

# 07. Flash File Systems

## Special Topics in Computer Systems:

Modern Storage Systems

(SE820-01)

**Instructor:**

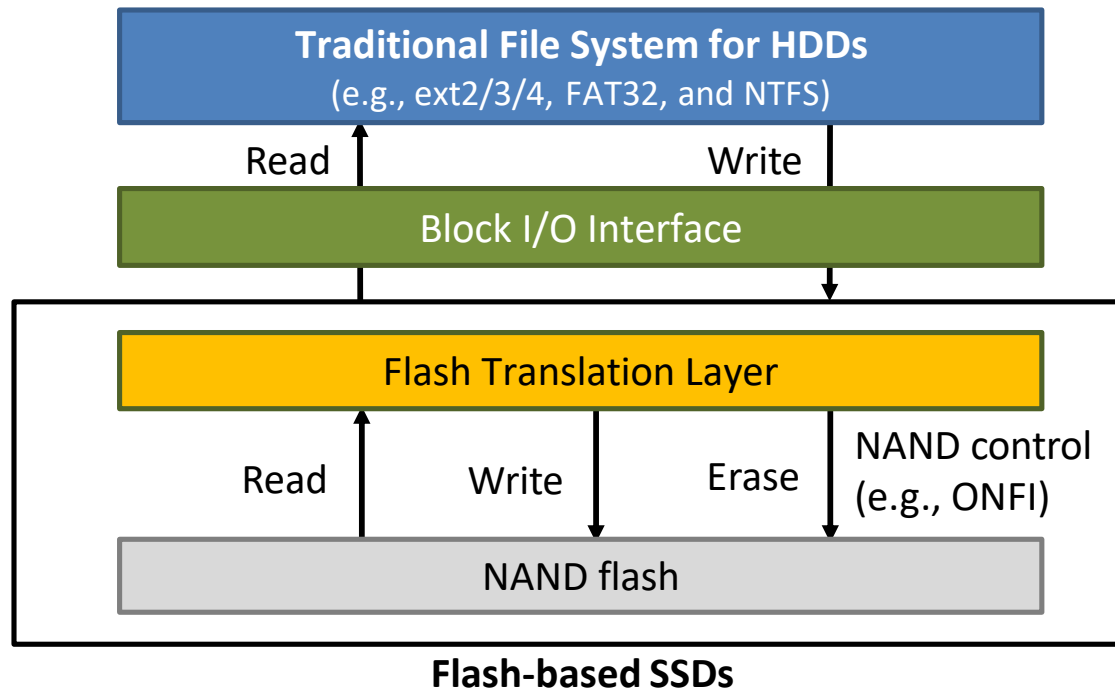
Prof. Sungjin Lee ([sungjin.lee@dgist.ac.kr](mailto:sungjin.lee@dgist.ac.kr))

# Outline

- **Traditional Flash File Systems**
- **SSD-Friendly Flash File Systems**
  - F2FS: Flash-friendly File System
- **Reference**

# Traditional File Systems for Flash

- Originally designed for block devices like HDDs
  - e.g., ext2/3/4, FAT32, and NTFS
- But, **NAND flash memory is not a block device**
  - The FTL provides block-device views outside, hiding the unique properties of NAND flash memory

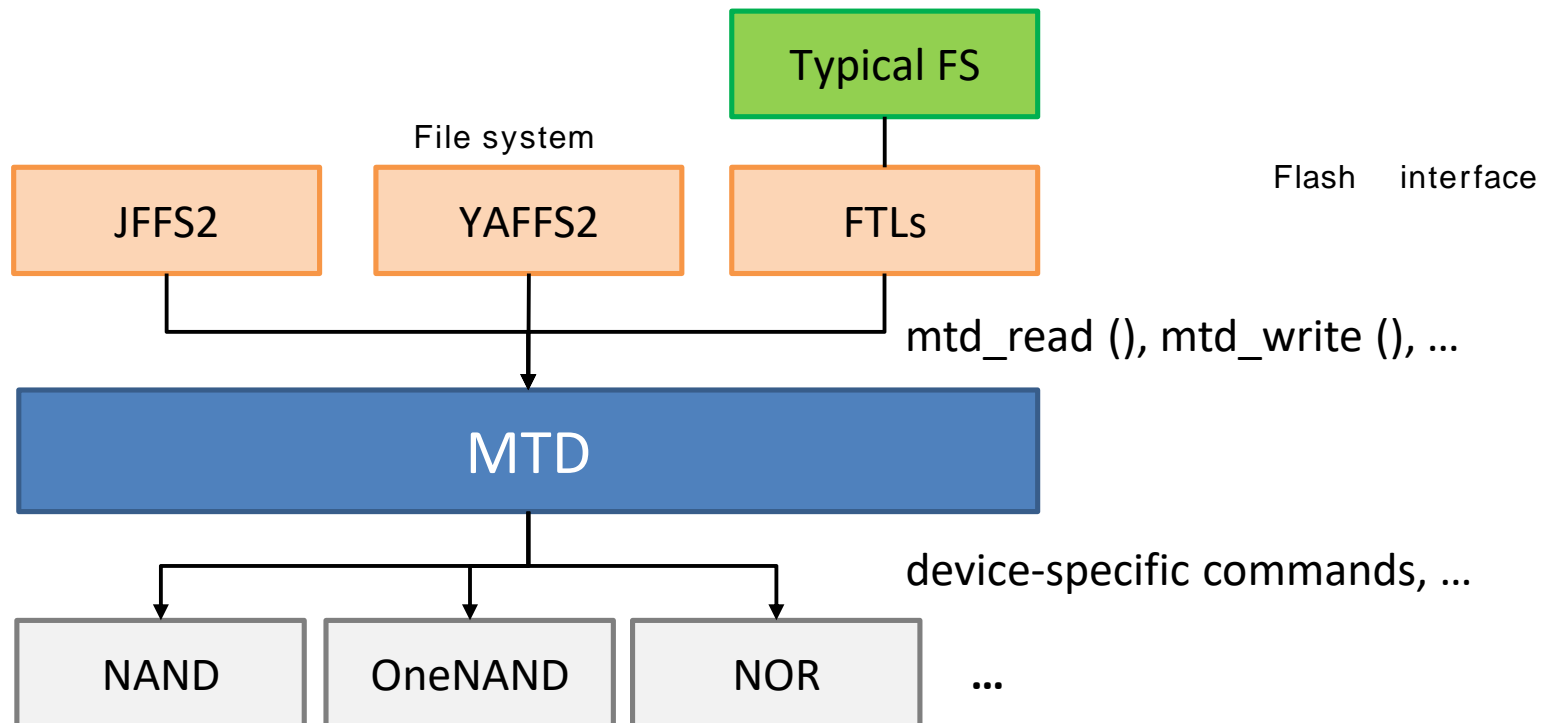




# Memory Technology Device (MTD)

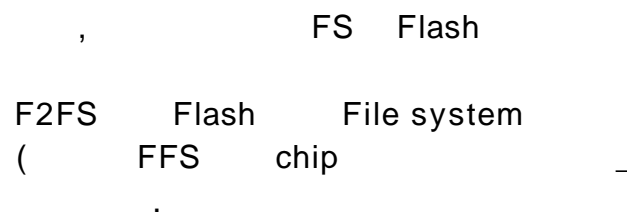
- **MTD is the lowest level for accessing flash chips**
  - Offer the same APIs for different flash types and technologies
    - e.g., NAND, OneNAND, and NOR

- **JFFS2 and YAFFS2 run on top of MTD**



# Traditional File Systems vs. Flash File Systems

	Flash	FFS
	File System + FTL	Flash File System
Method	- Access a flash device via FTL	- Access a flash device directly
Pros	<ul style="list-style-type: none"> <li>- High interoperability</li> <li>- No difficulties in managing recent NAND flash with new constraints</li> </ul>	<ul style="list-style-type: none"> <li>- High-level optimization with system-level information</li> <li>- Flash-aware storage management</li> </ul>
Cons	<ul style="list-style-type: none"> <li>- Lack of system-level information</li> <li>- Flash-unaware storage management</li> </ul>	<ul style="list-style-type: none"> <li>- Low interoperability</li> <li>- Must be redesigned to handle new NAND constraints</li> </ul>



Flash file systems now become obsolete because of difficulties for the adoption to new types of NAND devices

# JFFS2: Journaling Flash File System

## ■ A log-structured file system (LFS) for use with NAND flash

- Unlike LFS, however, it *does not allow* any in-place updates!!!

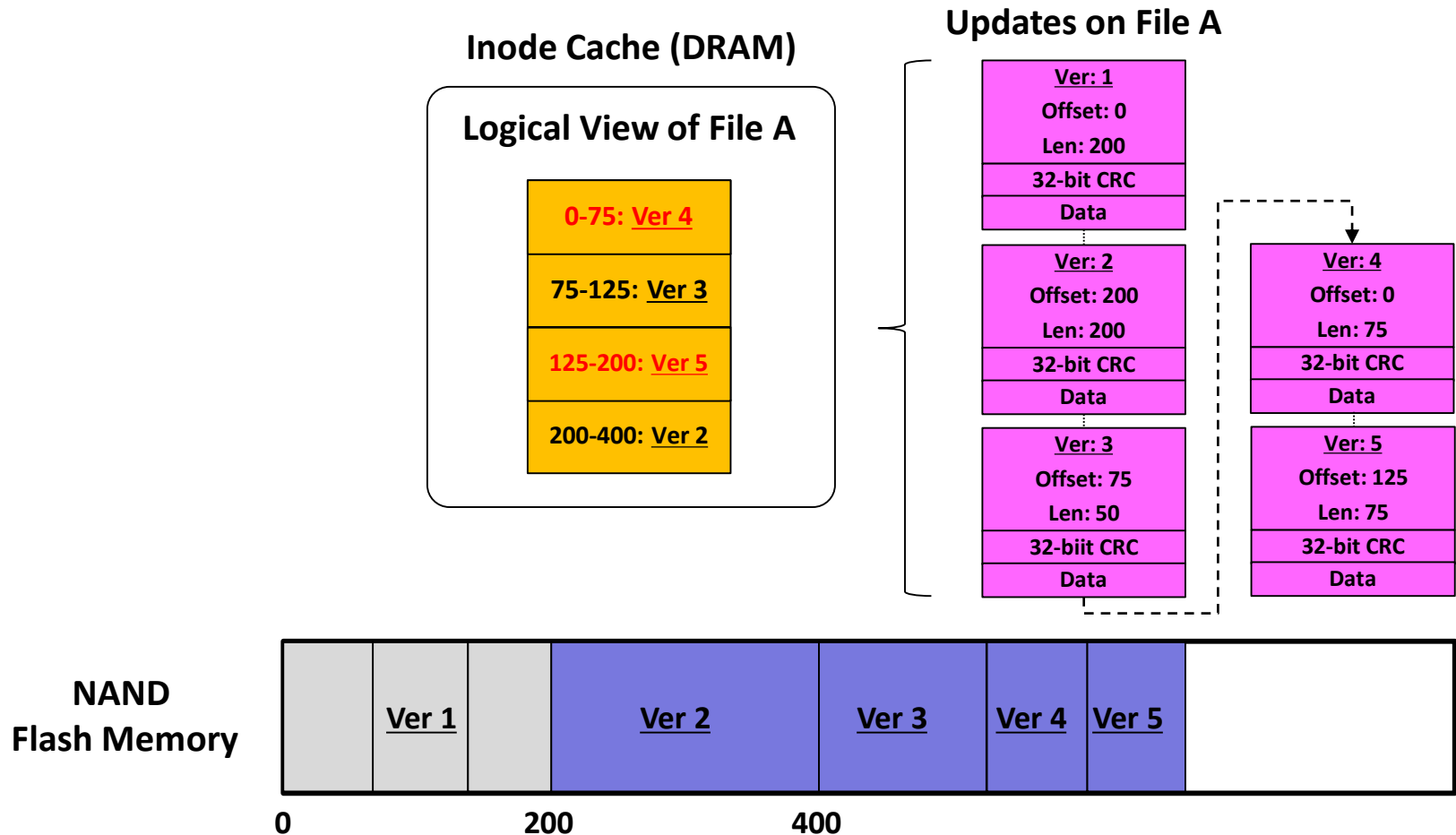
## ■ Main features of JFFS2

- *File data and metadata* stored as nodes in NAND flash memory
- Keep an inode cache holding the information of nodes in DRAM
- A greedy garbage collection algorithm
  - Select cheapest blocks as a victim for garbage collection
- A simple wear-leveling algorithm combined with GC
  - Consider the wearing rate of flash blocks when choosing a victim block for GC
- Optional data compression

victim block  
100      99      cost      가  
random      wear - leveling

# JFFS2: Write Operation

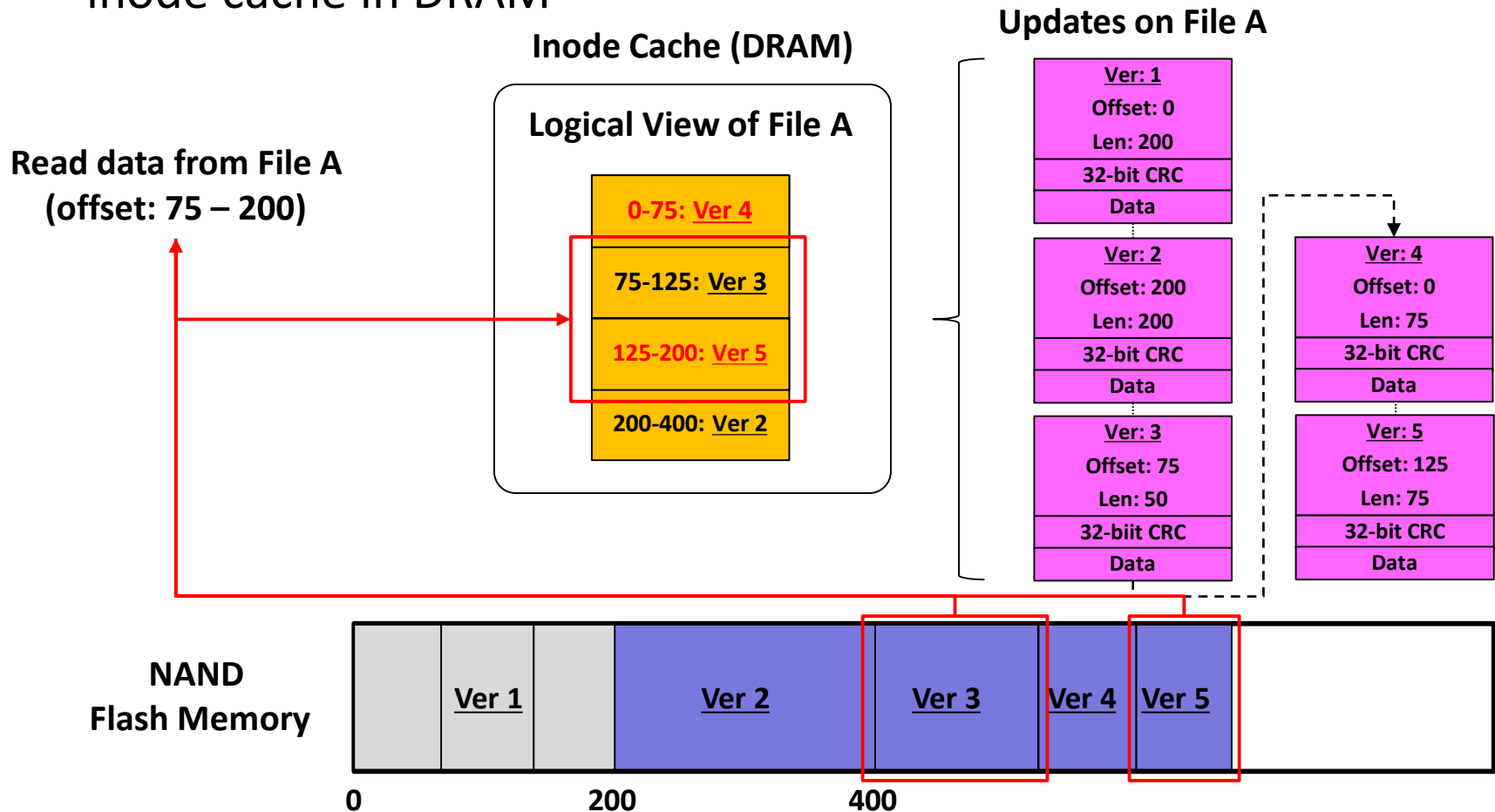
- All data are written sequentially to a log which records all changes





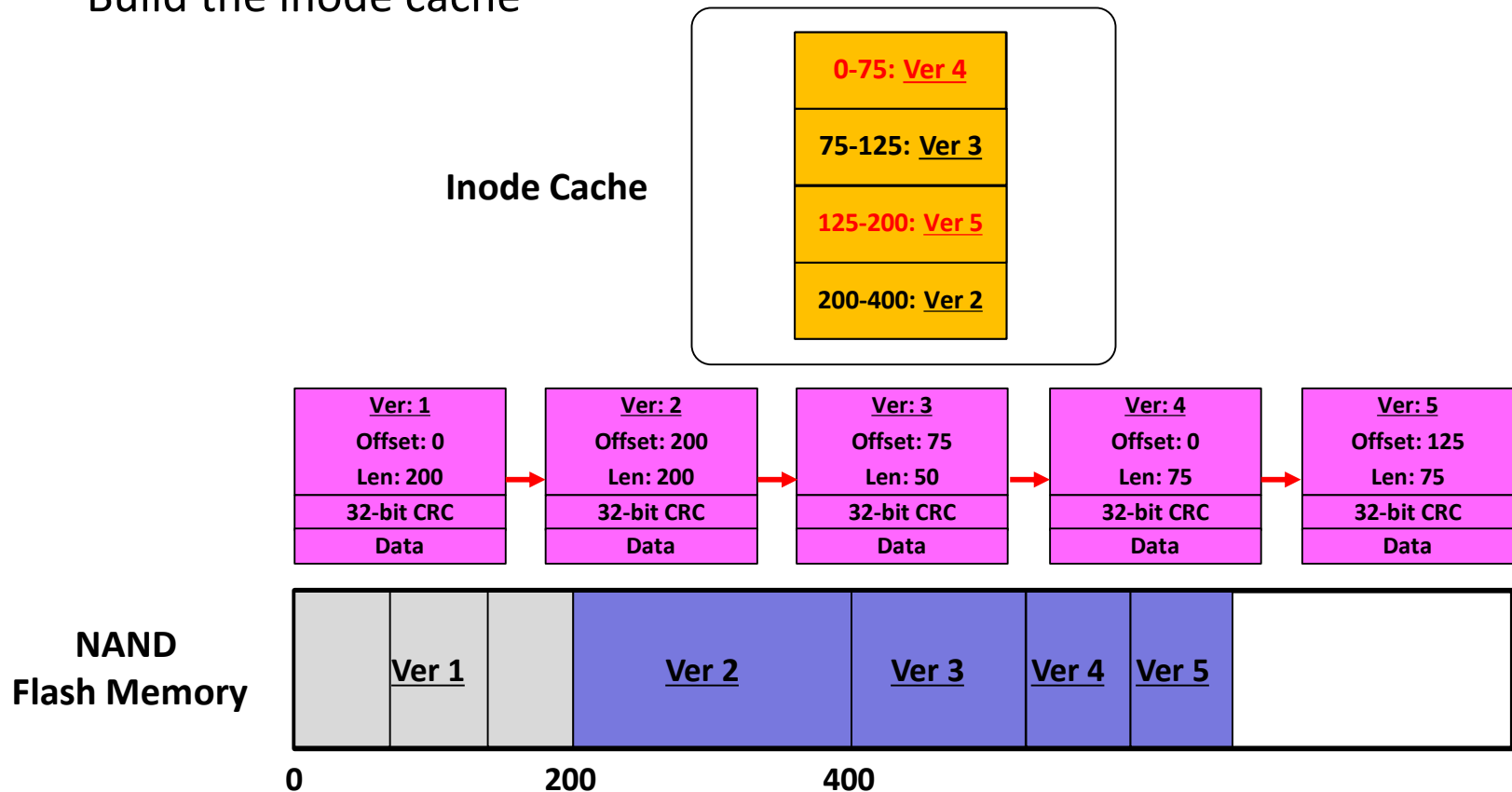
# JFFS2: Read Operation

- The latest data can be read from NAND flash by referring to the inode cache in DRAM



# JFFS2: Mount

- Scan the flash memory medium after rebooting
  - Check the CRC for written data and mark the obsolete data
  - Build the inode cache



# JFFS2: Problems

## ■ Slow mount time

inode cache

node scan

.

- All nodes must be scanned at mount time
- Mount time increases in proportion to flash size and file system contents

## ■ High memory consumption

node DRAM  
DRAM

.

- All node information must be maintained in DRAM
- Memory consumption linearly depends on file system contents

## ■ Low scalability

size가  
scalability 2가

.

- Infeasible for a large-scale flash device
- Mount time and memory consumption increase according to a flash size

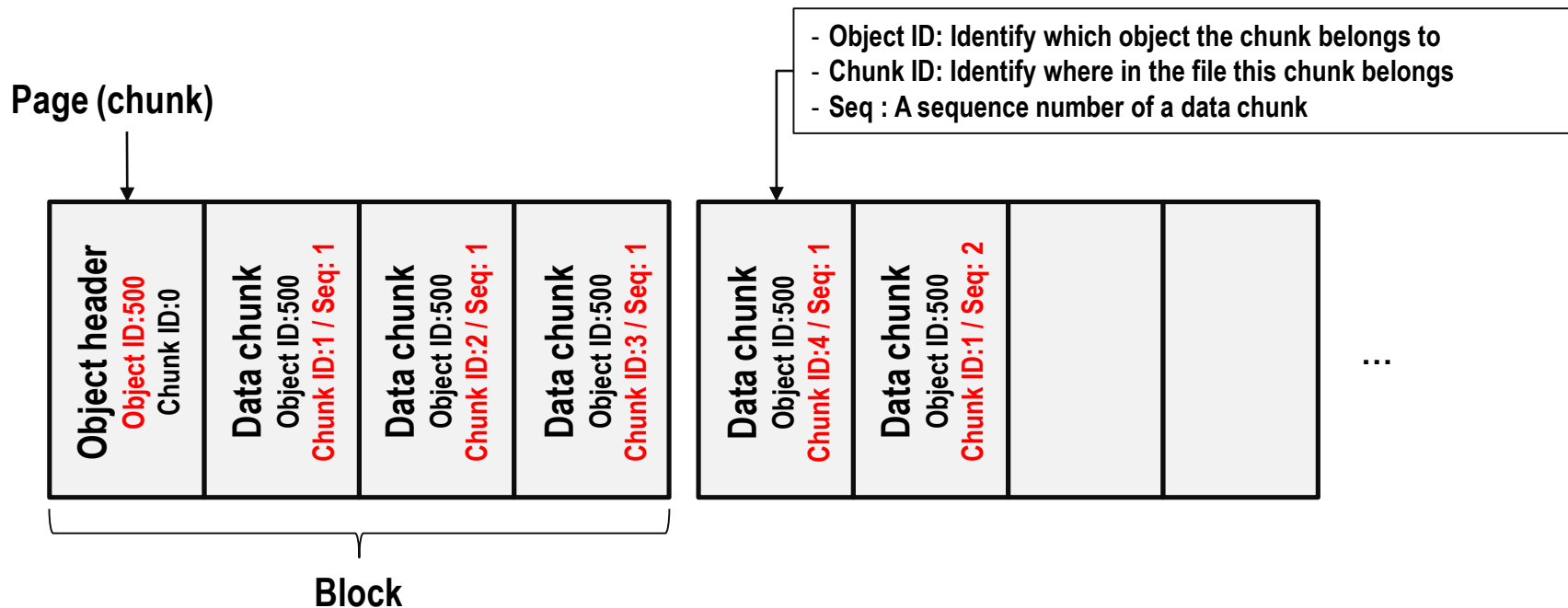
FS가 redhat

# YAFFS2: Yet Another Flash File System

- **Another log-structured file system for flash memory**
  - Store data to flash memory like a log with a unique sequence number like JFFS2
  - Reads and writes are performed similar to JFFS2
- **Mitigate the problems raised by JFFS2**
  - Relatively frugal with memory resource
  - Checkpoints for a fast file system mount
  - Dynamic wear-leveling
- **Support multiple platforms**
  - Linux, WinCE, RTOSs, etc

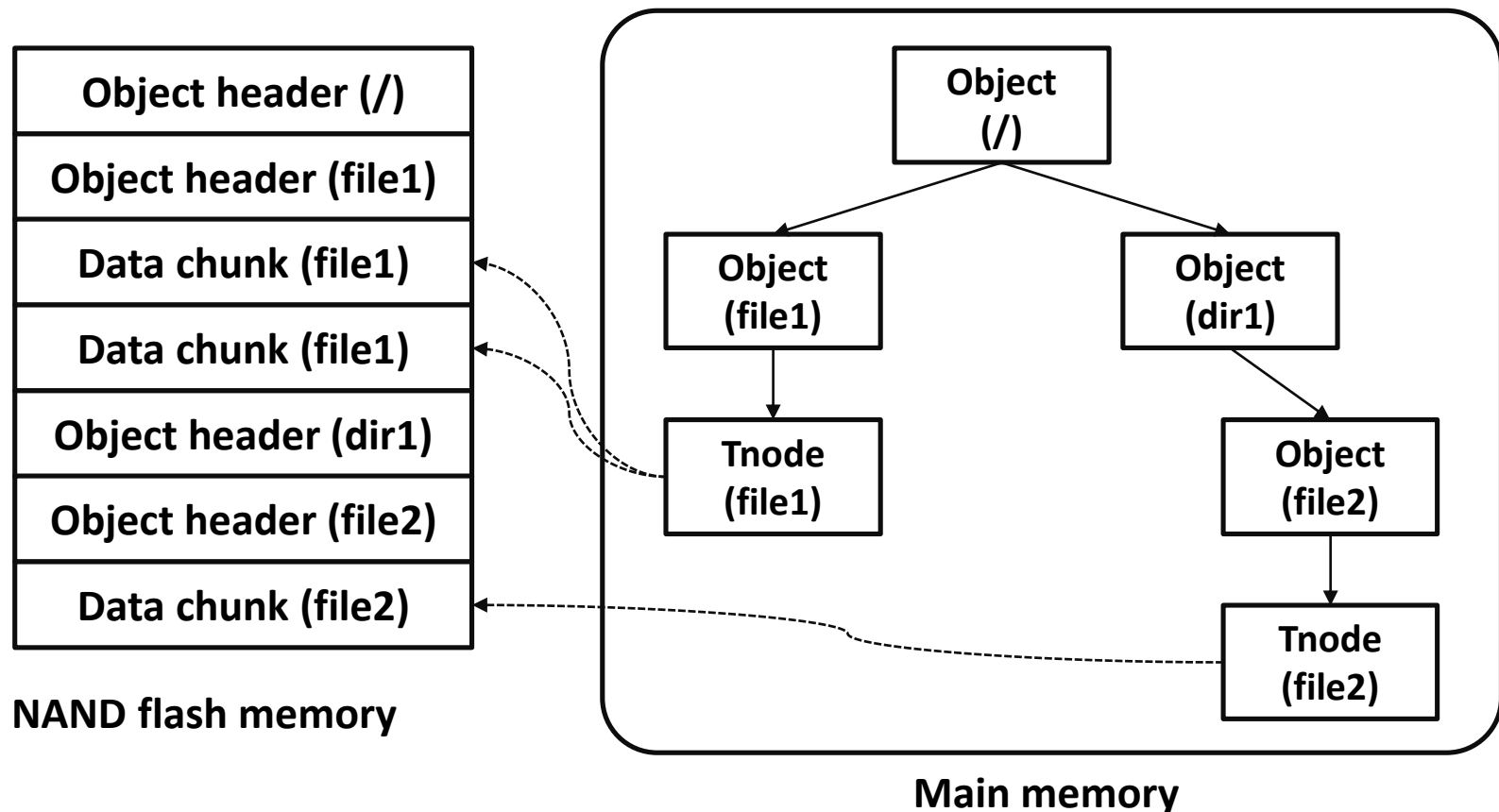
# YAFFS2: Physical Layout

- The entries in the log are all one chunk (one page) in size and can hold one of two types of chunk
  - Data chunk: A chunk holding regular data file contents
  - Object header: A descriptor for an object, such as a directory, a regular file, etc; similar to `struct stat` but include dentry



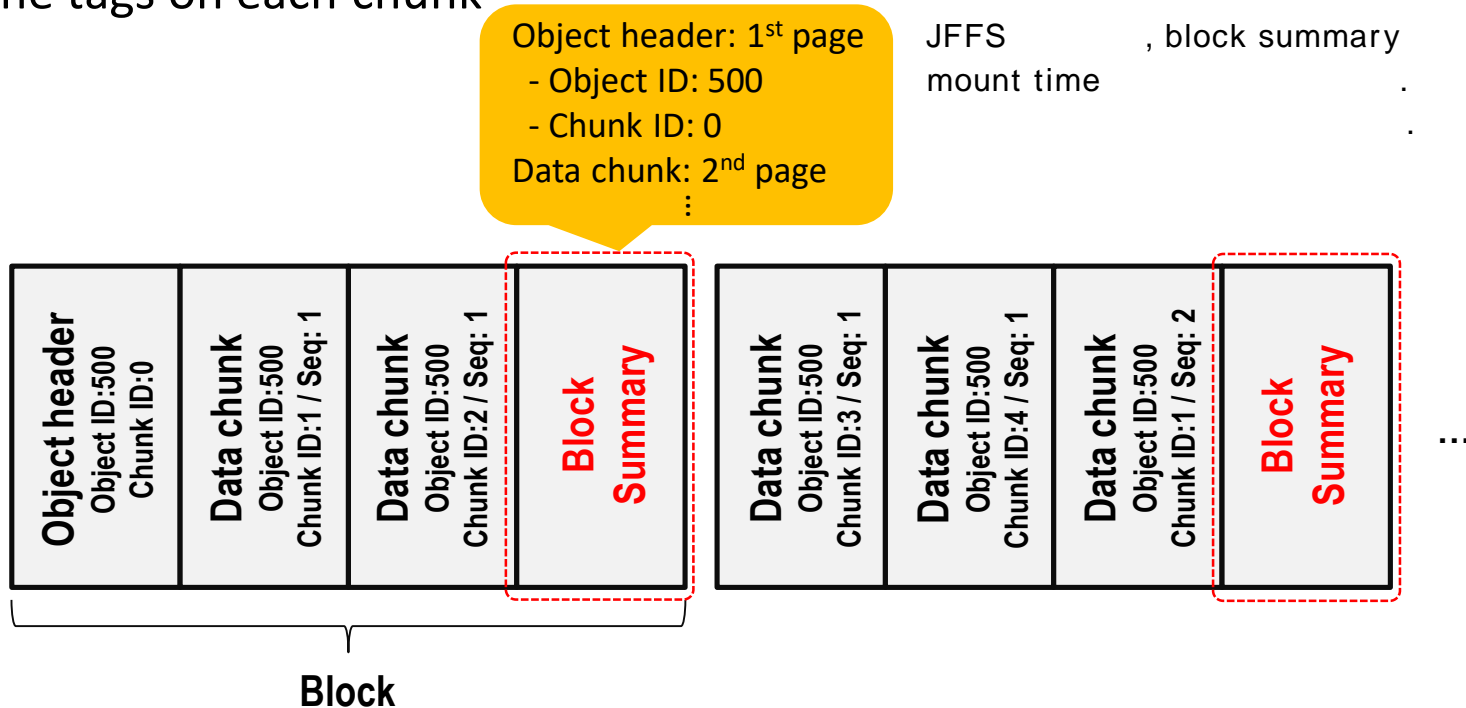
# YAFFS2: File System Layout

- Maintain the information about objects and chunks in DRAM



# YAFFS2: Block Summary

- A block summary (including object id and chunk id) for chunks in the block are written to the last chunk
- This allows all the tags for chunks in that block to be read in one hit, avoiding full disk scan
- If a block summary is not available for a block, then a full scan is used to get the tags on each chunk



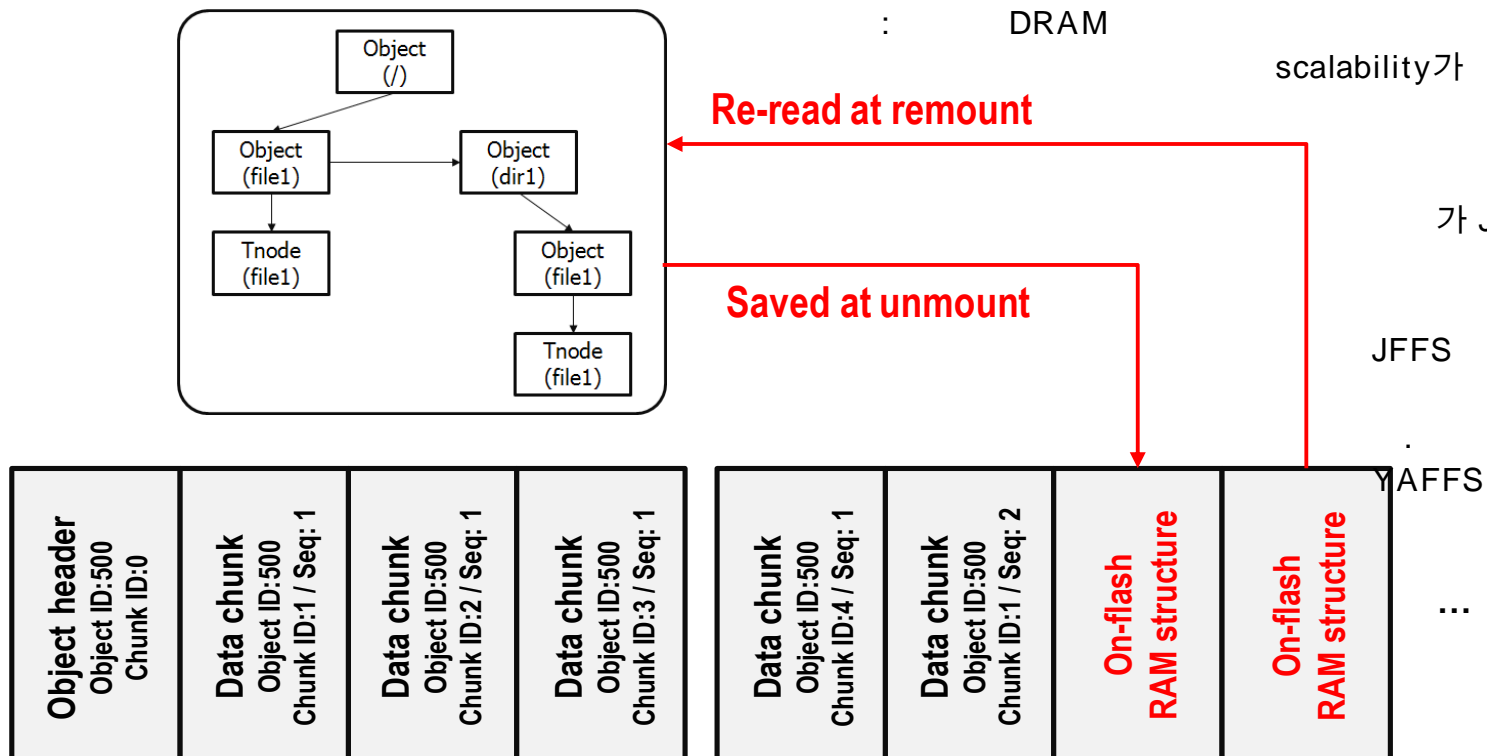
# YAFFS2: Checkpoints

- DRAM structures are saved on flash at unmount
- Structures re-read, avoiding boot scan
- Lazy loading also reduces mount time

DRAM  
boot scan

Flash

Failure ? summary block



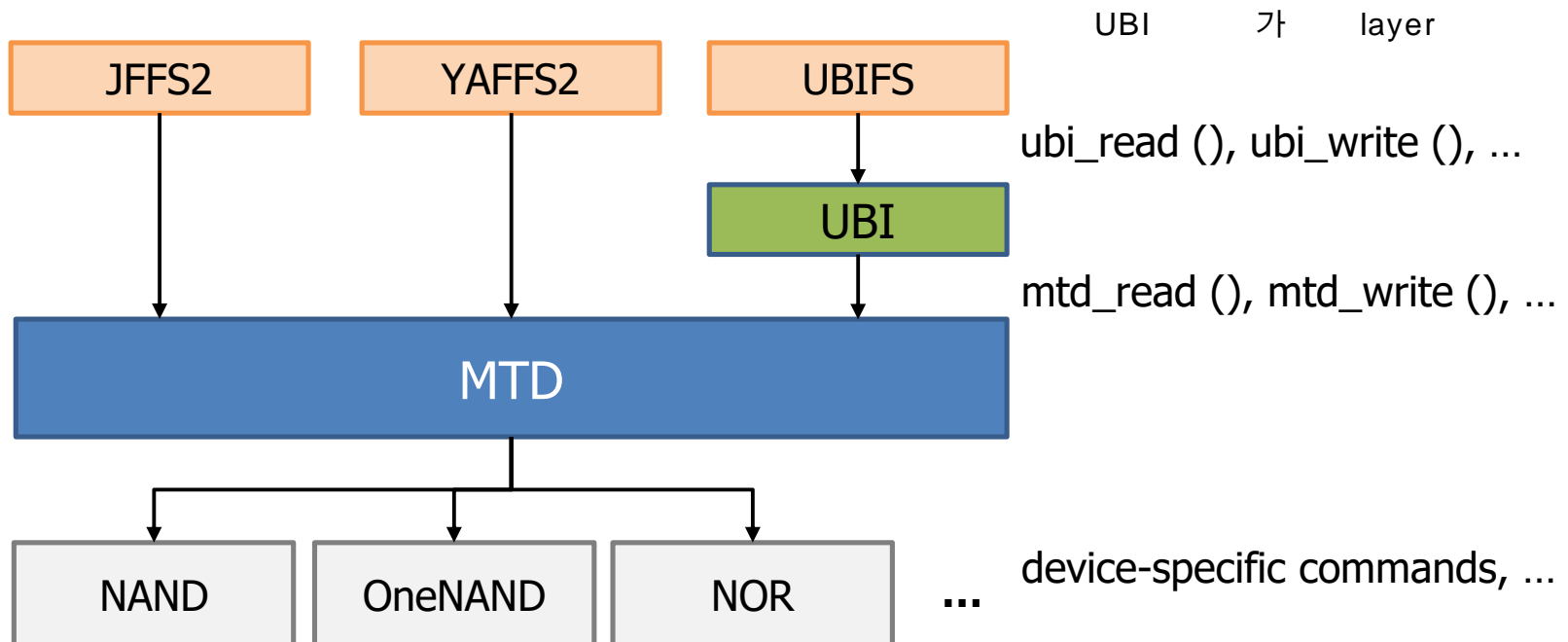


# UBIFS: Unsorted Block Image File System

- **A new flash file system developed by Nokia** <sup>FFS</sup>
  - Considered as the next generation of JFFS2
- **New features of UBIFS**
  - Scalability
    - Scale well with respect to flash size
    - Memory size and mount time do not depend on flash size
  - Fast mount
    - Do not have to scan the whole media when mounting
  - Write-back support
    - Dramatically improve the throughput of the file system in many workloads
  - ...

# UBI: Unsorted Block Image

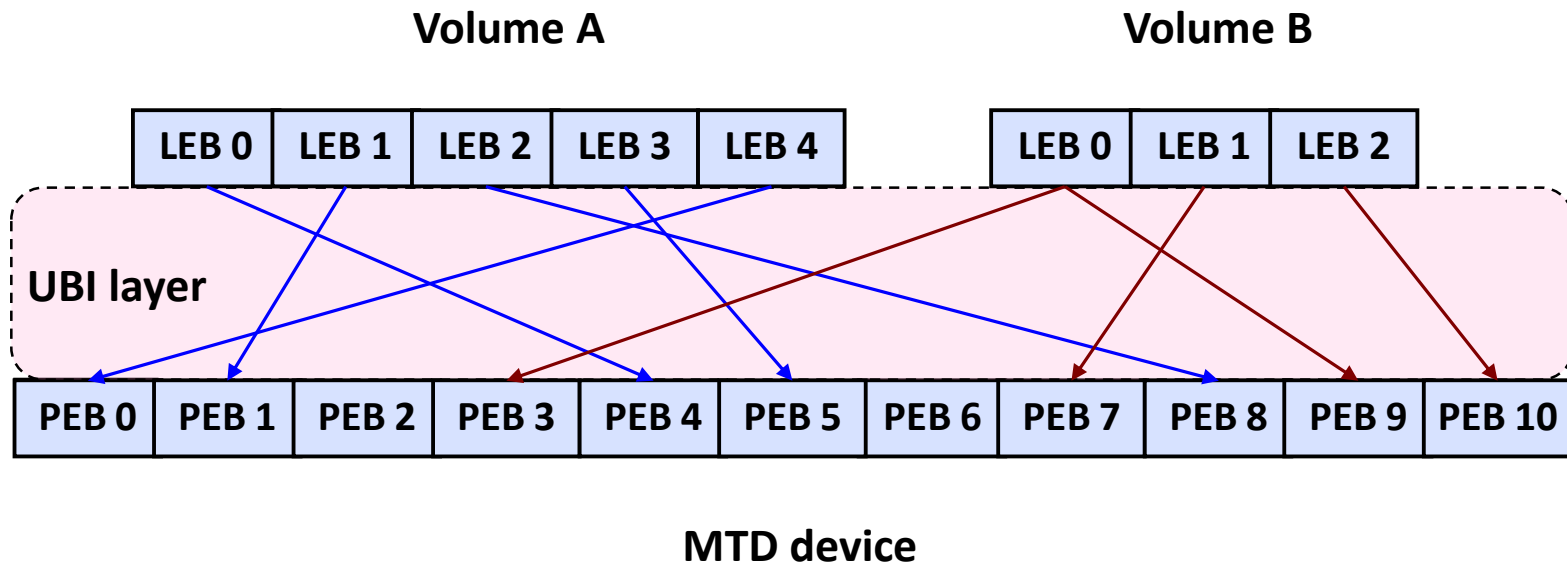
- **UBIFS runs on top of UBI volume**
  - UBI supports multiple volumes, bad block management, wear-leveling, and bit-flips error management
  - The upper level software can be simpler with UBI



# How UBI works

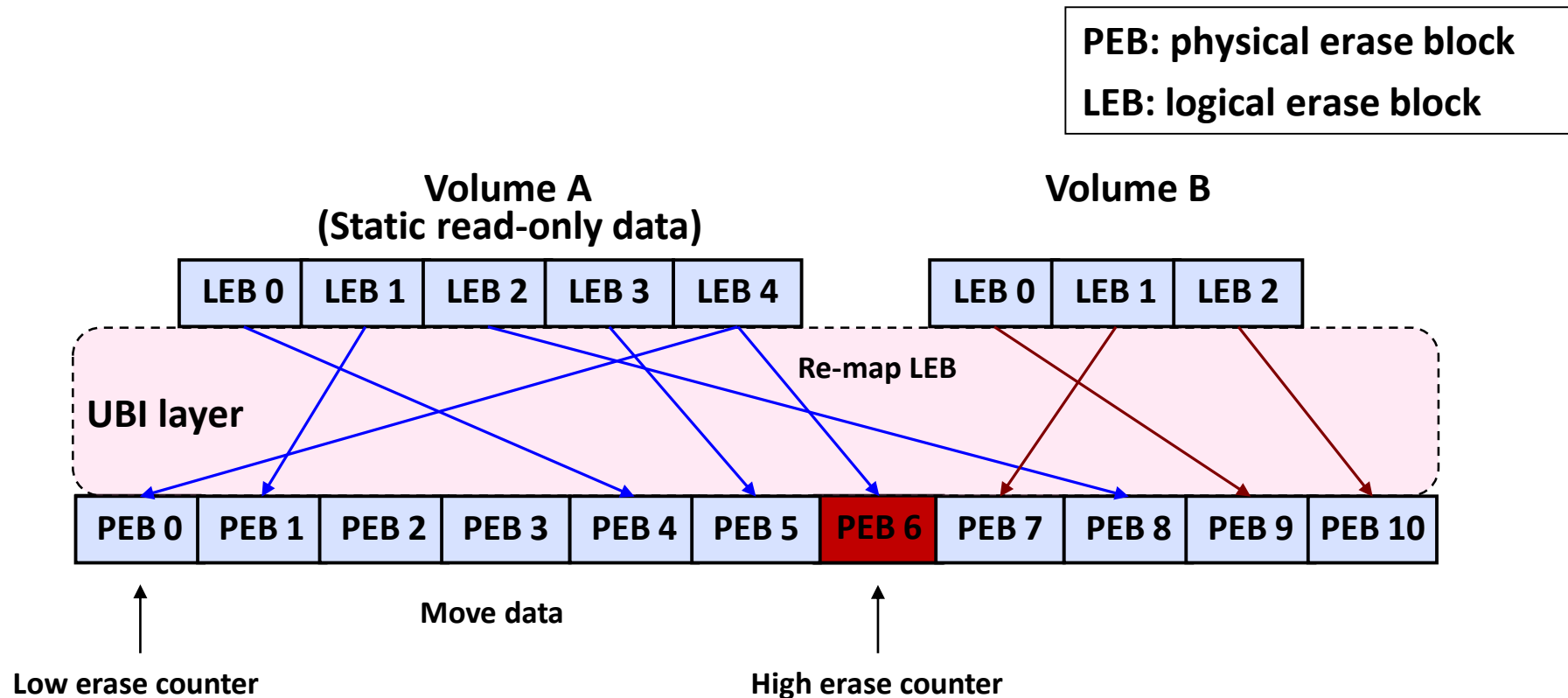
- Logical erase blocks (LEBs) are mapped to physical erase blocks (PEBs)
  - Any LEB can be mapped to any PEB

PEB: physical erase block  
LEB: logical erase block



# How UBI works

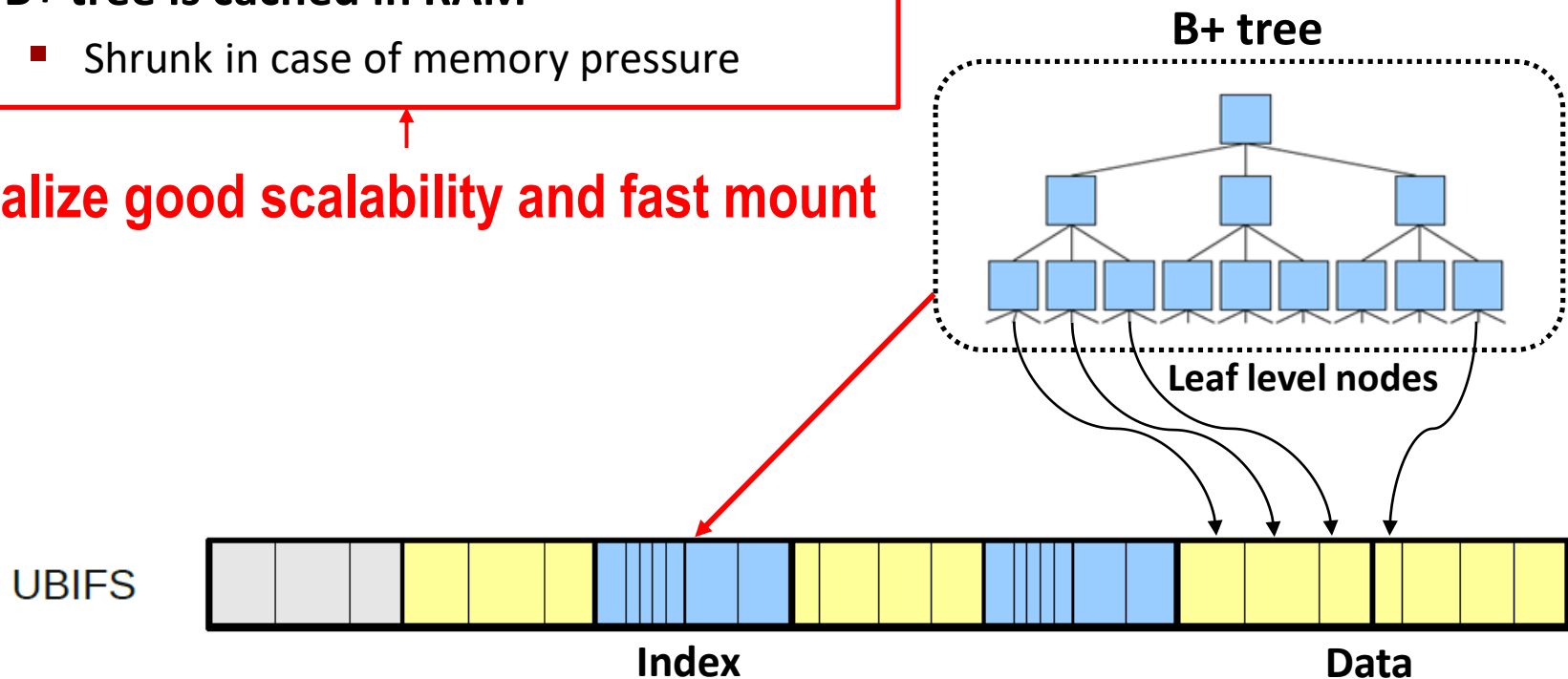
- UBI has its own wear-leveling algorithm that moves the data kept in highly erased blocks to lower one



# UBIFS: Indexing with B+ Tree

- UBIFS index is a B+ tree and is stored on NAND flash
  - c.f., JFFS2 does not store the index on flash
- Leaf level contains data
- Full scanning is not needed
- B+ tree is cached in RAM
  - Shrunk in case of memory pressure

Realize good scalability and fast mount



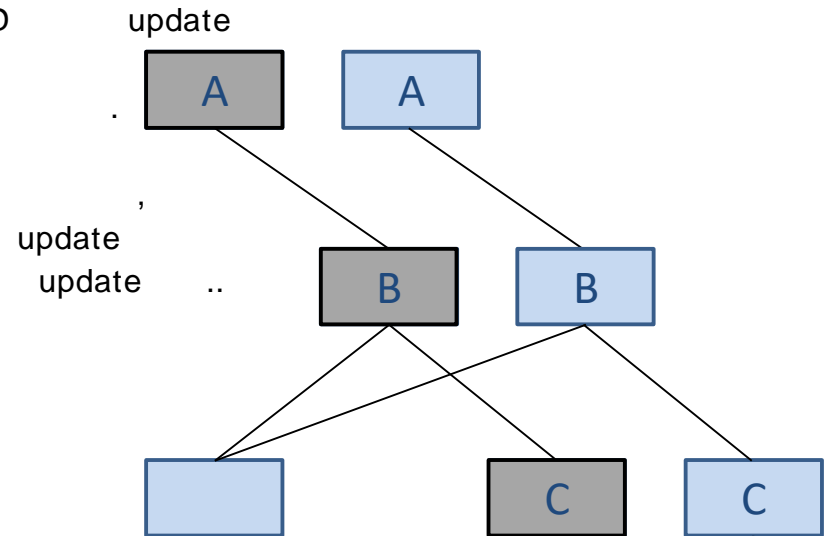
# UBIFS: Wandering Tree

## ■ How to find the root of the tree?

- Write data node "D"
- Old "D" becomes obsolete
- Write indexing node "C"
- Old "C" becomes obsolete
- Write indexing node "B"
- Old "B" becomes obsolete
- Write indexing node "A"
- Old "A" becomes obsolete

data node D  
D, C, B, A  
update

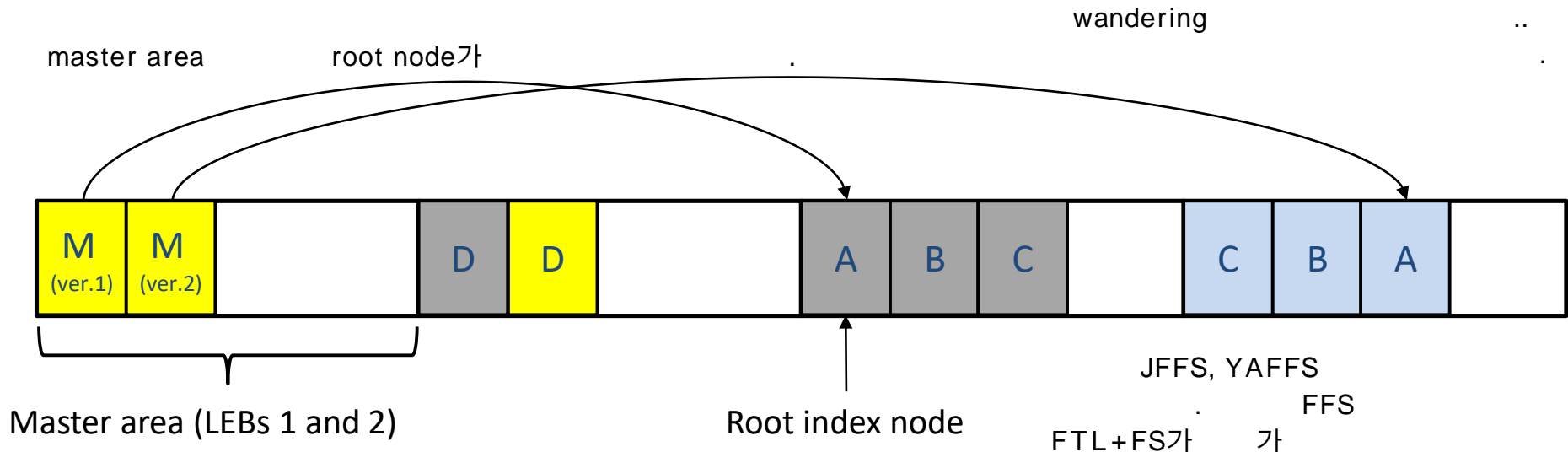
Tree  
node



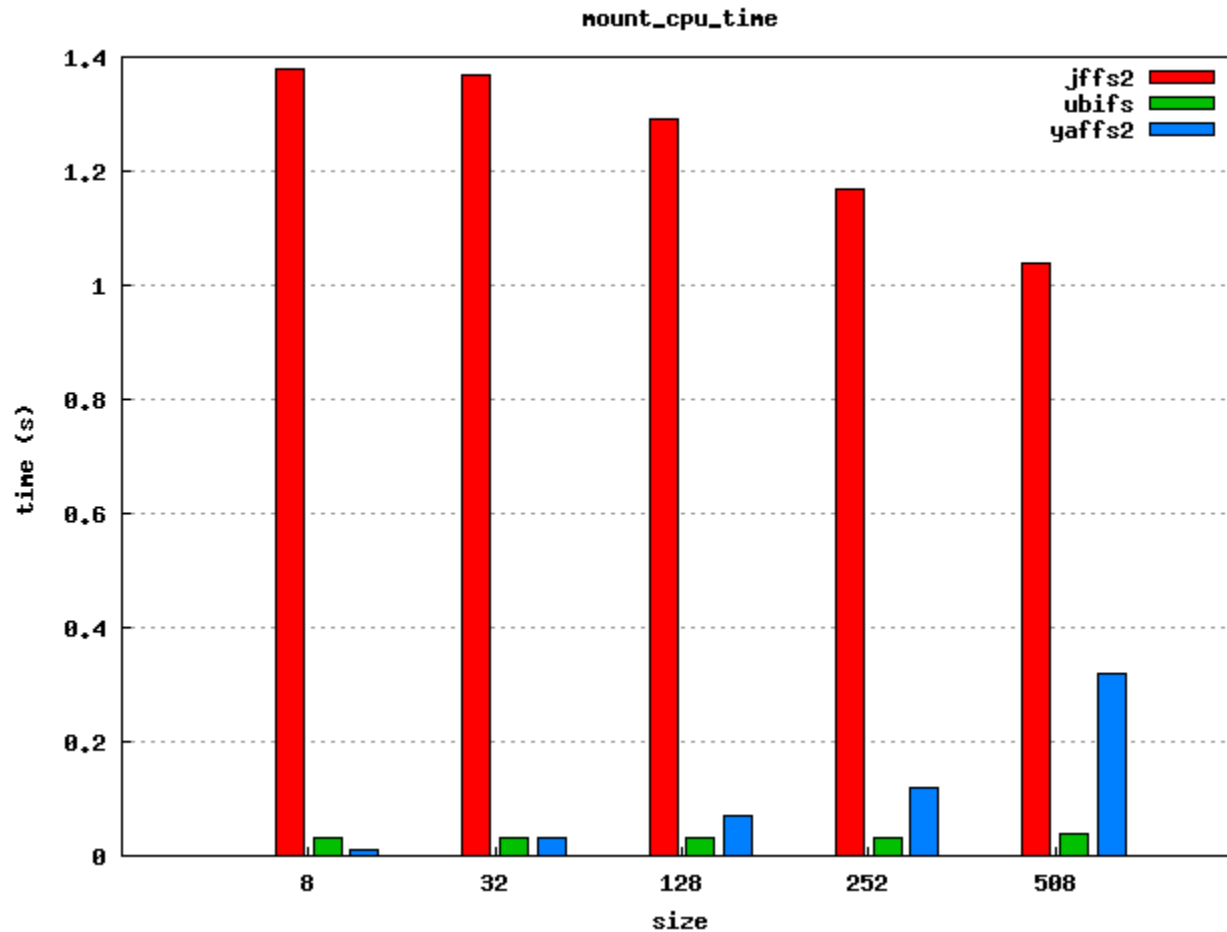
The position of the tree in flash is changed

# UBIFS: Master Area

- LEBs 1/2 are reserved for *a master area* pointing to *a root index node*
  - The master area can be quickly found on mount since its location is fixed (LEBs 1 and 2)
  - Using the root index node, B+ tree can be quickly constructed
- A mater area could have multiple nodes
  - A valid master node is found by scanning master area



# Mount Time Comparison





# Summary

- **JFFS2: Journaling Flash File System version 2**
  - Commonly used for low volume flash devices
  - Compression is supported
  - Long mount time & High memory consumption
  
- **YAFFS2: Yet Another Flash File System version 2**
  - Fast mount time with check-pointing
  
- **UBIFS: Unsorted Block Image File System**
  - Fast mount time and low memory consumption by adopting B+ tree indexing

# Outline

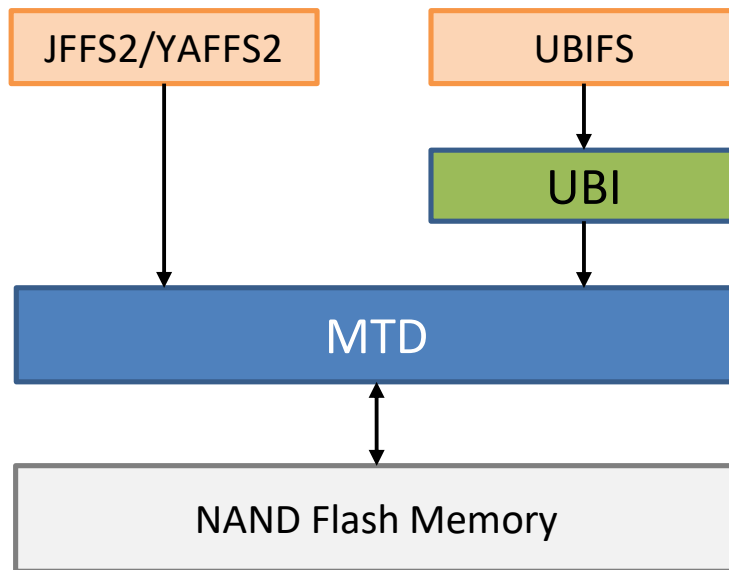
- Traditional Flash File Systems
- **SSD-Friendly Flash File Systems**
  - F2FS: Flash-friendly File System
- Reference

# F2FS: Flash-friendly File System

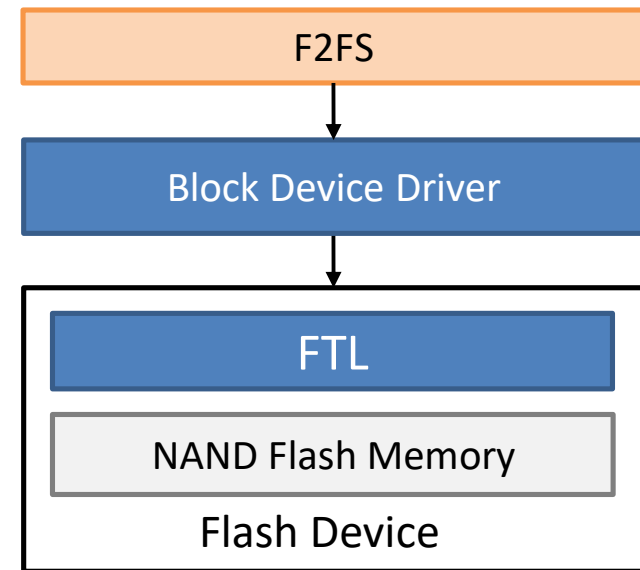
## ■ Log-structured file system for FTL devices

FFS, LFS  
block i/o subsystem

- Unlike other flash file systems, it runs atop FTL-based flash storage and is optimized for it
- Exploit system-level information for better performance and reliability (e.g., better hot-cold separation, background GC, ...)



Traditional flash file system



F2FS

# Design Concept of F2FS

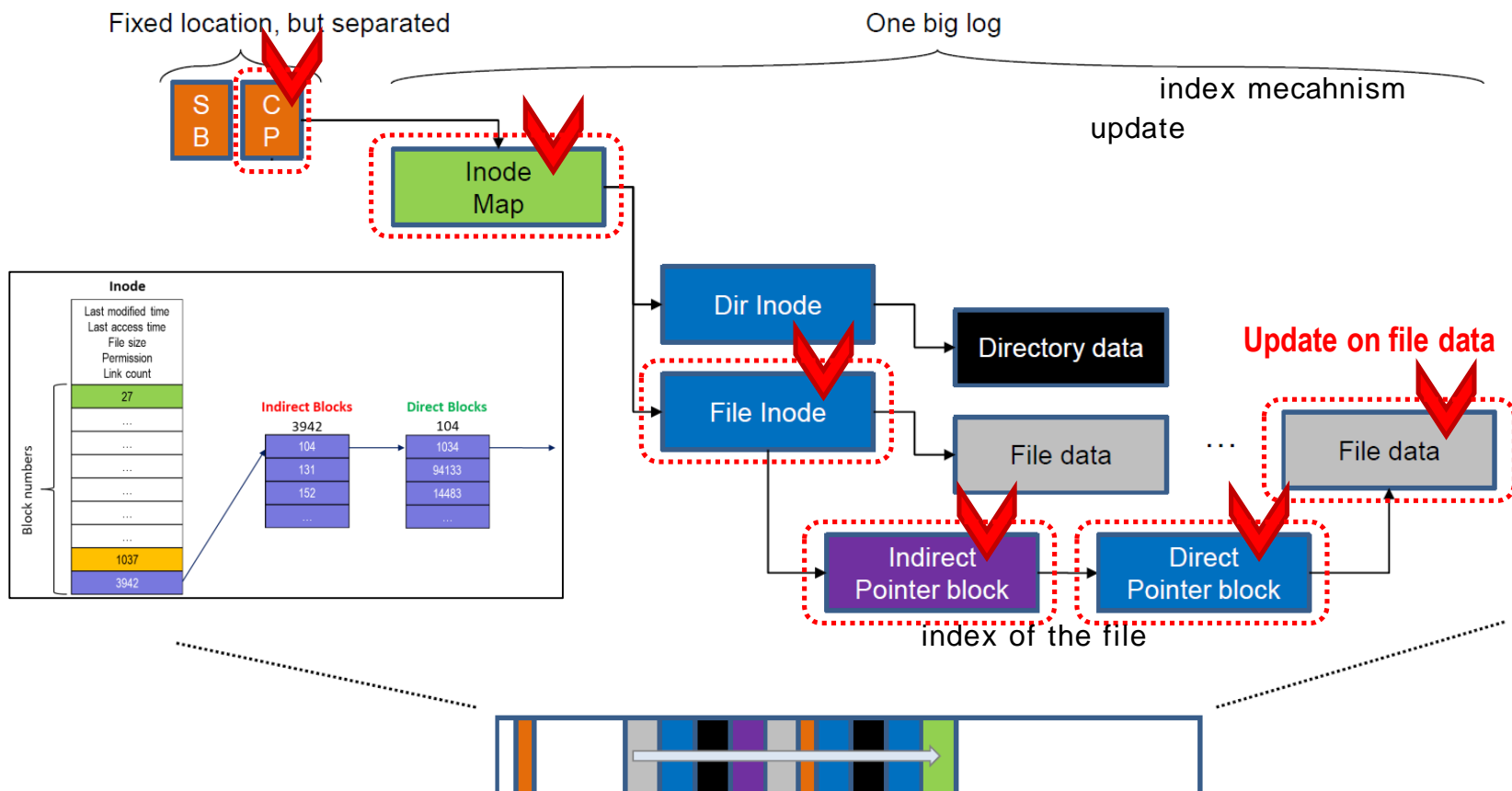
## ■ Designed to take advantage of both of two approaches

- Exploit high-level system information
- Flash-aware storage management
- Handle new NAND flash without redesign

	File System + FTL	Flash File System
Method	- Access a flash device via FTL	- Access a flash device directly
Pros	<ul style="list-style-type: none"> <li>- High interoperability</li> <li>- No difficulties in managing recent NAND flash with new constraints</li> </ul>	<ul style="list-style-type: none"> <li>- High-level optimization with system-level information</li> <li>- Flash-aware storage management</li> </ul>
Cons	<ul style="list-style-type: none"> <li>- Lack of system-level information</li> <li>- Flash-unaware storage management</li> </ul>	<ul style="list-style-type: none"> <li>- Low interoperability</li> <li>- Must be redesigned to handle new constraints</li> </ul>

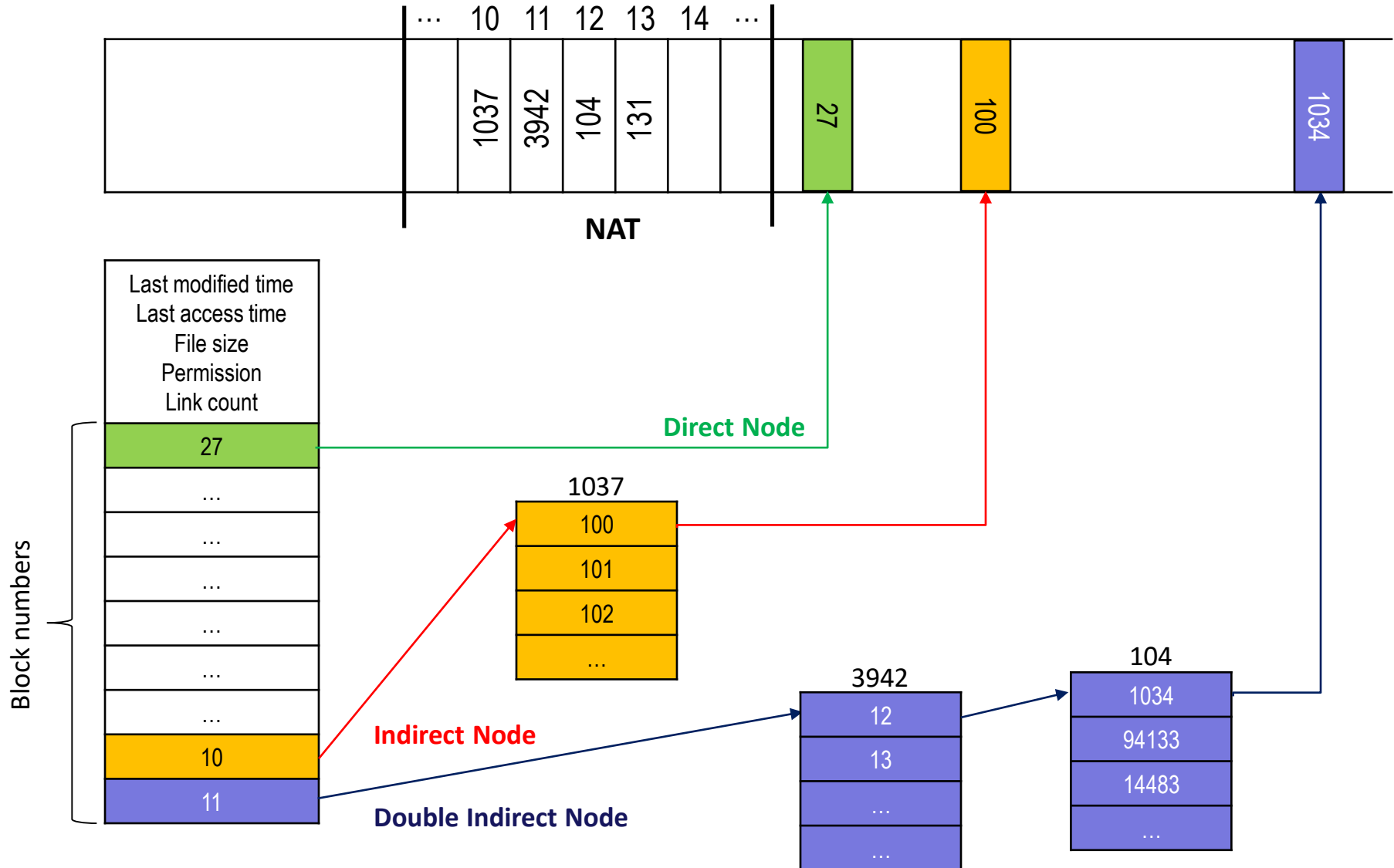
# Index Structure of LFS

- LFS and UBIFS suffer from the wandering tree problem
  - Update on a file causes several extra writes

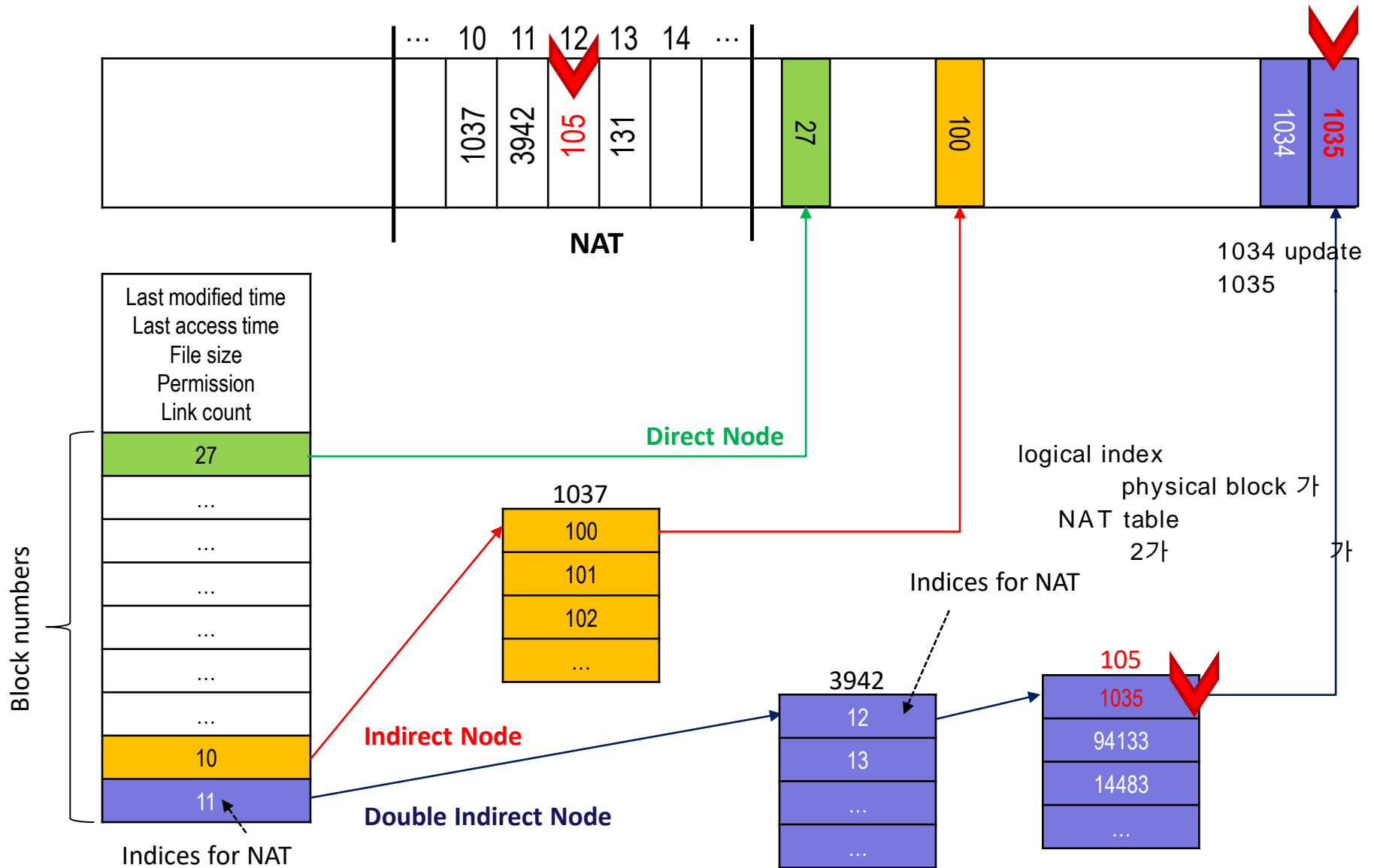


- 
- In-place update**  
Fixed location w/ locality
- Actual block address**
- wandering problem
- wandering problem
- SB CP NAT
- Dir Inode
- Directory data
- File Inode
- File data
- ...
- File data
- Indirect Node
- Offset in NAT
- Offset in NAT
- Direct Node

# Index Structure of F2FS (Cont)



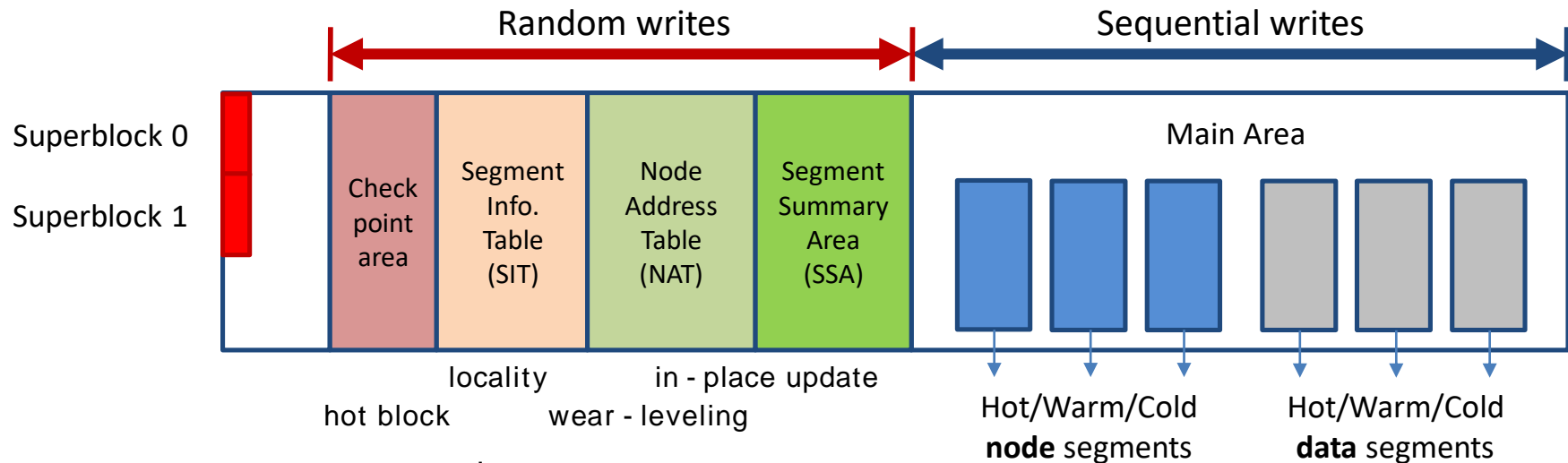
# Index Structure of F2FS (Cont)





# Logical Storage Layout

- **File-system's metadata is located together for locality**
  - Use an “in-place update” strategy for metadata
- **Files are written sequentially for performance**
  - Use an “out-of-place-update” strategy to exploit high throughput of multiple NAND chips
  - Six active logs for static hot and cold data separation



# Cleaning Process

## ■ Background cleaning process

- A kernel thread doing the cleaning job periodically at idle time

## ■ Victim selection policies

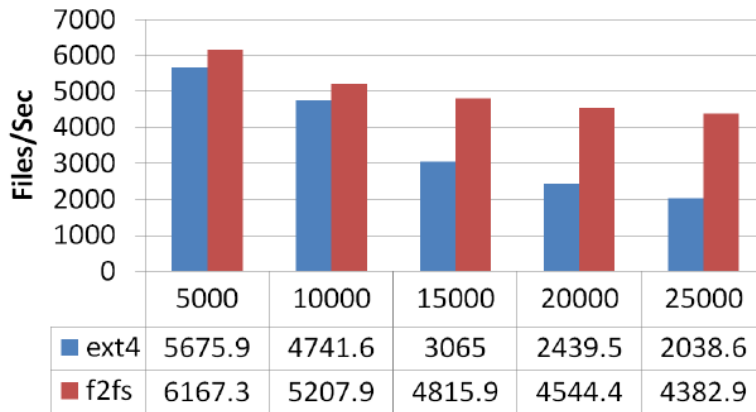
- Greedy algorithm for foreground cleaning job
  - Reduce the amount of data moved for cleaning
- Cost-benefit algorithm for background cleaning job
  - Reclaim obsolete space in a file system
    - Improve the lifetime of a storage device
    - Improve the overall I/O performance

greedy GC,  
cost - benefit

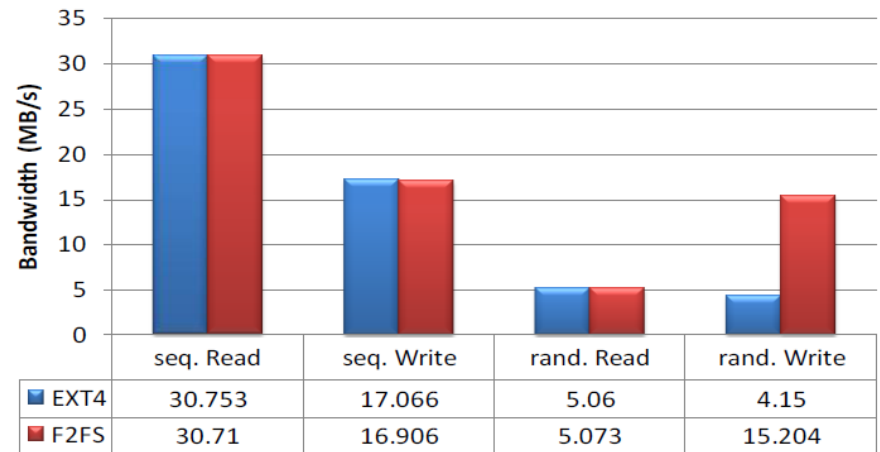
# F2FS Performance

[ System Specification ]

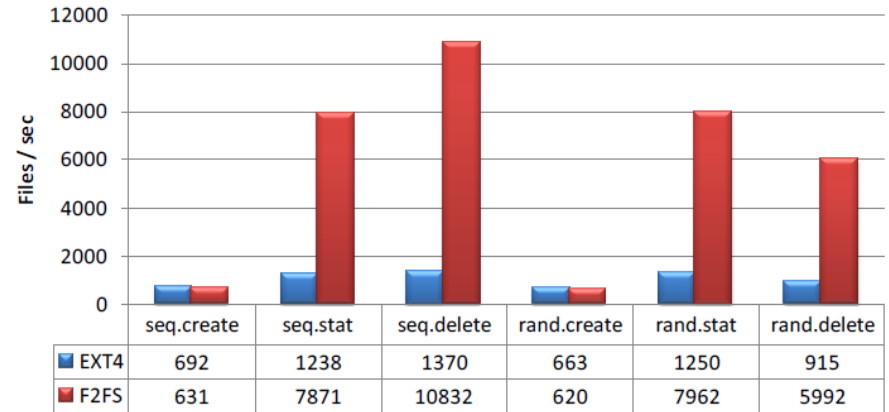
CPU	ARM Coretex-A9 1.2GHz
DRAM	1GB
Storage	Samsung eMMC 64GB
Kernel	<i>Linux 3.3</i>
Partition Size	<i>12 GB</i>



[ fs\_mark ]



[ iozone ]



[ bonnie++ ]

*End of Chapter 7*