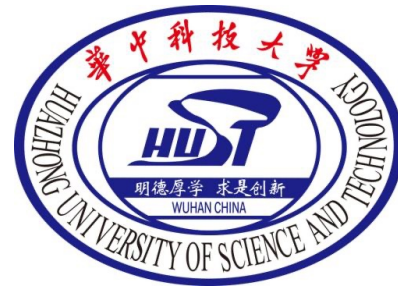


# FFCCD: Fence-Free Crash-Consistent Concurrent Defragmentation for Persistent Memory

Yuanchao Xu, Chencheng Ye, Yan Solihin, Xipeng Shen



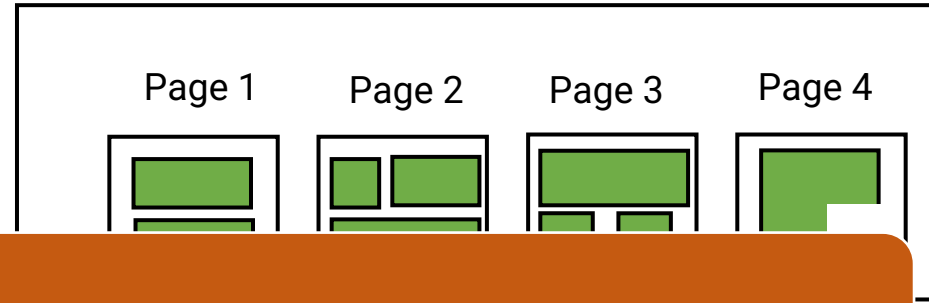
UNIVERSITY OF  
CENTRAL FLORIDA

# Persistent Fragmentation



Intel Persistent Memory

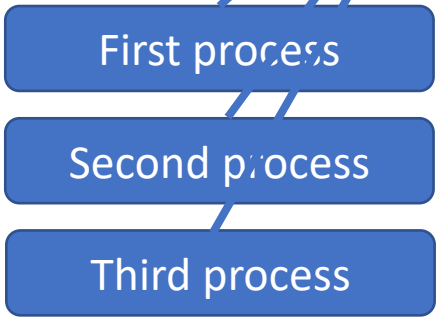
Persistent Memory Object Pool



Persistent fragmentation gets worsen with subsequent usages

Fragmentation

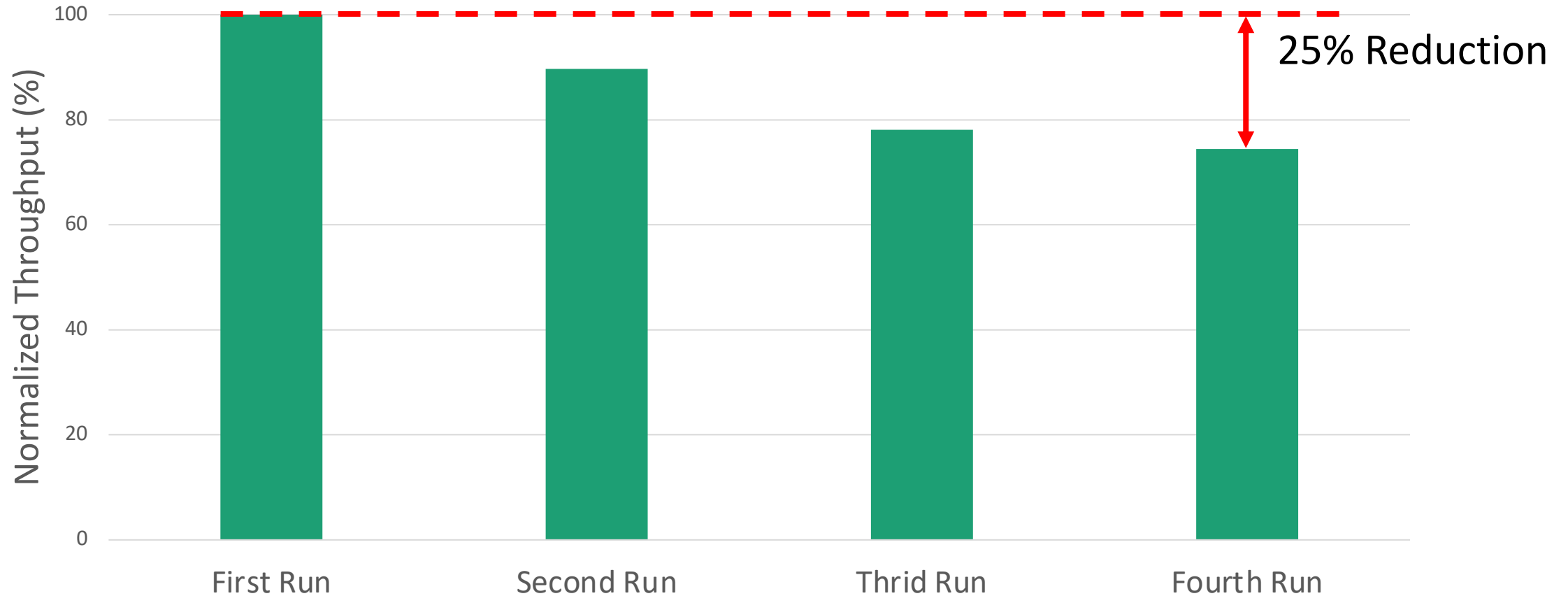
- Higher Density
- Byte-addressable Persistence
- DRAM-like Performance



 Live object

# Performance Degradation from Persistent Fragmentation

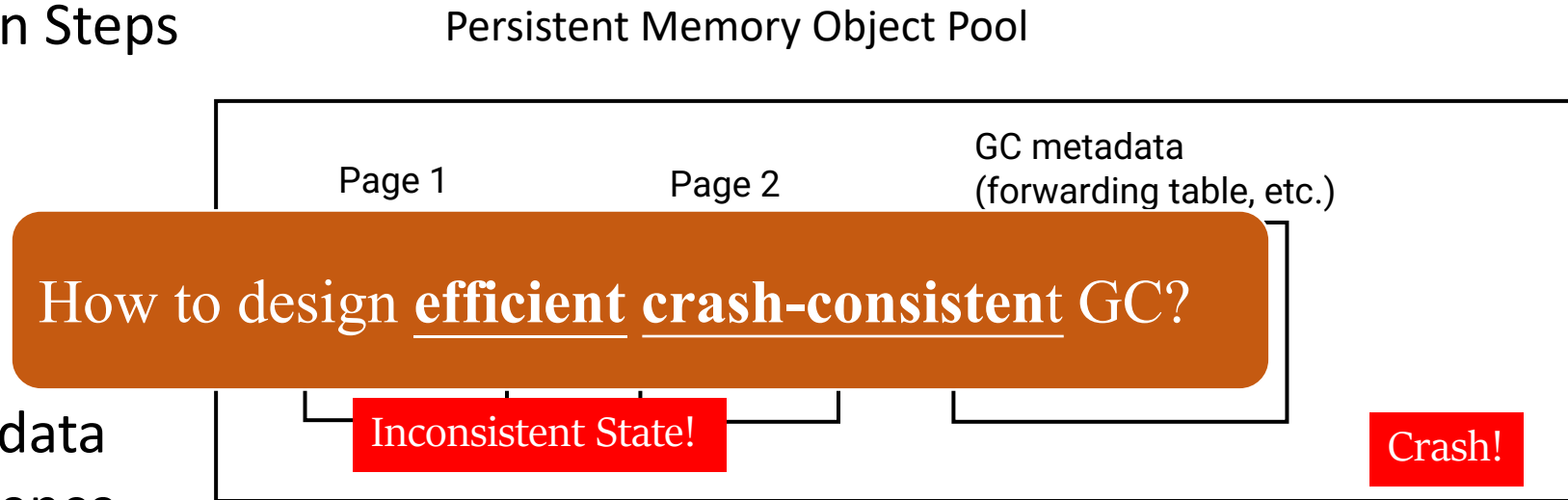
Redis throughput for 2M insertions and 1M deletions



# Garbage Collection (GC)

## Garbage Collection Steps

- 1) Marking
- 2) Summary
- 3) Compaction
  - Relocation
  - Update metadata
  - Update Reference



# Contributions

- First to analyze **PM fragmentation** systematically
  - Identify **sfences** as the key performance bottleneck
  - Analyze **post-crash states** to explore efficiency opportunities
- Design multiple concurrent GC solutions
  - Pure software **single-fence crash-consistent** design (SFCCD)
  - Architectural support for **fence-free crash-consistent** design (FFCCD)
- We evaluate designs and show its effectiveness

# Agenda

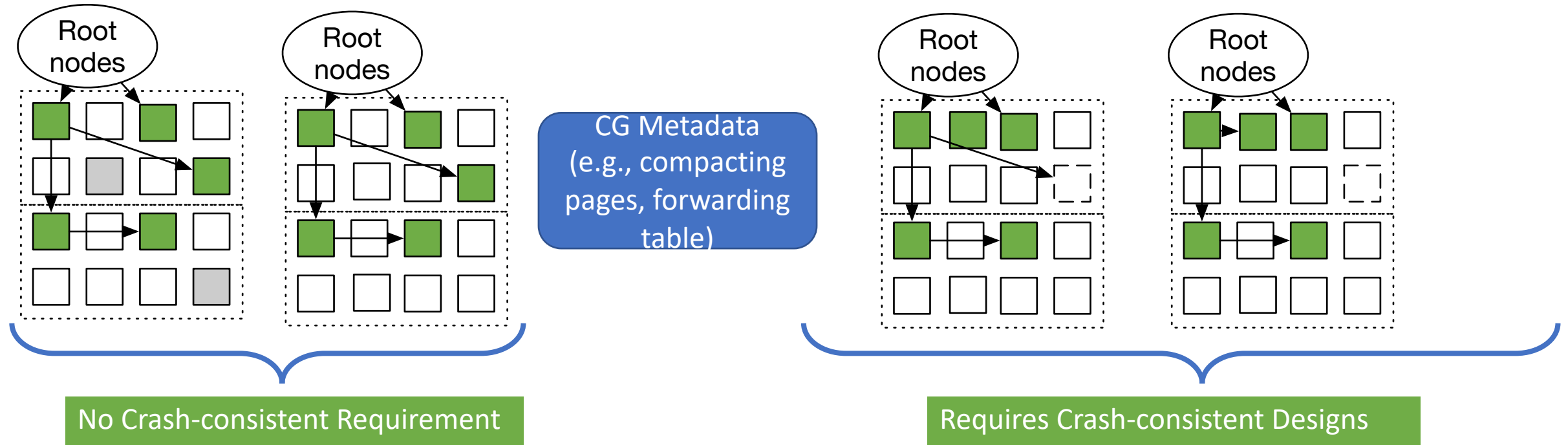
- Motivation
- Background
- Post-crash state exploration
- Architecture design
- Evaluation

# Concurrent Garbage Collection

Step 1: Marking

Step 2: Summary & Sweep

Step 3: Compaction



Free blocks Allocated object Reachable object One Page Moved object

# Concurrent Compaction: Read Barrier

```
p->first = input //p points to object A
```

```
<read barrier>
```

- 1) **Check** A is in compacting page
- 2) **Lookup** new address
- 3) **Check** A is not moved
  - 4) **Memcpy** A to the new address
  - 5) **Update** moved[A]=1
- 6) **Update** reference p to the new address



GC Metadata  
(e.g., compacting  
pages, forwarding  
table)



# Crash-consistent Concurrent Compaction Baseline[1]

```
p->first = input //p points to object A
```

```
<read barrier>
```

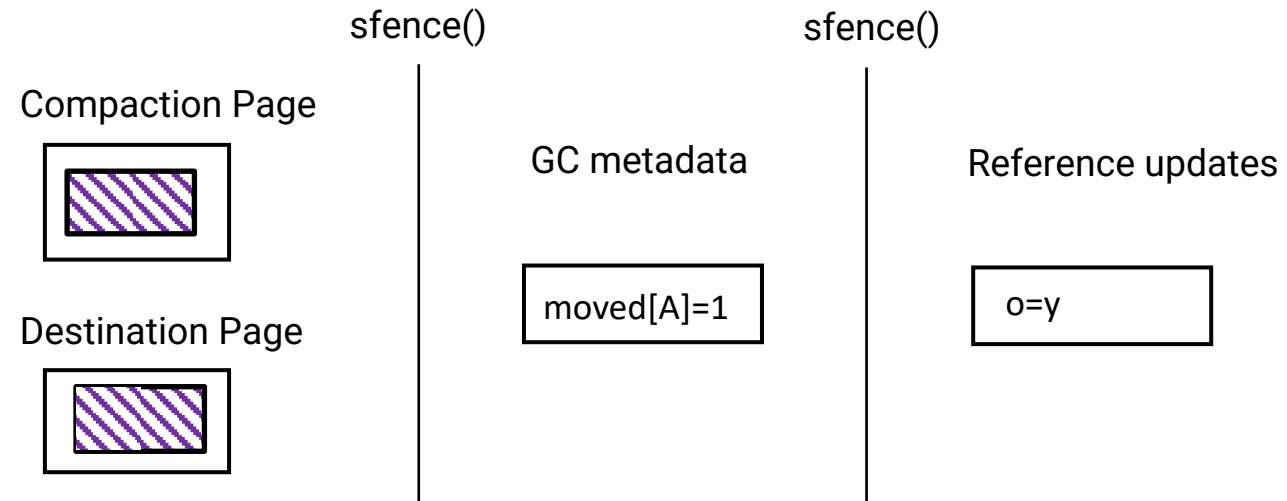
- 1) **Check** A is in compacting page
- 2) **Lookup** new address
- 3) **Check** A is not moved
- 4) **Memcpy** A to the new address
- 5) **Update** moved[A]=1
- 6) **Update** reference p to the new address

Persist GC metadata before compaction  
(e.g., compacting pages, forwarding table)

Crash-consistency relocation incurs about 50% overhead!

# Baseline Design

```
1 memcpy_nodrain(y,x,sizeof(A))
2 sfence();
3 moved[A]=1;
4 clwb(moved[A]);sfence();
5 o=y;
```



Recovery:

`moved[A]=0?`

`moved[A]=1? (no action)`

redo



Live object



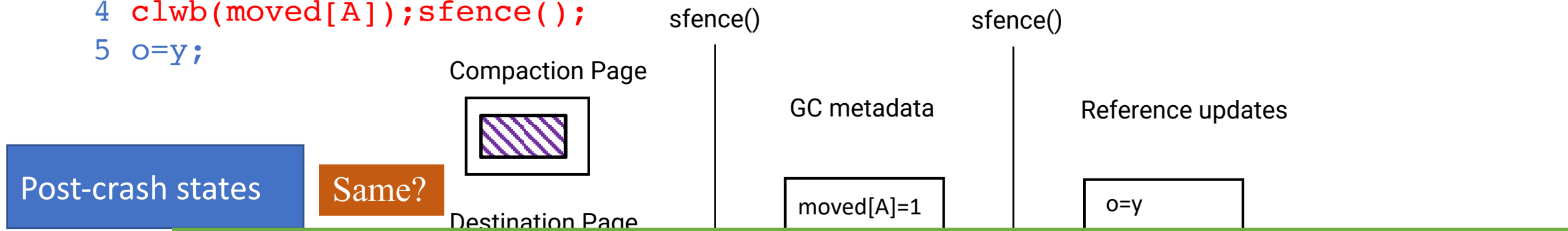
Cacheline

# Can we remove the first sfence?

```

1 memcpy_nodrain(y,x,sizeof(A))
2 sfence();
3 moved[A]=1;
4 clwb(moved[A]); sfence();
5 o=y;

```



Single-fence solution  
 Recovery: compare values with memcpy source for objects without reference update

Recovery:	Comparison Results	Recovery steps	moved[A]=1? (no action)	reference is updated? (no action)
	Same	set moved[A]=1		
	Partially Same	memcpy others & set moved[A]=1		
	Differ	unset moved[A]=1		

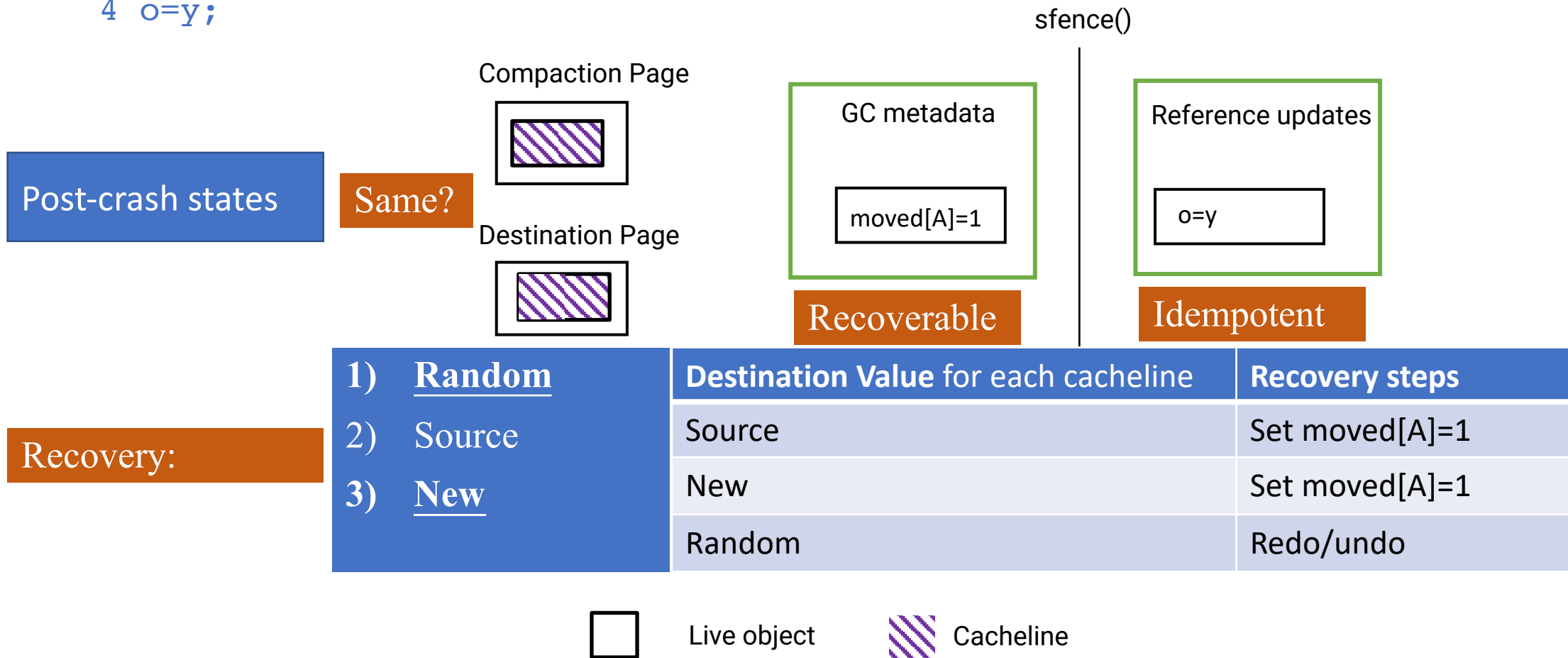


# Can we remove the second sfence?

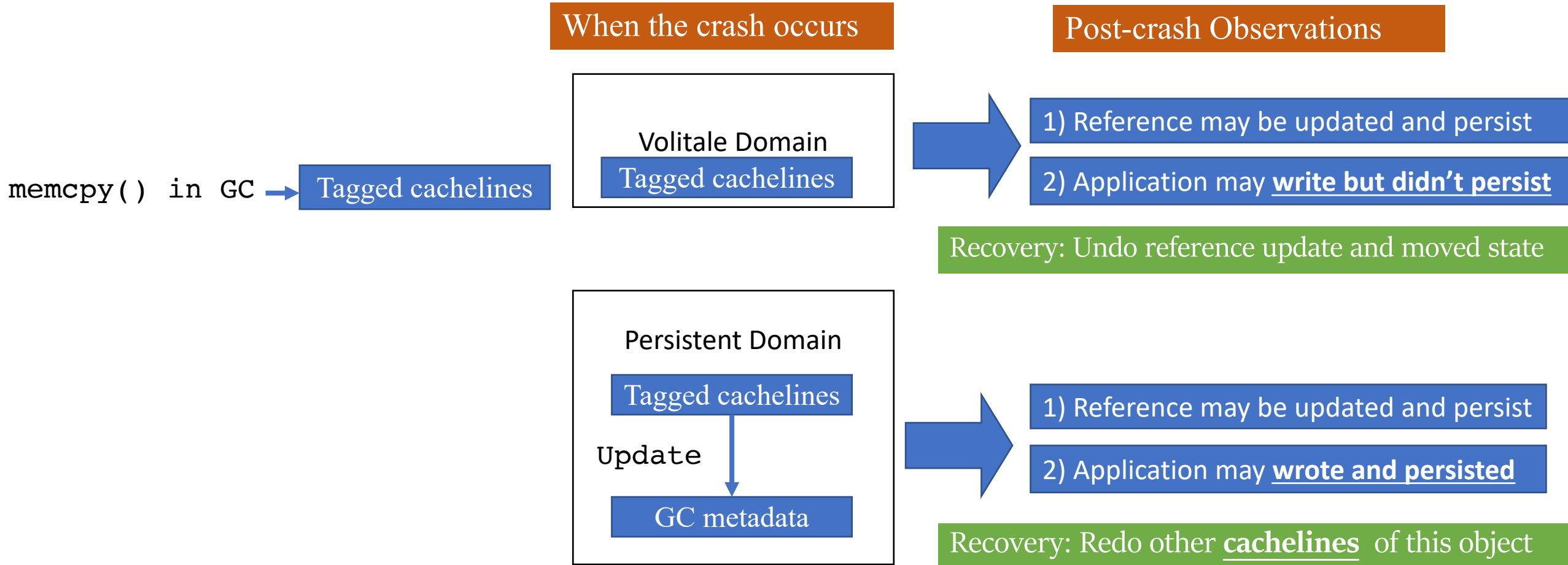
```

1 memcpy_nodrain(y,x,sizeof(A))
2 moved[A]=1;
3 clwb(moved[A]);sfence();
4 o=y;

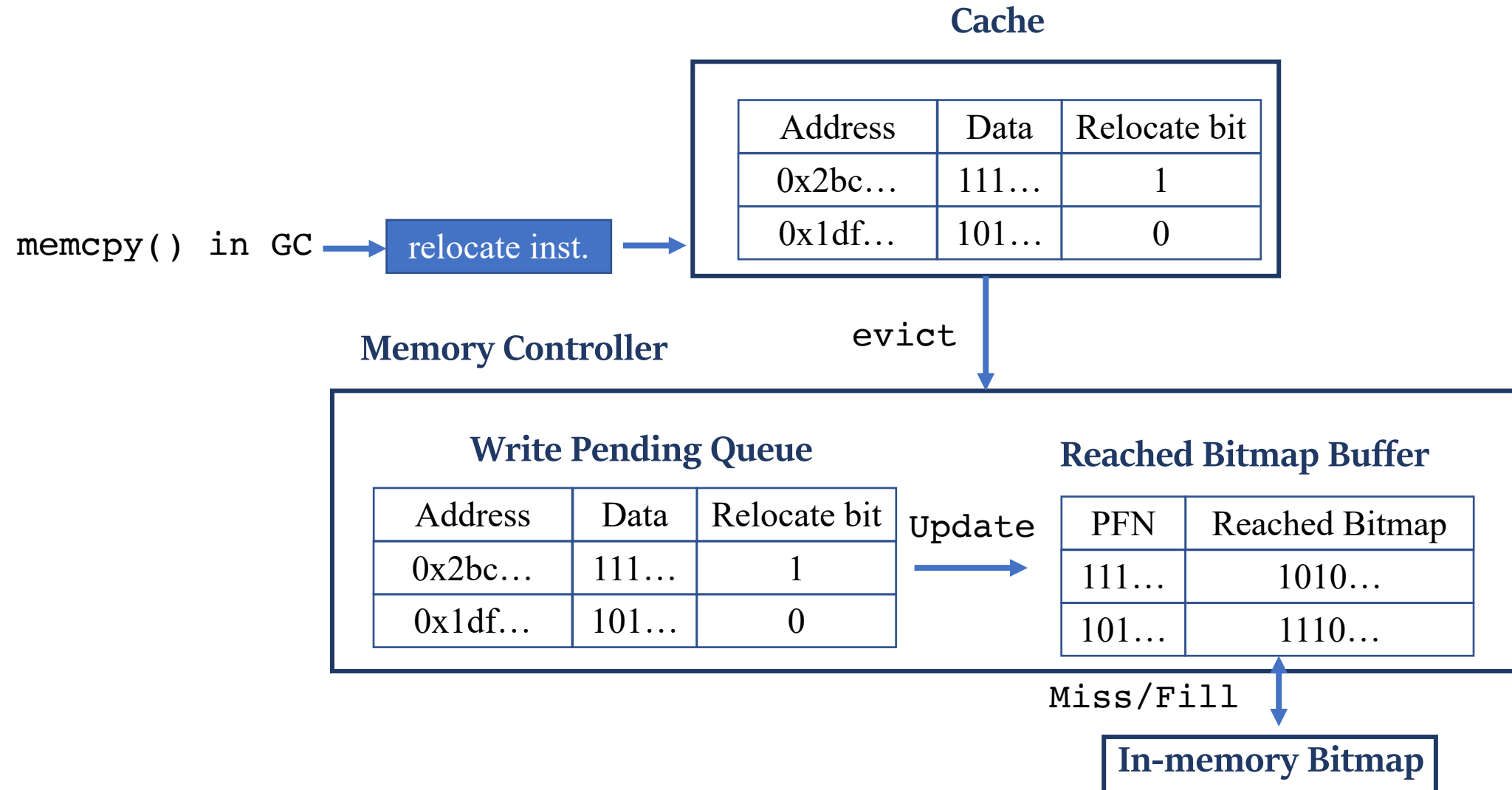
```



# Idea: Track Relocated Cachelines



# Architecture support



# Crash-consistent Concurrent Compaction Baseline

```
p->first = input //p points to object A
```

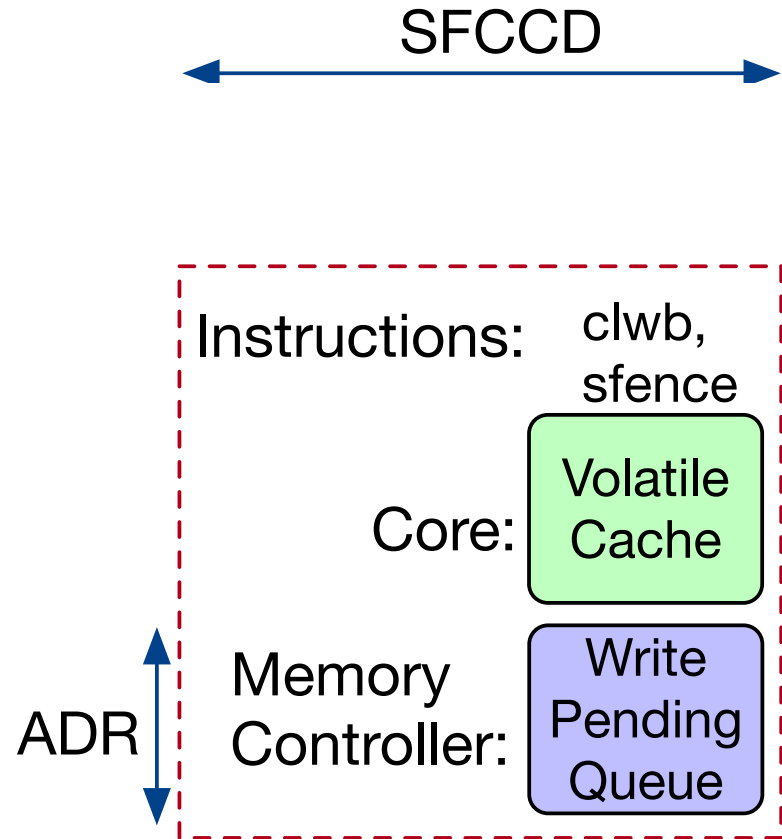
```
<read barrier>
```

- 1) **Check** A is in compacting page
- 2) **Lookup** new address
- 3) **Check** A is not moved
- 4) **Memcpy** A to the new address
- 5) **Update** moved[A]=1  
`clwb(moved[A]);sfence()`
- 6) **Update** reference p to the new address

More details in the paper

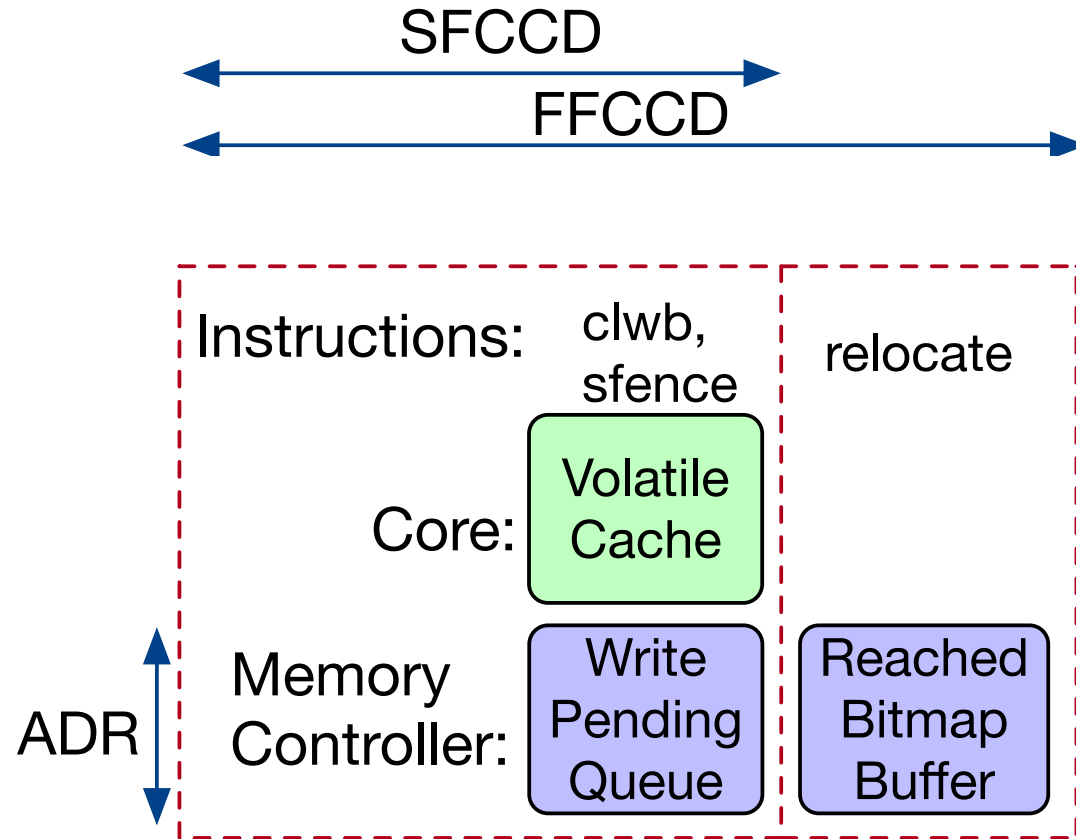
Fence-free design

# Design Summary

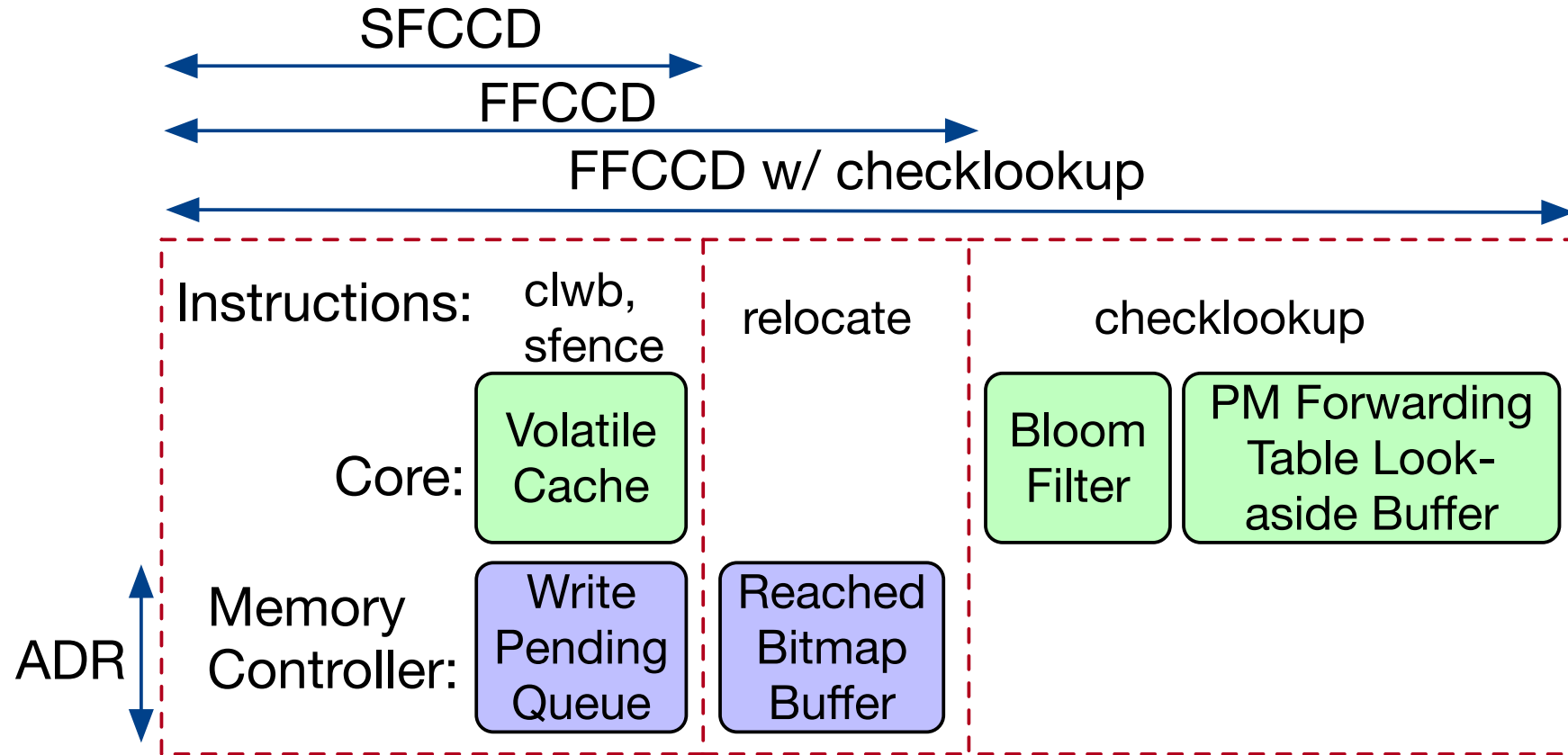




# Design Summary



# Design Summary



# Evaluation

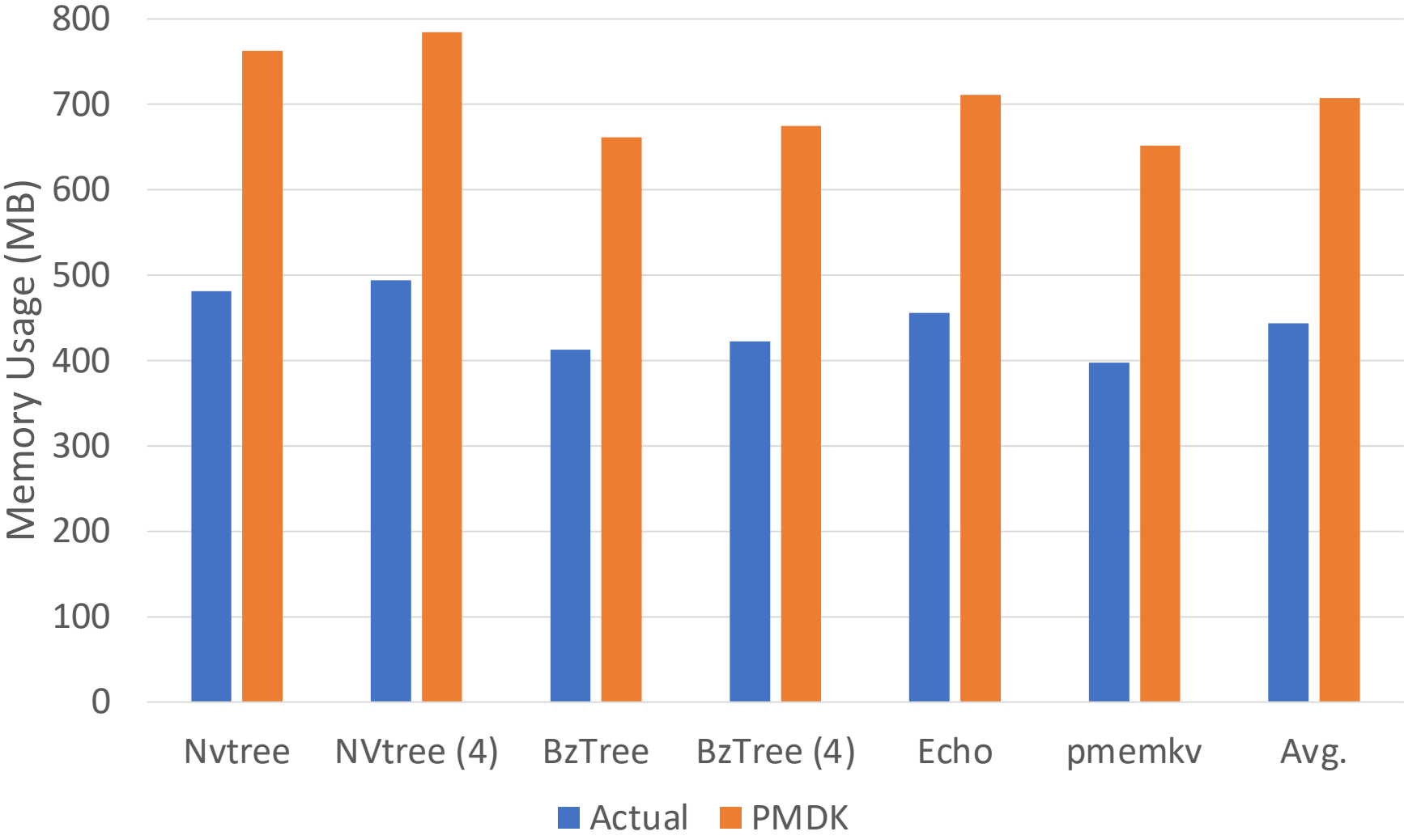
Benchmarks:

- PM data structures
- PM key-value stores

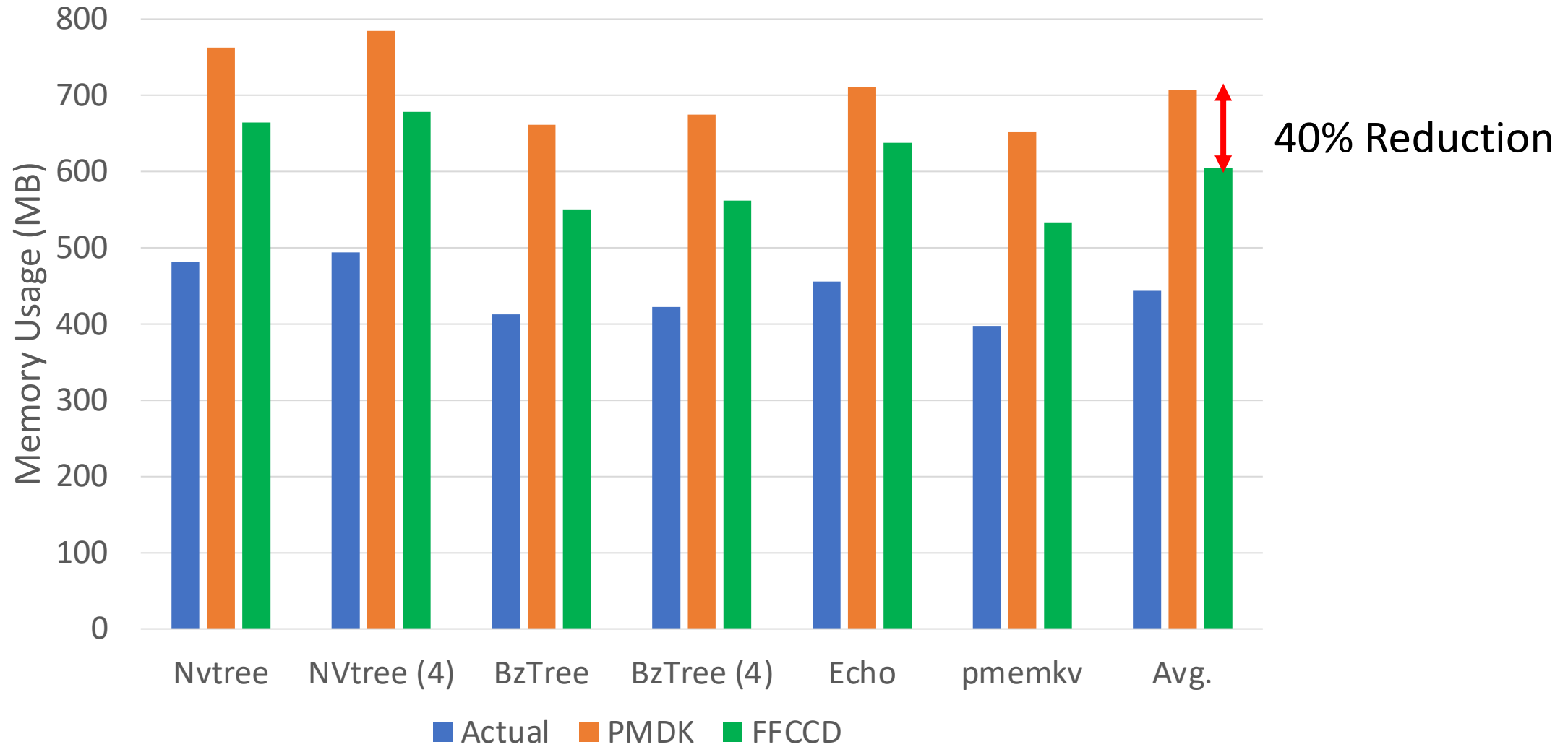
Workload:

- 1) Initialize data with 5M insertions
- 2) Delete 4M nodes
- 3) Insert 4M nodes
- 4) Delete 4M nodes

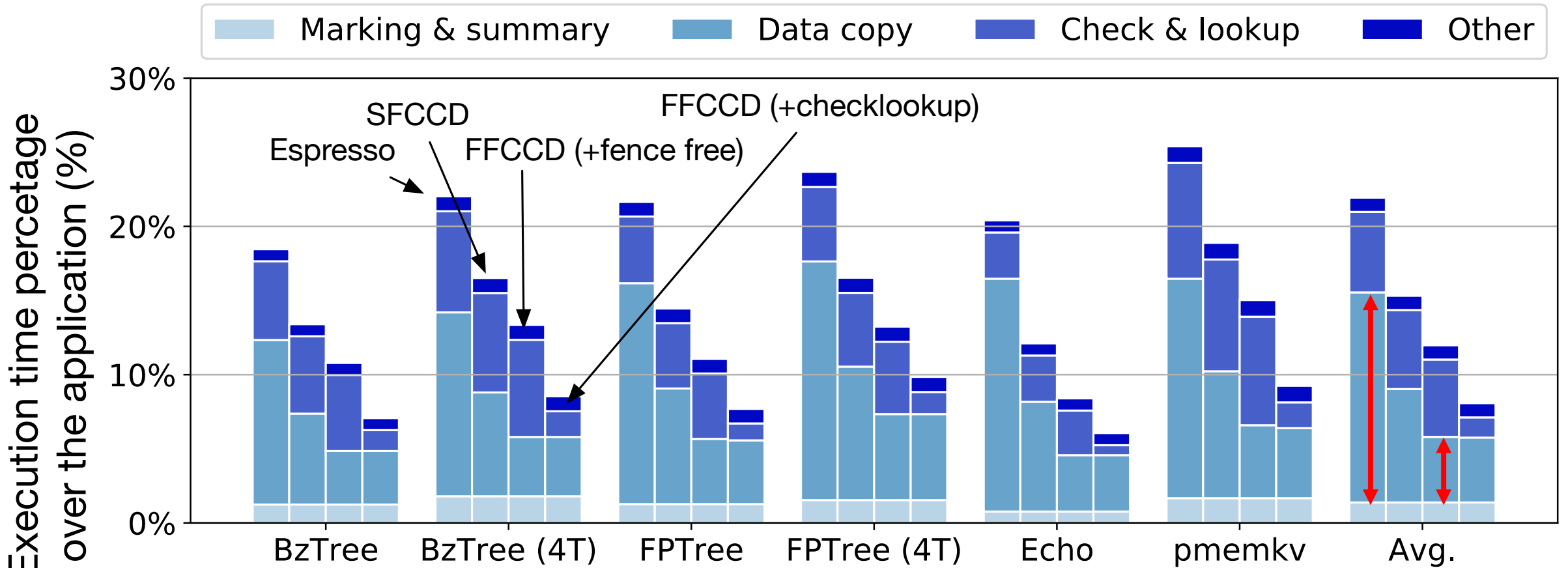
# Defragmentation



# Defragmentation

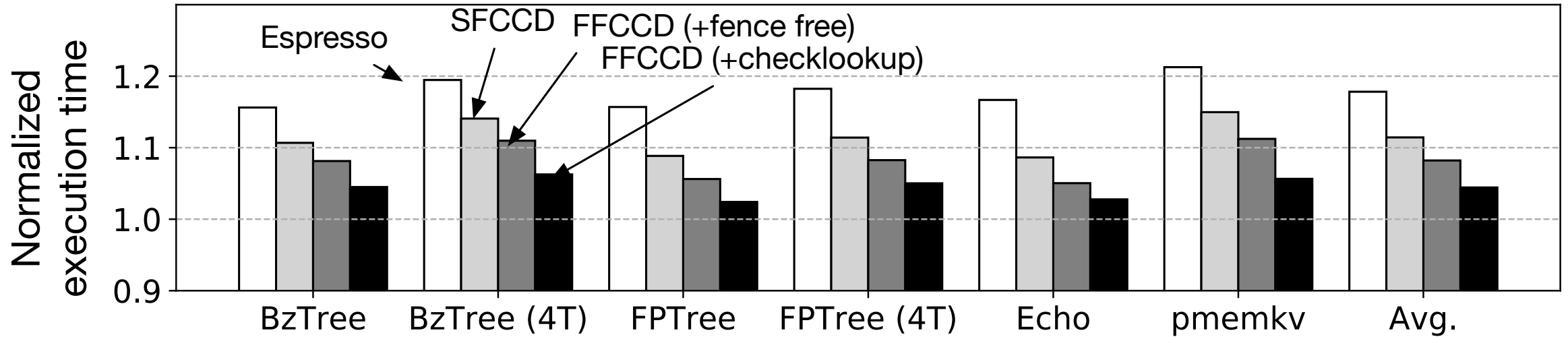


# GC Overhead Breakdown



FFCCD reduces 62% overhead in crash-consistent data copy

# Total Execution Time with GC



FFCCD incurs 4.1% overhead than non-defragmentation performance

# Conclusions

- We identify sfences as the key challenge to design efficient crash-consistent GC
- We design multiple solutions
  - SFCCD, software-only solution to remove 1 sfence
  - FFCCD, architectural-supportive solution to remove 2 sfences
- FFCCD provides 28–73% fragmentation reduction
- FFCCD incurs 4.1% overhead
  - Improve application performance due to better locality
  - Low overhead from GC

Thank you!  
Q&A