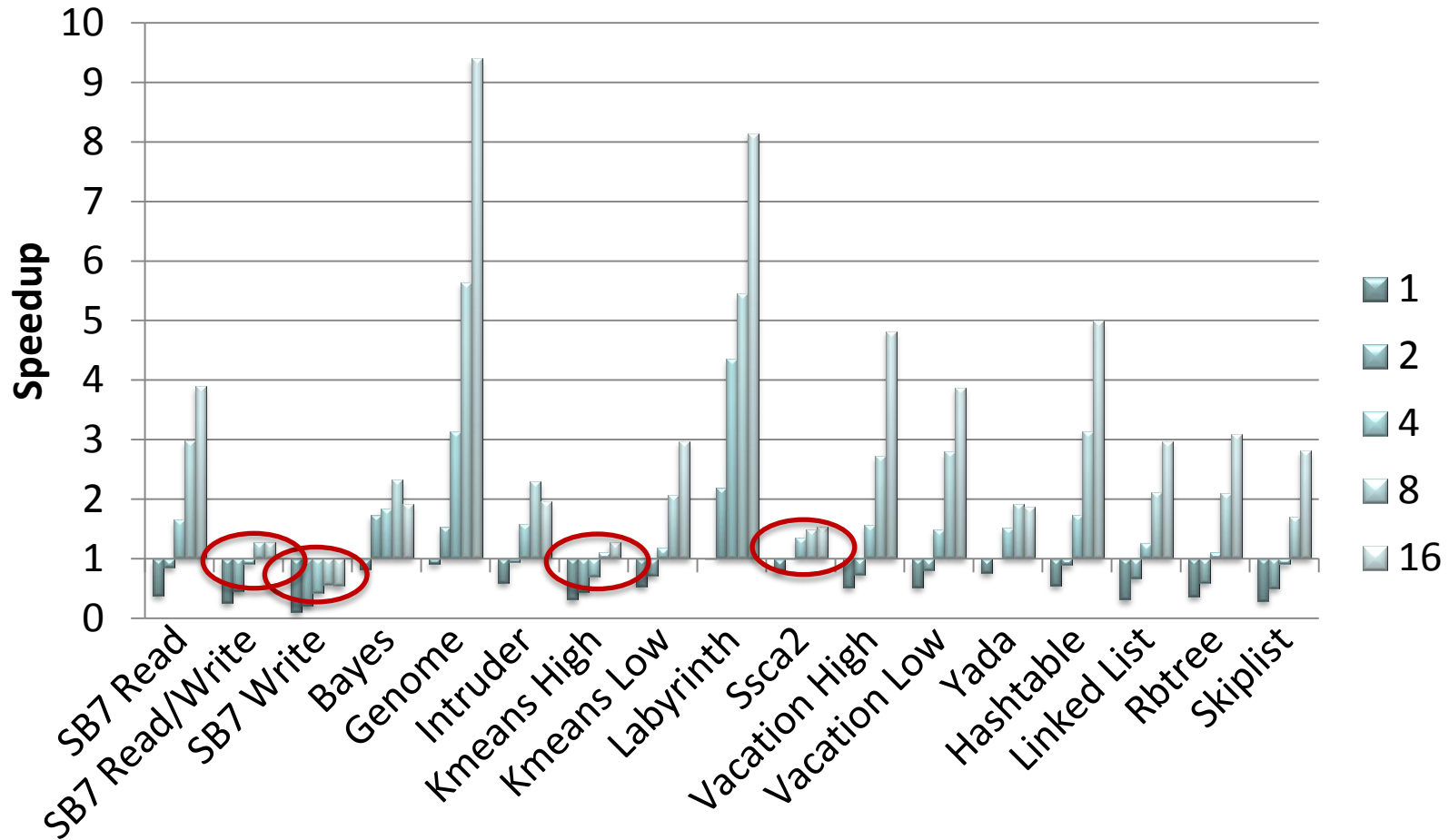
What is STM performance like?



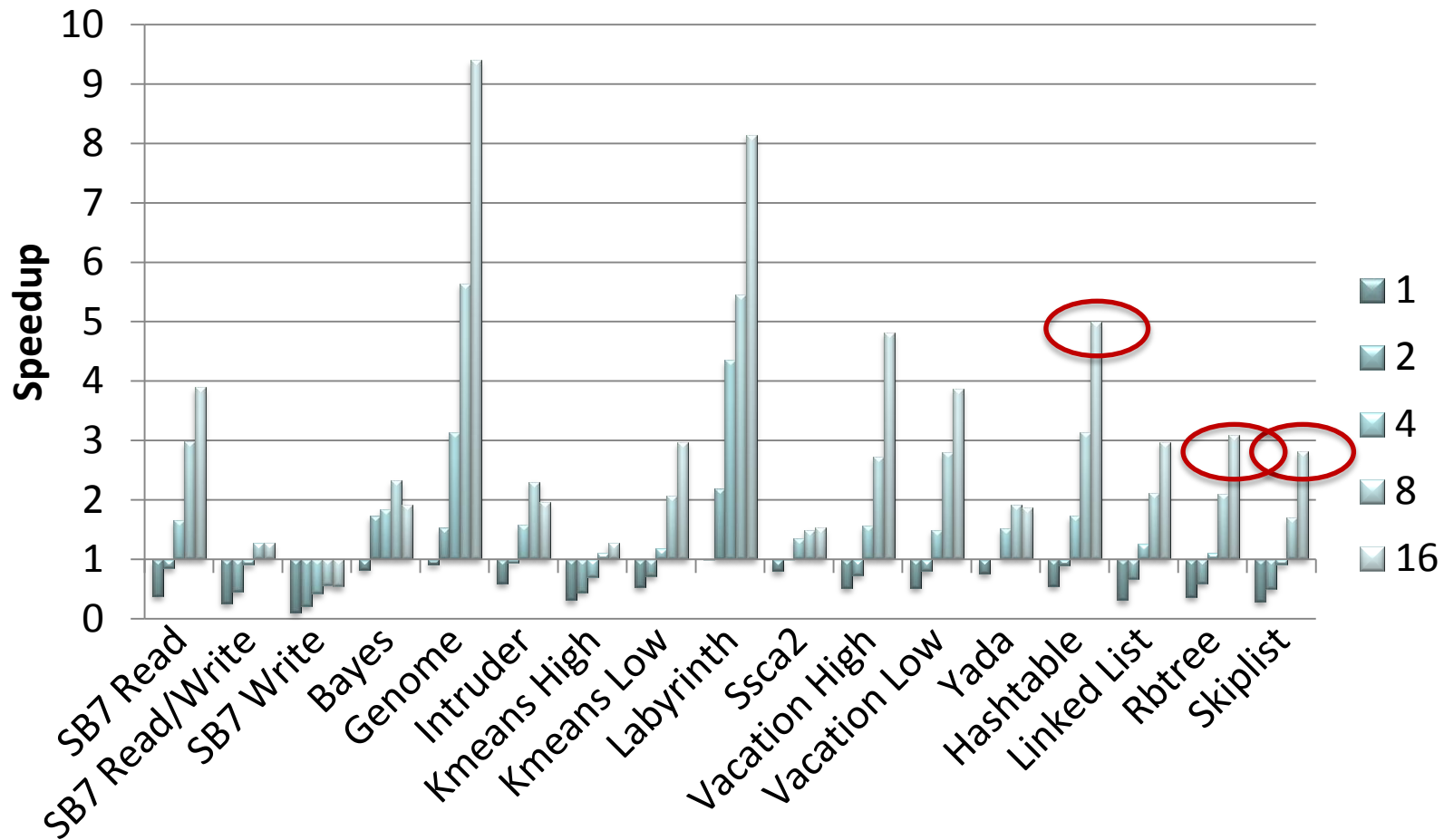
What is STM performance like?



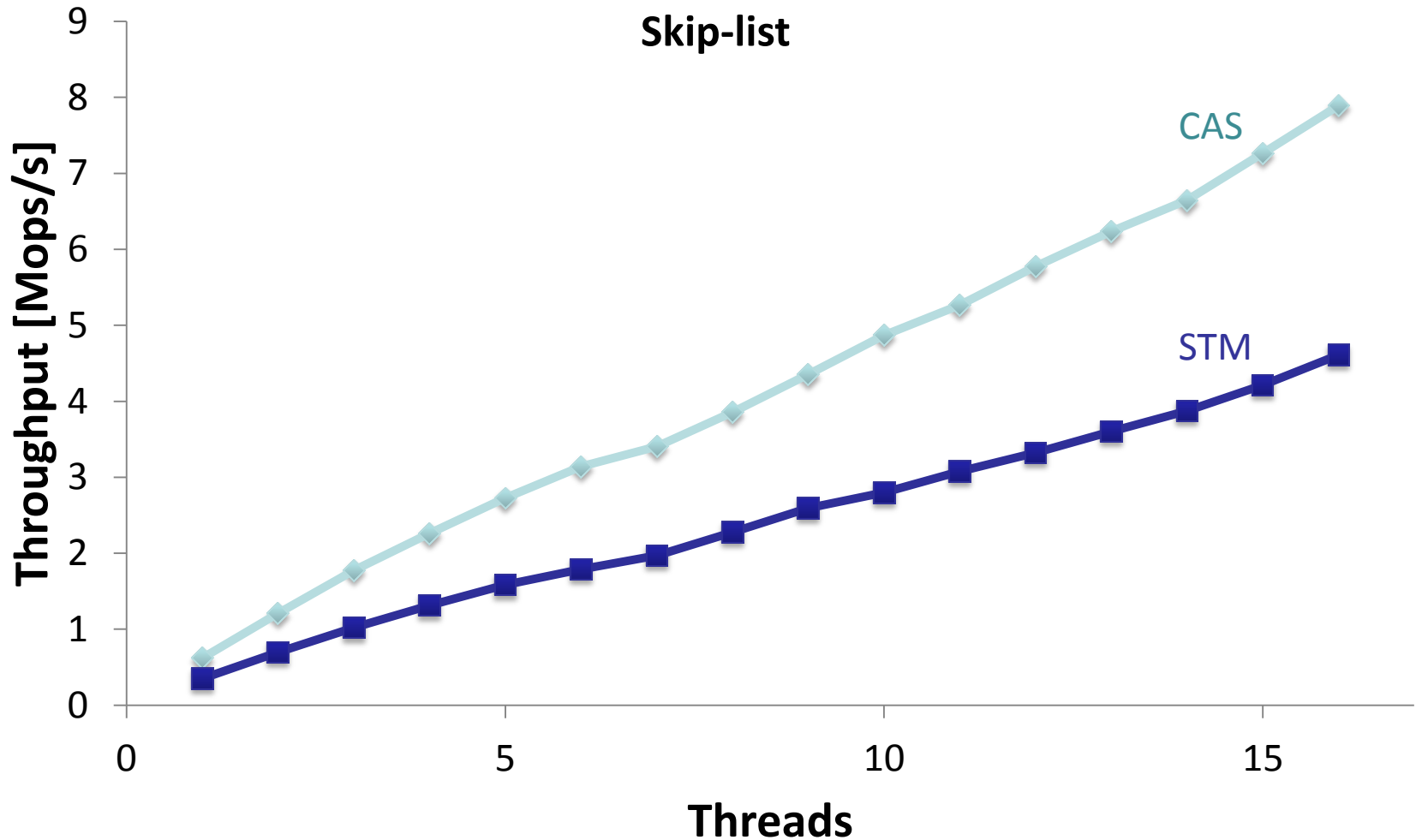
What is STM performance like?



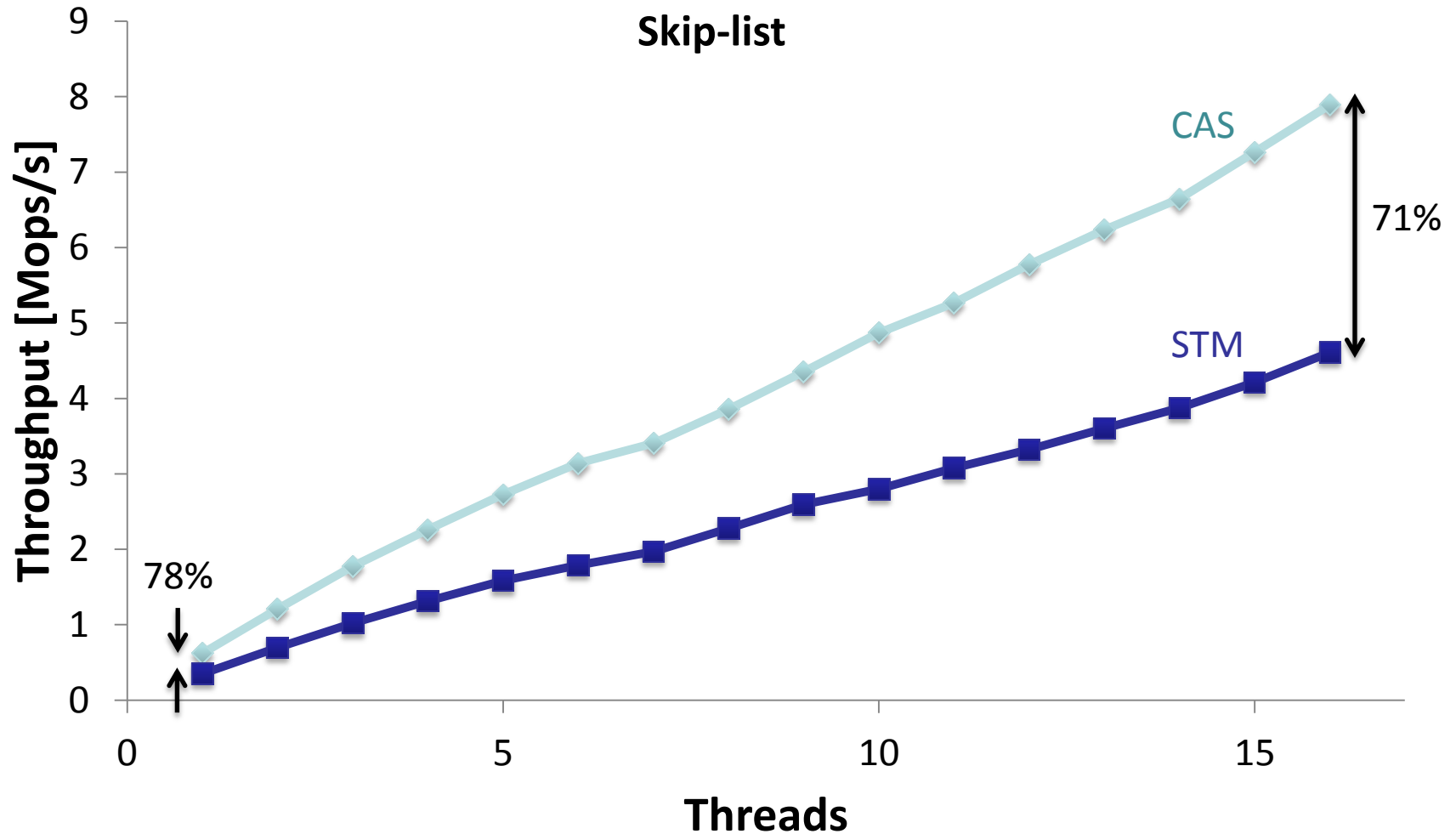
What is STM performance like?



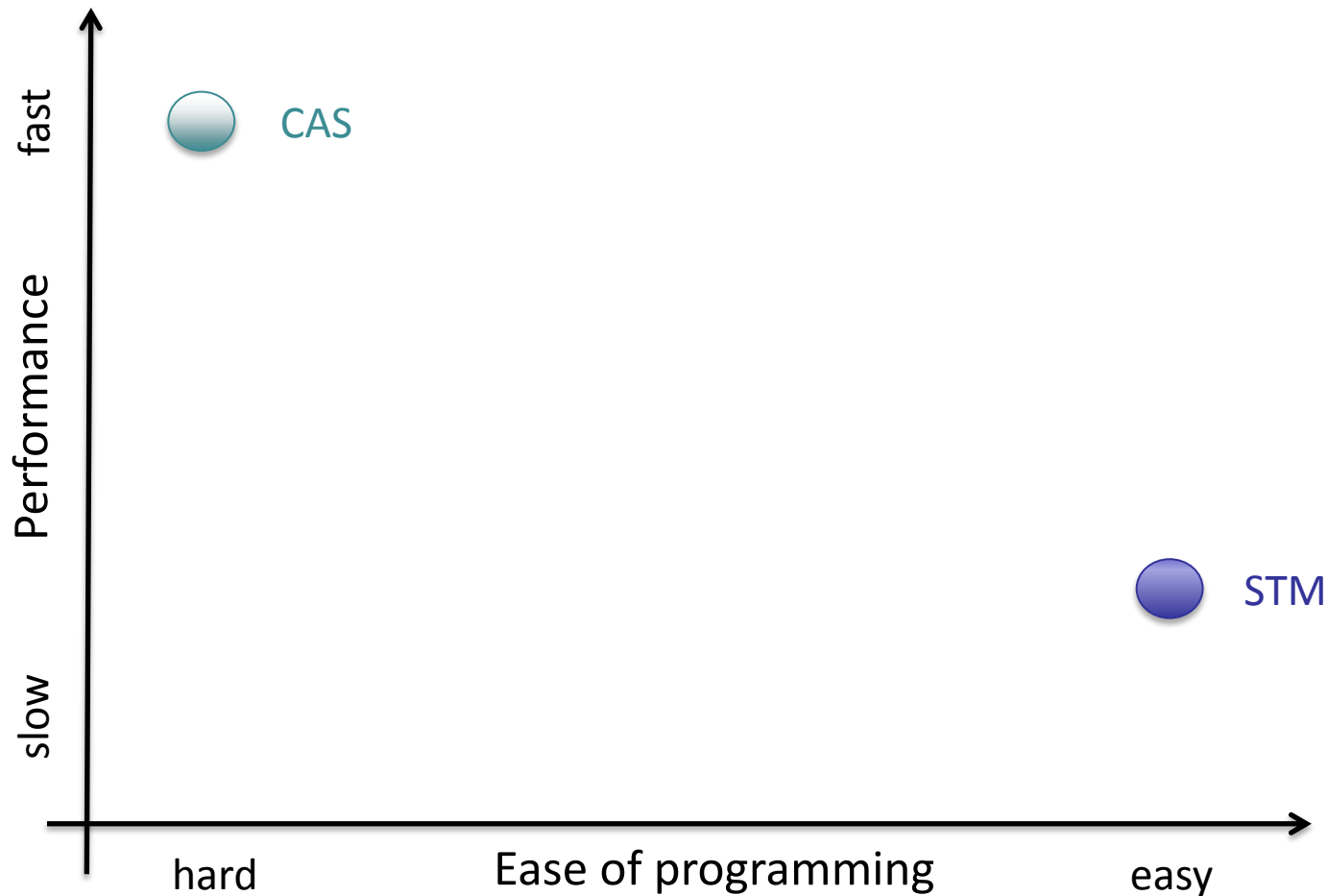
What is STM performance like?



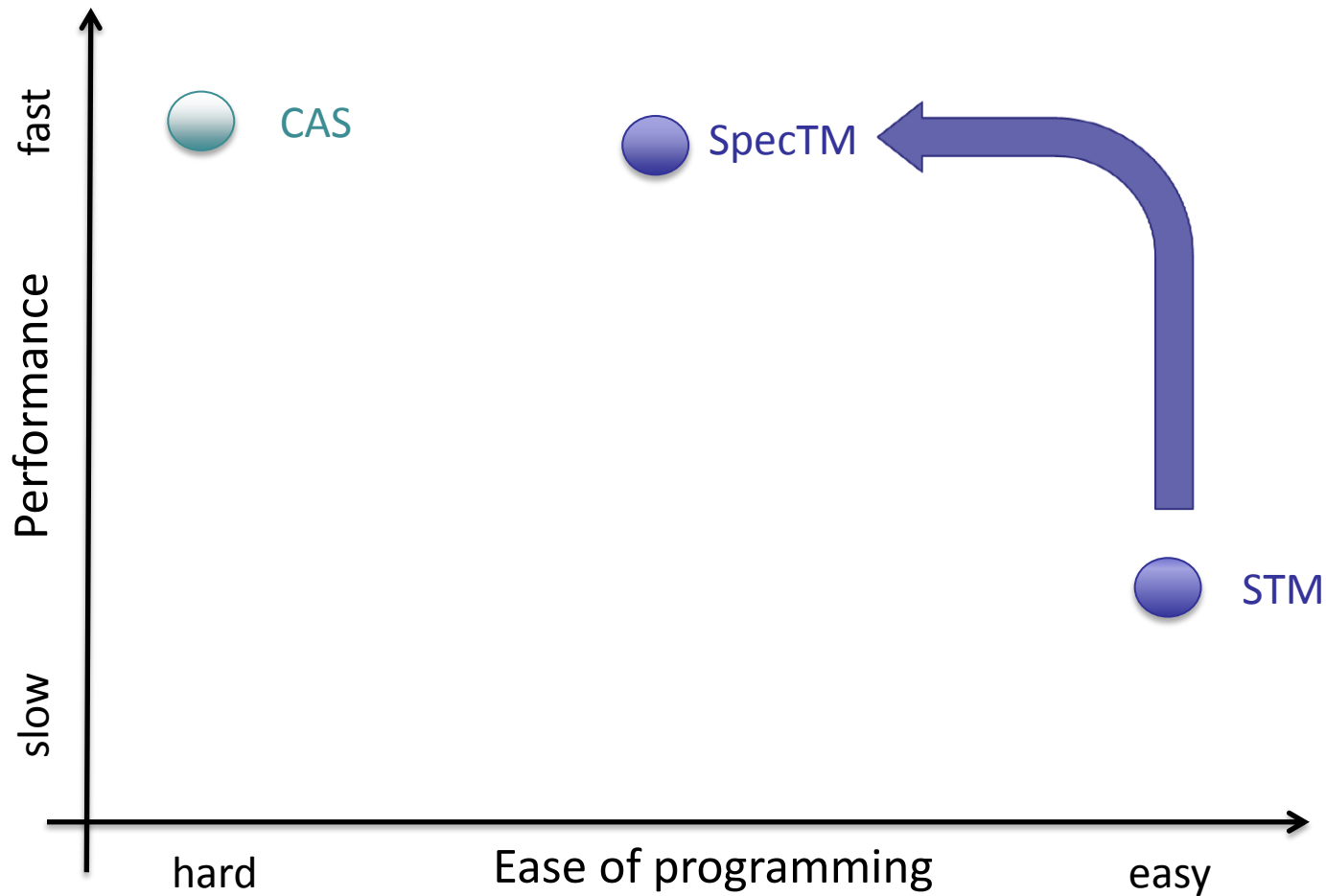
What is STM performance like?



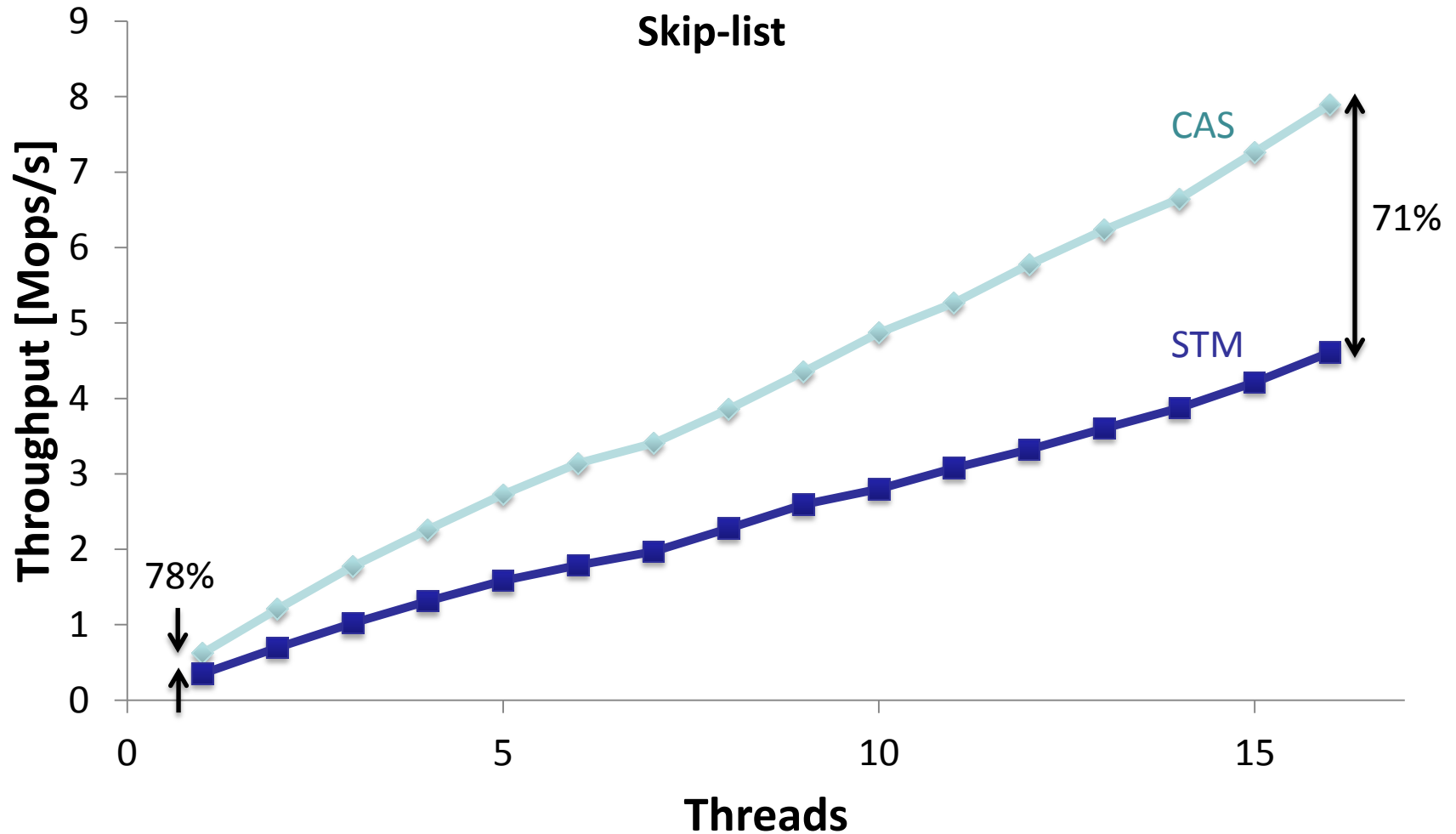
SpecTM: Specialized STM for Concurrent Data Structures



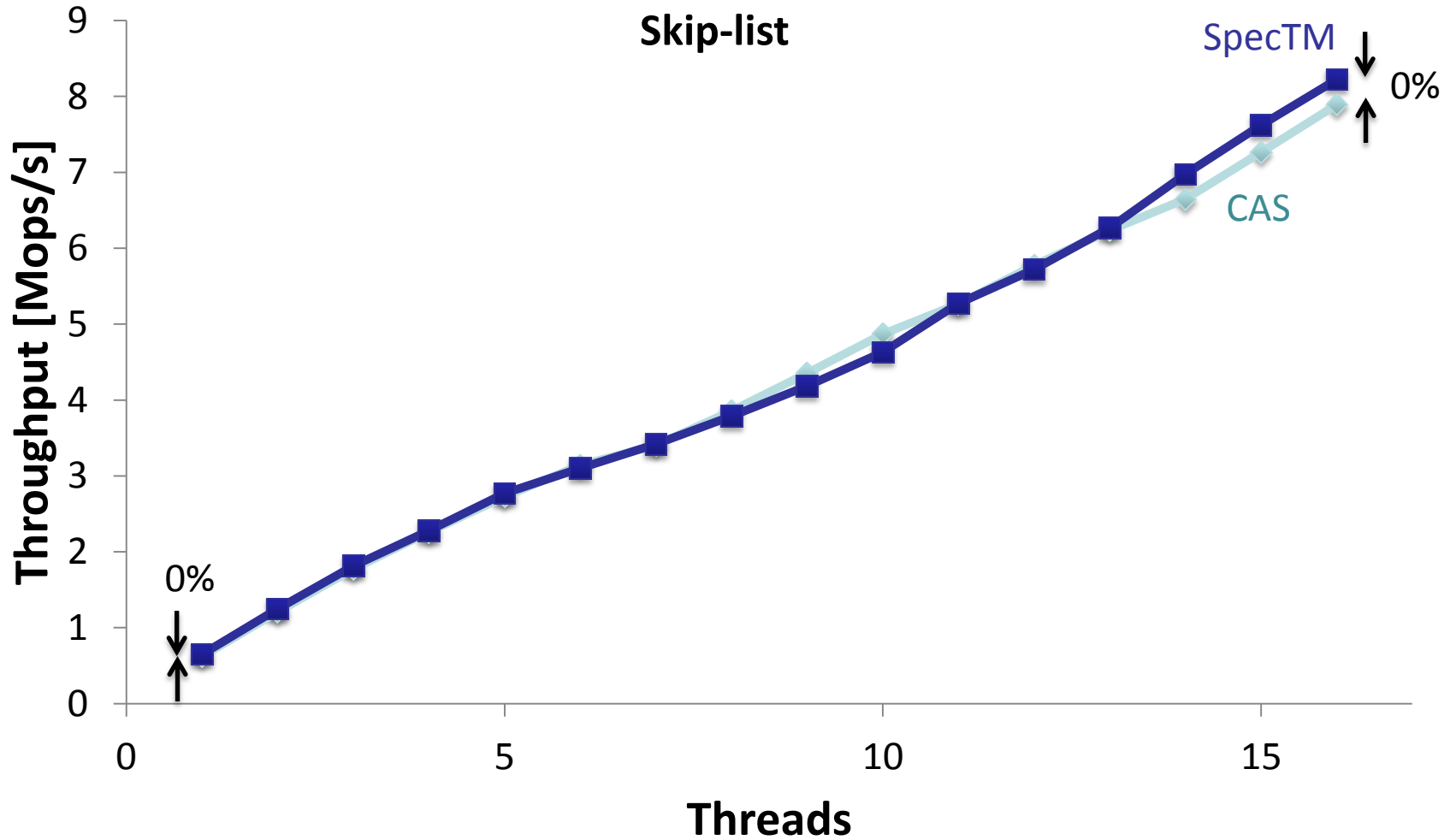
SpecTM: Specialized STM for Concurrent Data Structures



This talk



This talk



Overview

- STM overheads
- SpecTM
 - Short-transaction API
 - Collocating data and meta-data
 - In-word meta-data
- Evaluation

Overview

- STM overheads
- SpecTM
 - Short-transaction API
 - Collocating data and meta-data
 - In-word meta-data
- Evaluation

STM overheads

- Book-keeping done in software
- Memory accesses become STM calls
- Significant overheads³
 - With a single thread 50% overheads
 - Requires 4 threads to outperform sequential code in 75% of cases

[3] *Dragojević et al.* “Why STM Can Be more than a Research Toy” CACM Apr. 2011

STM API

```
void StartTx();
```

```
void CommitTx();
```

```
word_t ReadWord(word_t *addr);
```

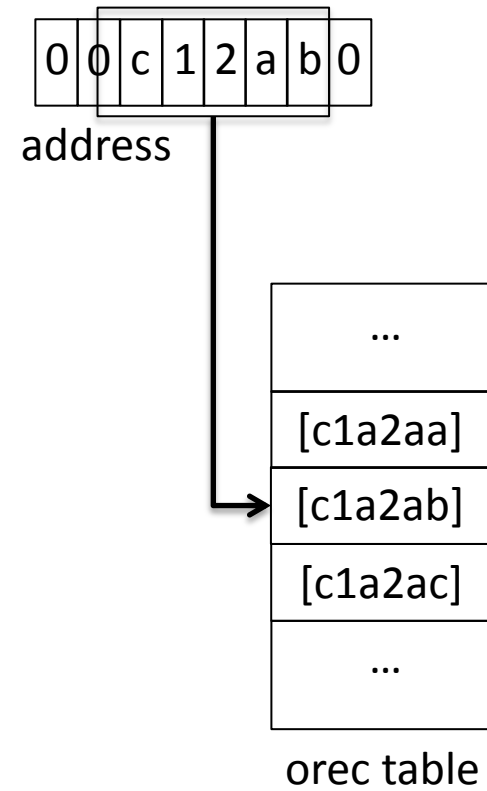
```
void WriteWord(word_t *addr, word_t val);
```

STM read

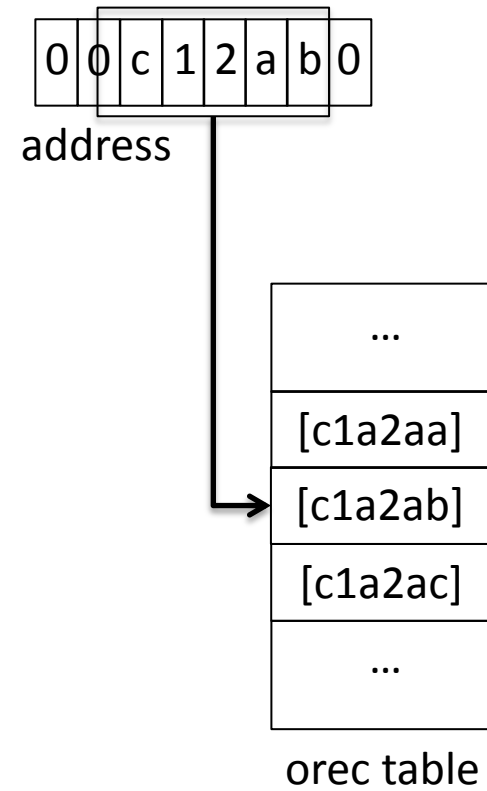
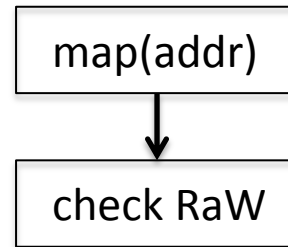
```
map(addr)
```

STM read

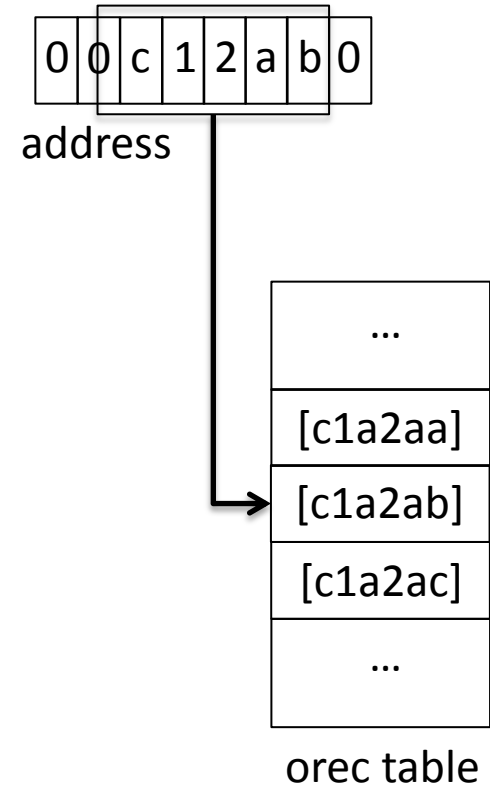
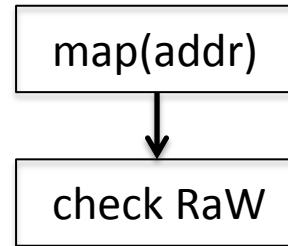
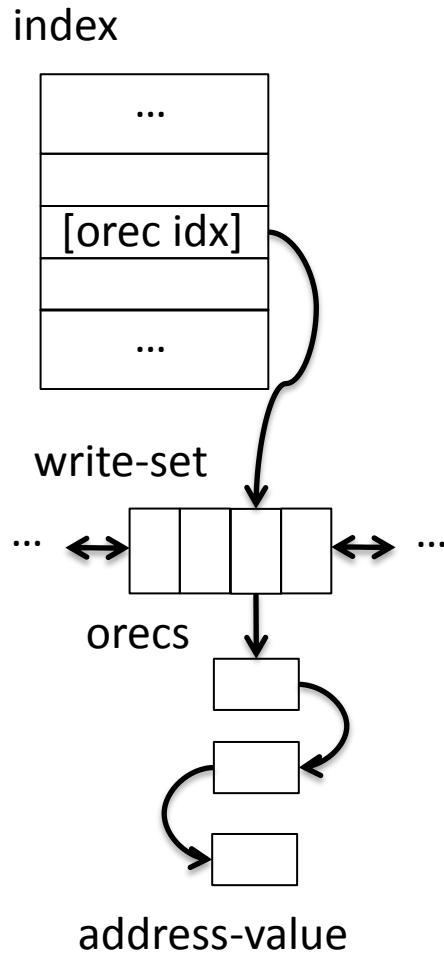
map(addr)



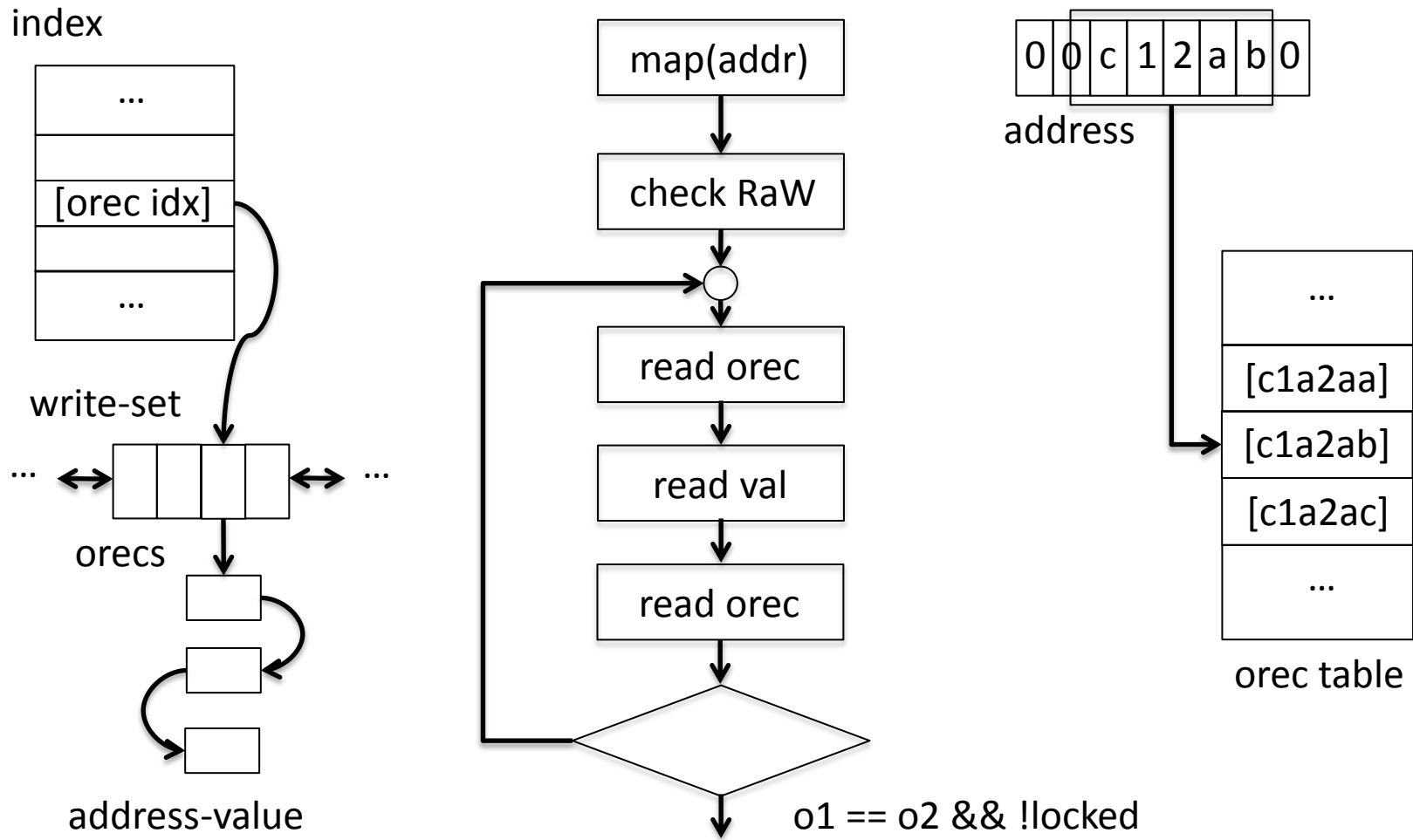
STM read



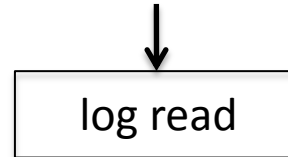
STM read



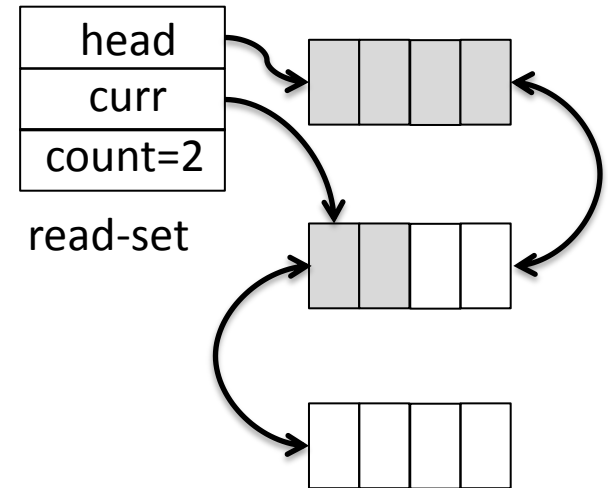
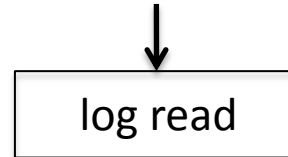
STM read



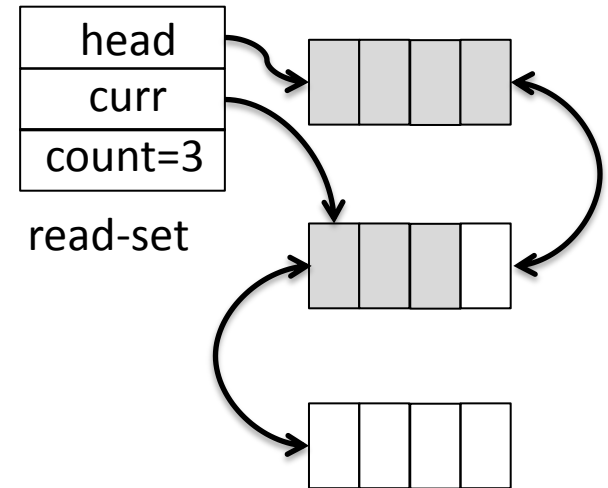
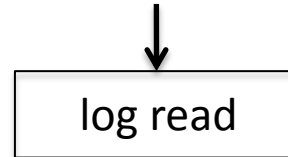
STM read (cont'd)



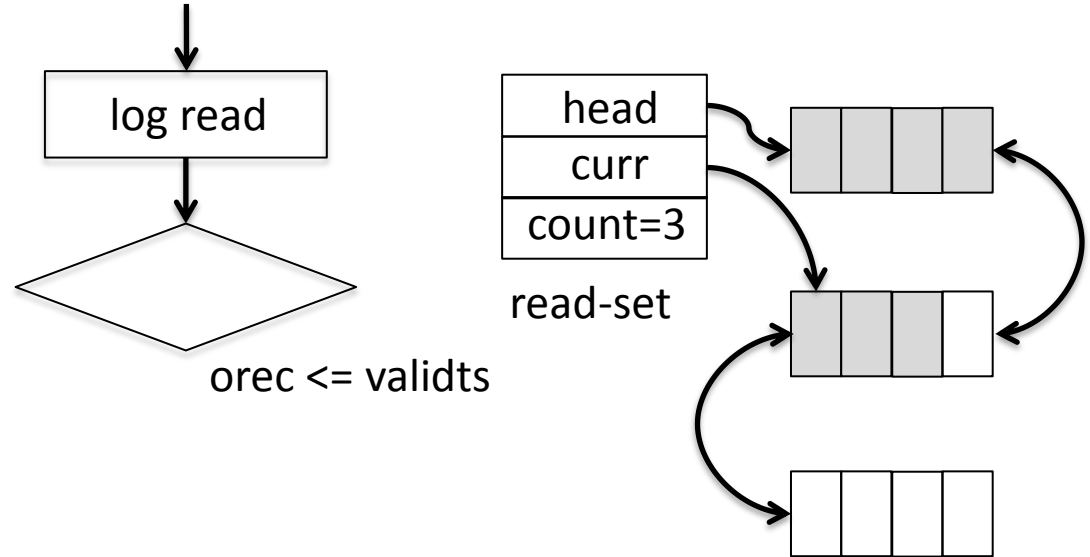
STM read (cont'd)



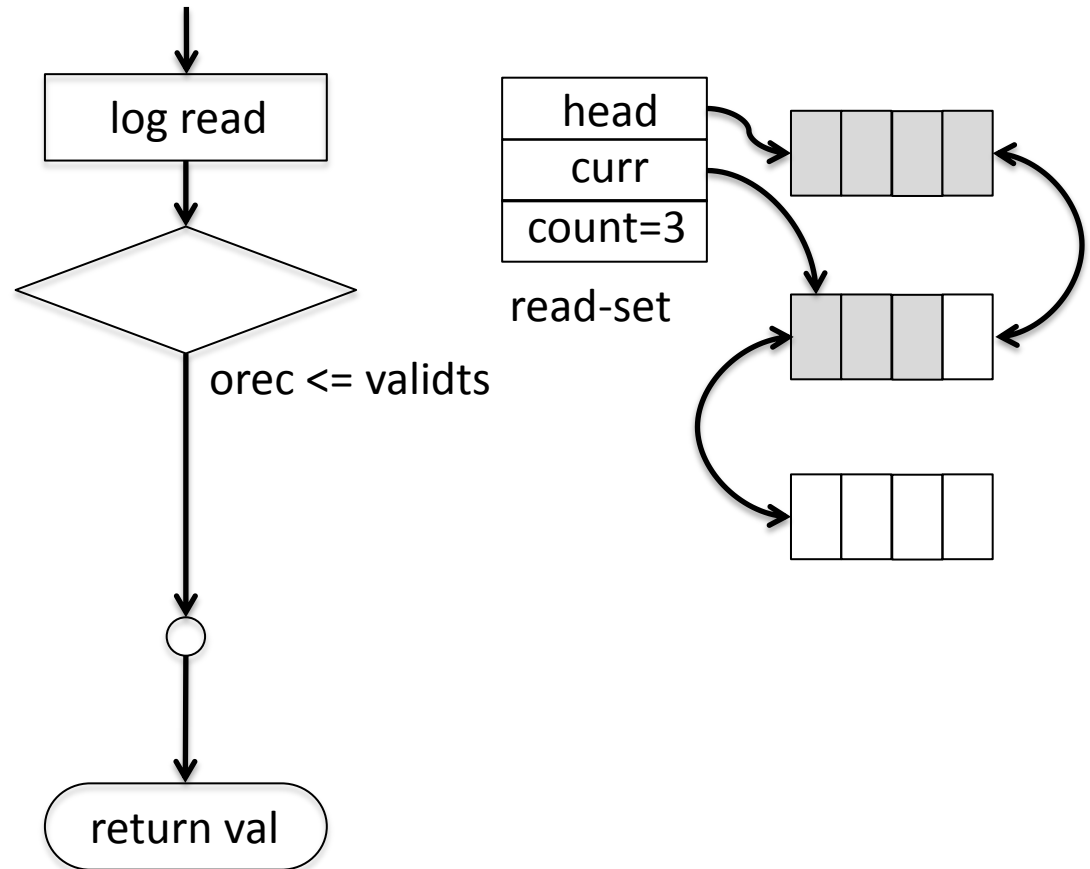
STM read (cont'd)



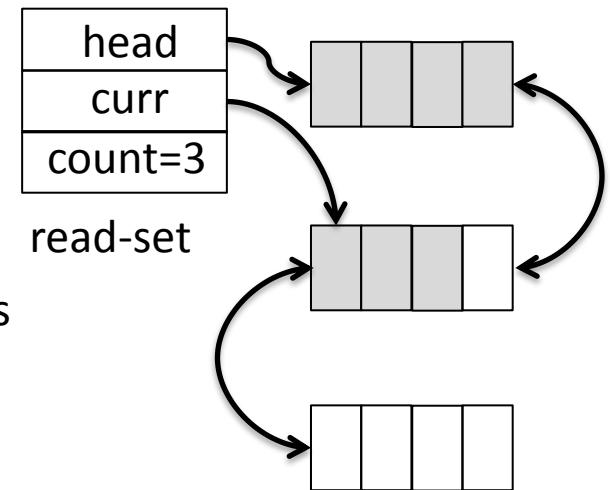
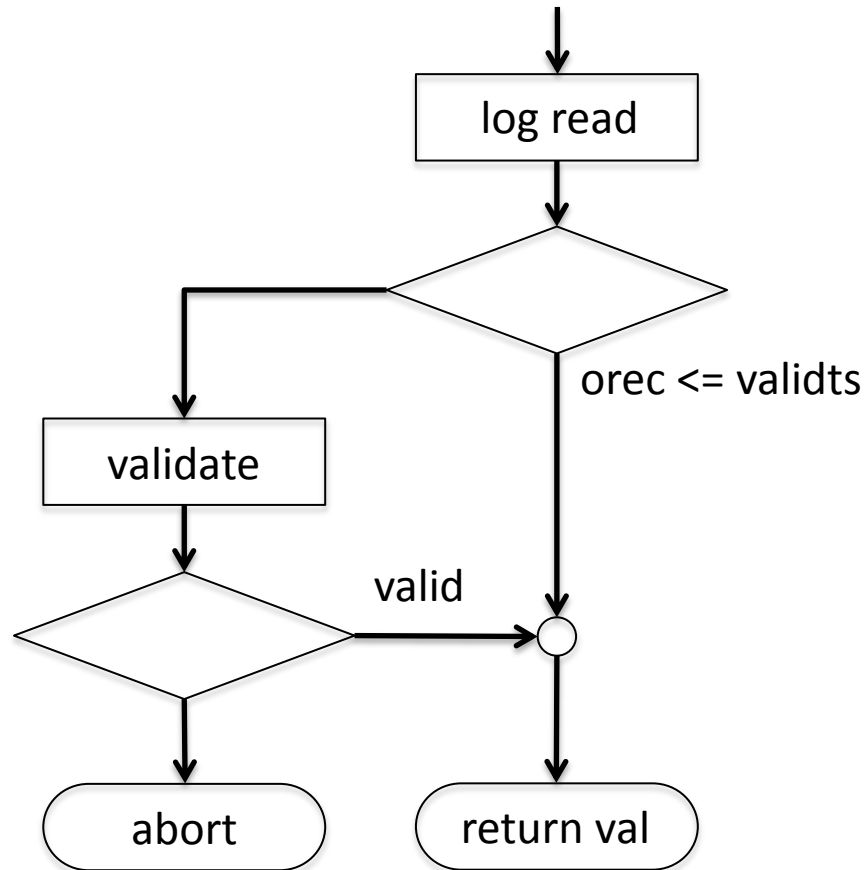
STM read (cont'd)



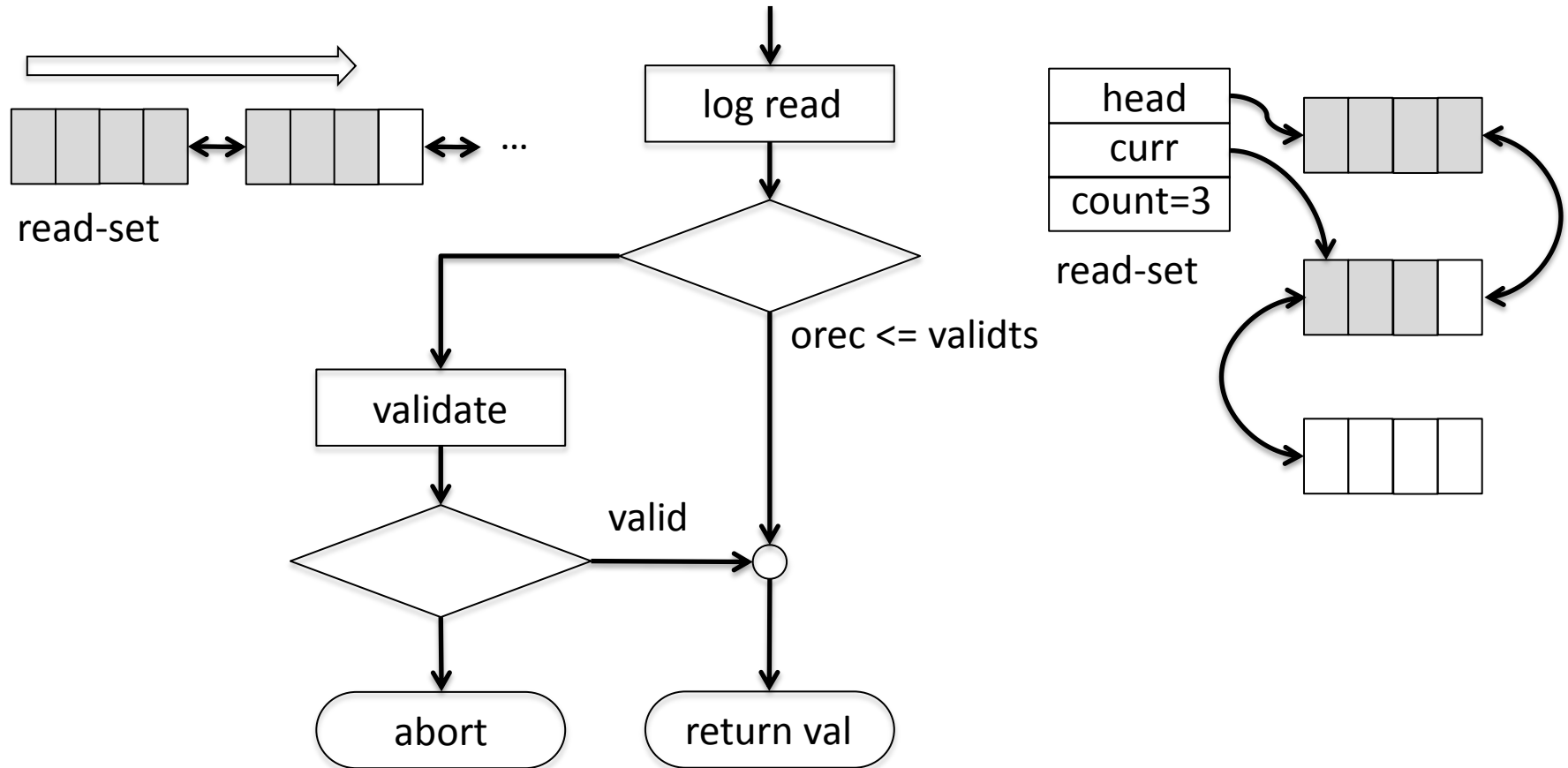
STM read (cont'd)



STM read (cont'd)



STM read (cont'd)



STM read fast-path

- Significantly more expensive than CPU read instructions
 - >10 instructions
- Writes have comparable costs
 - Incurs a compare-and-swap at commit time

Overview

- STM overheads
- SpecTM
 - Short-transaction API
 - Collocating data and meta-data
 - In-word meta-data
- Evaluation

Overview

- STM overheads
- SpecTM
 - Short-transaction API
 - Collocating data and meta-data
 - In-word meta-data
- Evaluation

Short transactions

- Short read-write transactions
- Short read-only transactions
- Single-access transactions
 - Read, Write, CAS

Short read-write transactions

```
word_t TxRWRead_1 (word_t *addr_1);  
word_t TxRWRead_2 (word_t *addr_2);  
word_t TxRWRead_3 (word_t *addr_3);  
...
```

```
void TxRWCommit_1 (word_t val);  
void TxRWCommit_2 (word_t v1, word_t v2);  
void TxRWCommit_3 (word_t v1,  
                   word_t v2,  
                   word_t v3);  
...
```

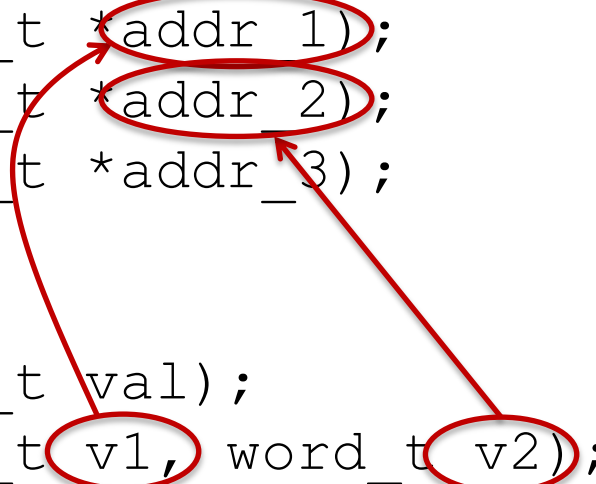
Short read-write transactions

```
word_t TxRWRead_1 (word_t *addr_1);  
word_t TxRWRead_2 (word_t *addr_2);  
word_t TxRWRead_3 (word_t *addr_3);  
...
```

```
void TxRWCommit_1 (word_t val);  
void TxRWCommit_2 (word_t v1, word_t v2);  
void TxRWCommit_3 (word_t v1,  
                   word_t v2,  
                   word_t v3);  
...
```

Short read-write transactions

```
word_t TxRWRead_1 (word_t *addr_1);  
word_t TxRWRead_2 (word_t *addr_2);  
word_t TxRWRead_3 (word_t *addr_3);  
...  
  
void TxRWCommit_1 (word_t val);  
void TxRWCommit_2 (word_t v1, word_t v2);  
void TxRWCommit_3 (word_t v1,  
                    word_t v2,  
                    word_t v3);  
  
...
```

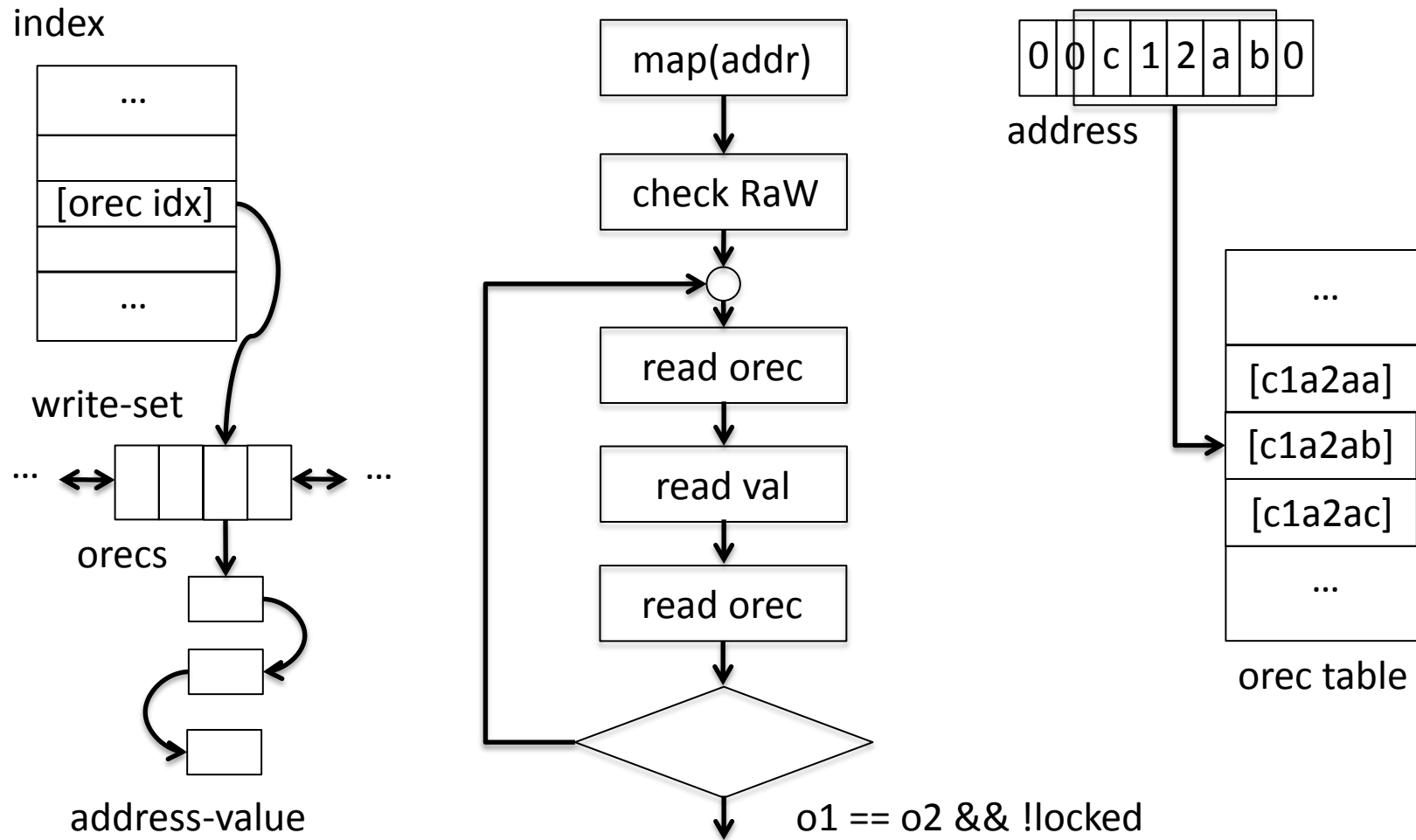


Short read-write transactions

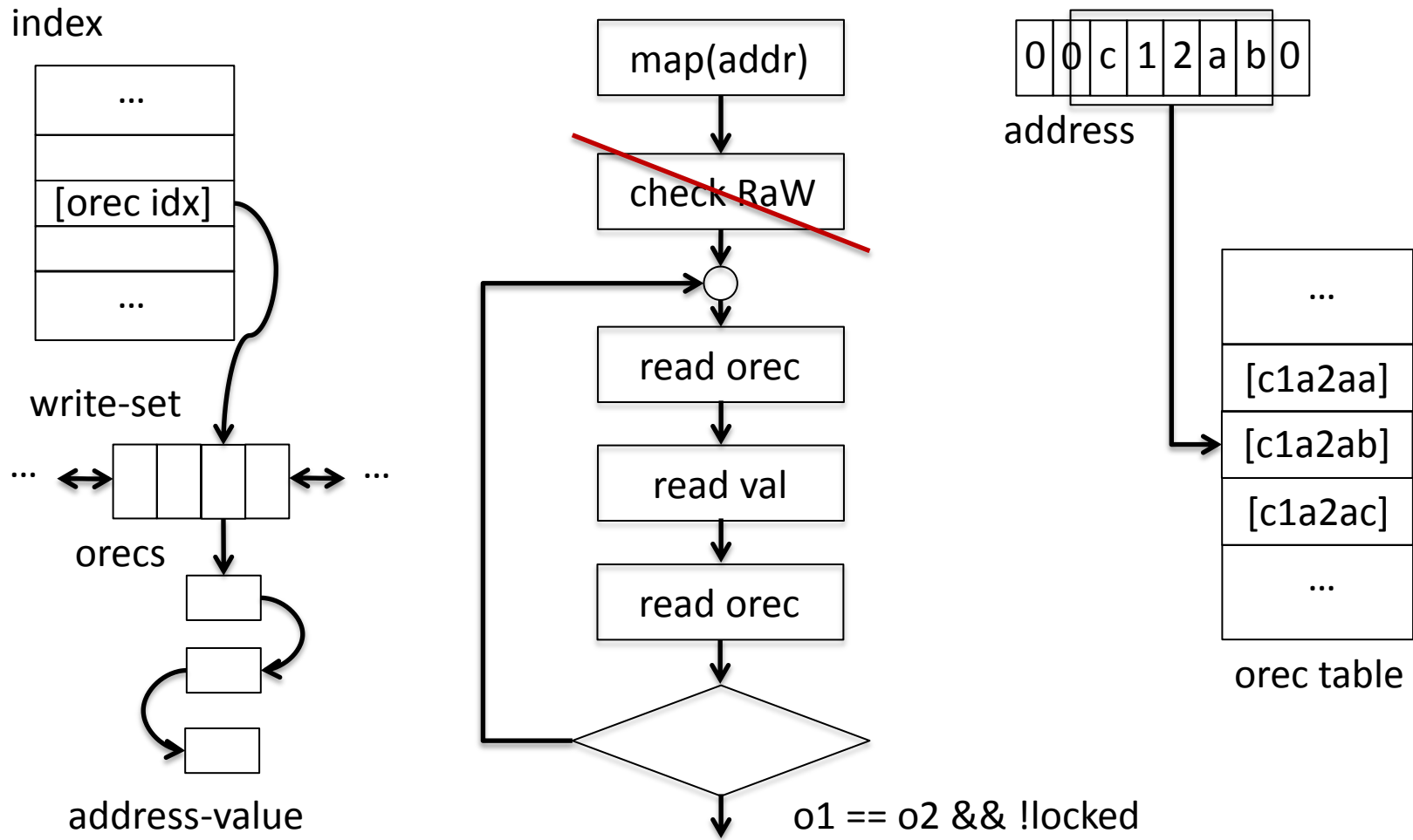
```
word_t TxRWRead_1 (word_t *addr 1);  
word_t TxRWRead_2 (word_t *addr 2);  
word_t TxRWRead_3 (word_t *addr 3);  
...  
  
void TxRWCommit_1 (word_t val);  
void TxRWCommit_2 (word_t v1, word_t v2);  
void TxRWCommit_3 (word_t v1,  
                    word_t v2,  
                    word_t v3);  
...
```

The diagram illustrates the flow of data in short read-write transactions. It shows three read functions (TxRWRead_1, TxRWRead_2, TxRWRead_3) and three commit functions (TxRWCommit_1, TxRWCommit_2, TxRWCommit_3). The pointer variables (*addr 1), (*addr 2), and (*addr 3) in the read functions are circled in red. The value variables (v1, v2, v3) in the commit functions are also circled in red. Red arrows point from the circled value variables in the commit functions back to the circled pointer variables in the read functions, indicating that the commit functions are writing back to the memory locations specified by the read functions.

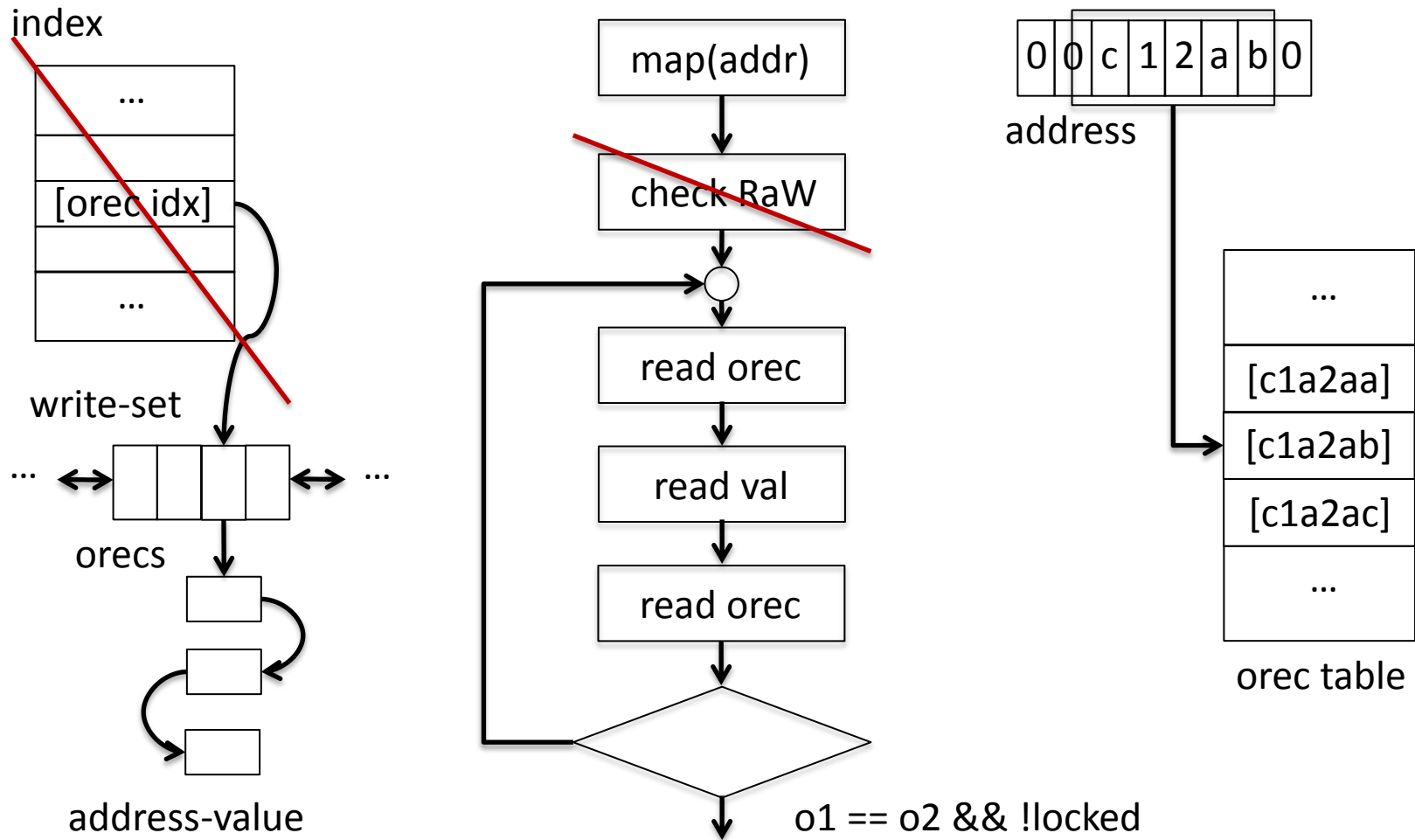
Short transaction optimizations



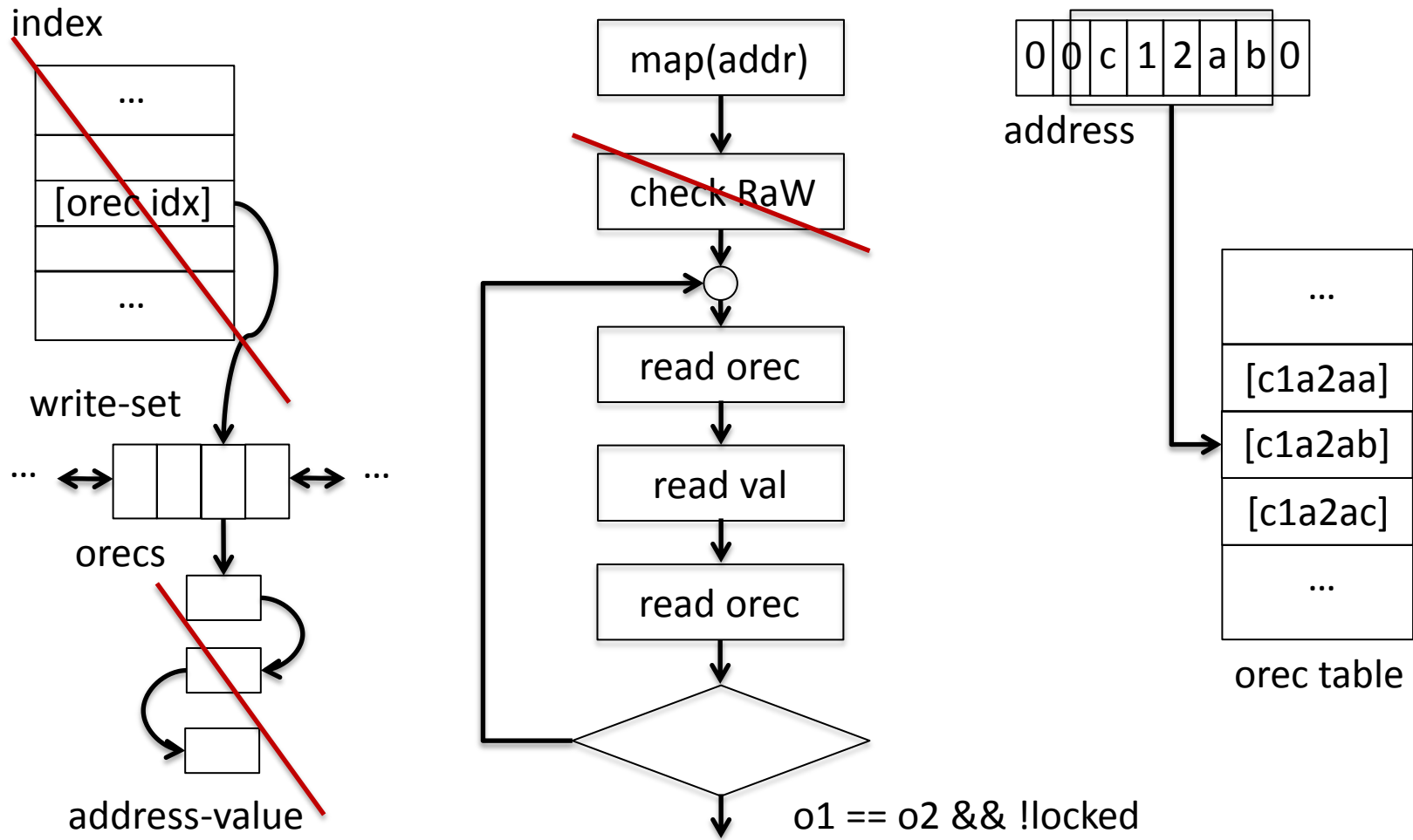
Short transaction optimizations



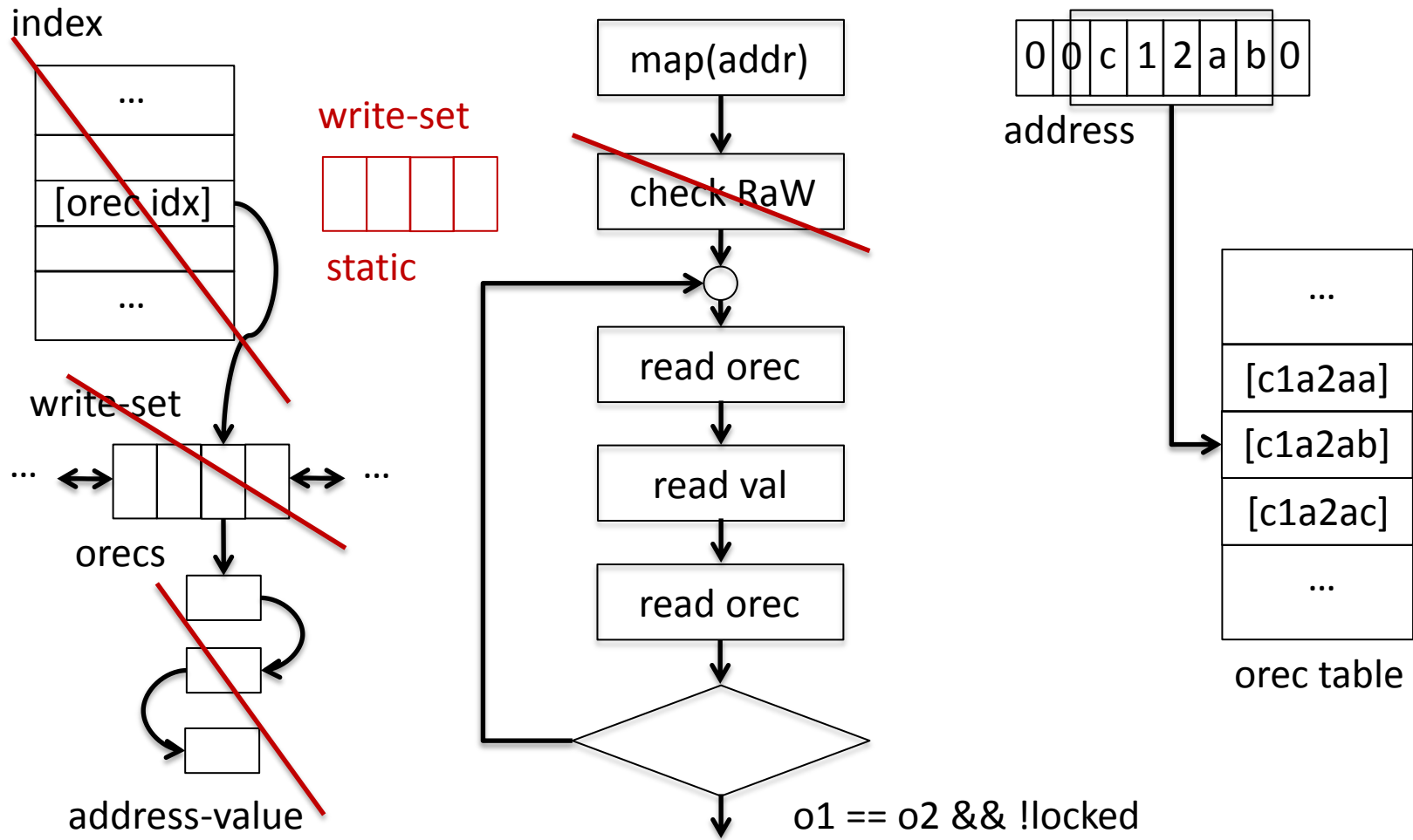
Short transaction optimizations



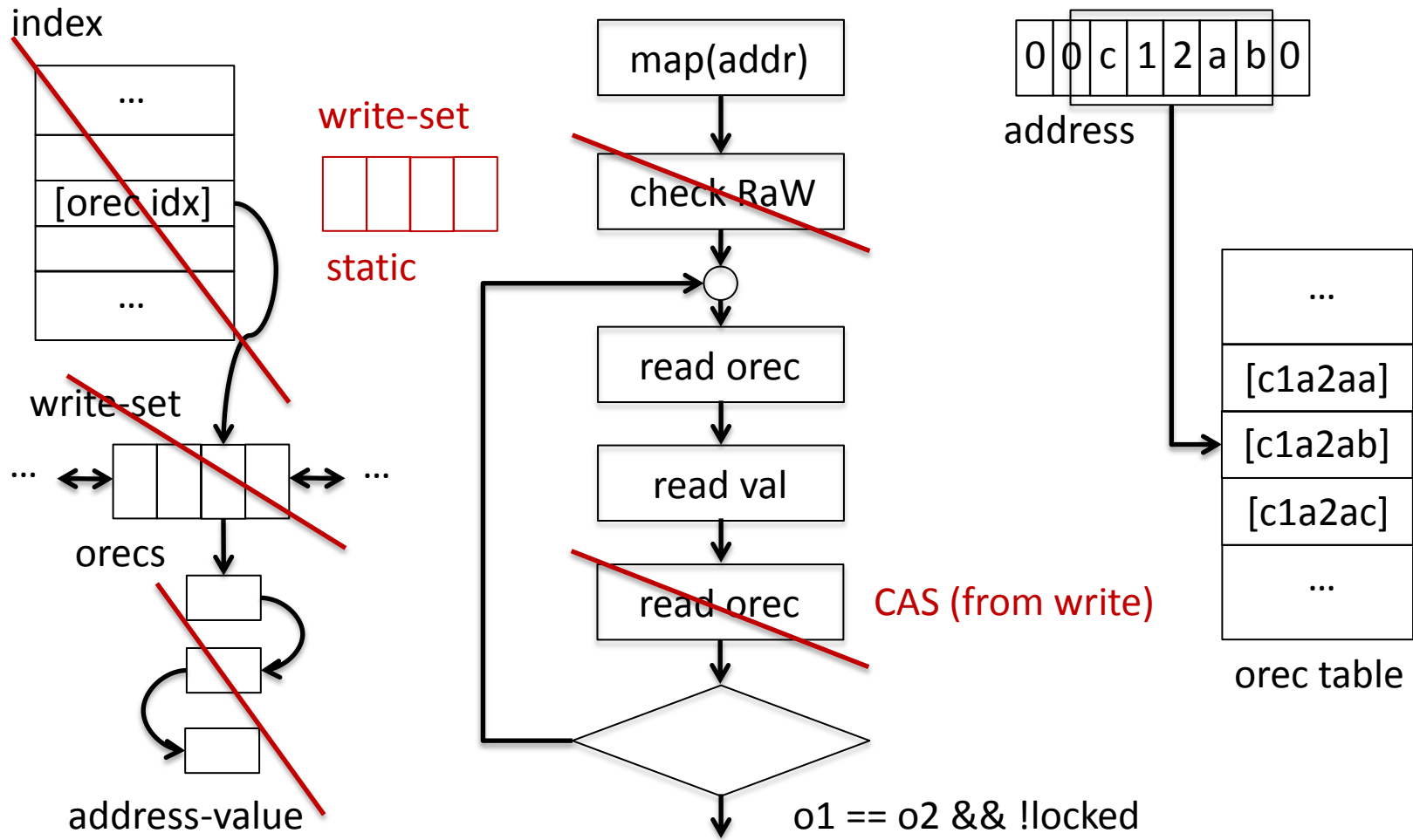
Short transaction optimizations



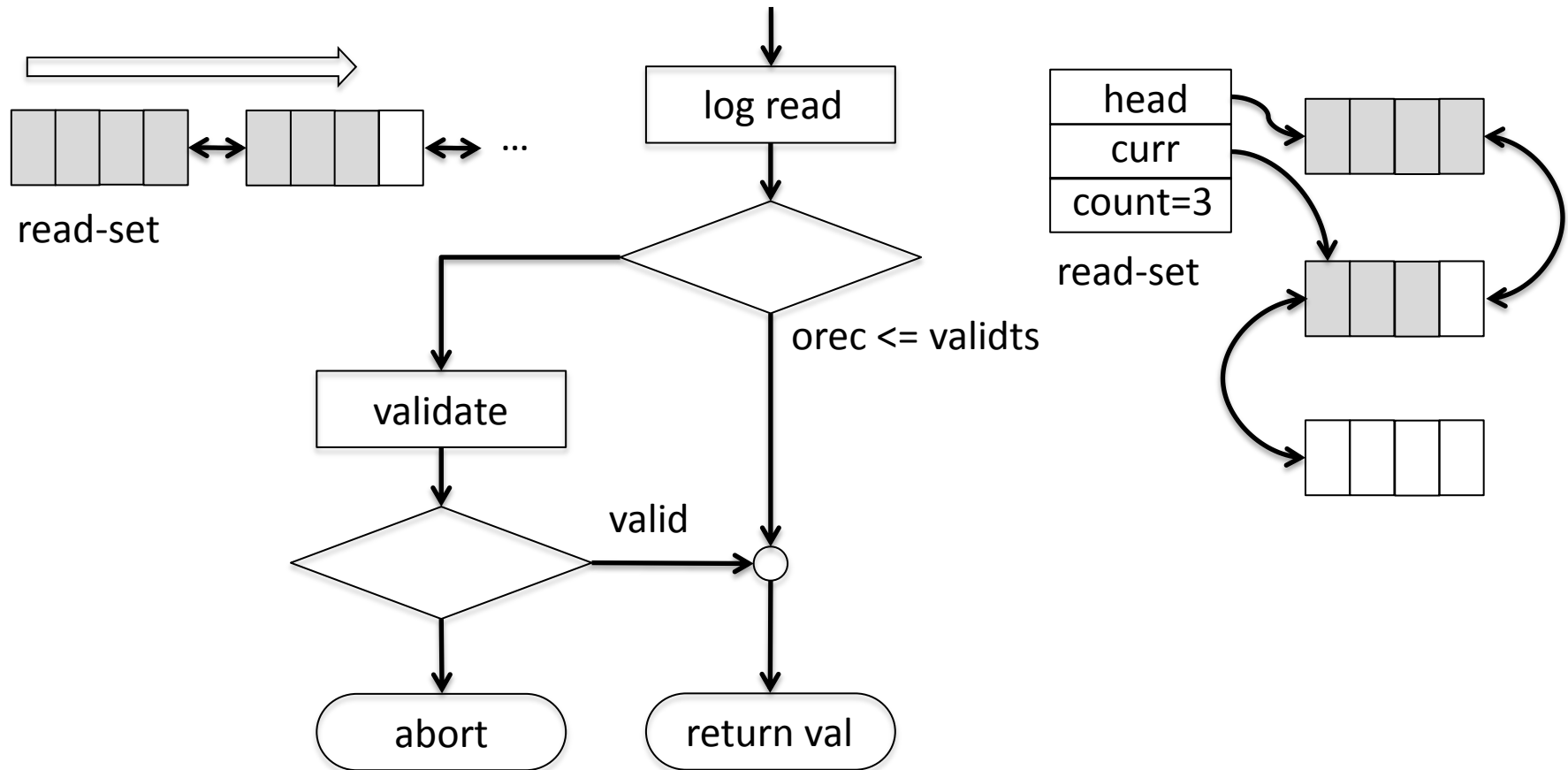
Short transaction optimizations



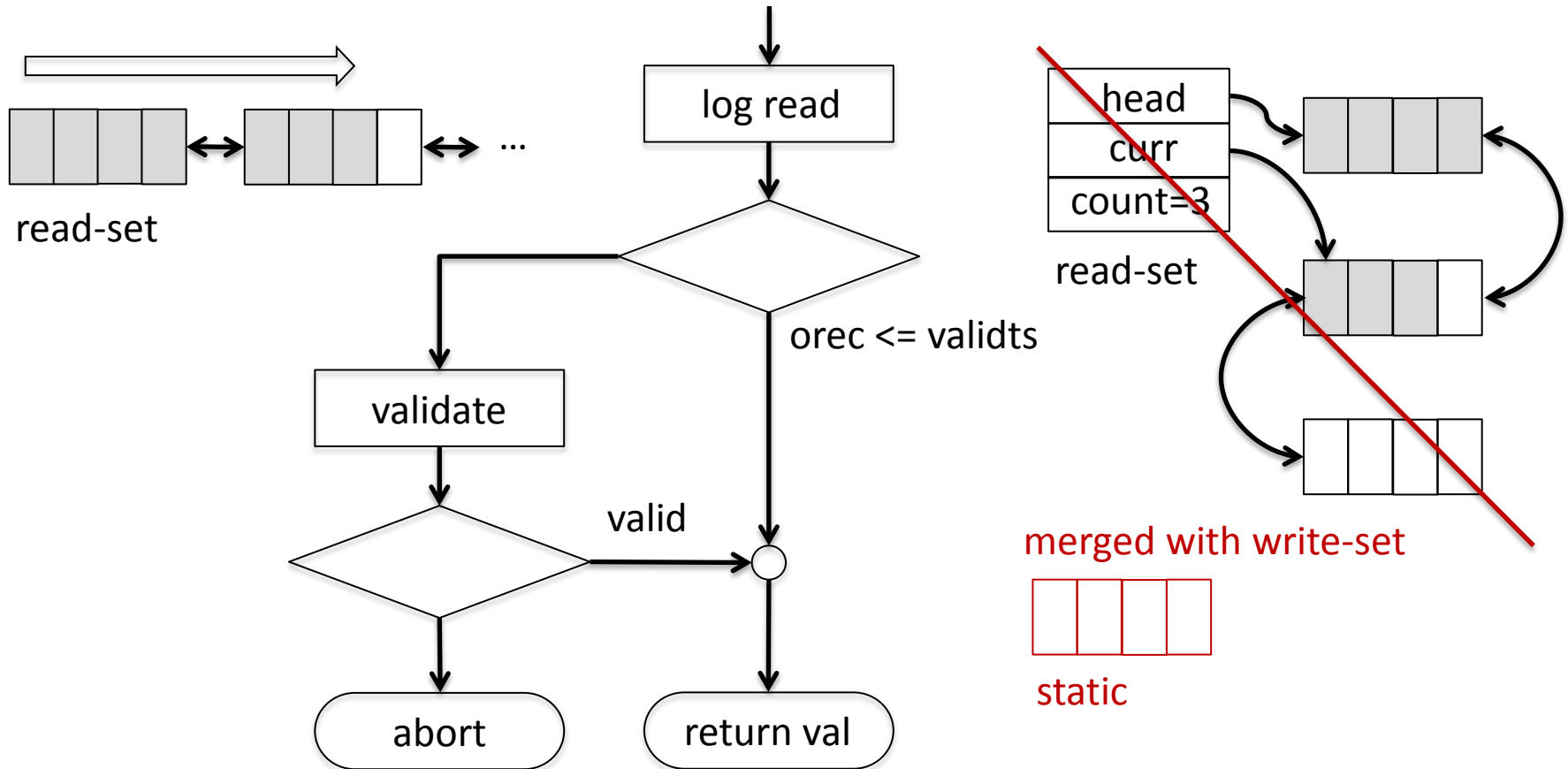
Short transaction optimizations



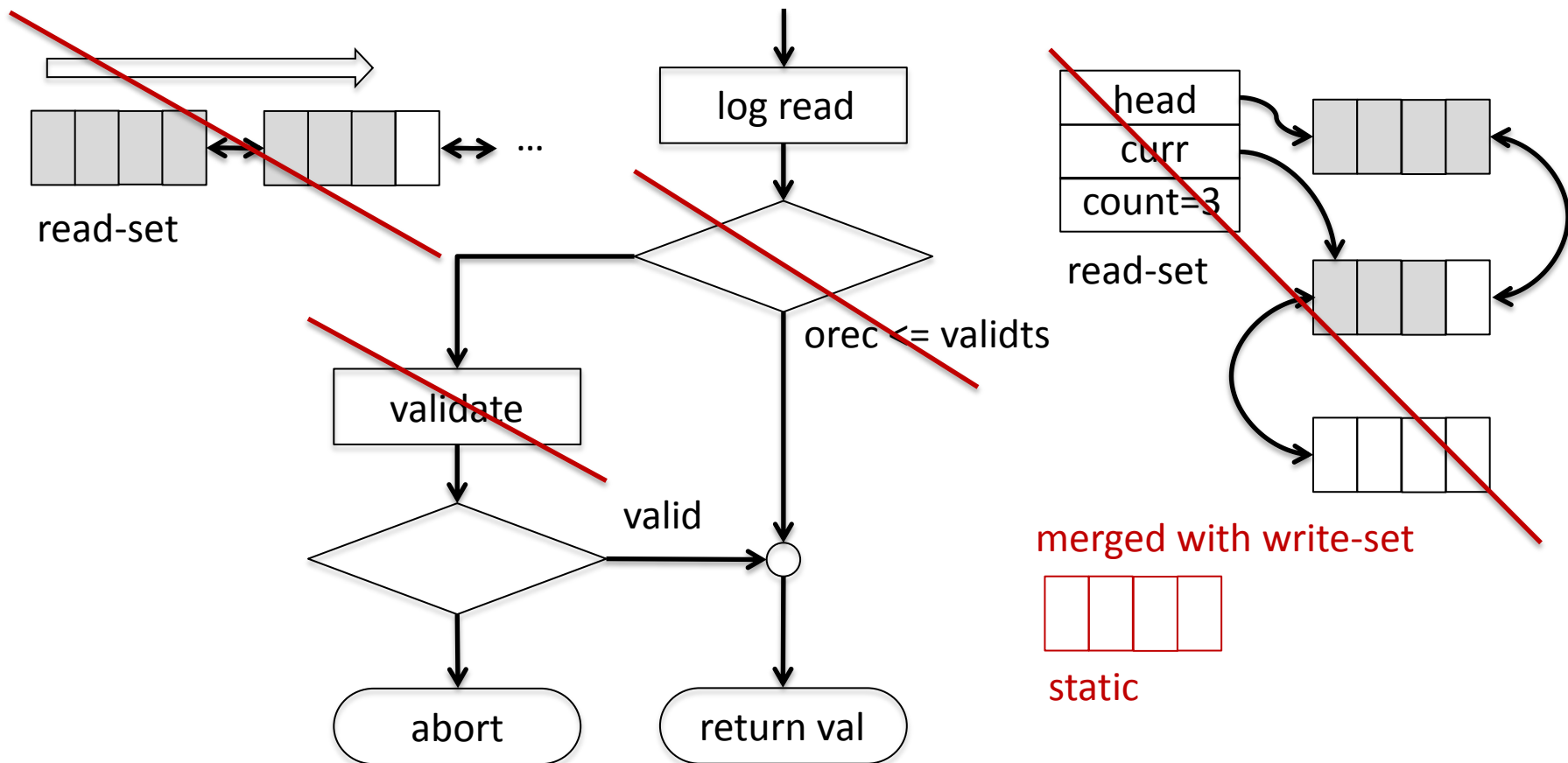
Short transaction optimizations (cont'd)



Short transaction optimizations (cont'd)



Short transaction optimizations (cont'd)

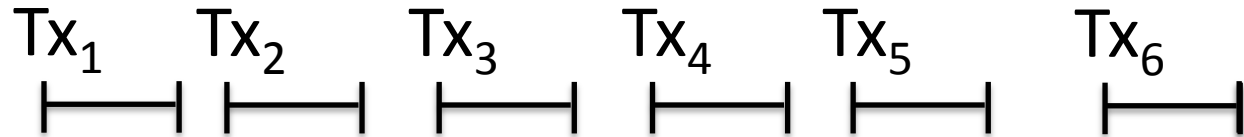


Using short transactions



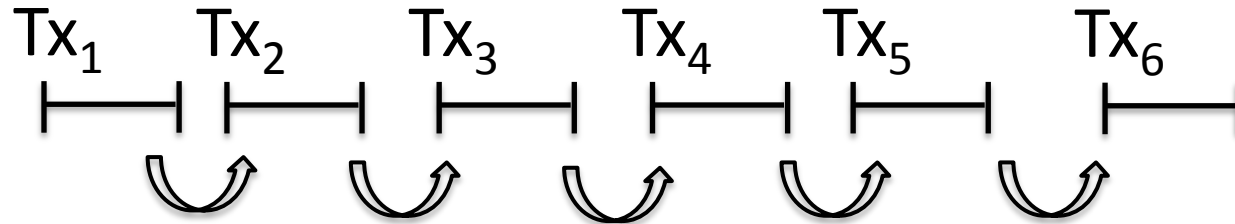
With “traditional” transactions the whole operation is a single transaction

Using short transactions



Split operation into a sequence of short transactions

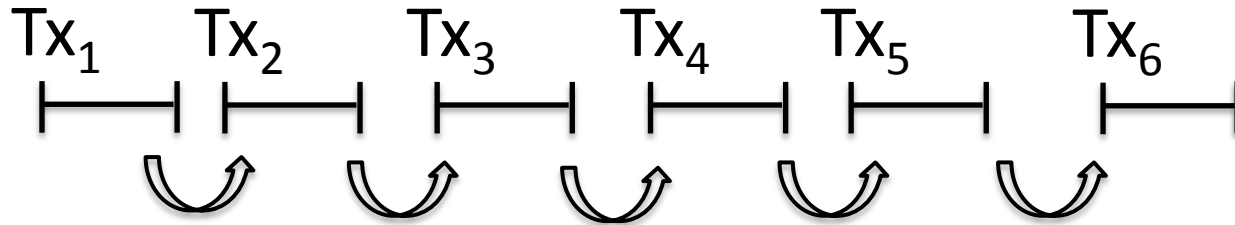
Using short transactions



“Glue” them together using techniques similar to lock-free algorithms (e.g. pointer marking)

Using short transactions

Is this simpler than using CAS directly?

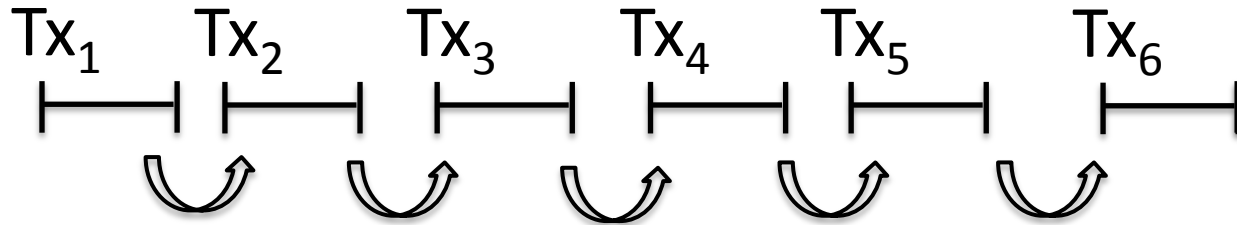


“Glue” them together using techniques similar to lock-free algorithms (e.g. pointer marking)

Using short transactions

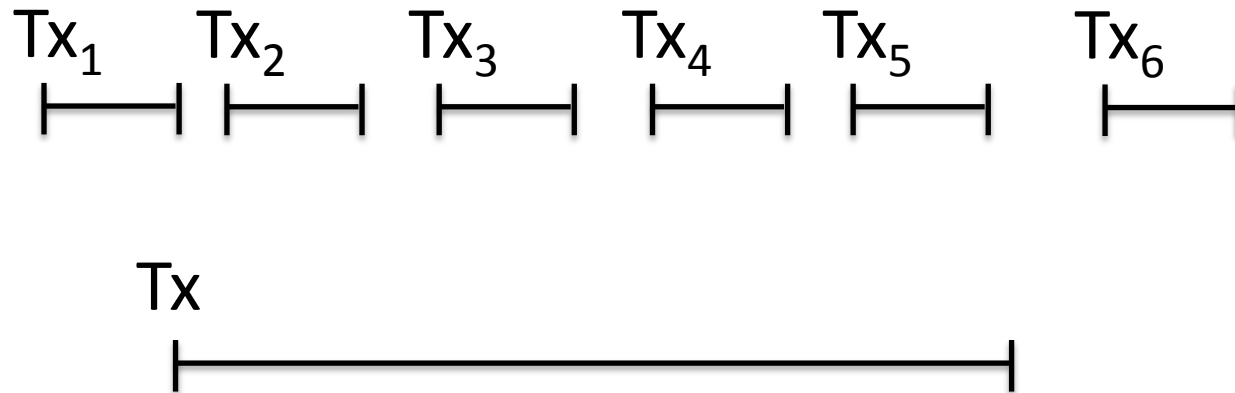
Is this simpler than using CAS directly?

Yes: transactions eliminate tricky races



“Glue” them together using techniques similar to lock-free algorithms (e.g. pointer marking)

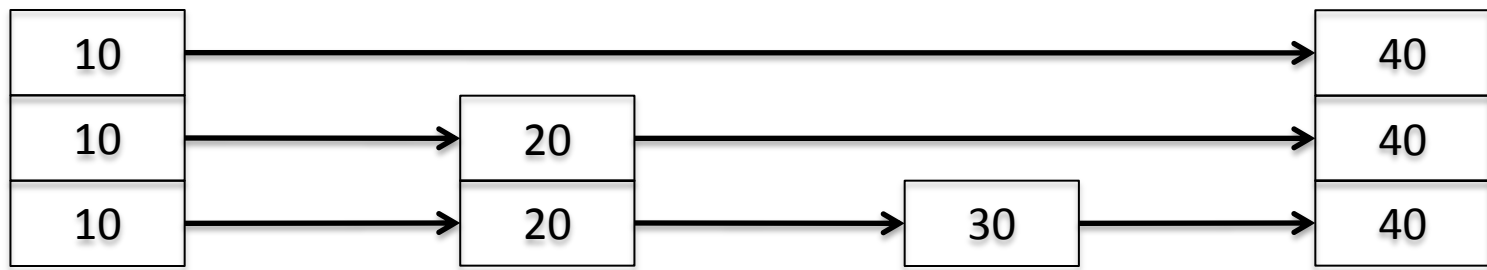
Mix short and ordinary transactions



Implement corner cases using ordinary transactions

Example

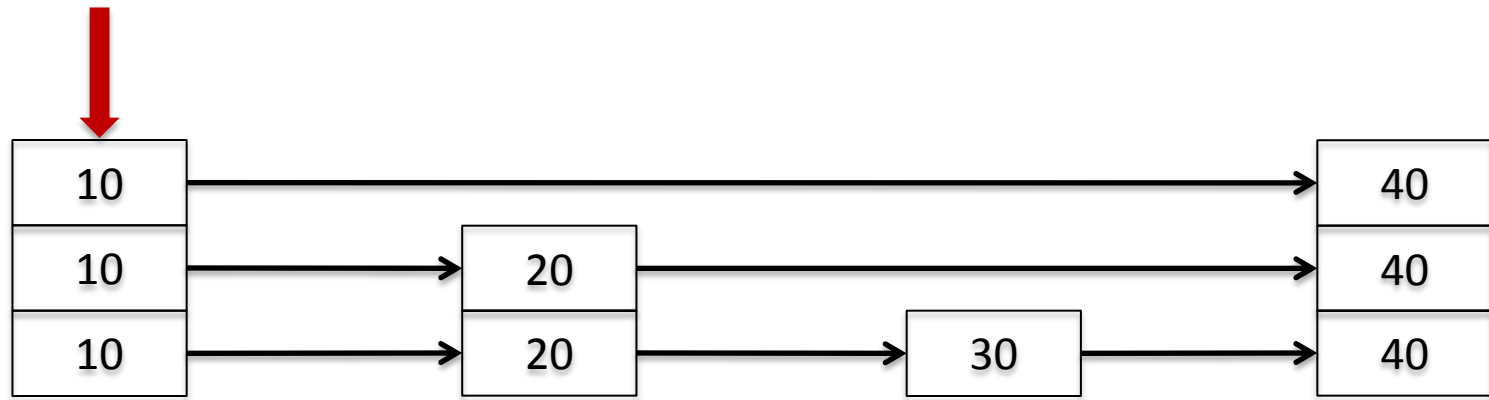
Remove 20



Example

Remove 20

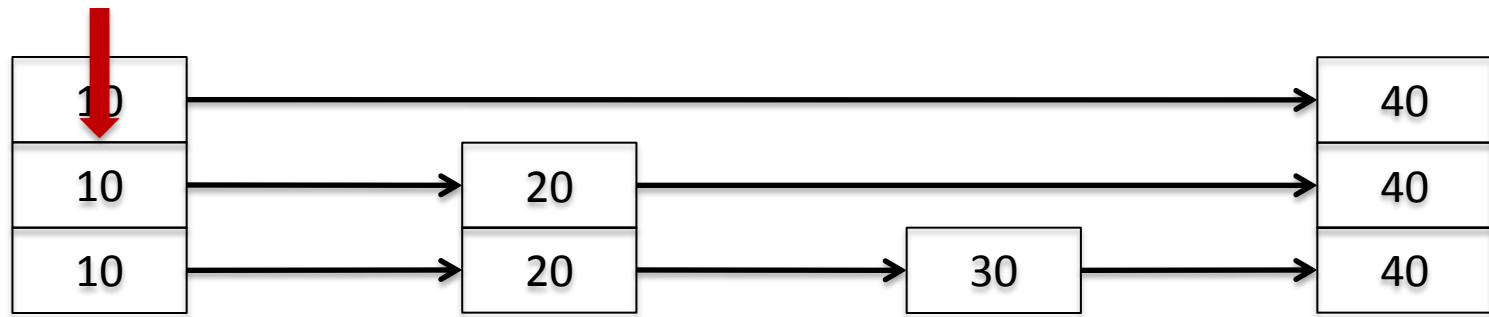
TxSingleRead



Example

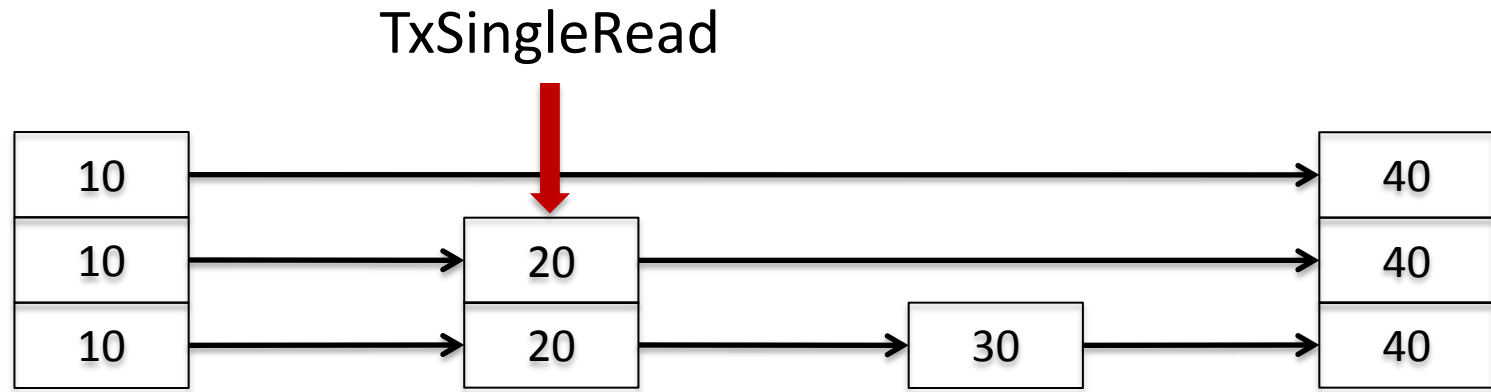
Remove 20

TxSingleRead



Example

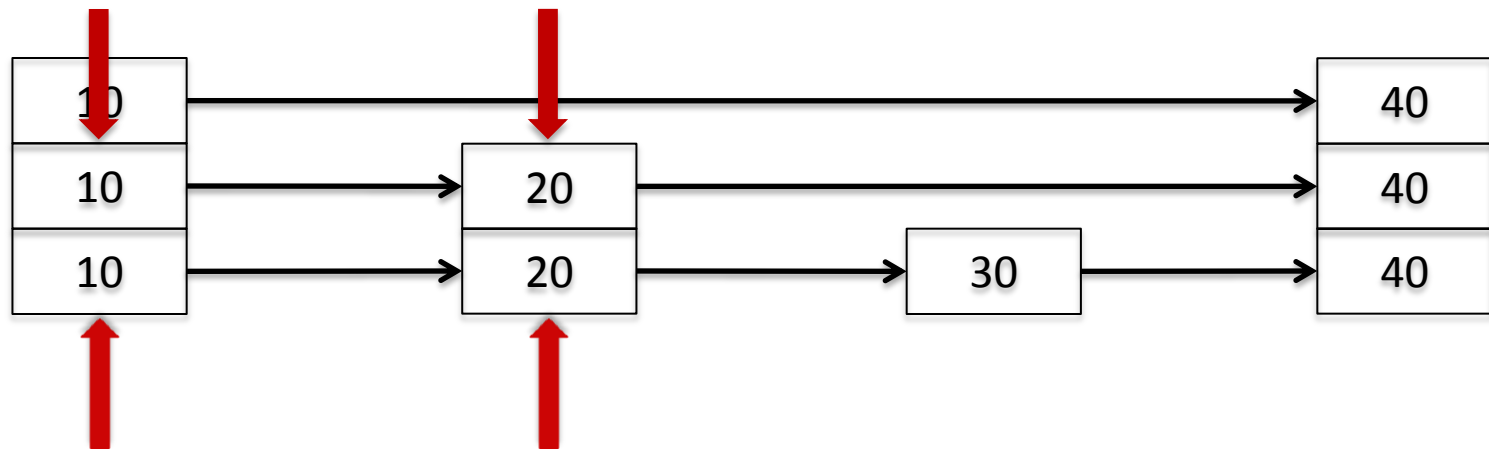
Remove 20



Example

Remove 20

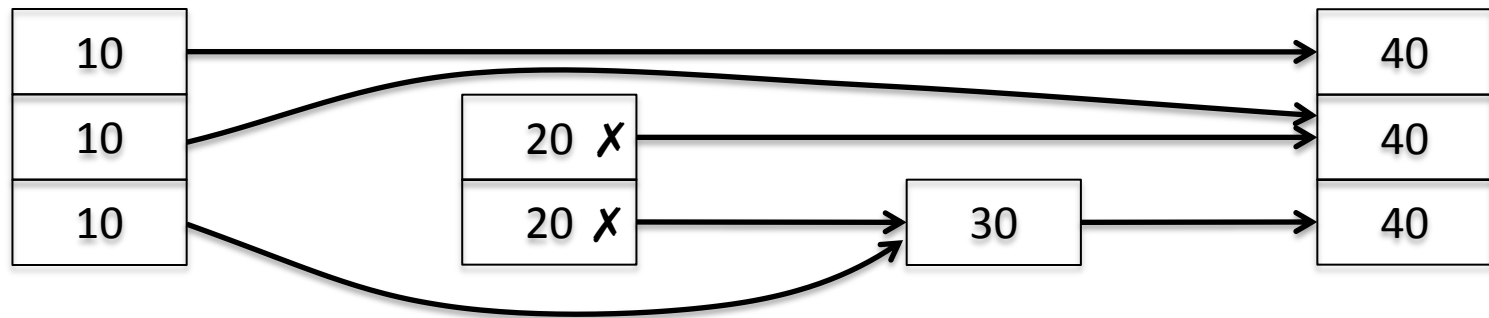
Short read-write transaction to mark the node and unlink it



Example

Remove 20

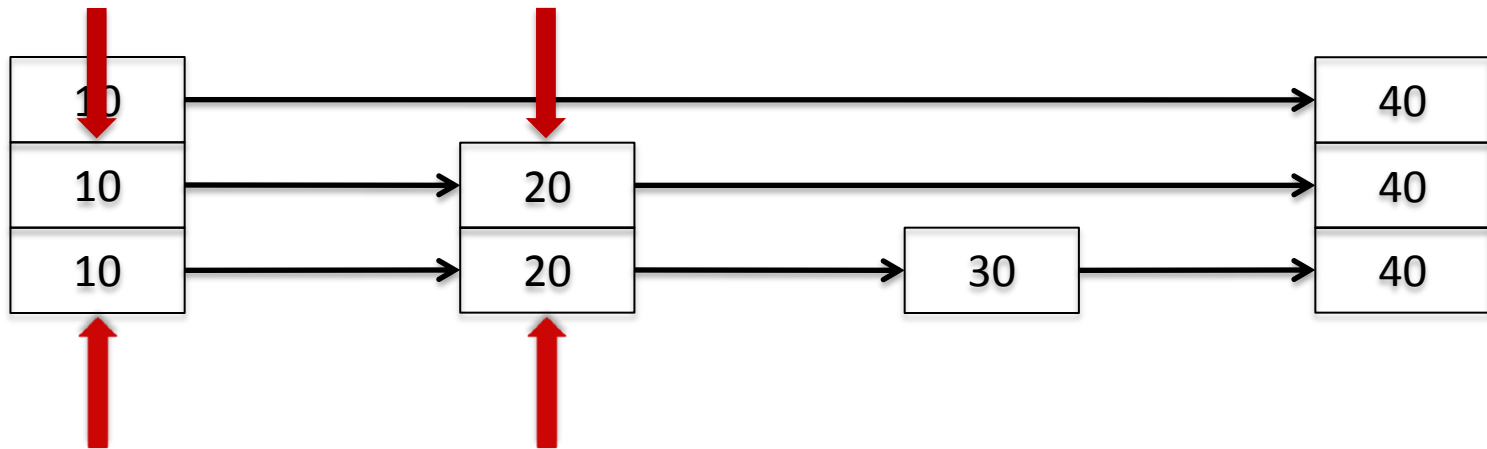
Short read-write transaction to mark the node and unlink it



Example

Remove 20

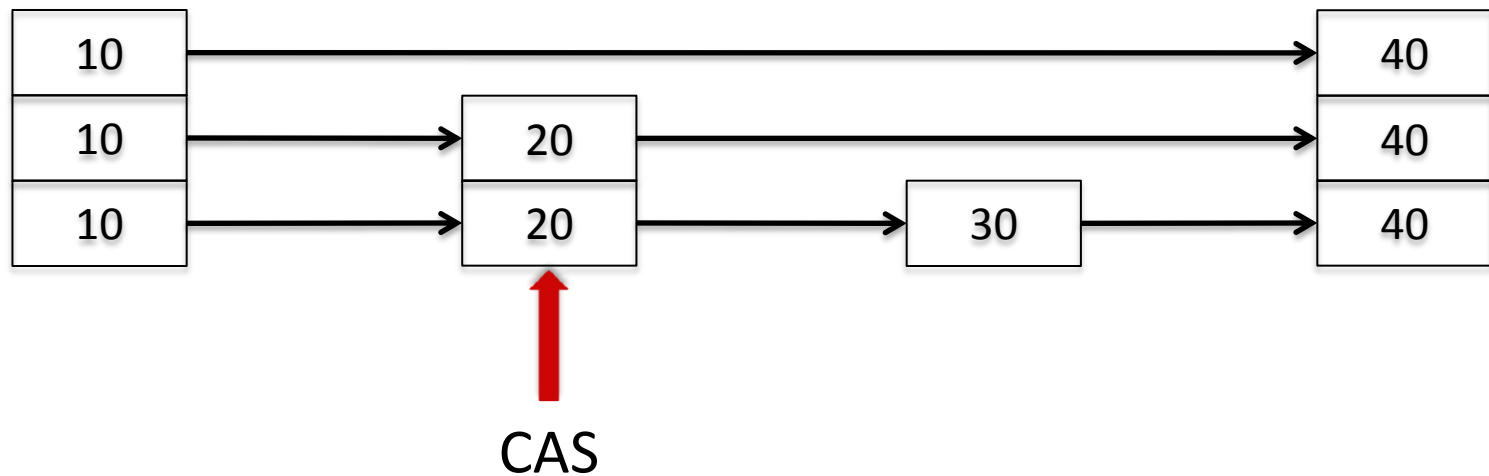
This becomes trickier when using CAS directly



Example

Remove 20

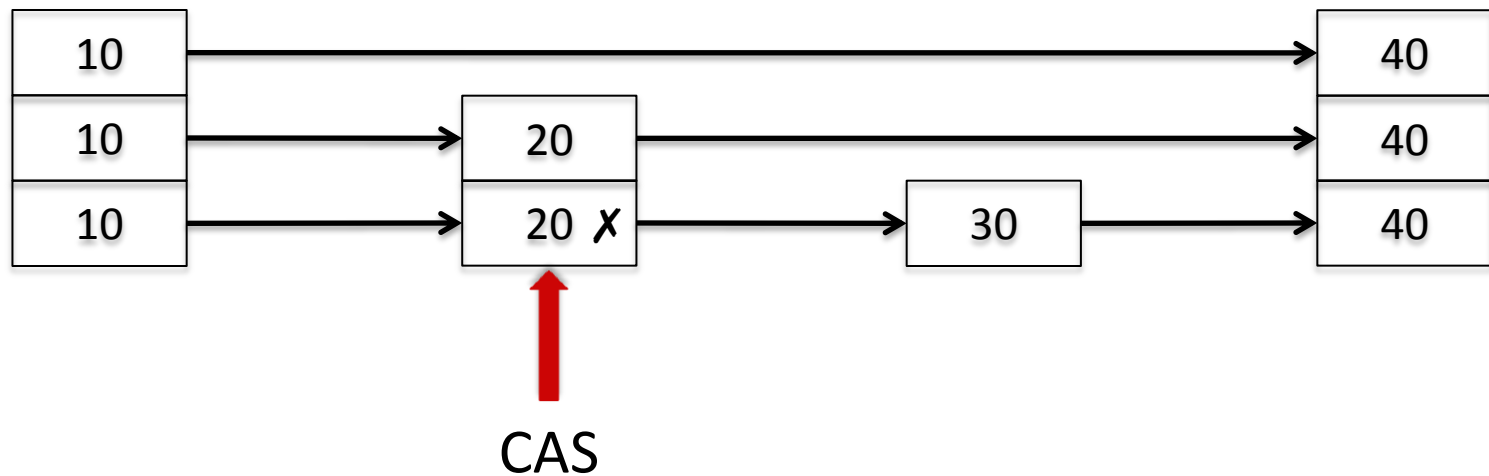
We have to use several CASes



Example

Remove 20

We have to use several CASes

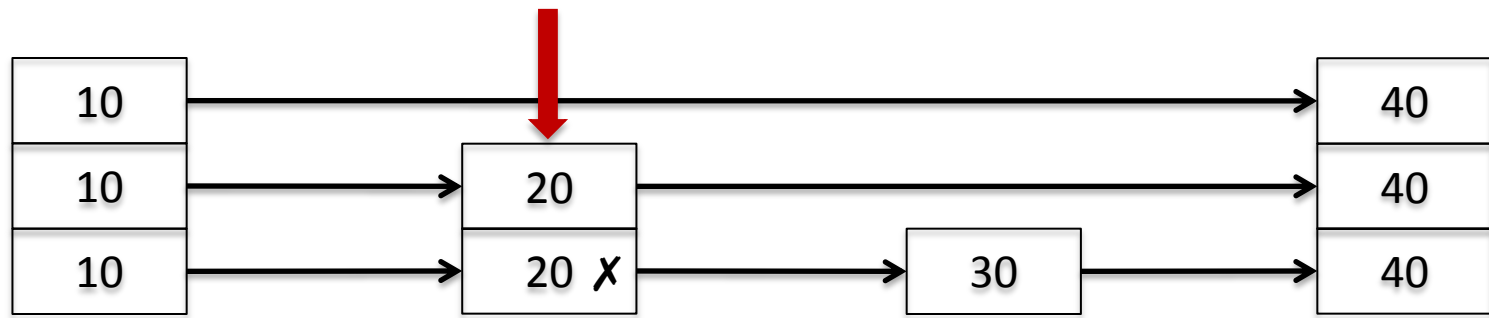


Example

Remove 20

We have to use several CASes

CAS

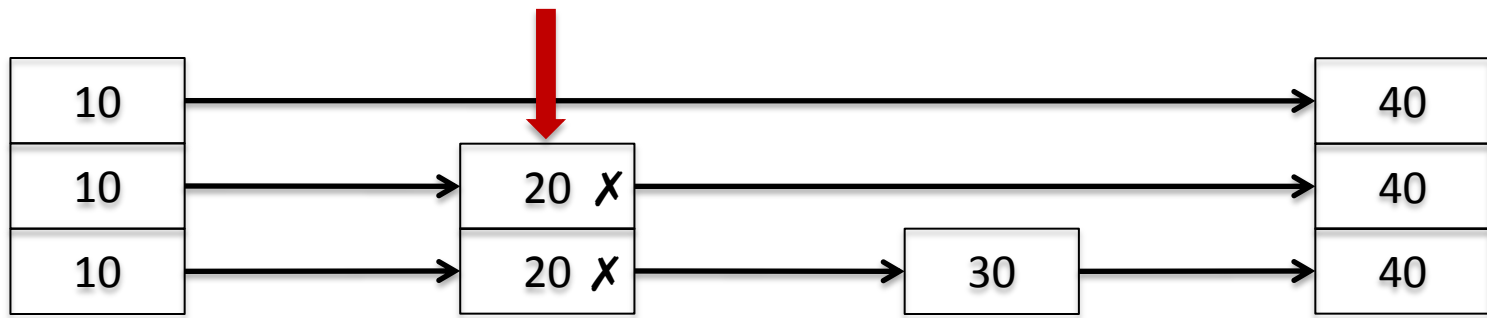


Example

Remove 20

We have to use several CASes

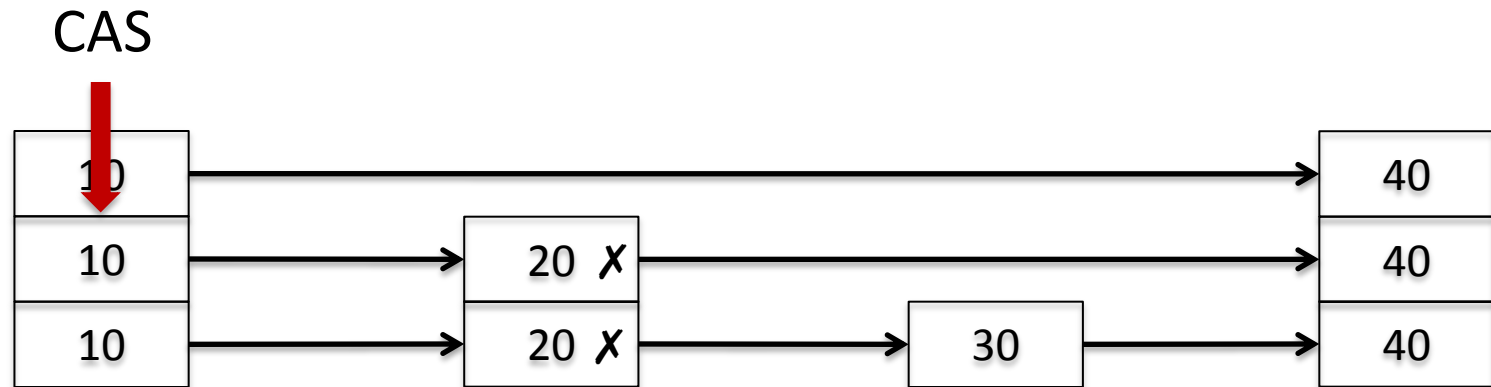
CAS



Example

Remove 20

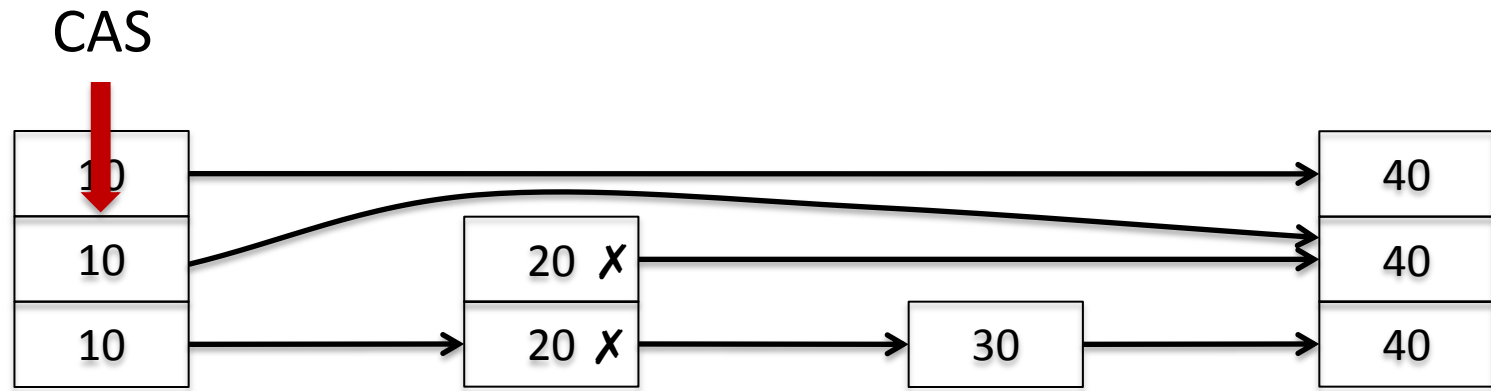
We have to use several CASes



Example

Remove 20

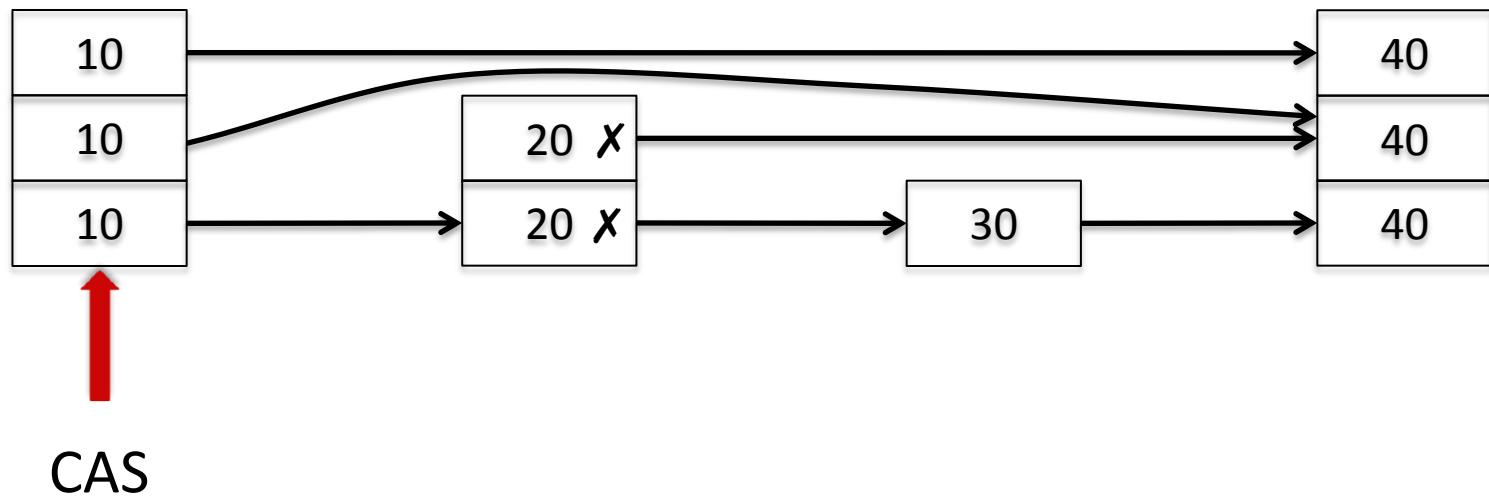
We have to use several CASes



Example

Remove 20

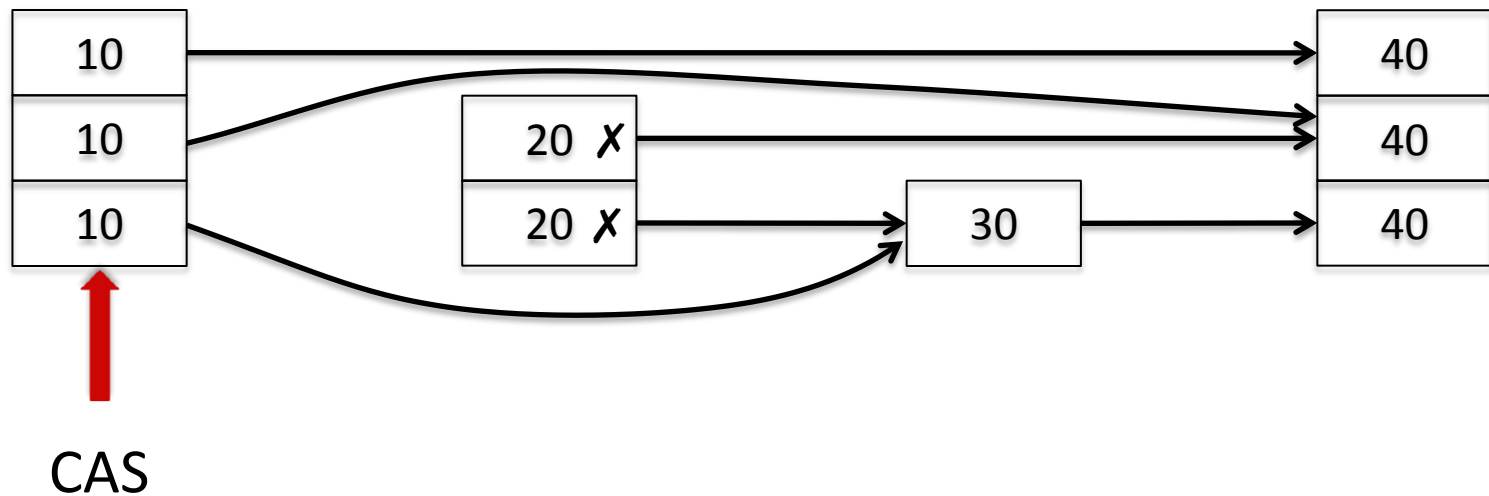
We have to use several CASes



Example

Remove 20

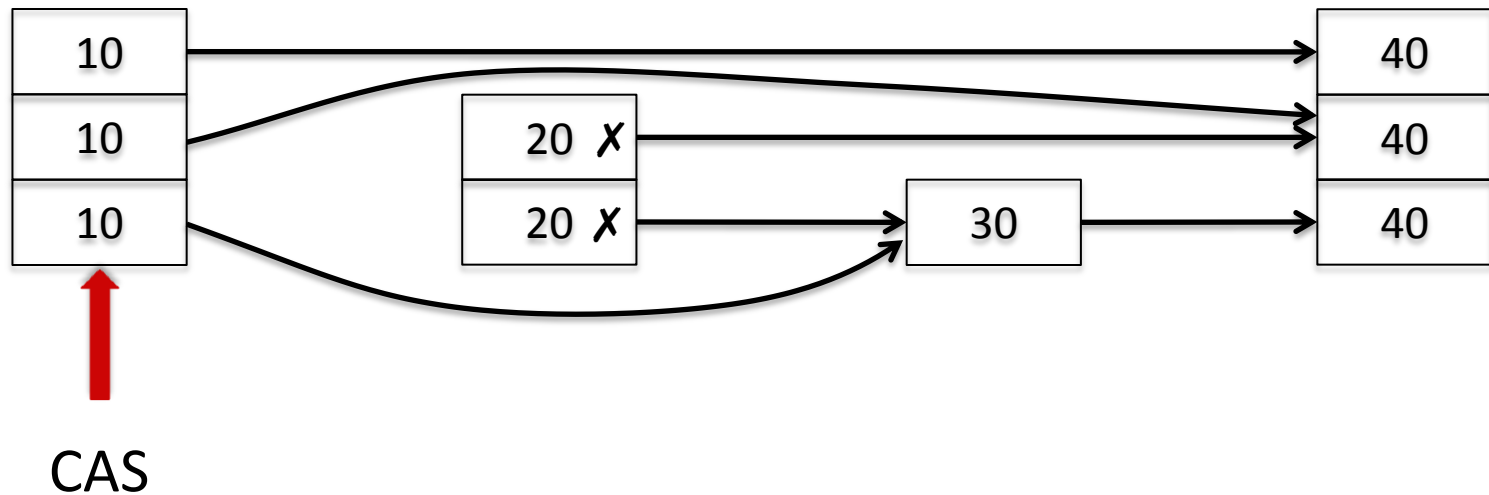
We have to use several CASes



Example

Remove 20

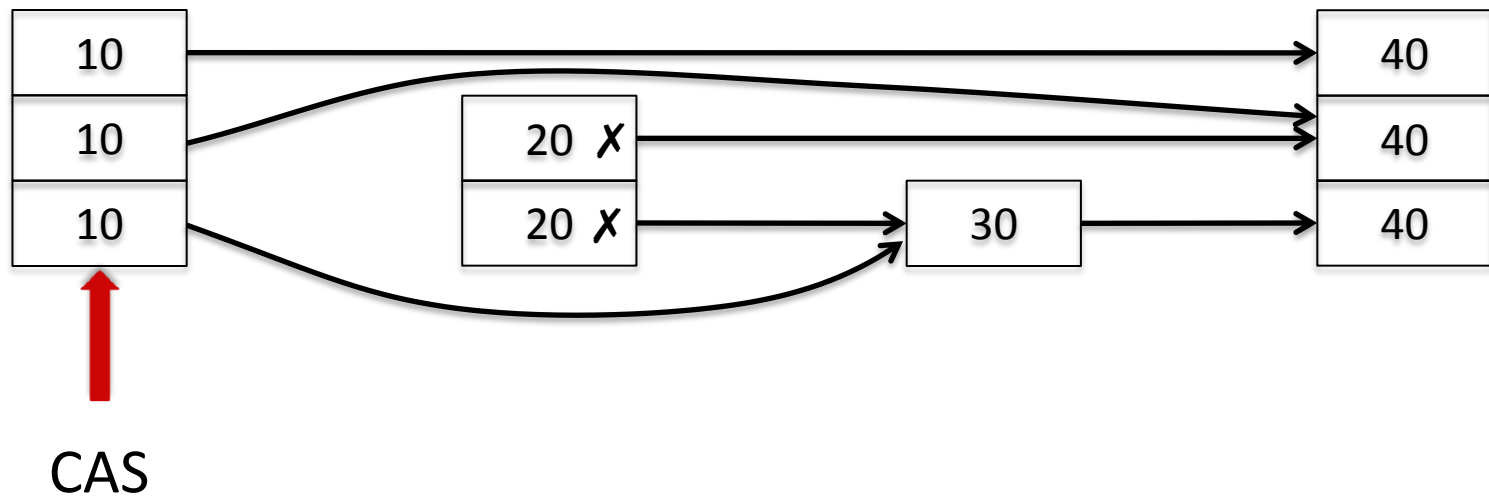
There might be a race at each of these steps



Example

Remove 20

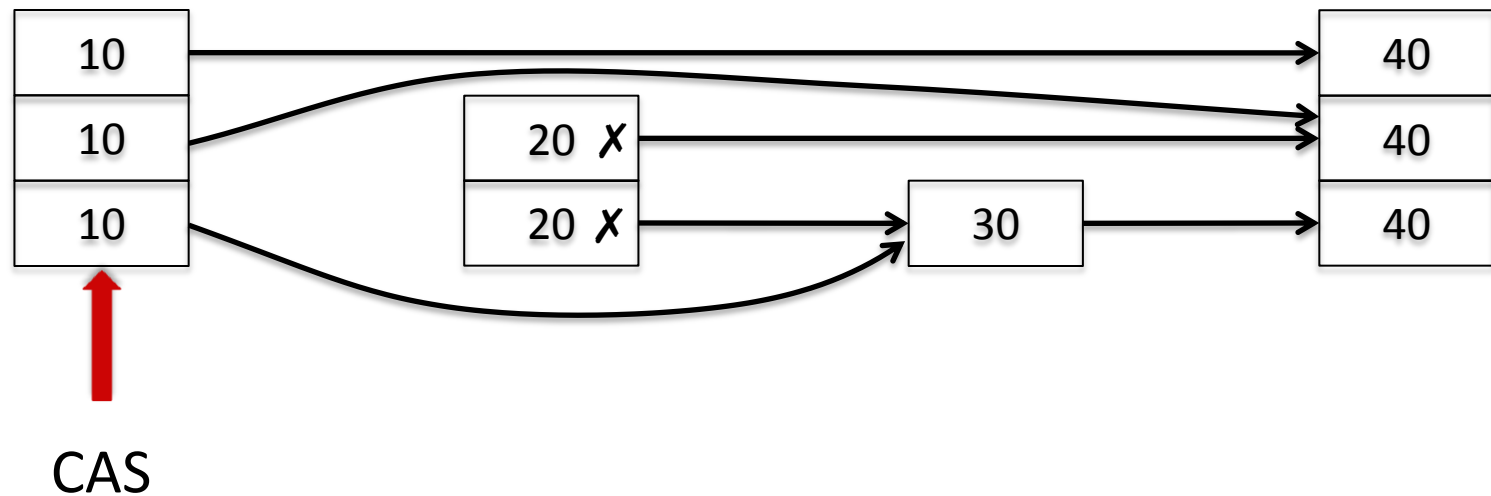
We need to make sure other operations clean-up for us



Example

Remove 20

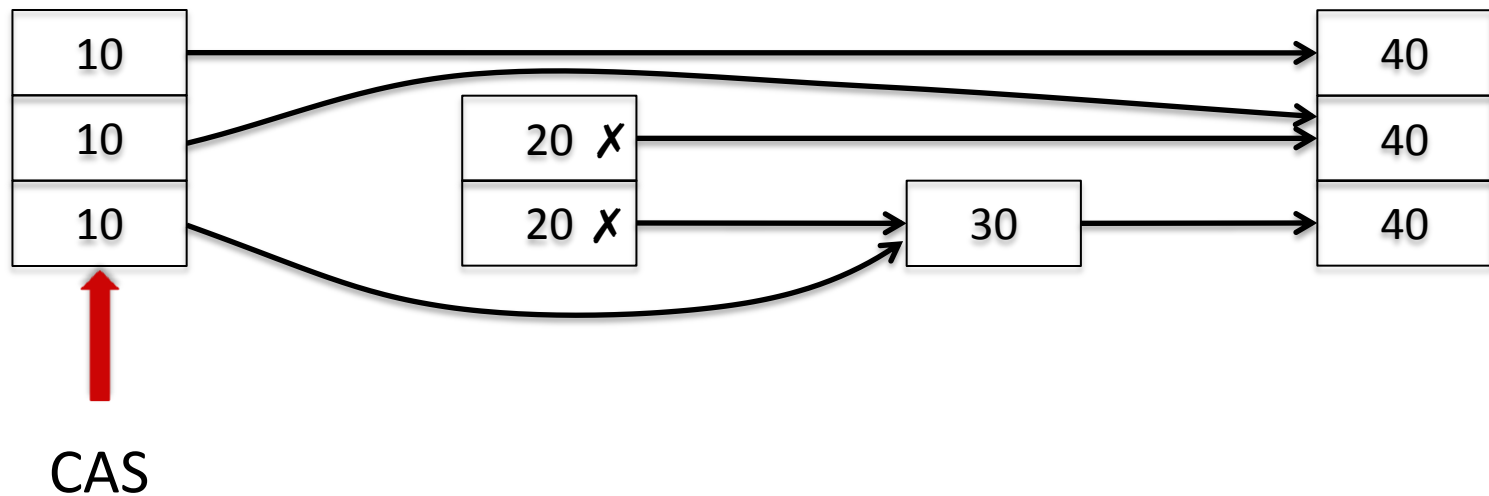
All interactions between operations have to be handled correctly



Example

Remove 20

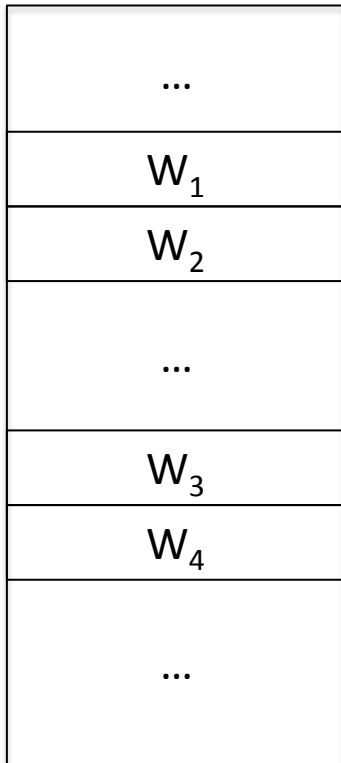
E.g. removal of an element that is still being added to the list



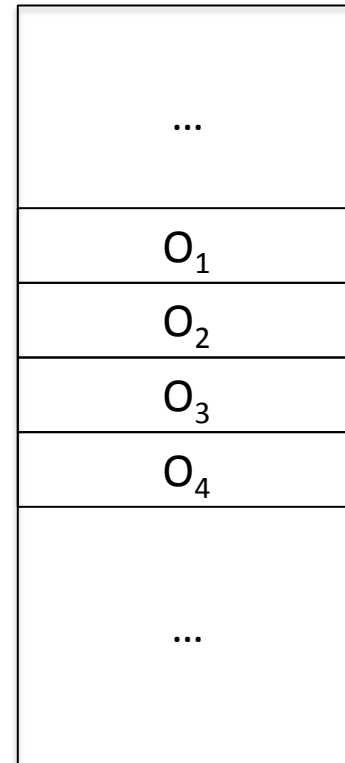
Overview

- STM overheads
- SpecTM
 - Short-transaction API
 - Collocating data and meta-data
 - In-word meta-data
- Evaluation

“Traditional” STM data layout

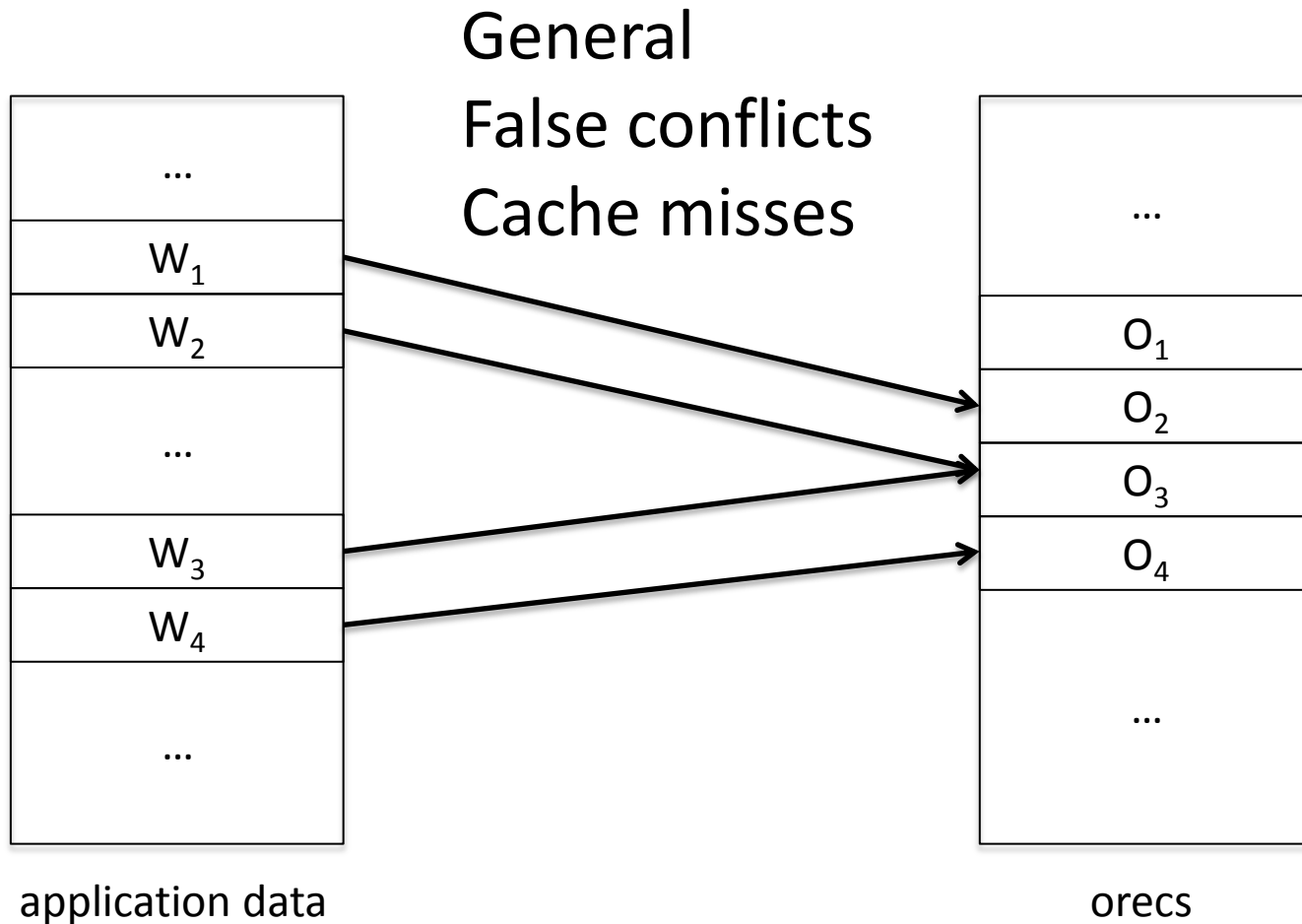


application data

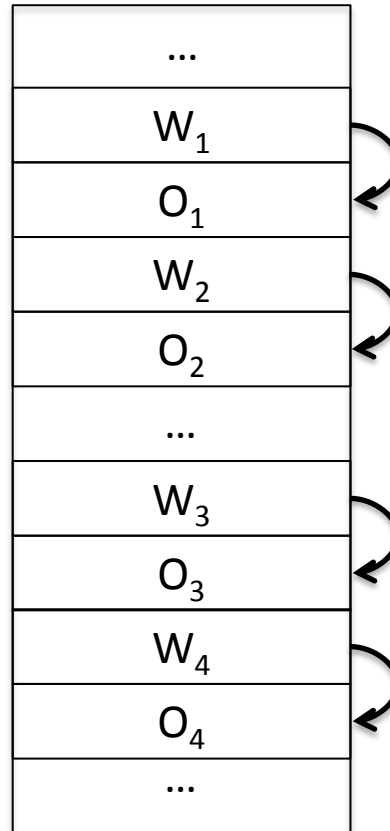


orecs

“Traditional” STM data layout

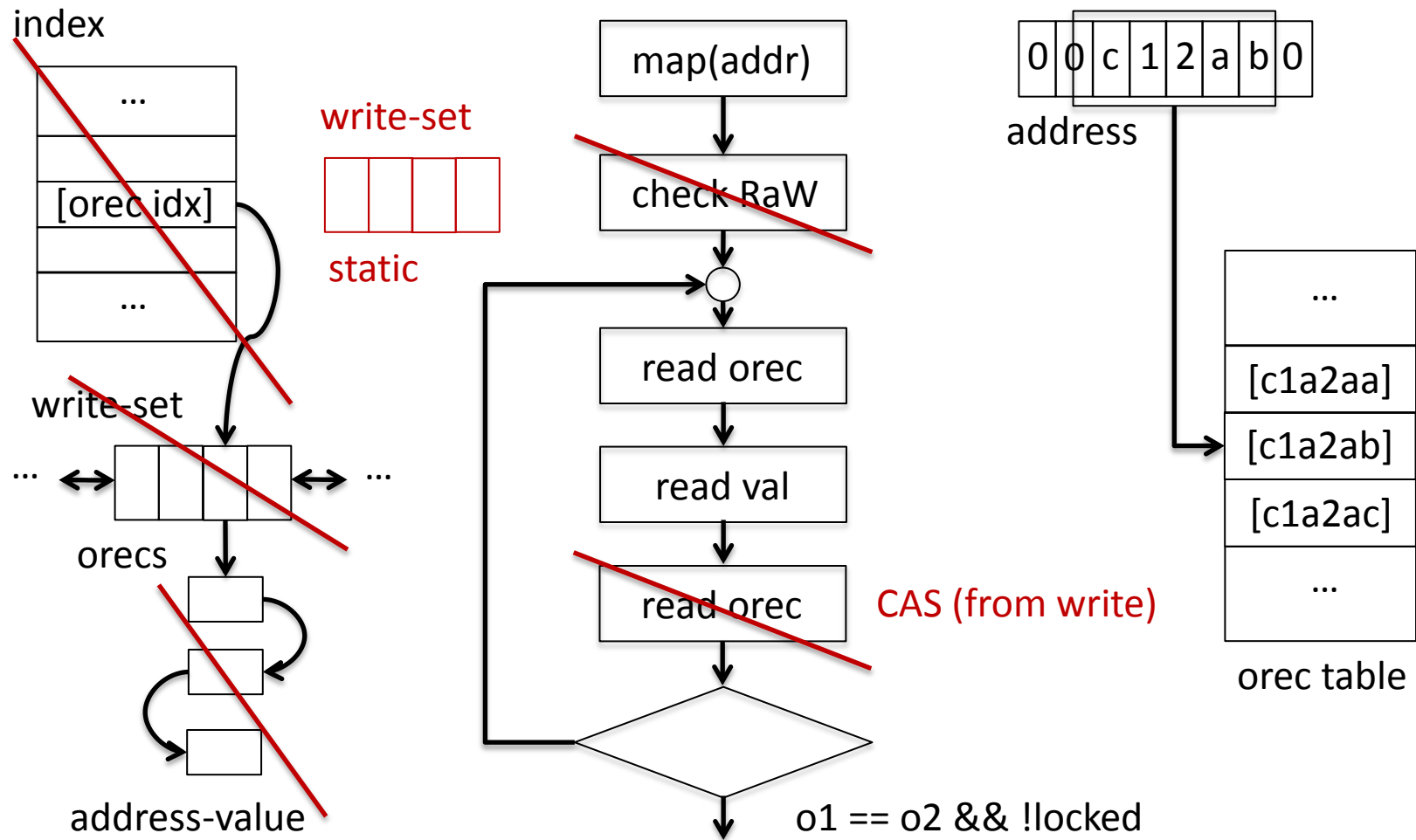


Collocate data and meta-data

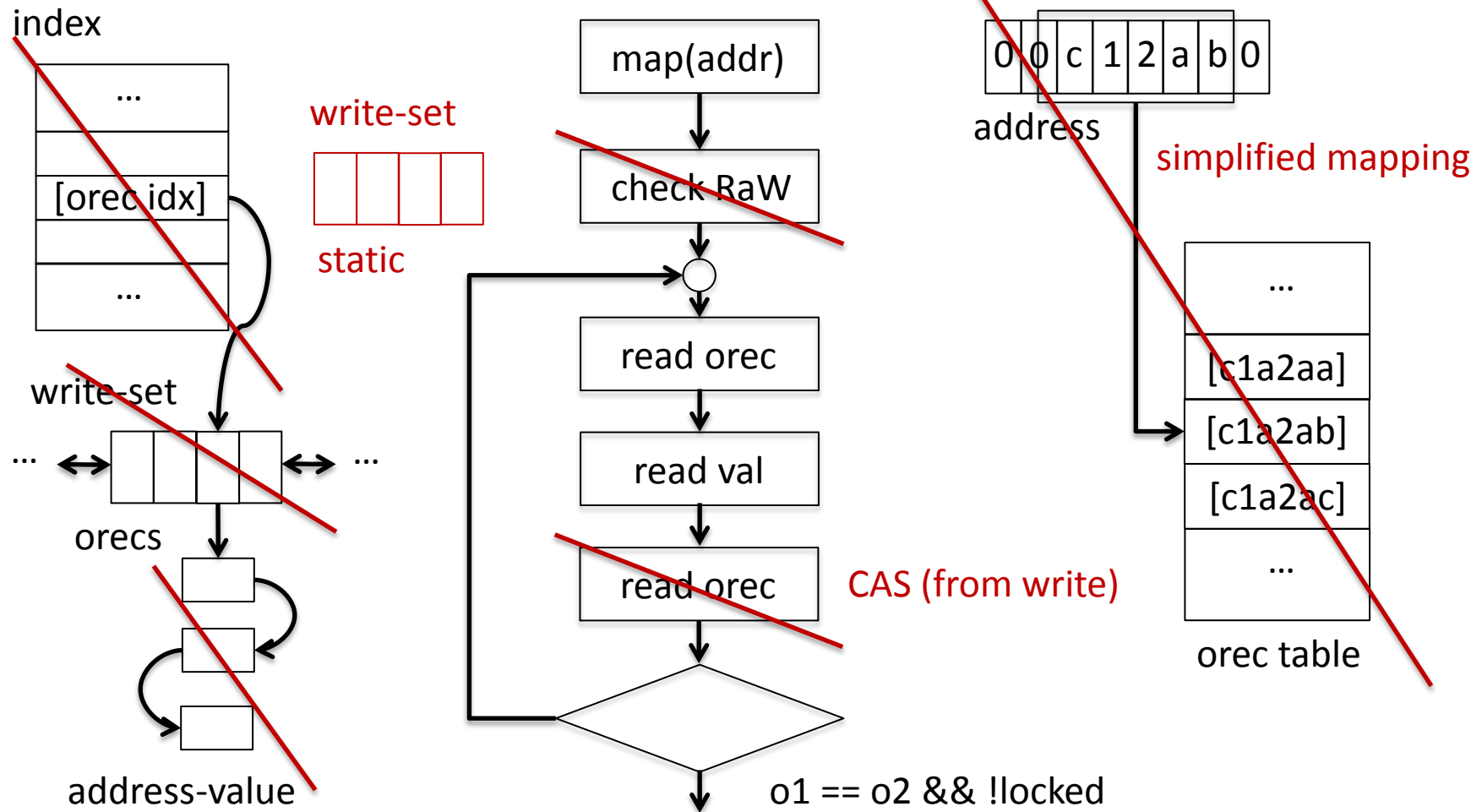


Simplify mapping
No false conflicts
Same cache line

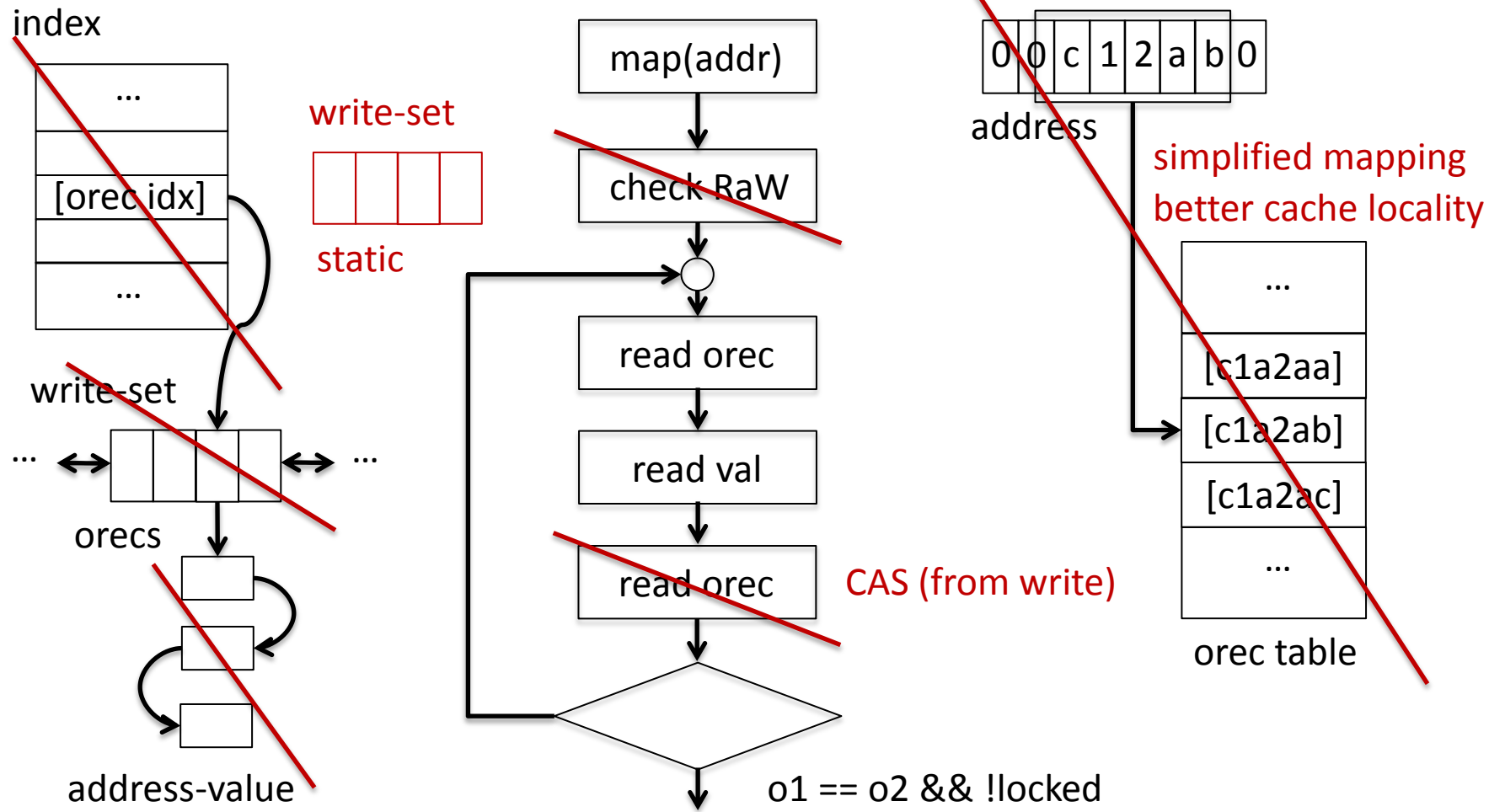
Collocated meta-data optimizations



Collocated meta-data optimizations



Collocated meta-data optimizations

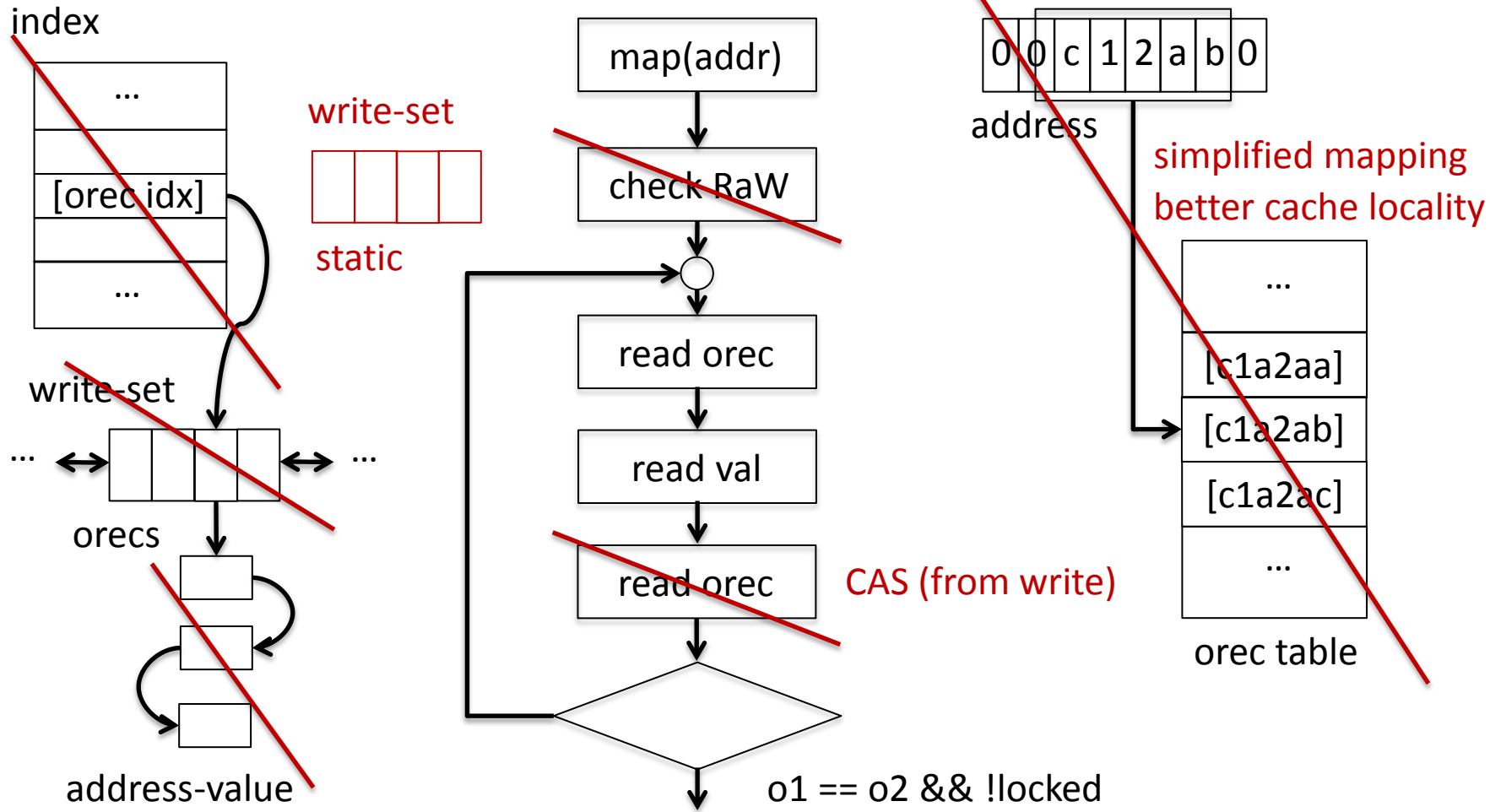


Pure value-based validation

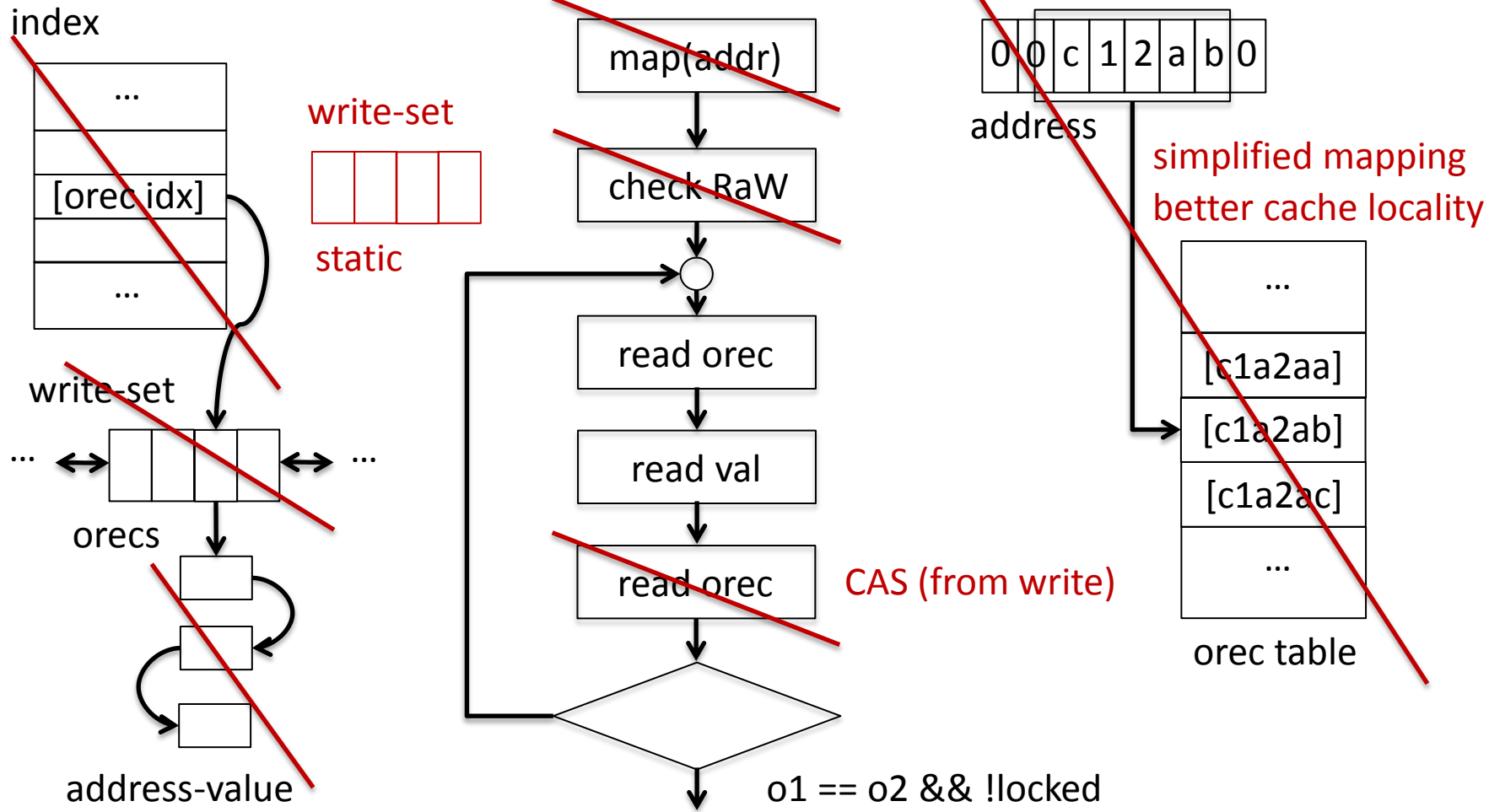
...	
W_1	
W_2	
...	
W_3	
W_4	
...	

Restrict values
Restrict access patterns
Single-word operations

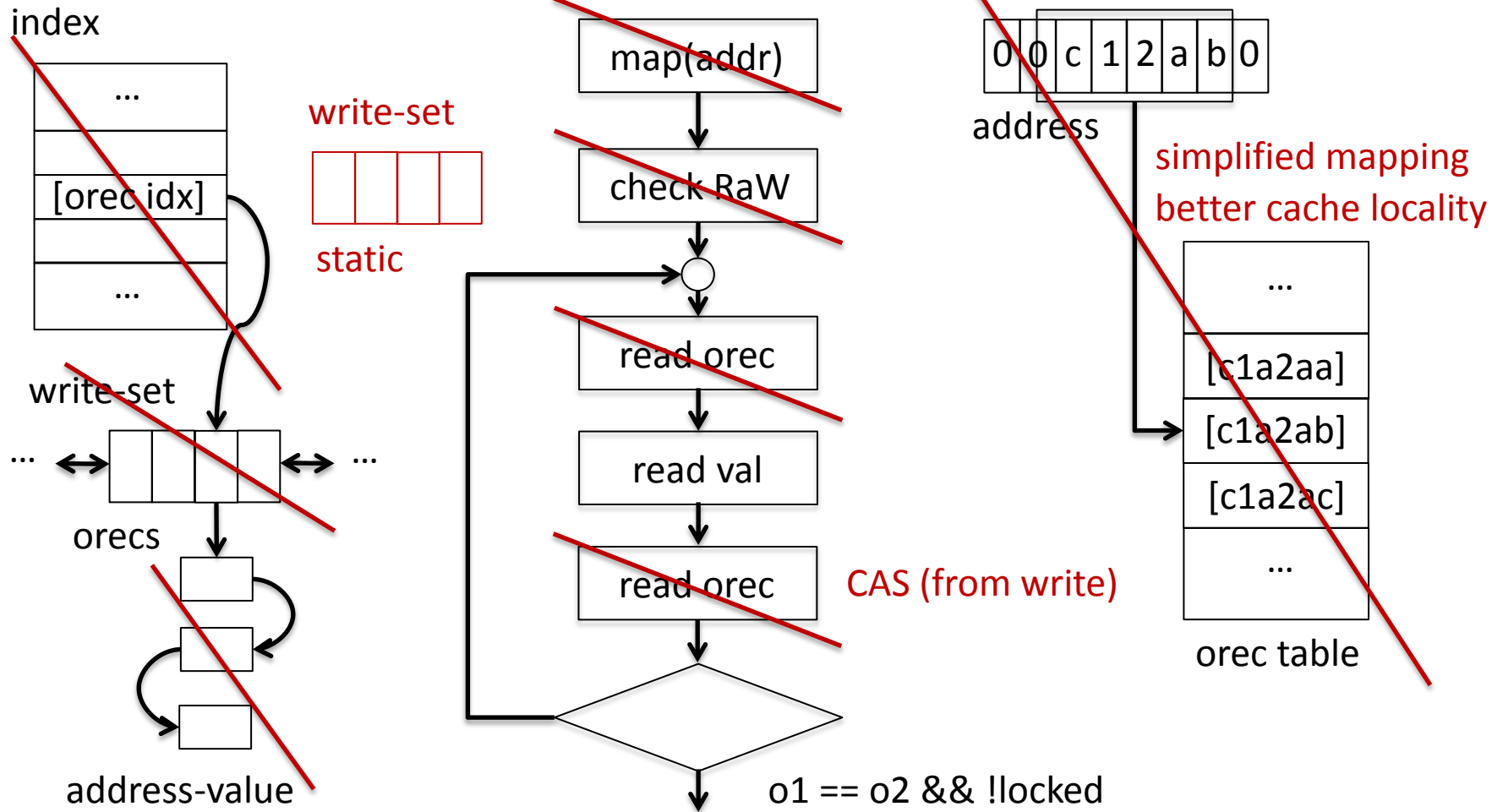
Value-based optimizations



Value-based optimizations



Value-based optimizations



API changes

API changes

```
word_t ReadWord(TVar *addr);
```

```
void WriteWord(TVar *addr, word_t val);
```

```
word_t TVarToWord(TVar addr);
```

```
TVar WordToTVar(word_t val);
```

API changes

Transactions access TVar, not word_t

```
word_t ReadWord(TVar *addr);
```

```
void WriteWord(TVar *addr, word_t val);
```

```
word_t TVarToWord(TVar addr);
```

```
TVar WordToTVar(word_t val);
```


API changes

Transactions access TVar, not word_t

```
word_t ReadWord(TVar *addr);
```

```
void WriteWord(TVar *addr, word_t val);
```

Calls to update TVar non-transactionally

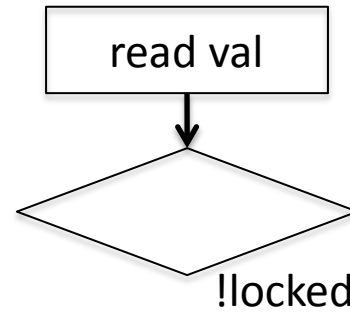
```
word_t TVarToWord(TVar addr);
```

```
TVar WordToTVar(word_t val);
```

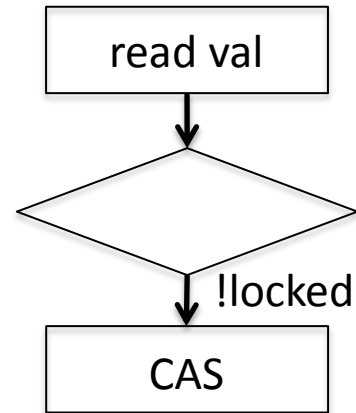
Pure value-based short read-write

read val

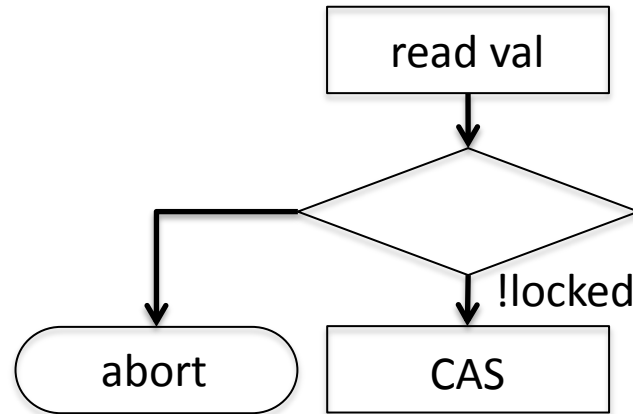
Pure value-based short read-write



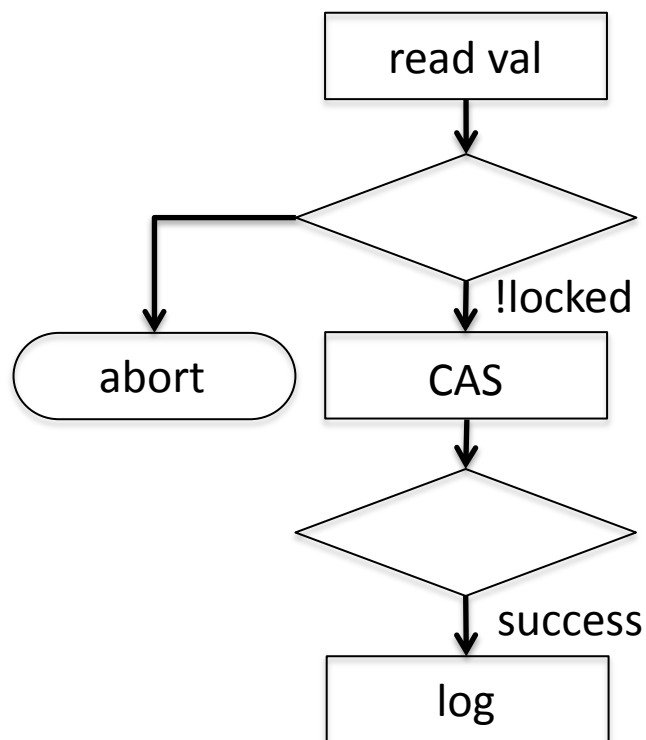
Pure value-based short read-write



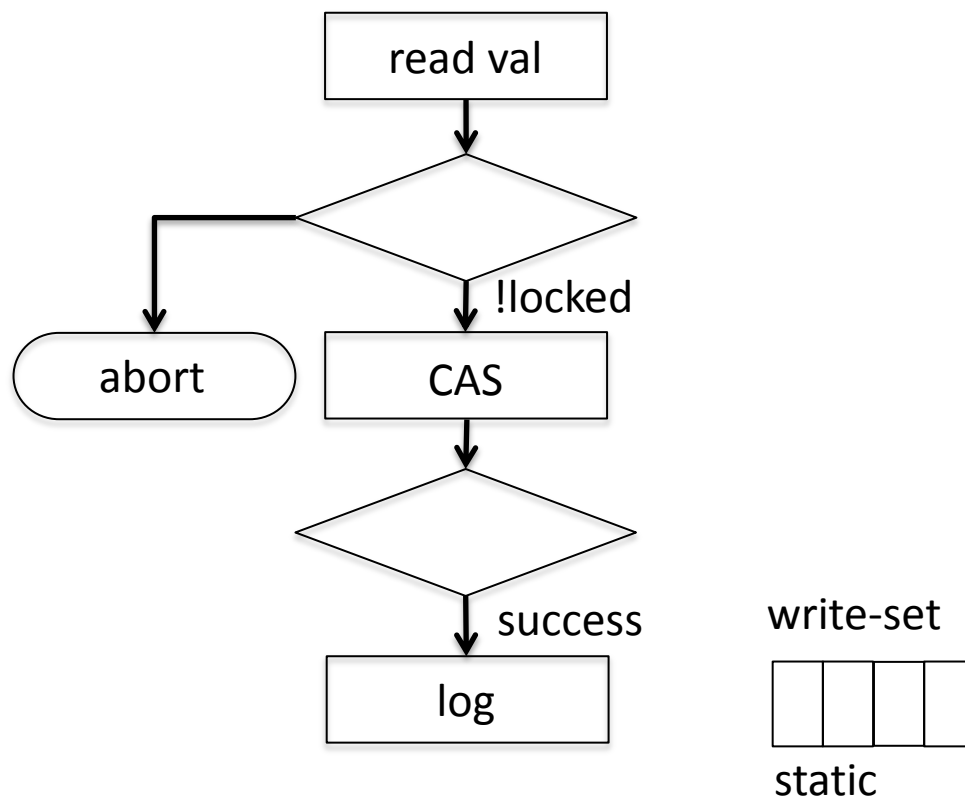
Pure value-based short read-write



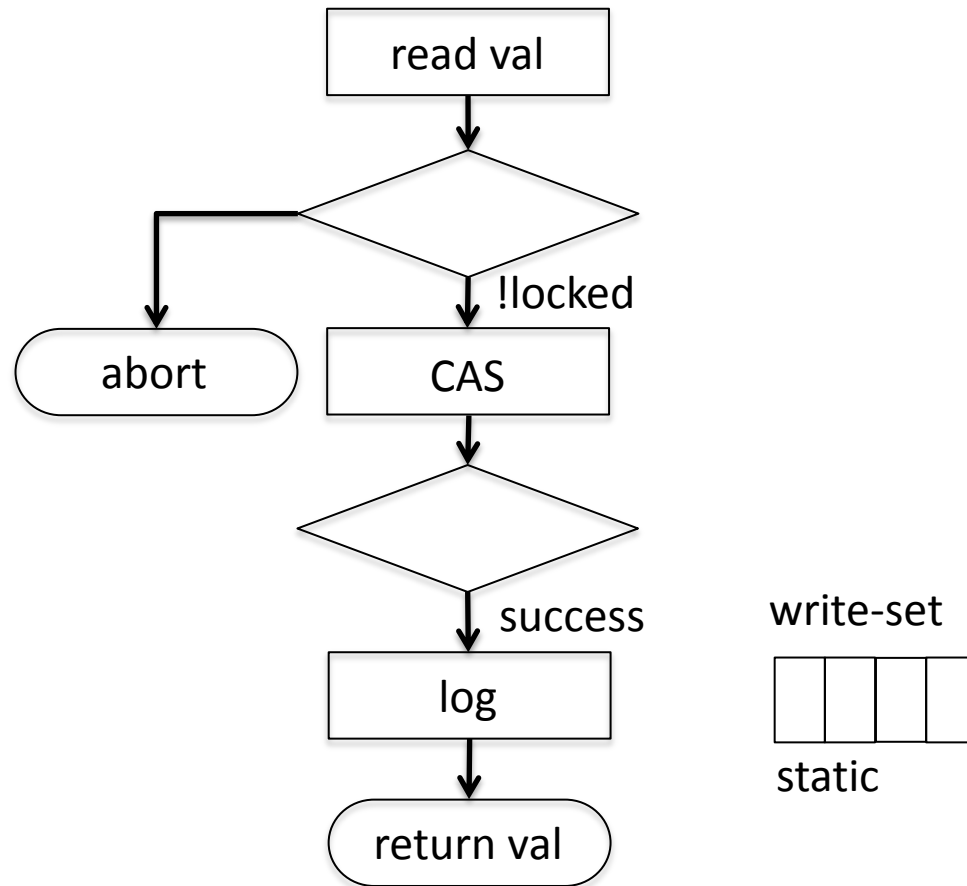
Pure value-based short read-write



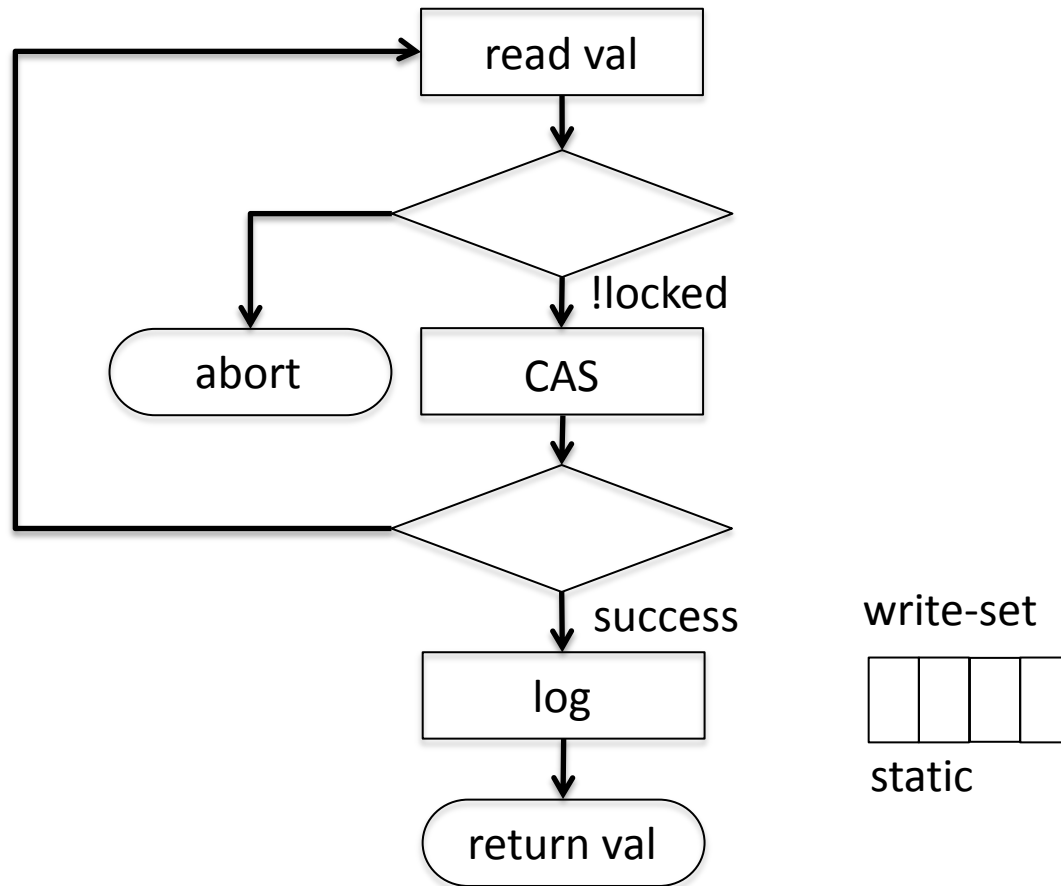
Pure value-based short read-write



Pure value-based short read-write

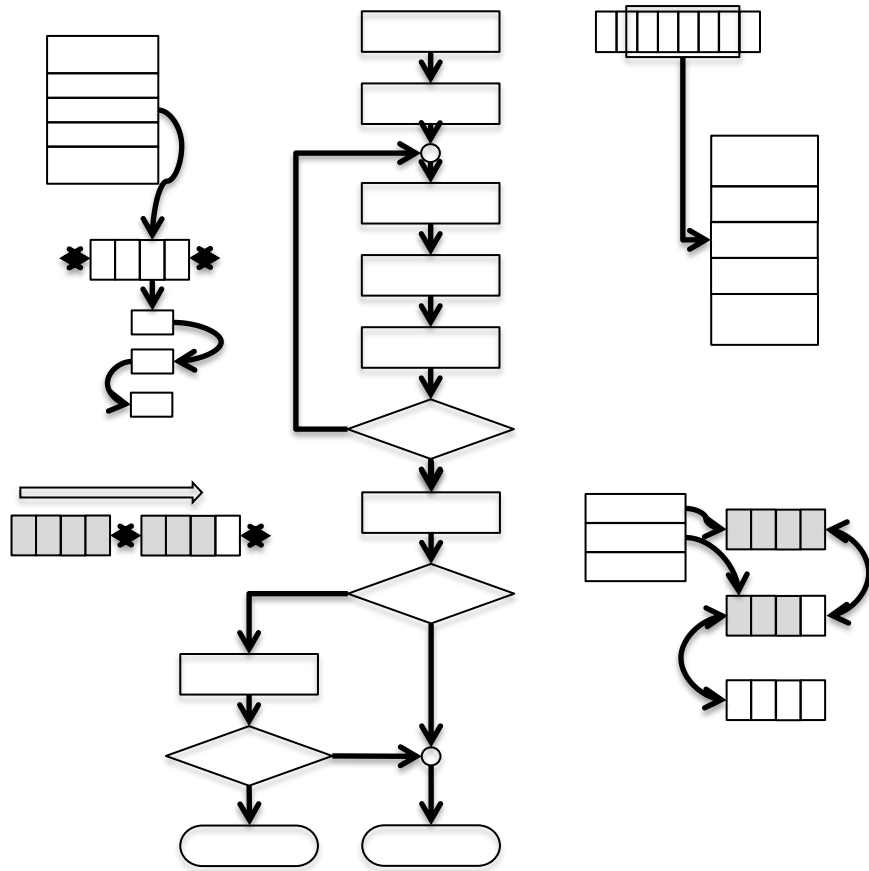


Pure value-based short read-write

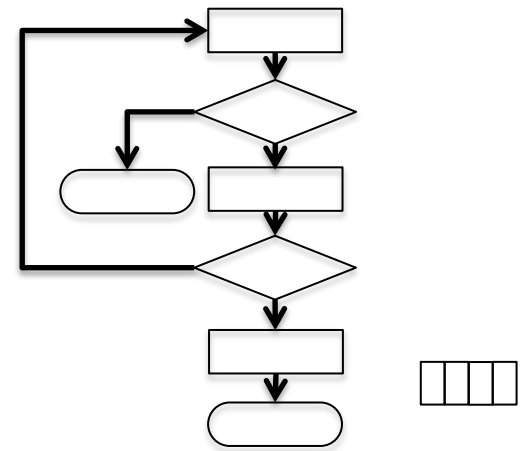


STM vs. SpecTM

STM



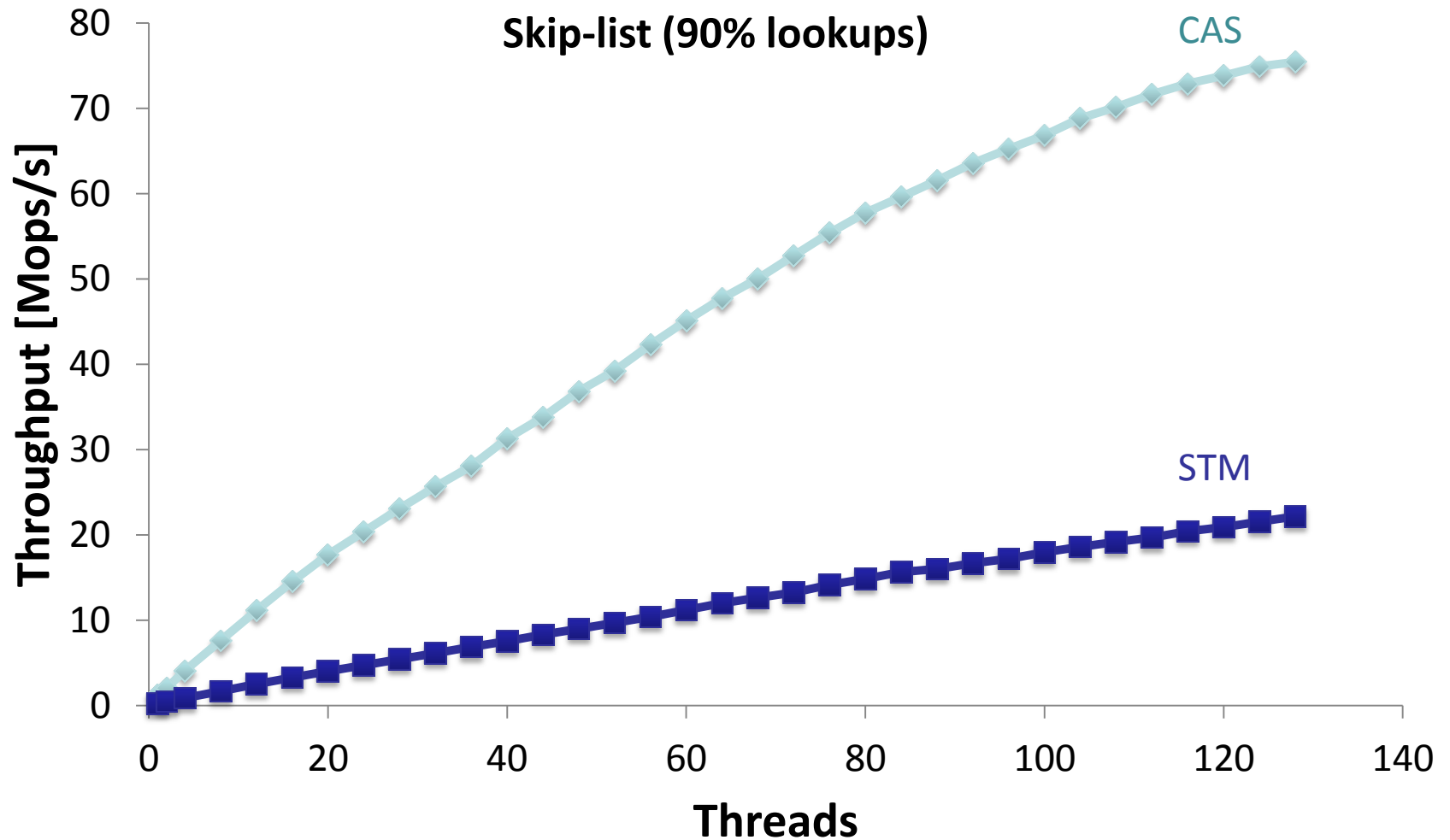
SpecTM



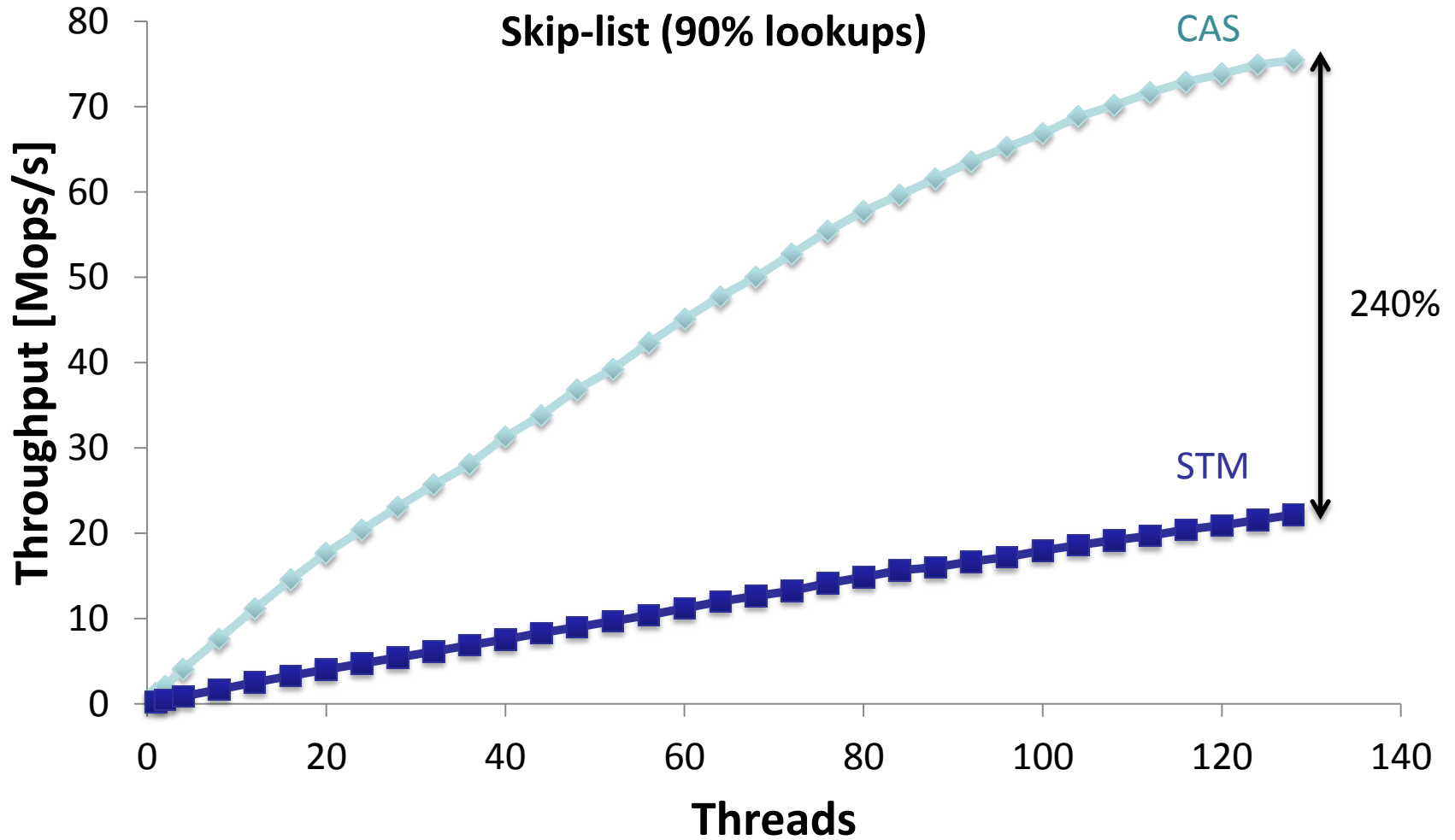
Overview

- STM overheads
- SpecTM
 - Short-transaction API
 - Collocating data and meta-data
 - In-word meta-data
- Evaluation

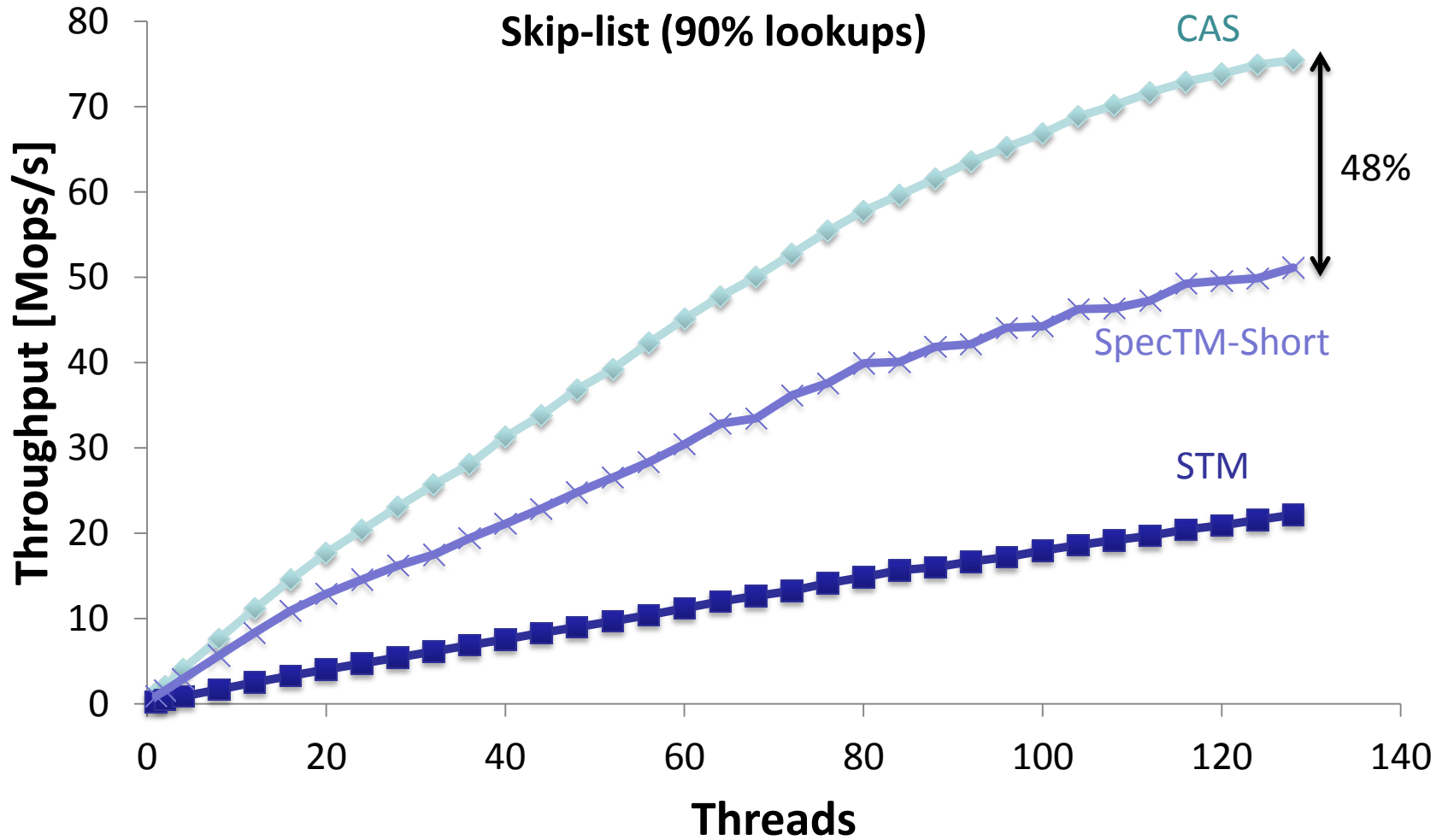
Evaluation



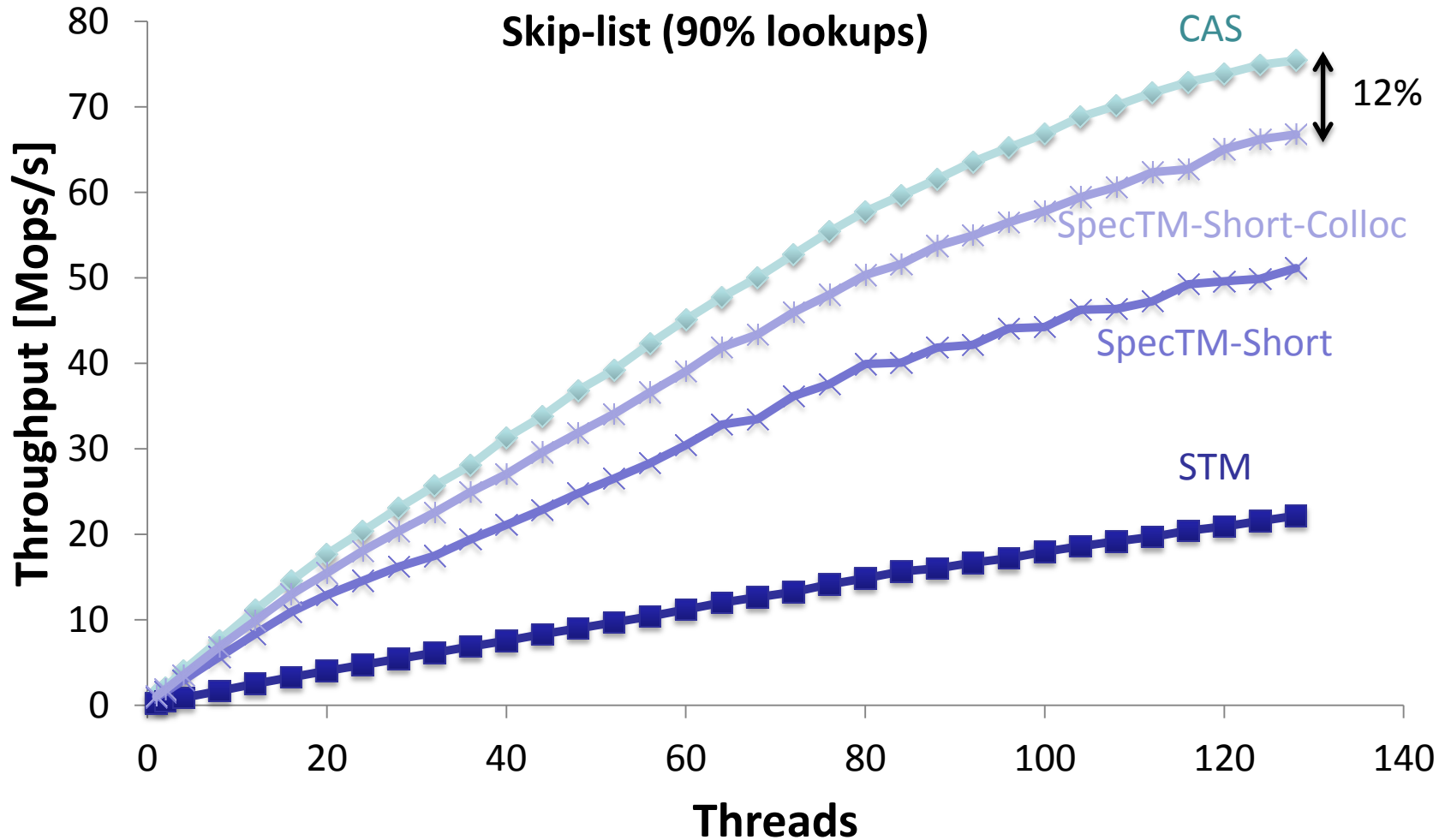
Evaluation



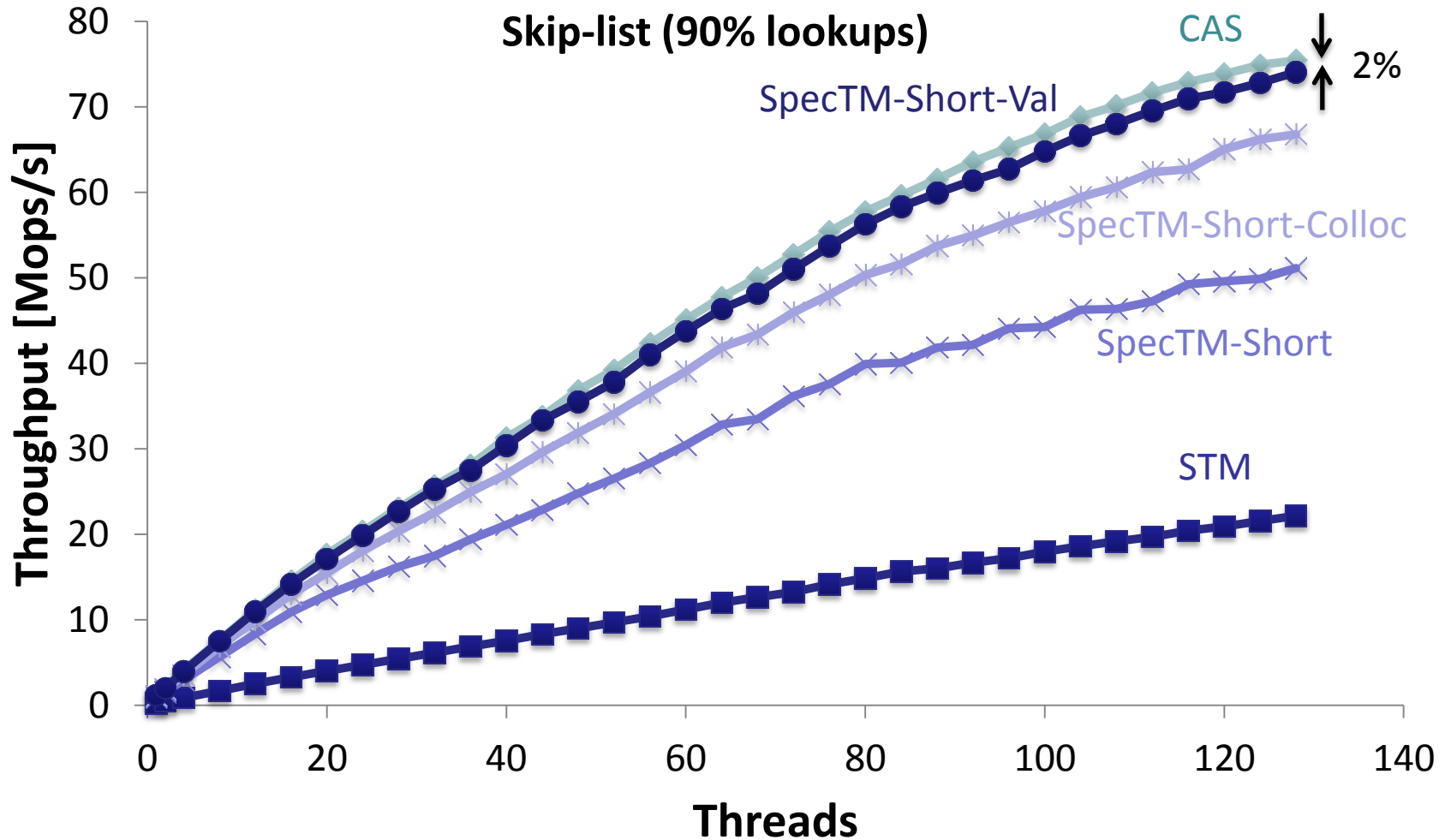
Evaluation



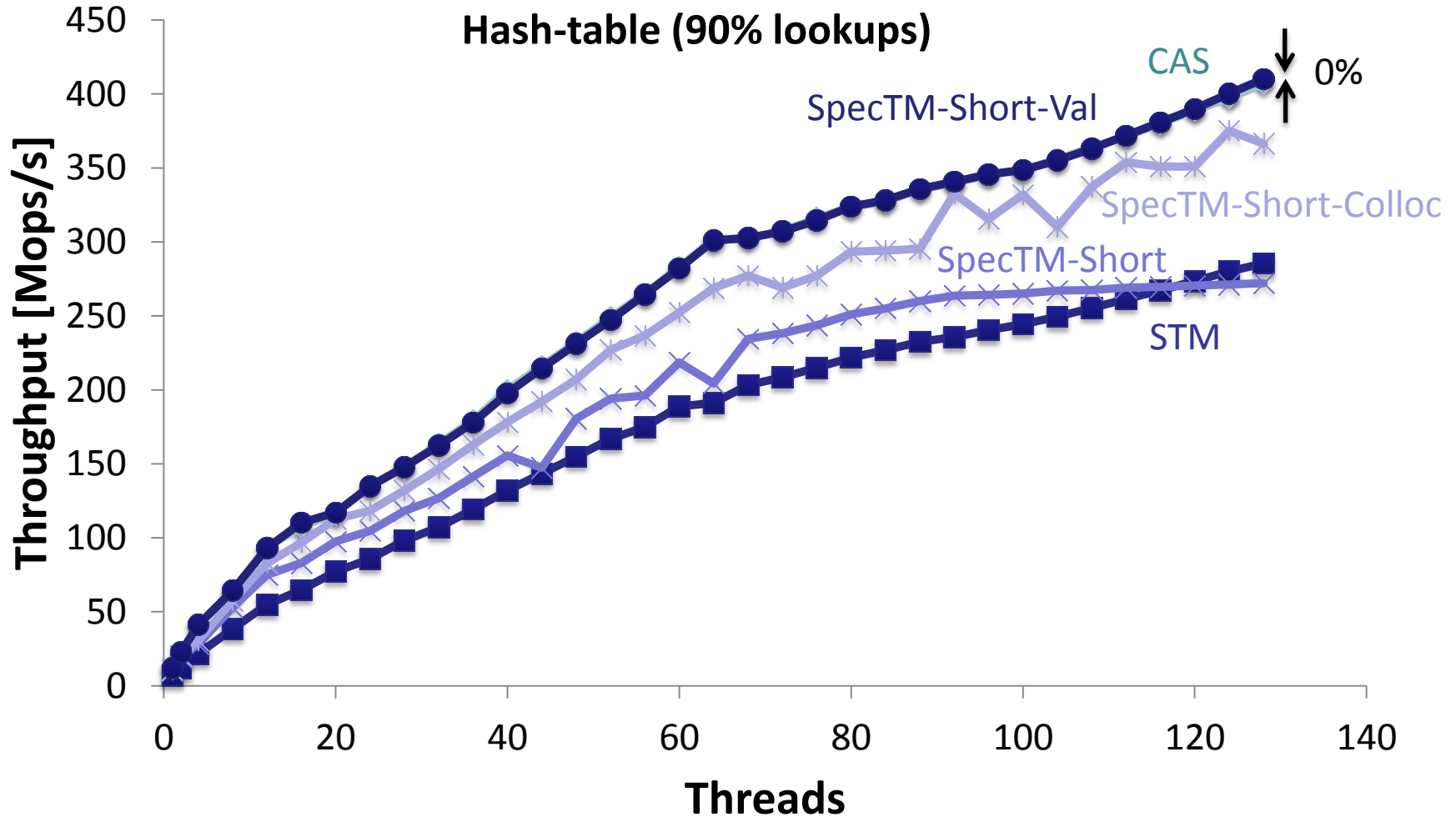
Evaluation



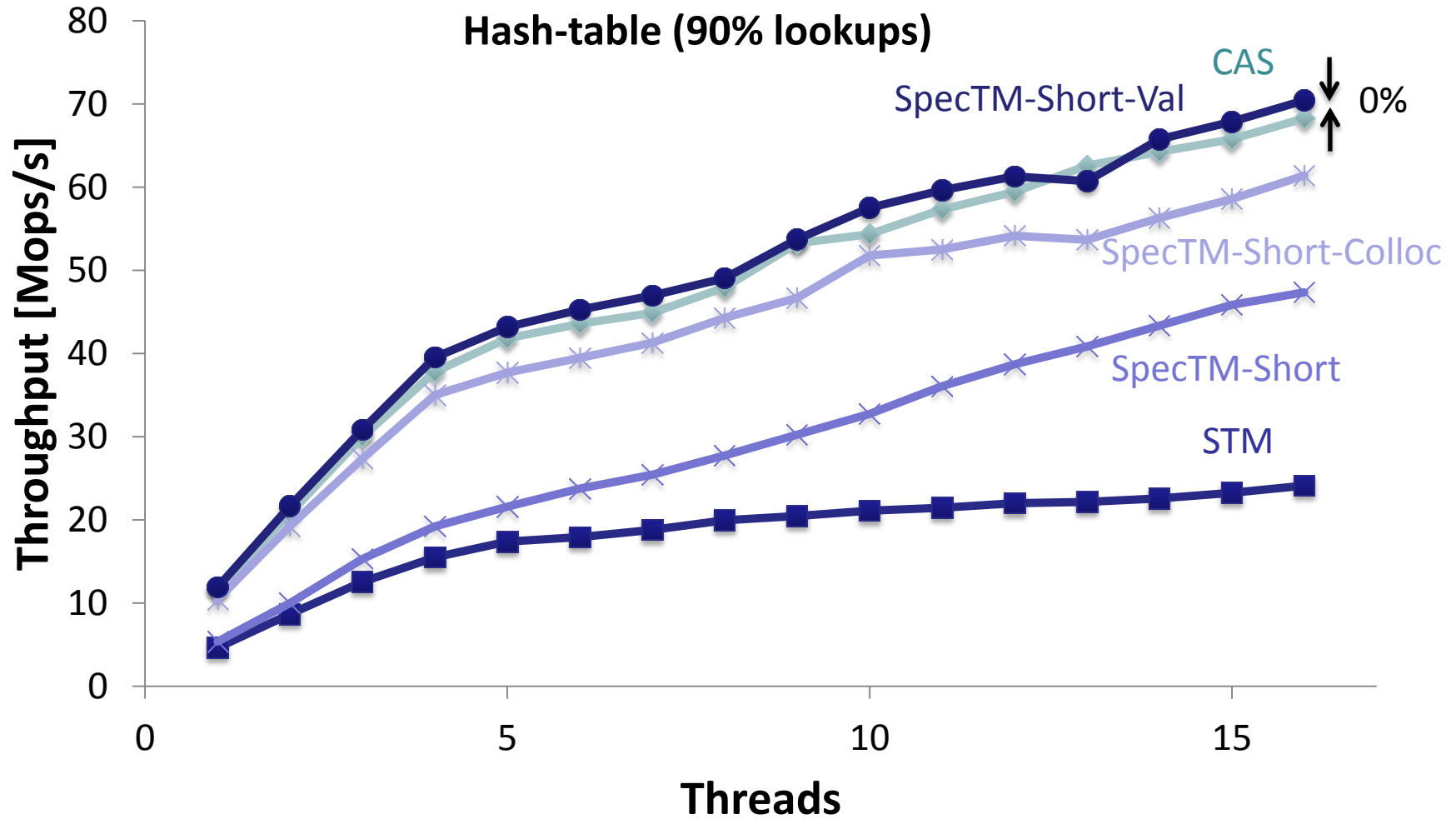
Evaluation



Evaluation



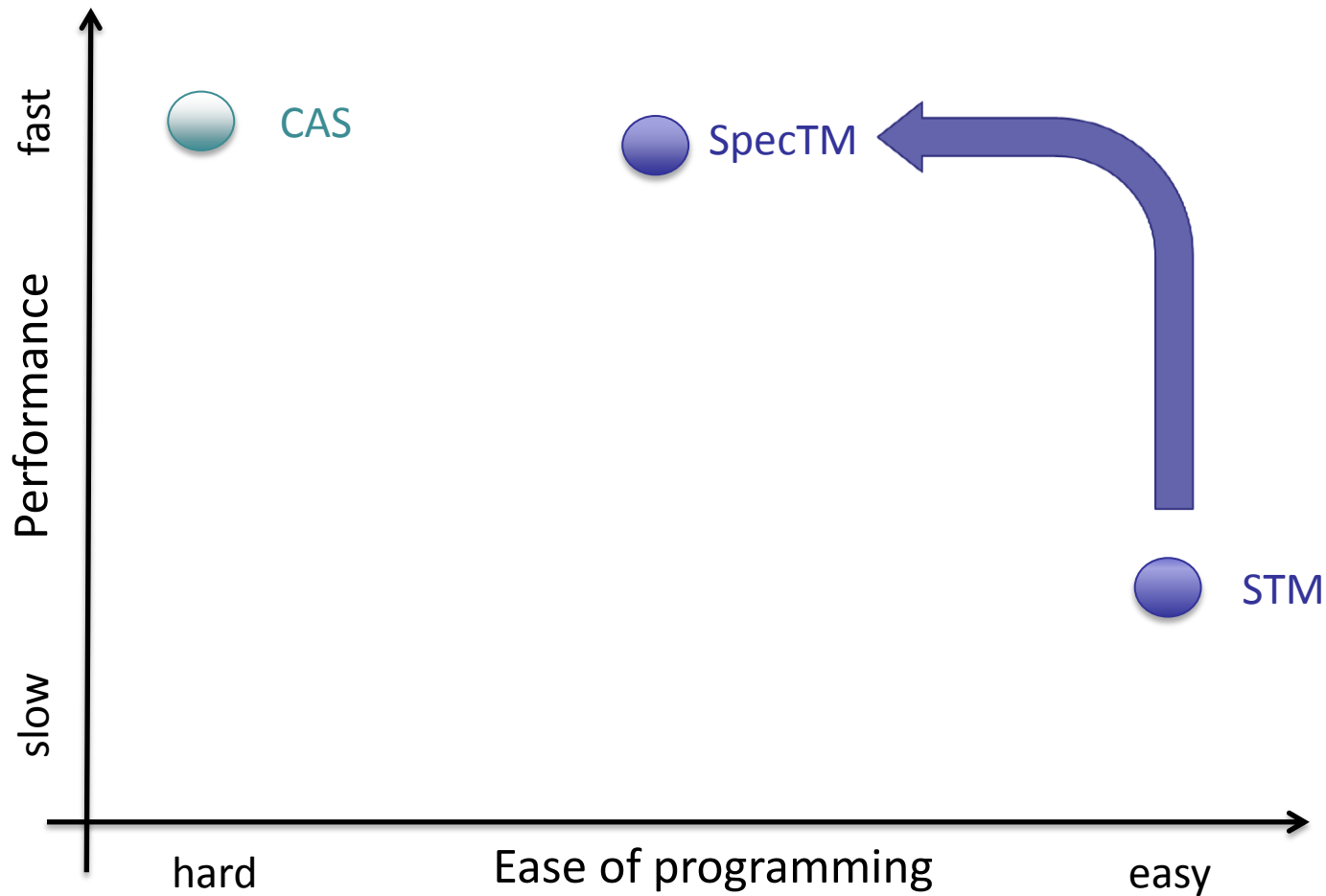
Evaluation



Conclusions

- Trade generality for performance in SpecTM
 - Restrict STM API
 - Control data layout
- STM-based data structures perform as well as the ones built directly from CAS
- Do we need SpecTM with upcoming Intel TSX?
 - Best-effort limited-size transactions

SpecTM: Specialized STM for Concurrent Data Structures



The Microsoft logo is displayed in a bold, italicized, black sans-serif font. The word "Microsoft" is followed by a registered trademark symbol (®). The background features a blue-to-white gradient with abstract, glowing white circular patterns that create a sense of motion and depth.

Microsoft[®]