# Supporting Block Device Abstraction on Storage Class Memory

**Youjip Won**

Hanyang University
Seoul, Korea

Hanyang University
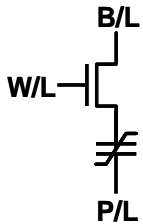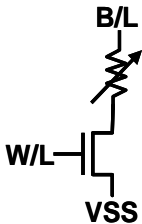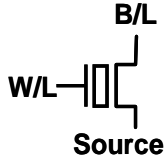**Distributed Multimedia Computing Laboratory**

# Background

# Comparison of storage class memory

| Items | FRAM | PRAM | NOR | NAND |
|---|---|---|---|---|
| Byte Addressable | Yes | Yes | Yes (read only) | No |
| Non-volatility | Yes | Yes | Yes | Yes |
| Read | 85ns | 62ns | 85ns | 16us |
| Write/Erase | 85ns / none | 300ns / none | 6.5us / 700ms | 200us / 2ms |
| Power Consumption | Low | High | High | High |
| Capacity | Low | Middle | Middle | High |
| Endurance | 1E15 | >1E7 | 100K | 100K |
| Unit Cell | | | | |

# Storage Class Memory Trend

- Expectation of the growth of Non-volatile Memory: NEDO(Japen. New Energy and Industrial Technology Development

**NVRAM Technology Trend**



0.5 GByte/Chip

"Integration of chip increases to marketable level", ByteandSwitch, 2006/6

# Media Speed and Operating System Technologies

# Byte Addressable NVRAM

**DRAM**

Fast

Byte Addressable

**NVRAM**

Non-volatile
Large size
Fast
Byte Addressable

**Flash/HDD**

Block Addressable

Large Size

## Storage Class Memory

Flash/HDD

**Non-volatile
memory storage**

# Computer with Storage Class Memory



Computer System: Current

CPU

H/W

MM

HDD

O/S

Process Mgmt

Memory mgmt | I/O buffer

MM

VM | FS

HDD

Computer System: Future

CPU

NVRAM Module

Process mgmt

Memory mgmt and file system

NVRAM

# Advantage: More Robust System



Current System

Future System

Process

Buffered data

write()

Delayed write

Process

write()

# Advantage: Maintaining Context

## Current System

**Proces s**

PCB
Code & Data

Register

**Fail !!**

Context is lot.

## Future System

**Proces s**

PCB
Code & Data

Register

**Fail !!**

Context is preserved.

# Timeliness of OS research for storage class memory

**File System for HDD**

Hard Disk

**HDD**

**Hybrid Disk**

10 yrs

10 yrs

Flash Memory

**USB drive**

**Solid State Disk**

**Flash F/S** 연구

**FTL research**

**SSD/Hybrid Disk** research

7 yrs

5 yrs

**NVRAM module**

Byte addressable NVRAM

**NVRAM F/S** research

OS technology for SCM

1970    1980    1990    2000    2007    2012

# Computer System with Storage Class Memory

# Computer for Storage Class Memory

**Memory Space**

**User Process**

load **0x0f8a0b**

**Page Table**

0

∞

Kernel Space(Static)

Kernel Space(Dynamic)

Memory space

Storage Space

Memory-storage Space

# Computer for Storage Class Memory

**Memory Space** **File system space**

**User Process**

**read (fp, offset)**

File space → Memory space

**Page Table**

| Kernel space(static) | Memory Space |
| Kernel space(dynamic) | Storage Space |
| | Memory-Storage Space |

0

∞

**Storage Space** Memory-Storage Space

**User Process**

load **0x0f8a0b**

read (**fp, offset**)

File offset  memory  **Memory address**

**Page Table**

0  ∞

Kernel space(static)   Memory space

Kernel space(dynamic)  Storage space

Memory−storage space

# Operating System for SCM



**Today's Computer**

Memory access → Memory Management → Memory → Virtual Memory → Virtual Memory Management → Storage

File access → I/O Buffer Management → I/O Buffer → File System → Disk Management → Storage

**Computer for Storage Class Memory**

Memory Access, File I/O → ??? → Memory + Storage

# Memory management

**Current Memory System**

**Access Request**

- Address Translation
- Page Table Structure
- Page Cache Management
- Page Replacement
- Page Protection
- Page Allocation Policy for Kernel/User
- Swap Area Management

**Swap Area**

**Memory System for Storage Class Memory**

**I/O Request**

**Consistency Support**

**Memory Manag**

- Address Translation
- Page Table Structure
- ~~Page Cache Management~~
- ~~Page Replacement~~
- ~~Page Protection~~
- Page Allocation Policy for Kernel/User
- ~~Swap Area Management~~
- **Management for Memory/Storage Area**

**Memory Storage**

# File System

**Current File System**

**I/O Request**

Block Device Management

Buffer

Consistency Support

Disk Management

HDD

- Buffer Allocation/Replacement
- Data Cache
- Meta Data Cache
- Read-ahead
- Write-behind
- Cache Replacement Policy
- Allocating Meta-data
- Storing Meta-data
- Defining Meta-data
- File System Structure
- Disk Usage Management

**File system for Storage Class Memory**

**I/O Request**

Consistency Support

Memory Management

Memory Storage

- Buffer Allocation/Replacement
- Data cache
- Meta-data Cache
- Read-ahead
- Write-behind
- Cache Replacement
- Allocating Meta-data
- Storing Meta-data
- Defining Meta-data
- File System Structure
- Disk Usage Management

# File System Management

## Current File System

**Recovery : Journaling Policy**

**Block Device Management**

**Buffer**

**Consistency Support**

**Disk Management**

**HDD**

- Journal Management on Memory
- Transaction Commit Algorithm
- Journal Commit on Disk
- Check Pointing
- Redo
- Undo

## File System for Storage Class Memory

**Recovery: Journaling /Transaction**

**Consistency Support**

**Memory Management**

**Memory -Storage**

- ~~Journal Management on Memory~~
- Transaction Commit Algorithm
- ~~Journal Commit on Disk~~
- Check Pointing
- ~~Redo~~
- Undo

# Atomicity of Block write for file system

# What makes memory-storage file system possible?

- We need to impose block device abstraction on Storage Class Memory Region!

- Key Ingredient in Imposing Block Device Abstraction on Memory region

We need atomicity guarantee on Block I/O!

# The notion atomicity

- Atomicity in block device level

  *Make 4KByte I/O atomic!*

- Atomicity in file system

  *Make system call atomic!*

  create()
  - Allocate free inode
  - Set inode bitmap
  - Allocate directory entry
  - Update directory block

22

# Block I/O and Partial Write

# Partial write in Storage Class Memory

- When partial write occurs?

    - Code Crash, Power outage, System failure

- Partial write and Storage Class memory

    - NVRAM can still preserve the result of incomplete I/O

    *Can affect the file system consistency.*

    *Can corrupt the data.*

# Problem of Partial Write

| New Data | Old Data |
|---|---|

Partially written data

Write Unit of Storage

Storage space

Wrong Data

# Atomic I/O in legacy Storage

- HDD

    - ECC for sector

    - Capacitor Assisted Write

    - Journaling file system

- Flash/SSD

    - ECC for every page(spare area)

    - Atomicity support by Flash Translation layer

# Atomic I/O in main memory file system

- Representative RAM based file systems
  - RAM disk, ramfs, tmpfs

- Characteristics
  - All data is located in volatile memory
  - File systems for Temporary file: naturally expected to be non-durable

## Unnecessary to consider atomicity of write

# Existing approaches storing data in SCM

PRAMFS [Longerbeam:2004]
- metadata
- data

MRAMFS [Edel:2004]
- Compressed metadata
- Compressed data

Conquest [Wang:2002]
- metadata
- Small file data

PRIMS [Greenan:2007]
- metadata log
- Software ECC

High complexity protection method

There is no partial write protection.

# Atomic Block I/O Interface for Storage Class Memory

**Application** ⇕ Read/write **File system** ⇕ Device driver **Block Device**

**File system for Hard disk**

**Application** ⇕ Read/write **File system** ⇕ Memory Copy **DRAM**

**File system for RAM**

**Application** ⇕ Read/write **File system** ⇕ Atomic Write Interface **Storage Class Memory**

**File system for Storage Class Memory**

# Log-based Block Mapping

# System Environment

Processor

DRAM

Storage Class Memory

0xE0000

0xE8000

Virtual Storage Class Memory Region

**Page Table**

0x08000

0x10000

Storage Class Memory

**Virtual Address Space**

# Organization of Storage Class Memory Region

## Storage Class Memory

### Management Area

### Page Area

**Management Info.**

| | |
|---|---|
| Update loc | NULL |
| Tmp loc | NULL |

**Mapping Table**

| | |
|---|---|
| Target | 1 |
| 1 | 3 |
| 2 | N |
| 3 | 0 |
| ⋮ | ...... |
| K | L |
| ⋮ | ....... |
| N | 1 |

Page Area (4KB pages):
- 0th page
- 1th page
- 2th page
- 3th page
- ...... ⋮
- Lth page
- ...... ⋮
- Nth page

➤ Entries of mapping table have physical page address
➤ Size of data in the Management Area is same as word size of processor

# Atomicity and Block I/O

- Read(LBA, buffer)
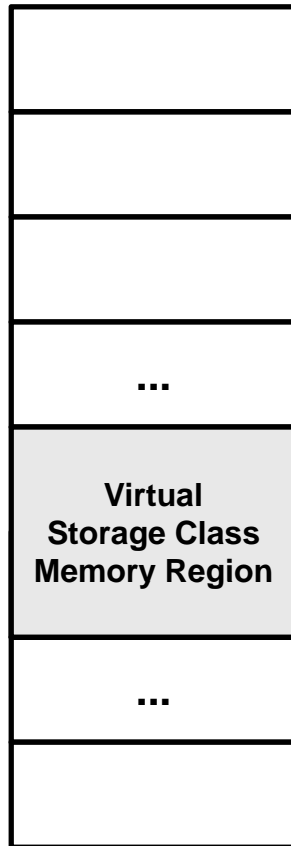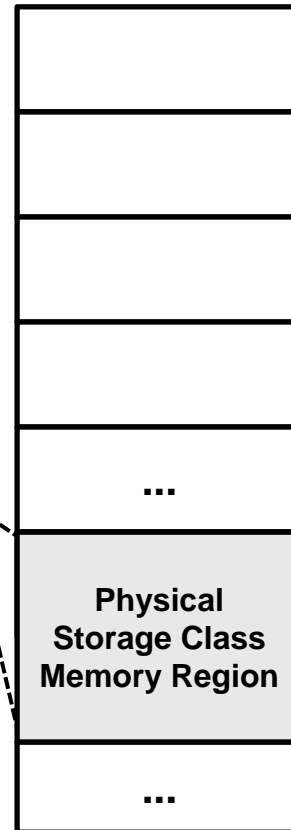  - Read 4KB page from storage class memory
  - Parameter
    - Buffer address
    - Location of data to read(Page number)
- Write(LBA, Buffer)
  - Write 4KB page to storage class memory
  - Parameter
    - Pointer of source data
    - Page number to write

# System architecture

**Virtual Address Space**

**Physical Address Space**

**Storage Class Memory**

| | |
|---|---|
| ... | |
| **Virtual Storage Class Memory Region** | |
| ... | |

**Page Table**

| |
|---|
| |
| |

| | |
|---|---|
| ... | |
| **Physical Storage Class Memory Region** | |
| ... | |

| | |
|---|---|
| Management Info. | **Management Area** |
| Mapping Table | |
| *Page* | **Page Area** |
| *Page* | |
| ..... | |
| *Page* | |

☐ **Conventional Memory Space**

▨ **Storage Class Memory Space**

# Read (LBA, Buffer)

1. Obtain MAP[LBA]

2. Copy from storage class memory to buffer

# Read Request

Read refer to mapping table and obtain physical address of requested page

**Page Area**

0th page

**buf**

1th page

2th page

**Management Area**

3th page

Target block | NULL

Read(1, buf)  1 | 3

......

......  ...

N | 2

L th page
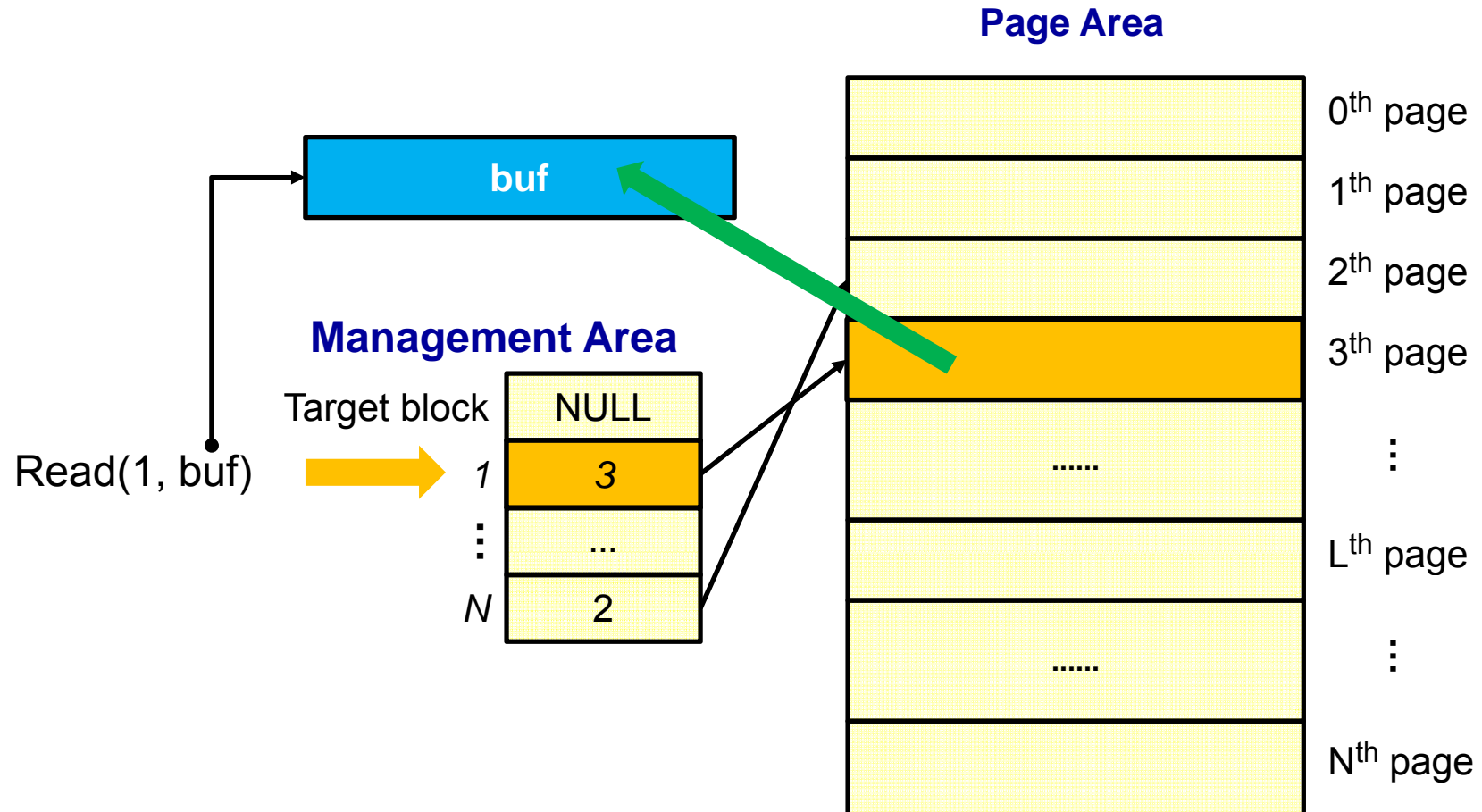
......

N th page

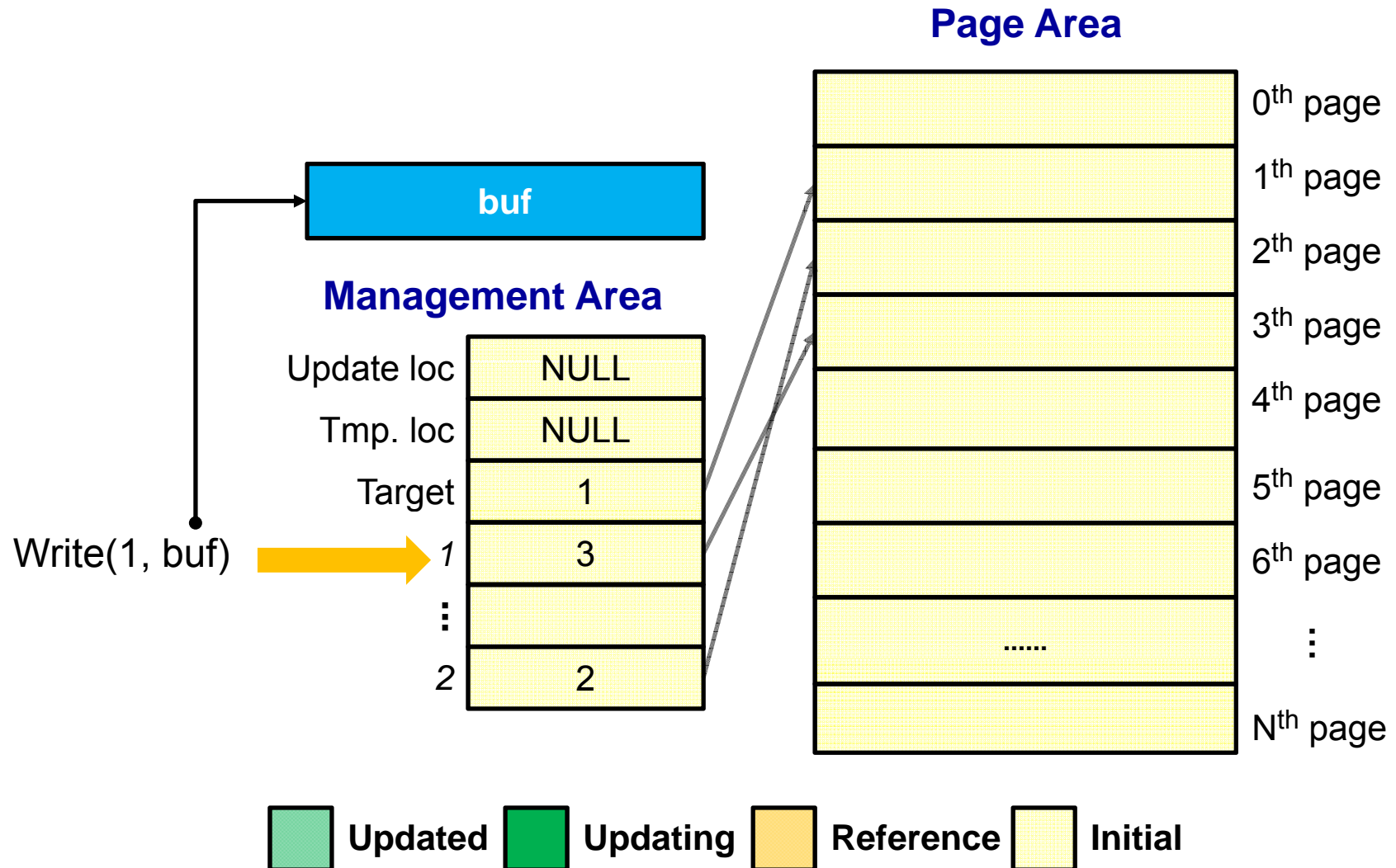# Write(LBA, Buffer)

Copy  data to MAP[target block]  and

swap(MAP[target block], MAP[LBA])

1. MAP[update loc]  ← LBA.

2. MAP[tmp loc] ← MAP[LBA].

3. Copy  data to MAP[target block].

4. MAP[LBA] ← MAP[target block]

5. MAP[target block ] ← MAP[tmp loc]

6. MAP[update loc] ← NULL

7. MAP[tmp loc] ← NULL

# Write Request

Write request is occurred

## Page Area

**buf**

## Management Area

| | |
|---|---|
| Update loc | NULL |
| Tmp. loc | NULL |
| Target | 1 |
| *1* | 3 |
| ⋮ | |
| *2* | 2 |

Write(1, buf)

Page Area:
- 0th page
- 1th page
- 2th page
- 3th page
- 4th page
- 5th page
- 6th page
- ⋮ ......
- Nth page

**Updated**  **Updating**  **Reference**  **Initial**

# Write Request Step1

Update "update location" to point to requested block number address

**Page Area**

# Write Request data transfer

Transfer data to physical location where target block points to

**Page Area**

| Area | |
|---|---|
| | 0th page |
| 1th page | 1th page |
| | 2th page |
| | 3th page |
| | 4th page |
| | 5th page |
| ...... | 6th page |
| | Nth page |

**buf**

**Management Area**

| | |
|---|---|
| Update loc | 1 |
| Tmp. loc | 3 |
| Target | 1 |
| *1* | 3 |
| ⋮ | |
| *2* | 2 |

Write(1, buf)

**Updated**  **Updating**  **Reference**  **Initial**

# Write Request Step3

✦ Update "Map[1]" that is mapping table entry of requested block address(swap "Map[1]" and "target block")

**Page Area**

**buf**

**Management Area**

| | |
|---|---|
| Update loc | 1 |
| Tmp. loc | 3 |
| Target | 1 |
| *1* | 1 |
| ⋮ | |
| *2* | 2 |

Write(1, buf)

$0^{th}$ page
$1^{th}$ page
$2^{th}$ page
$3^{th}$ page
$4^{th}$ page
$5^{th}$ page
$6^{th}$ page
⋮
......
$N^{th}$ page

| | | | | | | |
|---|---|---|---|---|---|---|
| ▦ | **Updated** | ▦ | **Updating** | ▦ | **Reference** | ▦ Initial |

# Write Request Step4

➕Update "target block" to point to old physical address of requested block stored in "temporary space"

**Page Area**
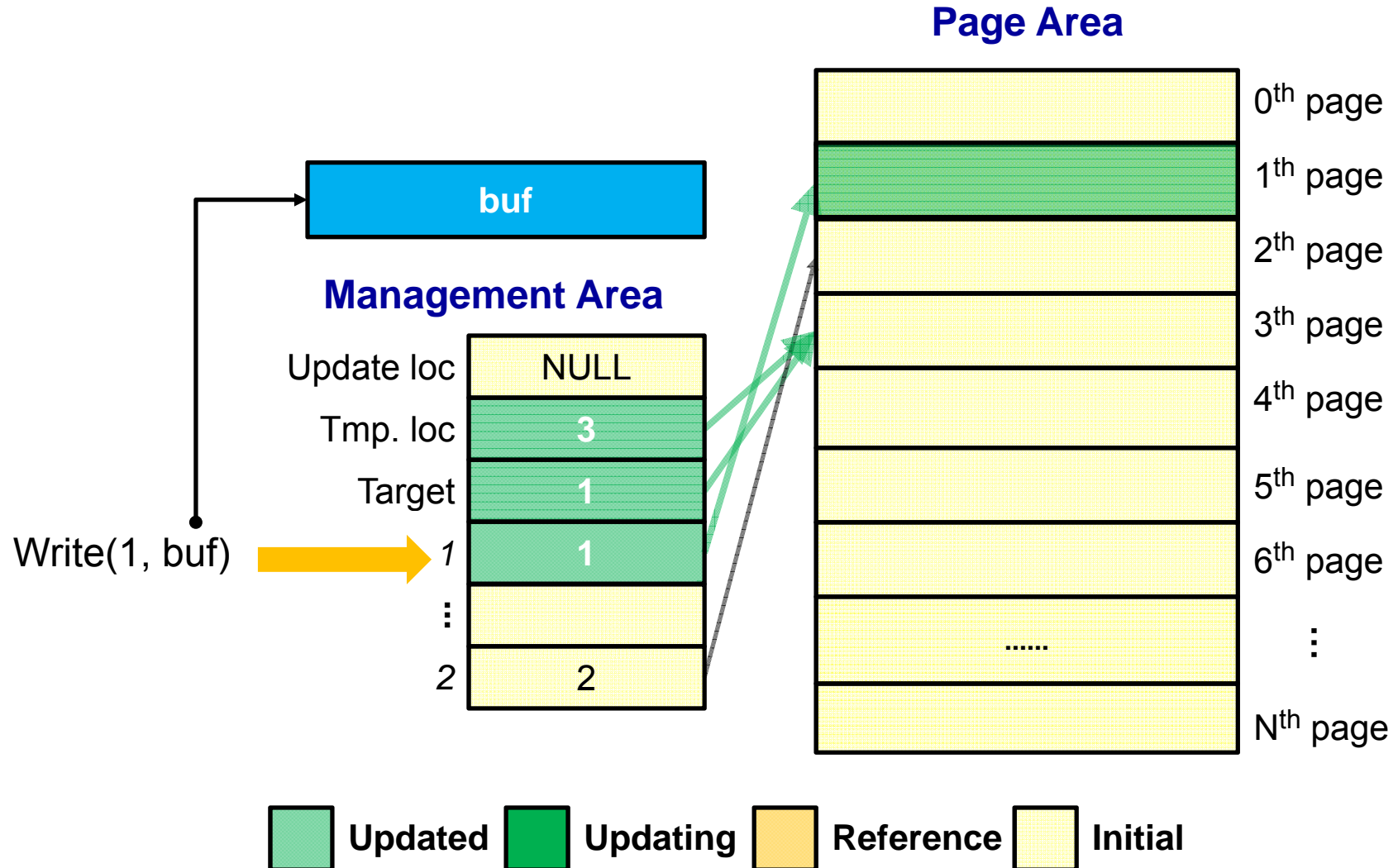
**buf**

**Management Area**

| | | |
|---|---|---|
| Update loc | 1 | |
| Tmp. loc | 3 | |
| Target | 1 | |
| 1 | 1 | |
| ⋮ | | |
| 2 | 2 | |

Write(1, buf)

0th page
1th page
2th page
3th page
4th page
5th page
6th page
⋮
......
Nth page

**Updated**  **Updating**  **Reference**  **Initial**

# Write Request Step5

Initialize "update location" as NULL value

**Page Area**

**buf**

**Management Area**

| | |
|---|---|
| Update loc | NULL |
| Tmp. loc | 3 |
| Target | 1 |
| 1 | 1 |
| ⋮ | |
| 2 | 2 |

Write(1, buf)

0th page
1th page
2th page
3th page
4th page
5th page
6th page
⋮
......
Nth page

■ **Updated**　■ **Updating**　■ **Reference**　□ **Initial**

# Write Request Step6

Initialize "temporary space" as NULL value

**Page Area**

**buf**

**Management Area**

| | |
|---|---|
| Update loc | NULL |
| Tmp. loc | NULL |
| Target | 3 |
| 1 | 1 |
| ⋮ | |
| 2 | 2 |

Write(1, buf)

0th page
1th page
2th page
3th page
4th page
5th page
6th page
⋮
......
Nth page

■ **Updated**　■ **Updating**　■ **Reference**　□ **Initial**

# Log-based history management

| Update location |
|---|
| **Temporary space** |

→ Step1

| Update location |
|---|
| Temporary space |
| *Map*[update location] |

→ Step2

page

| Update location |
|---|
| Temporary space |
| Target block |
| *Map*[upda... |

→ Step3

page

| Update location |
|---|
| Temporary space |

→ Step4

page

| Upda... |
|---|
| Temporary space |

→ Step5

**All step can be identified.**
**Proceed restore from completed step.**

Null value
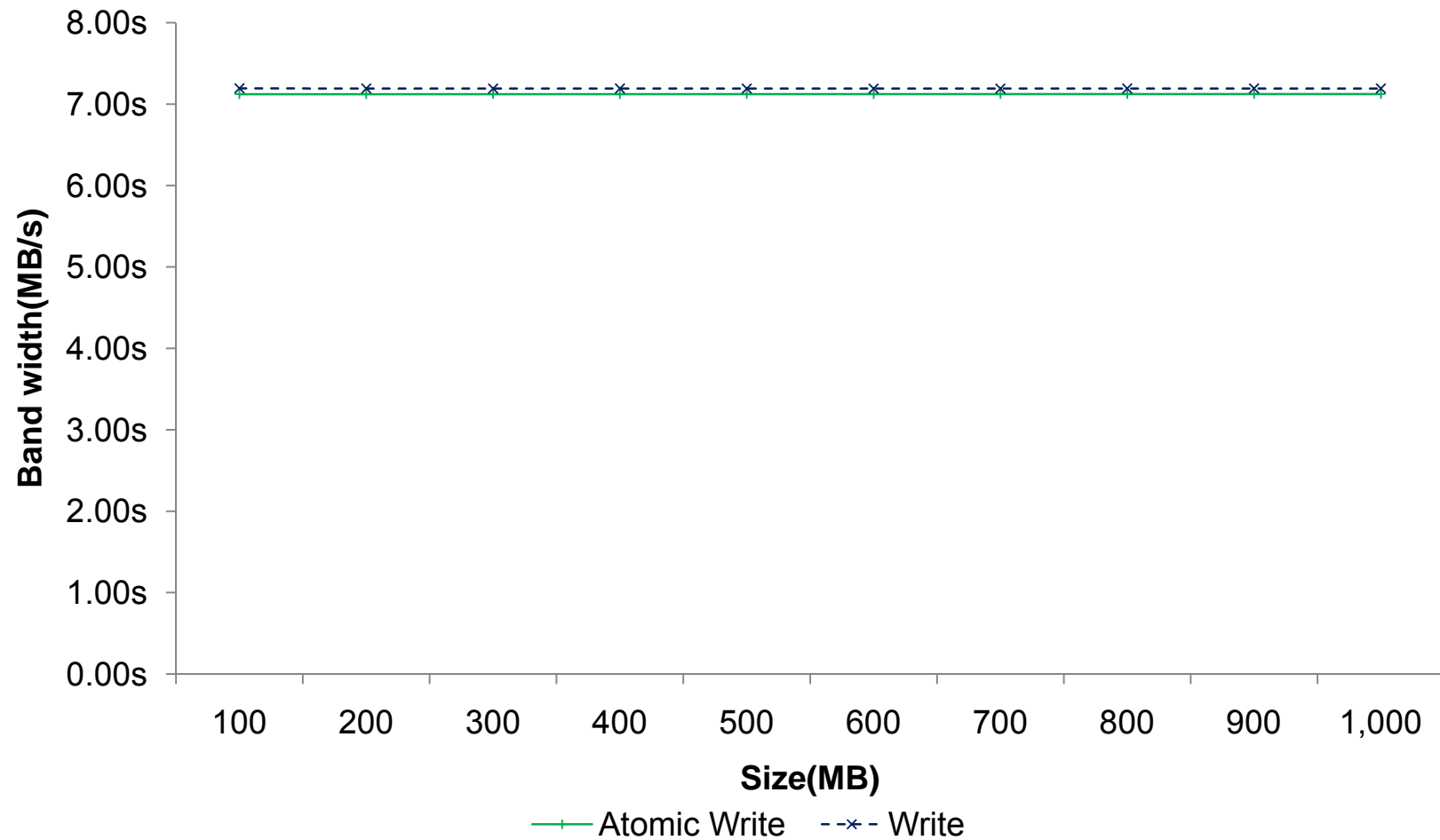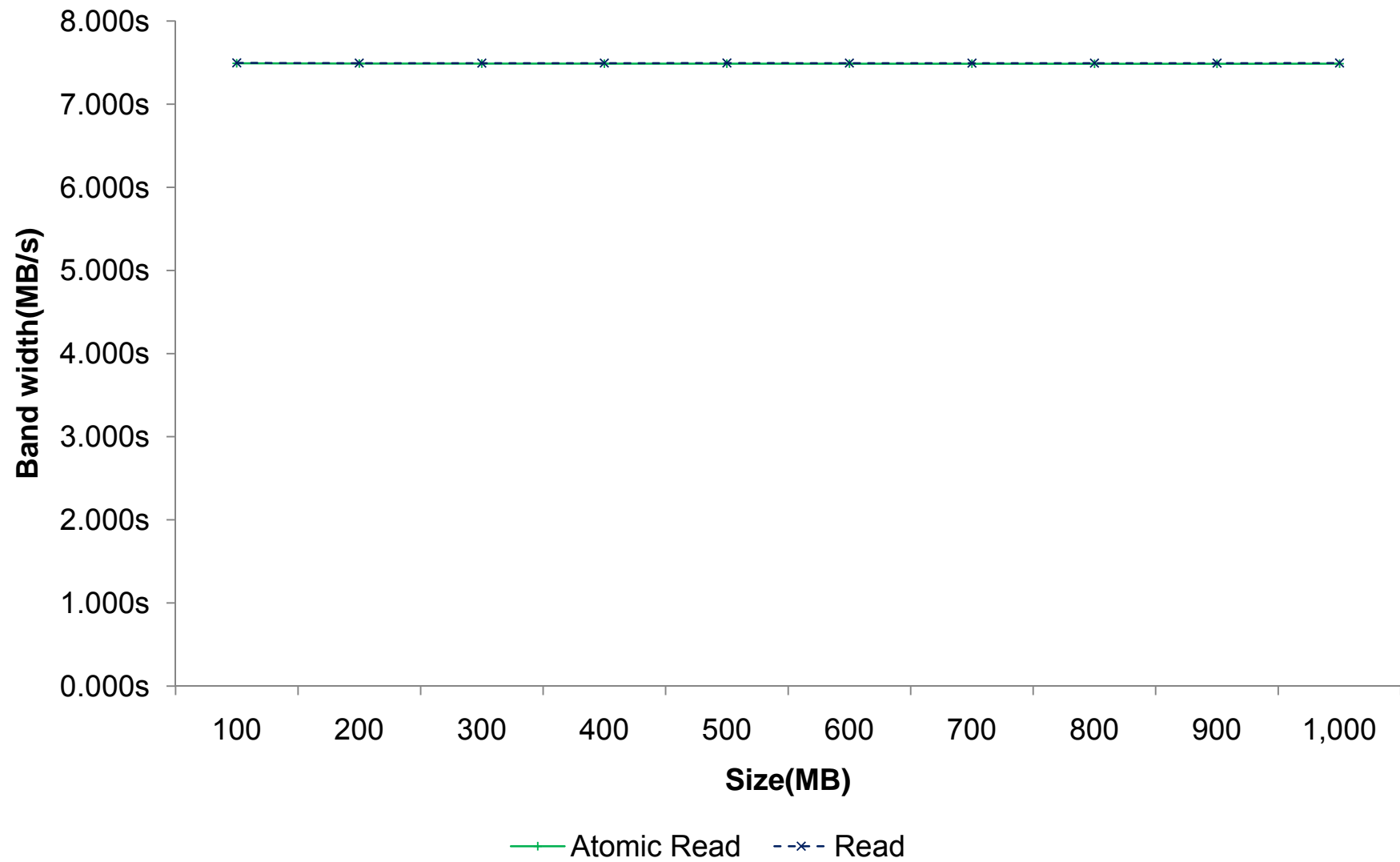
# Experiment

# Environment

- SMDK2440 Board
  - S3C2440
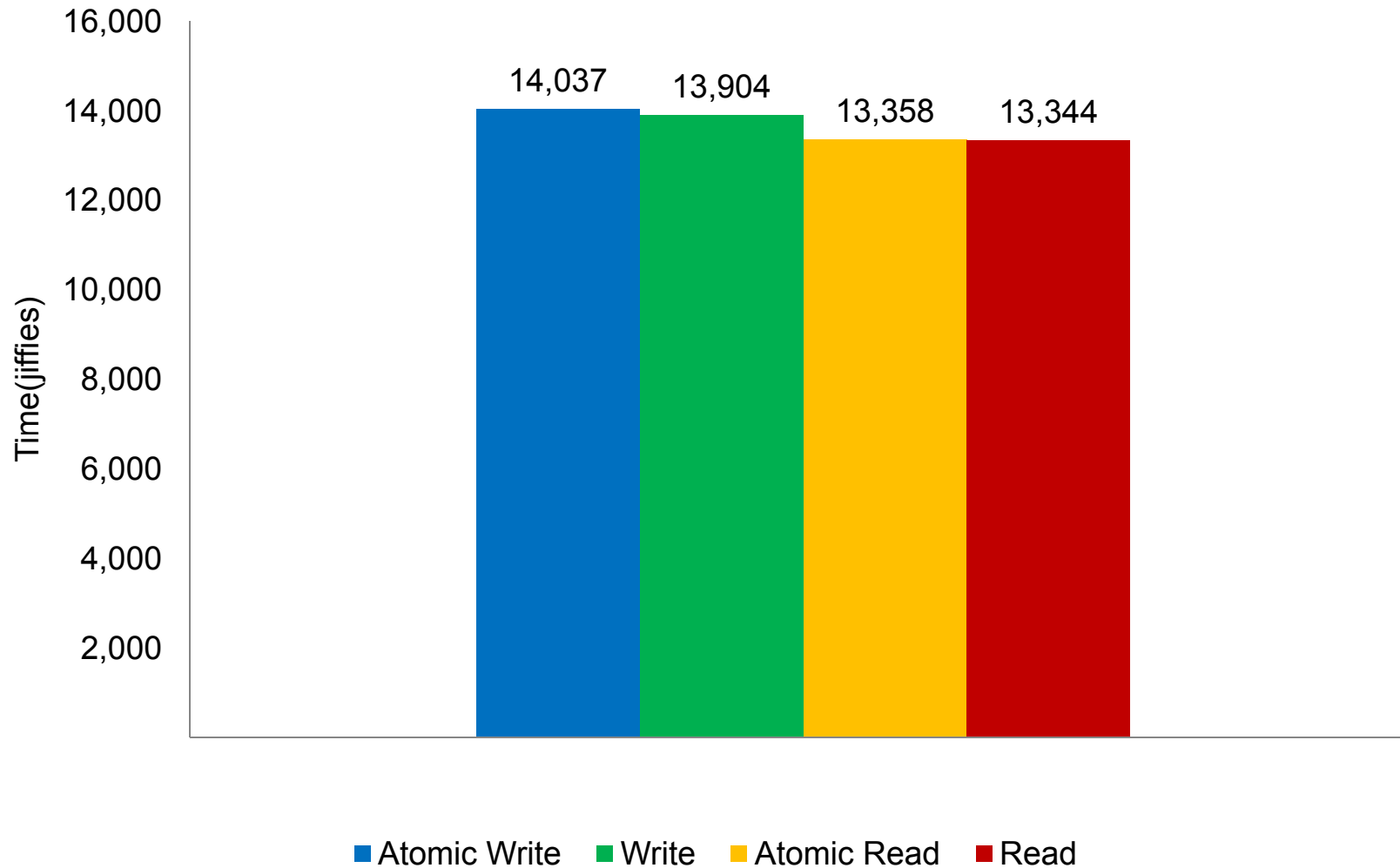  - 64MB DRAM
  - 8MB FRAM

# Write Performance

# Read Performance

# Overhead Comparison

# Conclusion & Future work

- Supporting atomicity of I/O operation
  - Overhead in providing atomicity is not significant

- Log-based block mapping mechanism is proposed to support atomicity

- Effect of the cache of processor

- Target block can be a bottleneck point