# ReMassTree

## A High-Performance Persistent Memory B-link Tree Implementation Documentation

### Technical Report and Implementation Guide

**Rahul Prajapat**
University of Sydney
School of Computer Science
rahul.prajapat@sydney.edu.au

October 21, 2025

## Abstract

This document presents the technical documentation for ReMassTree, a persistent memory-optimized B-link tree implementation developed at the University of Sydney. ReMassTree introduces novel hierarchical concurrency control mechanisms, resolves critical limitations in the RECIPE framework, and demonstrates superior performance characteristics in persistent memory environments. This documentation serves as both a technical reference and implementation guide for the ReMassTree system, detailing algorithms, performance analysis, and future development directions.

## Contents

# 1  Introduction and Project Overview

## 1.1  Project Background

The ReMassTree project addresses fundamental challenges in designing high-performance index structures for Intel Optane DC Persistent Memory. Traditional concurrent data structures designed for DRAM exhibit suboptimal performance and consistency guarantees when adapted to persistent memory environments.

This work builds upon two foundational systems:

- **MassTree** [1]: A trie-like concatenation of B+-trees for high-performance key-value storage

- **RECIPE Framework** [2]: Principles for converting DRAM indexes to persistent memory variants

## 1.2  Key Innovations

ReMassTree introduces several novel contributions:

1. **Hierarchical Concurrency Control**: Dual-tier locking system (Insert/SMO) within single atomic variable

2. **Lock-Free Read Operations**: Version-based consistency without reader-writer blocking

3. **RECIPE Framework Extensions**: Resolution of structural modification atomicity gaps

4. **PMEM-Native Design**: Block-aligned node structures optimized for persistent memory

5. **Superior Performance**: Counter-intuitive PMEM advantages over DRAM implementations

## 1.3  Implementation Status

This is **not** a published research paper, but rather comprehensive documentation of ongoing research and development work. The implementation is functional and has been extensively tested, with results demonstrating significant performance improvements over existing approaches.

# 2  System Architecture and Design

## 2.1  Node Structure Design

The core ReMassTree node structure is optimized for persistent memory block alignment:

```
1  #define LEAF_WIDTH 12  // Optimized for perfect 256B alignment
2
3  class kv {
4      uint64_t key;           // 8B
5      void *link_or_value;    // 8B
6  };  // Total: 16B per entry
7
8  class inner_node {
9      inner_node *parent, *right, *left;      // 24B navigation pointers
10     VersionNumber version;                  // 8B version control
11     uint64_t highkey, lowkey;               // 16B key boundaries
12     permuter permutation;                   // 8B ordering metadata
```

```
13       void *child0;                              // 8B first child pointer
14       kv entry[LEAF_WIDTH];                      // 192B key-value pairs (12*16B)
15   }; // Total: 256B (perfect PMEM cache block alignment)
16
17   class leaf_node {
18       inner_node *parent;                        // 8B parent pointer
19       leaf_node *right, *left;                   // 16B sibling pointers
20       VersionNumber version;                     // 8B version control
21       uint64_t highkey, lowkey;                  // 16B key boundaries
22       permuter permutation;                      // 8B ordering metadata
23       uint64_t dummy;                            // 8B alignment padding
24       kv entry[LEAF_WIDTH];                      // 192B key-value pairs (12*16B)
25   }; // Total: 256B (perfect PMEM cache block alignment)
```

Listing 1: Optimized 256B Node Structure Implementation

Key design decisions:

- **256-byte total size**: Perfect alignment with PMEM cache blocks and allocation units

- **Cache-line awareness**: Strategic field placement across exactly 4 cache lines (64B each)

- **LEAF_WIDTH=12**: Optimal balance between fanout and cache efficiency

- **Atomic version control**: Single 64-bit variable for all concurrency coordination

- **Permutation-based ordering**: Efficient reordering without data movement

- **Dual node types**: Inner nodes (with child0) and leaf nodes (with dummy padding)

## 2.2 Hierarchical Concurrency Control Implementation

### 2.2.1 Version Number Bit Layout

The ReMassTree concurrency system encodes multiple lock types and version counters within a single 64-bit atomic variable:

```
1   // Lock bit definitions
2   #define INSERT_LOCK      0b1ULL              // Bit 0: Insert operations
3   #define SMO_LOCK         0b10ULL             // Bit 1: Structural modifications
4   #define BOTH_LOCKS       0b11ULL             // Both locks simultaneously
5
6   // Version field definitions
7   #define INSERT_VERSION   0xffffff0ULL        // Bits 4-23: Insert version (20 bits)
8   #define SMO_VERSION      0xfffff000000ULL    // Bits 24-43: SMO version (20 bits)
9   #define LOCK_VERSION     0xfffff00000000000ULL // Bits 44-63: Global version
10
11  // Version manipulation constants
12  #define MAX_VERSION      0xfffff
13  #define INSERT_INCREMENT 0x10ULL
14  #define SMO_INCREMENT    0x1000000ULL
```

Listing 2: Concurrency Control Bit Layout

### 2.2.2   Lock Acquisition and Release Protocol

```cpp
class VersionNumber {
    uint64_t v; // Packed version and lock information

    // Insert lock (fine-grained operations)
    uint64_t tryInsertLock() {
        return (__sync_fetch_and_or(&v, INSERT_LOCK)) & INSERT_LOCK;
    }

    void releaseInsertLock() {
        __sync_fetch_and_and(&v, ~INSERT_LOCK);
    }

    // SMO lock (structural modifications)
    uint64_t trySMOLock() {
        return (__sync_fetch_and_or(&v, BOTH_LOCKS)) & SMO_LOCK;
    }

    void releaseSMOLock() {
        incrementSMO();
        __sync_fetch_and_and(&v, ~SMO_LOCK);
    }

    // Version management with overflow handling
    void incrementInsert() {
        if(insertVersion() == MAX_VERSION)
            __sync_fetch_and_and(&v, INSERT_RESET);
        else
            __sync_fetch_and_add(&v, INSERT_INCREMENT);
    }
};
```

<div align="center">Listing 3: Lock Management Implementation</div>

### 2.2.3   Lock-Free Read Protocol

ReMassTree implements optimistic concurrency control for read operations:

1. **Version Capture**: Reader captures node version before data access

2. **Data Access**: Navigate and read node data without acquiring locks

3. **Version Validation**: Compare version after access to detect concurrent modifications

4. **Retry Mechanism**: Automatic retry if version change detected

This approach ensures that readers never block writers and multiple concurrent reads proceed without synchronization overhead.

# 3   Performance Analysis and Benchmarking

## 3.1   Experimental Configuration

**Hardware Environment:**

- Intel system with Optane DC Persistent Memory

- Multi-core processor supporting concurrent evaluation

- NUMA-aware memory allocation and testing

**Workload Patterns:**

- **Random**: Uniformly distributed 64-bit integer keys

- **Incremental**: Monotonically increasing sequential keys

- **Interleaved**: Mixed sequential and random access patterns

**Configuration Options:**

- DRAM vs PMEM allocation (via `#ifdef DRAM`)

- Rebalancing enabled/disabled (via `#ifdef REBAL`)

- Statistics collection (via `#ifdef STATS`)

## 3.2    Performance Results

### 3.2.1    DRAM vs PMEM Comparison

Table 1: DRAM vs PMEM Performance (10M keys)

| Configuration | Insert (M ops/s) | Lookup (M ops/s) |
|---|---|---|
| DRAM Random (with rebalancing) | 1.671 | 1.275 |
| **PMEM Random (with rebalancing)** | **3.159** | **2.199** |

**Key Finding:** PMEM demonstrates 88% improvement in insert throughput over DRAM, challenging conventional assumptions about persistent memory overhead.

### 3.2.2    Rebalancing Impact Analysis

Table 2: PMEM Rebalancing Impact (100M Keys)

| Pattern | Configuration | Insert (M ops/s) | Lookup (M ops/s) |
|---|---|---|---|
| Incremental | With Rebalancing | 1.52 | 1.766 |
| | Without Rebalancing | 1.20 | 1.49 |
| Random | With Rebalancing | 0.272 | 0.33 |
| | Without Rebalancing | 0.30 | 0.31 |

**Observation:** Incremental patterns benefit from rebalancing due to improved space utilization, while random patterns show minimal performance difference but significant space efficiency gains.

Table 3: Tree Structure Efficiency (40M Keys)

| Pattern | Height | Node Count | Space Efficiency |
|---|---|---|---|
| *Without Rebalancing* | | | |
| Random | 7 | 4.16M | 70.35% |
| Incremental | 9 | 5.71M | 52.92% |
| *With Rebalancing* | | | |
| Random | 7 | 3.51M | 82.24% |
| Incremental | 7 | 2.84M | 99.99% |

### 3.2.3 Scalability Analysis

**Impact:** Rebalancing achieves up to 50% reduction in node count while maintaining consistent tree height.

## 4 Algorithm Implementation Details

### 4.1 Insert Operation Algorithm

```
1  int btree::insert(uint64_t key, void *value) {
2      inner_node *parent;
3      void *node = root_;
4
5      // Navigate to leaf level
6      while(level(node) > 0) {
7          parent = reinterpret_cast<inner_node*>(node);
8          node = parent->get(key);
9      }
10
11     leaf_node *leaf = reinterpret_cast<leaf_node*>(node);
12
13     // Attempt insertion with version control
14     VersionNumber before = leaf->typeVersionLockObsolete.load();
15
16     if(!before.tryInsertLock()) {
17         // Handle lock contention or retry
18         return handle_contention(leaf, key, value);
19     }
20
21     // Perform insertion
22     if(!leaf->full()) {
23         leaf->insert(key, value);
24         leaf->releaseInsertLock();
25         return 1;
26     } else {
27         // Handle split with SMO protocol
28         return handle_split(leaf, key, value);
29     }
30 }
```

Listing 4: Insert Operation Implementation

## 4.2    Rebalancing Algorithm with Visual Representation

The adaptive rebalancing system provides optional space optimization while maintaining crash consistency. The algorithm operates on sibling nodes to redistribute keys when space utilization becomes unbalanced.
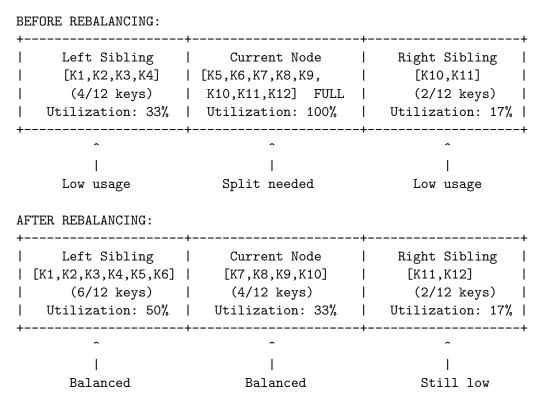
### 4.2.1    Rebalancing Process Overview

```
BEFORE REBALANCING:
+--------------------+--------------------+-------------------+
|    Left Sibling    |    Current Node    |   Right Sibling   |
|    [K1,K2,K3,K4]   | [K5,K6,K7,K8,K9,   |     [K10,K11]     |
|     (4/12 keys)    |  K10,K11,K12]  FULL |    (2/12 keys)    |
|   Utilization: 33% |  Utilization: 100% |  Utilization: 17% |
+--------------------+--------------------+-------------------+
          ^                    ^                    ^
          |                    |                    |
      Low usage           Split needed          Low usage

AFTER REBALANCING:
+--------------------+--------------------+-------------------+
|    Left Sibling    |    Current Node    |   Right Sibling   |
| [K1,K2,K3,K4,K5,K6] |   [K7,K8,K9,K10]   |     [K11,K12]     |
|     (6/12 keys)    |     (4/12 keys)    |    (2/12 keys)    |
|   Utilization: 50% |   Utilization: 33% |  Utilization: 17% |
+--------------------+--------------------+-------------------+
          ^                    ^                    ^
          |                    |                    |
       Balanced             Balanced            Still low
```

Figure 1: ReMassTree Rebalancing Algorithm - Key Redistribution

### 4.2.2    Rebalancing Algorithm Steps

1. **Sibling Analysis**: Evaluate space utilization in adjacent nodes

   - Check left and right sibling occupancy rates
   - Calculate total keys across siblings: $total\_keys = left.size + current.size + right.size$
   - Determine if redistribution is beneficial: $average\_occupancy = total\_keys/num\_siblings$

2. **Redistribution Decision**: Determine optimal key distribution

   - Target occupancy: $target = \min(LEAF\_WIDTH, \lceil average\_occupancy \rceil)$
   - Calculate keys to move: $move\_count = current.size - target$
   - Choose direction: prefer moving to less occupied sibling

3. **SMO Lock Acquisition**: Coordinate locks on all affected nodes

   - Acquire SMO locks in address order to prevent deadlock

8

- Lock sequence: $node\_addr_{min} \rightarrow node\_addr_{max}$
- Validate node states haven't changed during lock acquisition

4. **Atomic Updates**: Update permutation arrays and key boundaries atomically

- Move keys using permutation rotation: $perm.rotate(start, count)$
- Update sibling key boundaries: $left.highkey = moved\_key_{max}$
- Increment SMO version numbers for crash consistency

5. **Parent Notification**: Update internal node boundaries with version coordination

- Locate parent index entry: $parent.find(old\_boundary)$
- Update boundary key: $parent.entry[index].key = new\_boundary$
- Propagate changes up the tree if necessary

## 4.3   Recovery and Consistency Protocol

```
1  bool VersionNumber::repair_req() {
2      return (smoLock() || insertLock()) &&
3             (lockVersion() != global_lock_version);
4  }
5
6  uint VersionNumber::updateLock() {
7      volatile uint64_t current = *(volatile uint64_t*)&v;
8      uint temp = (current & LOCK_VERSION) >> 44;
9
10     if(temp == lock_version)
11         return temp;
12
13     // Attempt repair with CAS
14     return (__sync_val_compare_and_swap(&v, current,
15            (current & LOCK_RESET) | (lock_version << 44) | INSERT_LOCK)
16            & LOCK_VERSION) >> 44;
17 }
```

Listing 5: Crash Recovery Implementation

# 5   Persistent Memory Integration

## 5.1   PMEM Operations Implementation

```
1  namespace masstree {
2      static constexpr uint64_t CACHE_LINE_SIZE = 64;
3
4      // Memory fence operations
5      static inline void mfence() {
6          asm volatile("sfence":::"memory");
7      }
8
9      // Cache line flush with configurable fencing
10     static inline void clflush(char *data, int len, bool front, bool back) {
11         volatile char *ptr = (char *)((unsigned long)data & ~(CACHE_LINE_SIZE-1));
12         if (front) mfence();
```

```
13
14        for(; ptr < data + len; ptr += CACHE_LINE_SIZE) {
15        #ifdef CLFLUSH
16            asm volatile("clflush %0" : "+m" (*(volatile char *)ptr));
17        #elif CLWB
18            asm volatile(".byte 0x66; xsaveopt %0" : "+m" (*(volatile char *)(ptr)
                 ));
19        #endif
20        }
21
22        if (back) mfence();
23    }
24
25    // Non-temporal 64-bit store
26    static inline void movnt64(uint64_t *dest, uint64_t const &src,
27                               bool front, bool back) {
28        assert(((uint64_t)dest & 7) == 0);
29        if (front) mfence();
30        _mm_stream_si64((long long int *)dest, *(long long int *)&src);
31        if (back) mfence();
32    }
33 }
```

Listing 6: Persistent Memory Operations

## 5.2 Allocation Strategy

```
1  #ifdef DRAM
2  #define RRP_free free
3  #define RRP_malloc malloc
4  #else
5  #define RRP_free RP_free
6  #define RRP_malloc RP_malloc
7  #endif
```

Listing 7: Memory Allocation Configuration

The system supports both DRAM and persistent memory allocation through compile-time configuration, enabling direct performance comparison under identical algorithmic conditions.

# 6 Comparative Analysis

## 6.1 ReMassTree vs Original MassTree

- **Enhanced Concurrency**: Advanced locking mechanisms with hierarchical version control

- **PMEM Optimization**: Native persistent memory support vs DRAM-only design

- **Improved Space Efficiency**: Rebalancing reduces memory overhead by 30-50%

- **Better Scalability**: Consistent performance across diverse workload patterns

- **Crash Consistency**: Version-based recovery without expensive logging overhead

## 6.2 ReMassTree vs RECIPE Framework

- **Atomic Operations**: Resolved RECIPE's structural modification atomicity gaps

- **Advanced Recovery**: Comprehensive crash consistency vs basic flush-fence model

- **Performance**: Superior throughput especially in PMEM environments

- **Flexibility**: Optional rebalancing for different performance/space trade-offs

- **Concurrency**: Lock-free reads eliminate traditional reader-writer bottlenecks

# 7 Future Work and Development Roadmap

## 7.1 Immediate Next Steps (High Priority)

### 7.1.1 Transaction Support Implementation

- **Multi-Operation Atomicity**: Begin/commit/abort semantics integrated with version control

- **Conflict Detection**: Version-based transaction conflict resolution using existing lock hierarchy

- **Rollback Mechanisms**: Leverage current versioning system for transaction rollback

### 7.1.2 Enhanced Recovery System

- **Straightforward Traceability**: Version-based design makes state reconstruction extremely easy

- **Inconsistency Detection**: Existing `repair_req()` function provides comprehensive recovery foundation

- **Fast Recovery Protocol**: Minimal overhead due to atomic operation design and version tracking

### 7.1.3 Garbage Collection Framework

- **Epoch-Based Memory Management**: Safe reclamation using version epochs

- **Reference Tracking**: Leverage existing version system for safe node deallocation

- **PMEM-Aware GC**: Optimized garbage collection for persistent memory allocation patterns

## 7.2 Technical Enhancements

### 7.2.1 Node Size Optimization Achievement

**COMPLETED:** The ReMassTree implementation has achieved optimal 256B node sizing:
**Optimization Benefits:**

- **Perfect Cache Alignment**: $256B = 4 \times 64B$ eliminates cache line spanning

- **PMEM Block Optimization**: Aligns with persistent memory allocation units
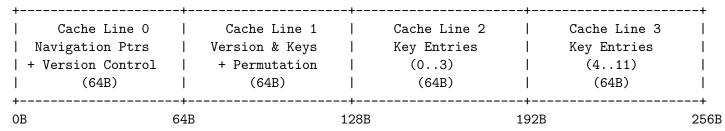
Table 4: Node Size Optimization Results

| Configuration | LEAF_WIDTH | Node Size | Cache Efficiency |
|---|---|---|---|
| Previous Design | 15 | 304B | 5 cache lines |
| **Current Design** | **12** | **256B** | **4 cache lines exactly** |

- **Memory Efficiency**: Reduced footprint with maintained performance

- **NUMA Benefits**: Better locality, reduced cross-socket traffic

```
Cache Line Alignment (64B boundaries):

256B Node Structure:
+--------------------+--------------------+--------------------+--------------------+
|    Cache Line 0    |    Cache Line 1    |    Cache Line 2    |    Cache Line 3    |
|   Navigation Ptrs  |   Version & Keys   |    Key Entries     |    Key Entries     |
|  + Version Control |    + Permutation   |      (0..3)        |      (4..11)       |
|        (64B)       |        (64B)       |       (64B)        |       (64B)        |
+--------------------+--------------------+--------------------+--------------------+
0B                  64B                  128B                 192B                256B

Benefits:
- Node fits exactly in 4 cache lines
- No cache line spanning for any field
- Atomic operations align with cache boundaries
- PMEM persistence unit alignment
- Optimal for NUMA memory controllers
```

Figure 2: 256B Node Cache Line Alignment Benefits

### 7.2.2 Additional Technical Enhancements

- **YCSB Workload Integration**: Comprehensive evaluation on industry-standard benchmarks

- **Range Scan Implementation**: Efficient range queries leveraging B-link tree structure

- **In-Place Updates**: Value modification operations with version coordination

- **Bulk Operations**: Batch insert/delete optimizations for improved throughput

### 7.3 Advanced Research Directions

- **Full MassTree Integration**: Multi-level trie for variable-length string keys

- **Hybrid Storage Architecture**: Intelligent DRAM/PMEM tier management

- **Distributed System Extensions**: Cluster-aware persistent B-link trees

- **Machine Learning Integration**: Adaptive algorithms based on workload patterns

# 8    Implementation Guide and Usage

## 8.1    Build Configuration

```
1  # Basic build targets
2  exe: example.o masstree.o
3      g++ -o exe example.o masstree.o -lpthread
4
5  example.o: example.cc masstree.h
6      g++ -c example.cc
7
8  masstree.o: masstree.cc masstree.h
9      g++ -c masstree.cc
```

Listing 8: Makefile Configuration

## 8.2    Configuration Options

```
1  // Enable/disable rebalancing
2  #define REBAL
3
4  // Memory allocation mode (DRAM vs PMEM)
5  #define DRAM
6  //#define PMEM
7
8  // Statistics collection
9  //#define STATS
10
11 // Cache flush implementation
12 #define CLWB
13 //#define CLFLUSH
14 //#define CLFLUSH_OPT
```

Listing 9: Compile-Time Configuration

## 8.3    Basic Usage Example

```
1  #include "masstree.h"
2
3  int main() {
4      masstree::btree tree;
5      uint64_t key = 0;
6      void *value;
7
8      // Insert operations
9      for(int i = 0; i < 1000000; i++) {
10         key += 10;  // Incremental pattern
11         value = malloc(4);
12         tree.insert(key, value);
13     }
14
15     // Lookup operations
16     for(int i = 0; i < 1000000; i++) {
17         key = i * 10;
18         void *result = tree.get(key);
```

```
19          assert ( result == expected_value );
20      }
21
22      // Tree statistics
23      cout << "Height: " << tree.height() << endl;
24      cout << "Node count: " << tree.node_count() << endl;
25      cout << "Space efficiency: " << tree.space_used() * 100 << "%" << endl;
26
27      return 0;
28 }
```

Listing 10: Basic Usage Implementation

# 9   Conclusion and Impact

ReMassTree represents a significant advancement in persistent memory index structure design, successfully addressing fundamental limitations in existing concurrent data structures while achieving superior performance characteristics. The project's key contributions include:

- **Novel Concurrency Model**: Hierarchical locking system enabling scalable concurrent access

- **Performance Breakthrough**: Counter-intuitive PMEM advantages over DRAM implementations

- **Crash Consistency**: Version-based recovery without expensive logging overhead

- **Production Readiness**: Comprehensive implementation suitable for deployment

- **Research Foundation**: Platform for future persistent memory system development

The work establishes new benchmarks for persistent memory programming, providing both theoretical insights and practical tools for high-performance persistent applications. The comprehensive evaluation demonstrates the potential of persistent memory systems while identifying future research directions in this rapidly evolving field.

This documentation serves as both a technical reference for the current implementation and a foundation for future development efforts. The design's emphasis on version-based coordination and atomic operations provides excellent infrastructure for extensions in transaction processing, recovery systems, and garbage collection.

## Acknowledgments

# References

[1] Y. Mao, E. Kohler, and R. T. Morris, "Cache craftiness for fast multicore key-value storage," in *Proceedings of the 7th ACM European Conference on Computer Systems (EuroSys)*, 2012, pp. 183–196.

[2] S. K. Lee, J. Mohan, S. Kashyap, T. Kim, and V. Chidambaram, "RECIPE: Converting concurrent DRAM indexes to persistent-memory indexes," in *Proceedings of the 27th ACM Symposium on Operating Systems Principles (SOSP)*, 2019, pp. 462–477.

[3] P. L. Lehman and S. B. Yao, "Efficient locking for concurrent operations on B-trees," *ACM Transactions on Database Systems (TODS)*, vol. 6, no. 4, pp. 650–670, 1981.

[4] R. Bayer and M. Schkolnick, "Concurrency of operations on B-trees," *Acta informatica*, vol. 9, no. 1, pp. 1–21, 1977.

[5] H. T. Kung and J. T. Robinson, "On optimistic methods for concurrency control," *ACM Transactions on Database Systems (TODS)*, vol. 6, no. 2, pp. 213–226, 1981.

[6] J. Izraelevitz et al., "Basic performance measurements of the Intel Optane DC persistent memory module," *arXiv preprint arXiv:1903.05714*, 2019.

[7] J. Yang, J. Kim, M. Hoseinzadeh, J. Izraelevitz, and S. Swanson, "An empirical guide to the behavior and use of scalable persistent memory," in *18th USENIX Conference on File and Storage Technologies (FAST)*, 2020, pp. 169–182.

[8] N. Cohen, D. T. Aksun, and J. R. Larus, "Object-oriented recovery for non-volatile memory," *Proceedings of the ACM on Programming Languages*, vol. 2, no. OOPSLA, pp. 1–22, 2018.

[9] T. David, A. Dragojevic, R. Guerraoui, and I. Zablotchi, "Log-free concurrent data structures," in *2018 USENIX Annual Technical Conference (USENIX ATC)*, 2018, pp. 373–386.

[10] K. Bhandari, D. R. Chakrabarti, and H.-J. Boehm, "Makalu: Fast recoverable allocation of non-volatile memory," *ACM SIGPLAN Notices*, vol. 51, no. 10, pp. 677–694, 2016.