# 6. Advanced File Systems

**Special Topics in Computer Systems:**
Modern Storage Systems
(SE820-01)

**Instructor:**

Prof. Sungjin Lee (sungjin.lee@dgist.ac.kr)

# What is a File System?

- **Provides a virtualized logical view of information stored on various storage media, such as disks, tapes, and flash-based SSDs**

- **Two key abstractions have developed over time in the virtualization of storage**
  - *File*: A linear array of bytes, each of which you can read or write
    - Its contents are defined by a creator (e.g., text and binary)
    - It is often referred to as its *inode* number
  - *Directory*: A special file that is a collection of files and other directories
    - Its contents are quite specific – it contains a list of (user-readable name, inode #) pairs (e.g., ("foo", 10))
    - It has a hierarchical organization (e.g., tree, acyclic-graph, and graph)
    - It is also identified by an inode number

# Operations on Files and Directories
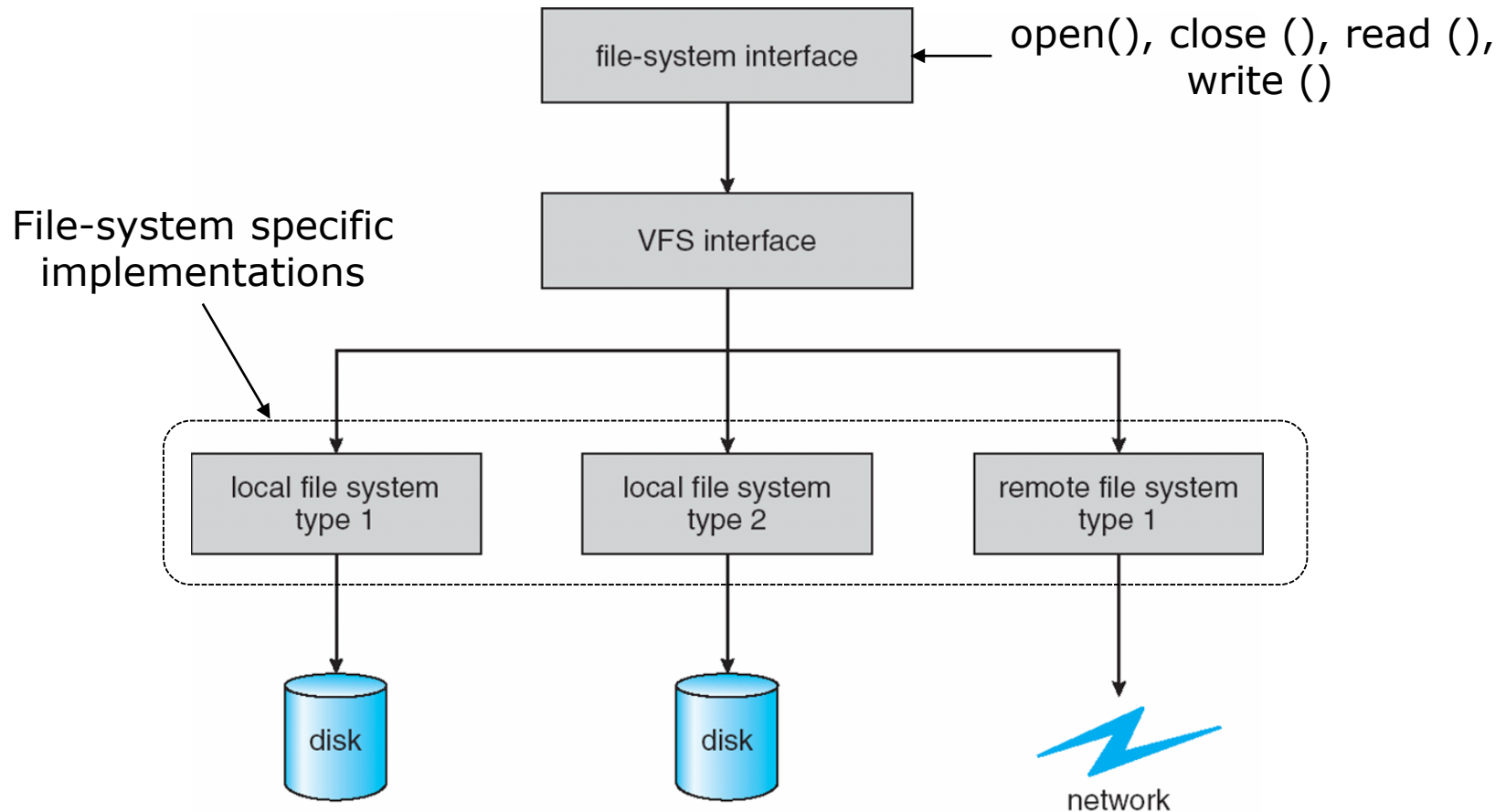
## POSIX Operations on Files

| POSIX APIs | Description |
|---|---|
| creat () | Create a file |
| open () | Create/open a file |
| write () | Write bytes to a file |
| read () | Read bytes from a file |
| lseek () | Move byte position inside a file |
| unlink () | Remove a file |
| truncate () | Resize a file |
| close () | Close a file |
| | … |

## POSIX Operations on Directories

| POSIX APIs | Description |
|---|---|
| opendir () | Open a directory for reading |
| closedir () | Close a directory |
| readdir () | Read one directory entry |
| rewinddir () | Rewind a directory so it can be reread |
| mkdir () | Create a new directory |
| rmdir () | Remove a directory |
| | … |

# Virtual File System

- **The POSIX API is to the VFS interface, rather than any specific type of file system**
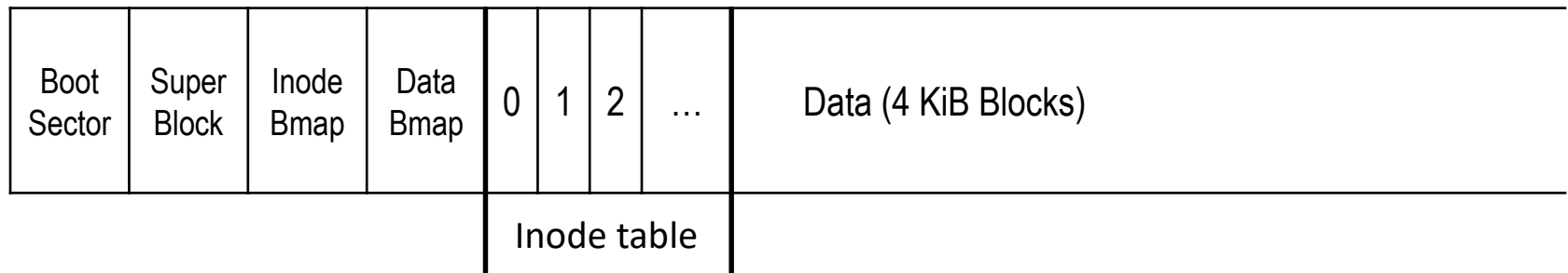


page_quality

4

# Evolution of File System

- **UNIX File System**

- **Fast File System**

- **Journaling File System**

- **Extents File System**

- **Log-structured File Systems**
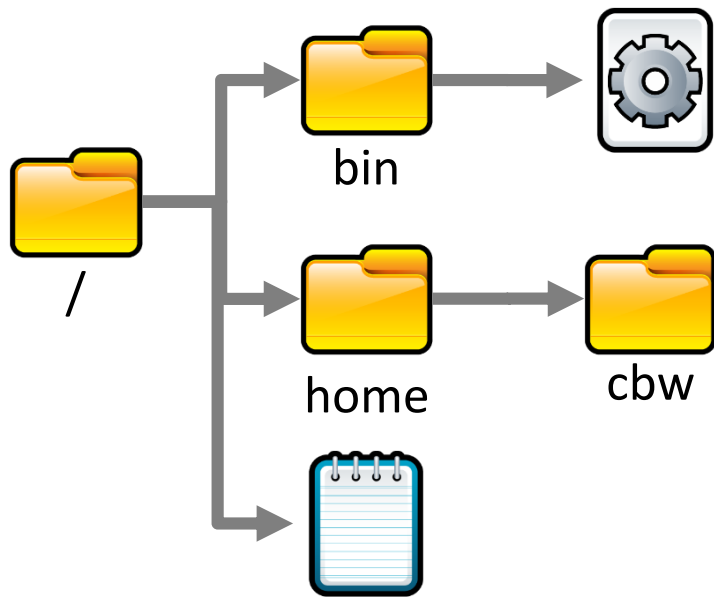
# Evolution of File System

- **UNIX File System ← ext1**
- **Fast File System ← ext2**
- **Journaling File System ← ext3**
- **Extents File System ← ext4**
- **Log-structured File Systems ← WAFL, F2FS, …**

# UNIX File System

■ A traditional file system first developed for UNIX systems

| Boot Sector | Super Block | Inode Bmap | Data Bmap | 0 | 1 | 2 | … | Data (4 KiB Blocks) |
|---|---|---|---|---|---|---|---|---|

Inode table
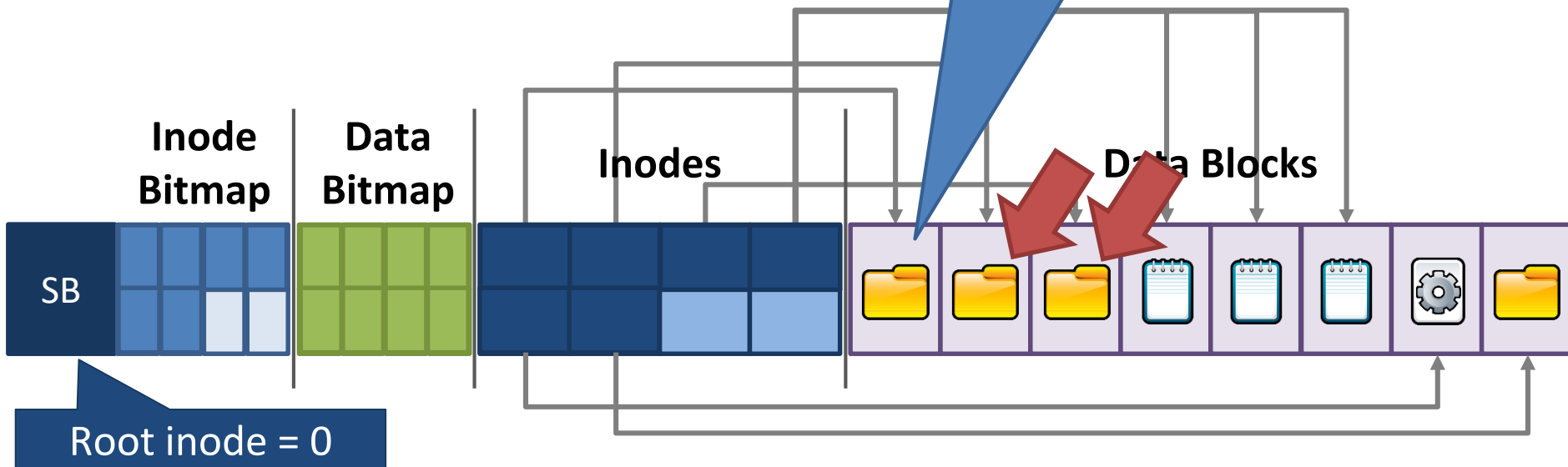
- **Boot sector**: Information to be loaded into RAM to boot up the OS
- **Superblock**: File system's metadata (e.g., file system type, size, …)
- **Inode & data Bmaps**: Keep the status of blocks belonging to an inode table and data blocks
- **Inode table**: Keep file's metadata (e.g., size, permission, …) and data block pointers
- **Data blocks**: Keep users' file data

- Directories are files
- Contains the list of entries in the directory

- Each inode can directly point to 12 blocks
- Can also indirectly point to blocks at 1, 2, and 3 levels of depth

**Inode Bitmap**

**Data Bitmap**

**Inodes**

**Data Blocks**

SB

Root inode = 0

# Inode & Block Pointers

# Good and Bad of UNIX File System

- **The Good – ext file system (inodes) support:**
  - All the typical file/directory features
  - Hard and soft links
- **The Bad: poor locality**
  - It is not optimized for spinning disks
  - inodes and associated data are far apart on the disk!

# Evolution of File System

- **UNIX File System ← ext1**
- **Fast File System ← ext2**
- **Journaling File System ← ext3**
- **Extents File System ← ext4**
- **Log-structured File Systems ← WAFL, F2FS, …**

# Fast File System (FFS)

- **Fast File System (FFS) developed at Berkeley in 1984**
  - First attempt at a disk aware file system
  HDD
  - Optimized for performance on spinning disks

- **Observation:**
  - Processes tend to access files that are in the same (or close) directories

- **Key idea:**
  - place groups of directories and their files into cylinder groups
  - Introduced into ext2, called block groups

# FFS Locality for SEER Traces



FFS Locality

**Observation 3:**
About 65% of file accesses
were shorter than
two distances

**Observation 2:**
40% of file accesses were to
either the same file or to
one in the same directory

**Observation 1:**
7% of file accesses
were to the same file

```
The distance of two file access is 1
proc/src/foo.c
proc/src/bar.c
```

```
The distance of two file access is 2
proc/src/foo.c
proc/obj/foo.o
```

# Block Groups

unix fs    layer
block                     .
directory          locality
block                file                          .

- ## In ext, there is a single set of key data structures

  - One data bitmap, one inode bitmap
  - One inode table, one array of data blocks

- ## In ext2, each block group contains its own key data structures

| Inode Bitmap | Data Bitmap | Inodes | Data Blocks |
|---|---|---|---|

| SB | Block Group 1 | Block Group 2 | Block Group 3 | Block Group 4 | Block Group 5 | Block Group 6 |
|---|---|---|---|---|---|---|

# Allocation Policy

- **ext2 attempts to keep related files and directories within the same block group**



| SB | Block Group 1 | Block Group 2 | Block Group 3 | Block Group 4 | Block Group 5 | Block Group 6 |

# The Good and the Bad of FFS

- **The good – ext2 supports:**
  - All the features of ext with even better performance because of increased spatial locality

- **The bad**

  - Large files must cross block groups
  - As the file system becomes more complex, the chance of file system corruption grows
    - E.g. invalid inodes, incorrect directory entries, etc.

# Consistent Update Problem

- What happens if sudden power loss occurs while writing data to a file

write ($\overset{\text{inode}}{10}$, "foo", strlen ("foo") );

| Boot Sector | Super Block | | Inode Bmap | | Data Bmap | 0 | 1 | 2 | ... | foo | Data (4 KiB Blocks) |
|---|---|---|---|---|---|---|---|---|---|---|---|

Inode table

ext2

The file system will be inconsistent!!!
→ **Consistent update problem**

# Evolution of File System

- **UNIX File System ← ext1**
- **Fast File System ← ext2**
- **Journaling File System ← ext3**
- **Extents File System ← ext4**
- **Log-structured File Systems ← WAFL, F2FS, …**

# How to Ensure Consistency after a Crash?

- **Strategy 1: Don't bother to ensure consistency**
  - Accept that the file system may be inconsistent after a crash
  - Run a program that fixes the file system during bootup
  - *File system checker (fsck)*     , fsck     inconsistent problem     .

- **Strategy 2: Use a transaction log to make multi-writes atomic**
  - Log stores a history of all writes to the disk
  - After a crash, the log can be "replayed" to finish updates
  - *Journaling file system*     DB     . WAL

Strategy #1:
# File System Checker

- **Key idea: fix inconsistent file systems during bootup**
  - Unix utility called *fsck* (*chkdsk* on Windows)
  - Scans the entire file system multiple times, identifying and correcting inconsistencies

- **Why during bootup?**
  - No other file system activity can be going on
  - After fsck runs, bootup/mounting can continue

# File System Checker Tasks

- **Superblock:**
  - Validate the superblock, replace it with a backup if it is corrupted
- **Free blocks and inodes:**
  - Rebuild the bitmaps by scanning all inodes
- **Reachability:**
  - Make sure all inodes are reachable from the root of the file system
- **Inodes:**
  - Delete all corrupted inodes, and rebuild their link counts by walking the directory tree
- **Directories:**
  - Verify the integrity of all directories
- … and many other minor consistency checks

# The Good and the Bad of fsck

- **Advantages of *fsck***
  - Doesn't require the file system to do any work to ensure consistency
  - Makes the file system implementation simpler

- **Disadvantages of *fsck***
  - Very complicated to implement the *fsck* program
    - Many possible inconsistencies that must be identified
    - Many difficult corner cases to consider and handle
  - *fsck* is ***super slow***
    - Scans the entire file system multiple times
    - Imagine how long it would take to fsck a 40 TB RAID array

Strategy #2:
# Journaling File System

- Journaling file systems address the consistent update problem by adopting an idea of *write-ahead logging (or journaling)* from database systems

- Ext3, Ext4, ReiserFS, XFS, and NTFS are based on journaling

write (10, "foo", strlen ("foo") );

Journal write & commit

TxB          TxE

| Boot Sector | Super Block | Inode Bmap | Data Bmap | 0 | 1 | 2 | … | foo | Data (4 KiB Blocks) | 0 | foo | |

Inode table

Journaling space

Checkpoint

# Write-Ahead Log (WAL)

- **Key idea: writes to disk are first written into a log**
  - After the log is written, the writes execute normally
  - In essence, the log records transactions

- **What happens after a crash…**
  - If the writes to the log are interrupted?
    - The transaction is incomplete
    - The user's data is lost, but the file system is consistent
  - If the writes to the log succeed, but the normal writes are interrupted?
    - The file system may be inconsistent, but…
    - The log has exactly the right information to fix the problem

# Data Journaling Example

- **Here, assume that we are appending to a file**
  - Three writes: inode v2, data bitmap v2, data $D_2$
- **Before executing these writes, first log them**

**Journal**

| TxB ID=1 | I v2 | B v2 | $D_2$ | TxE ID=1 | |
|---|---|---|---|---|---|

1. Begin a new transaction with a unique ID=k
2. Write the updated meta-data block(s)
3. Write the file data block(s)
4. Write an end-of-transaction with ID=k

# Commits and Checkpoints

- **A transaction is committed after all writes to the log are complete**
- **After a transaction is committed, the OS checkpoints the update**



**Committed!**

**Checkpointed!**

| Journal | TxB | I v2 | B v2 | $D_2$ | | TxE | |

| Inode Bitmap | Data Bitmap | Inodes | | Data Blocks | |
|---|---|---|---|---|---|

v2

$D_1$  $D_2$

checkout    failure
transaction
checkout           .

- **Final step: free the checkpointed transaction**

26

# The Good and the Bad of Journaling

- **Advantages of journaling**
  - Robust, fast file system recovery
    - No need to scan the entire journal or file system
  - Relatively straight forward to implement

- **Disadvantages of journaling**
  - Write traffic to the disk is ***doubled***
    - Especially the file data, which is probably large
  - Deletes are very hard to correctly log

# Making Journaling Faster

■ **Journaling adds a lot of write overhead**

■ **OSes typically batch updates to the journal**

  ▪ Buffer sequential writes in memory, then issue one large write to the log

    memory    buffering

  ▪ Example: ext3 batches updates for 5 seconds

                                    batch updates

■ **Tradeoff between performance and persistence**

  ▪ Long batch interval = fewer, larger writes to the log

    ▪ Improved performance due to large sequential writes

  ▪ But, if there is a crash, everything in the buffer will be lost

# Meta-Data Journaling

- **The most expensive part of data journaling is writing the file data twice**
  - Meta-data is small (~1 sector), file data is large
- **ext3 implements meta-data journaling**

**Journal**

| TxB | I v2 | B v2 | TxE | | | |
|-----|------|------|-----|---|---|---|

| **Inode Bitmap** | **Data Bitmap** | **Inodes** | **Data Blocks** |
|------------------|-----------------|------------|-----------------|

v2

$D_1$ $D_2$

D2                    I, B    checkout        .
                                              .

# Journaling Wrap-Up

- **Today, most OSes use journaling file systems**
  - ext3/ext4 on Linux
  - NTFS on Windows
- **Provides excellent crash recovery with relatively low space and performance overhead**
- **Next-gen OSes will likely move to file systems with copy-on-write semantics**
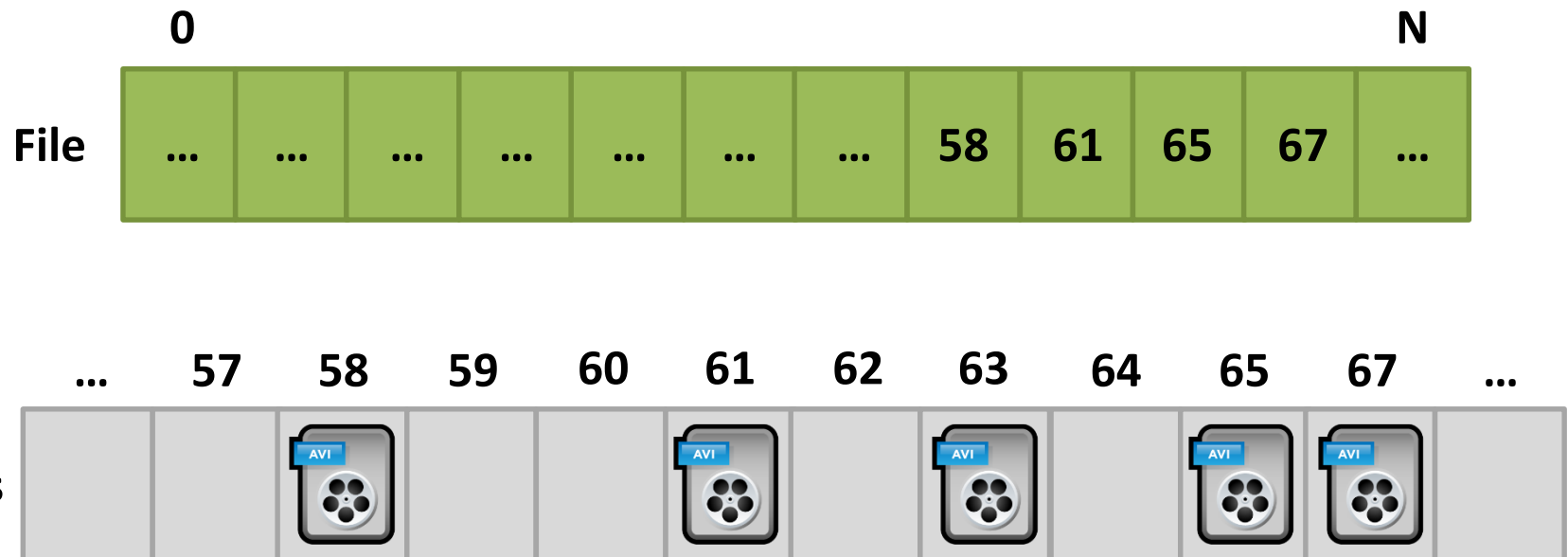  - btrfs and zfs on Linux

# Evolution of File System

- **UNIX File System ← ext1**

- **Fast File System ← ext2**

- **Journaling File System ← ext3**

- **Extents File System ← ext4**    ext3    directory data structure
                                                    .

- **Log-structured File Systems ← WAFL, F2FS, ...**

# Revisiting Inodes

logical      contiguous   , physical       .
ext3      Fragmentation         ,
inode         .

- **Inodes use indirection to acquire additional blocks of pointers**

- **Problem: inodes are not efficient for large files**
  - e.g., For a 100MB file, you need 25600 block pointers (assuming 4KB blocks)

- **This is unavoidable if the file is 100% fragmented**
  - However, what if large groups of blocks are contiguous?

| | 0 | | | | | | | | | | | N |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **File** | … | … | … | … | … | … | … | 58 | 61 | 65 | 67 | … |

| | … | 57 | 58 | 59 | 60 | 61 | 62 | 63 | 64 | 65 | 67 | … |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Blocks** | | | AVI | | | AVI | | AVI | | AVI | AVI | |

# From Pointers to Extents

- **Modern file systems try hard to minimize fragmentation**
  - Since it results in many seeks, thus low performance
- **Extents are better suited for contiguous files**

contiguous file          length          .

**inode**
block 1
block 2
block 3
block 4
block 5
block 6

**inode**
block 1
length 1
block 2
length 2
block 3
length 3

Each extent includes a block pointer and a length

# Revisiting Directories

- **Each directory is a file with a list of entries**
  - Entries are not stored in sorted order
  - Some entries may be blank, if they have been deleted

- **Problem: searching for files in large directories takes O(n) time**
  - Practically, you can't store >10K files in a directory
  - It takes way too long to locate and open files

file entry          linear search          .
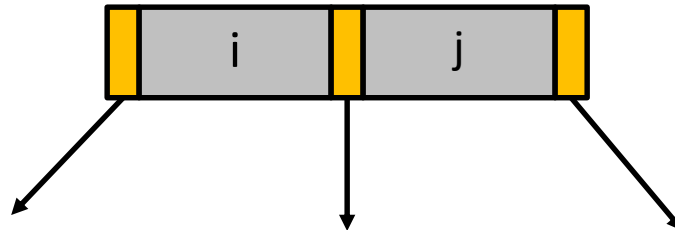                .

# From Lists to B-Trees

ext4                entry   B- tree         .              .

■ **ext4 and NTFS encode directories as B-Trees to improve lookup time to O(log N)**

■ **A B-Tree is a type of balanced tree optimized for disk**

  ▪ Items are stored in sorted order in blocks

■ **Suppose items *i* and *j* are in the root of the tree**

  ▪ The root must have 3 children, since it has 2 items

  ▪ The three child groups contain items *a < i*, *i < a < j*, and *a > j*
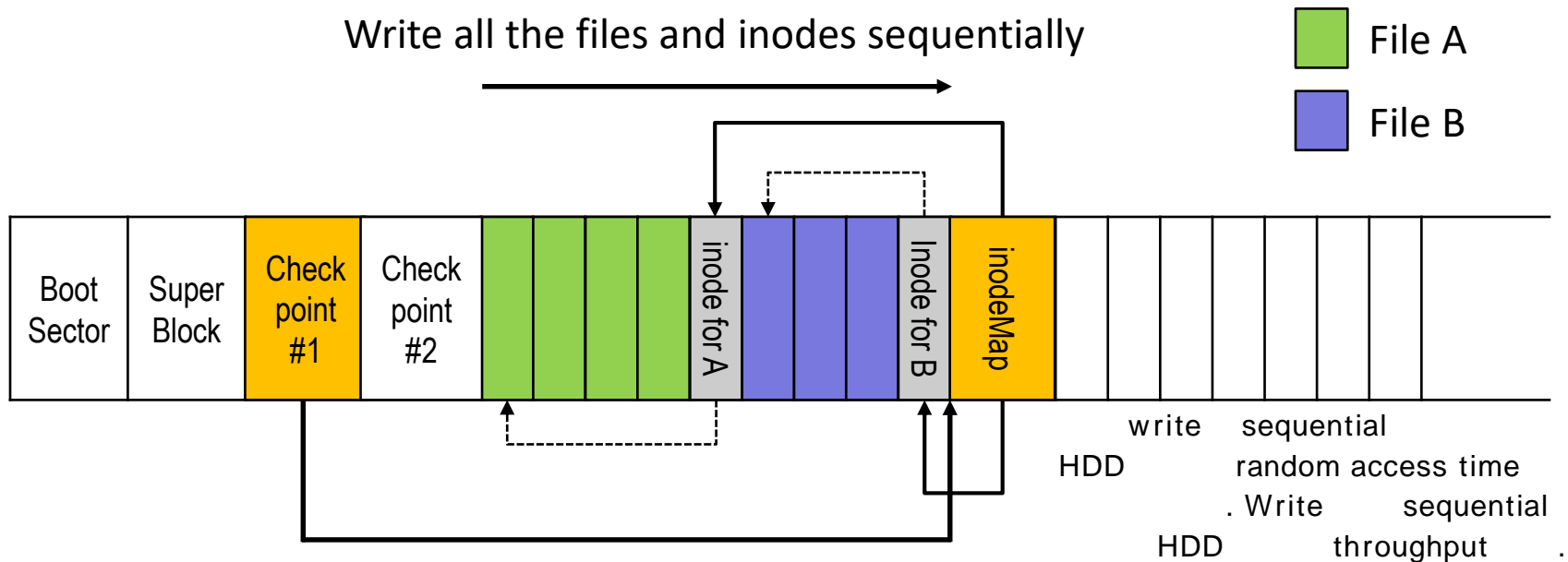
# Evolution of File System

- **UNIX File System ← ext1**
- **Fast File System ← ext2**
- **Journaling File System ← ext3**
- **Extents File System ← ext4**
- **Log-structured File Systems ← WAFL, F2FS, …**

Log- structured FS

90                            .

NetApp                  FS    CoW approach

# Log-structured File System

■ Log-structured file systems (LFS) treat a storage space as a huge log, appending all files and directories sequentially

■ The state-of-the-art file systems are based on LFS or CoW
  ▪ e.g., Sprite LFS, F2FS, NetApp's WAFL, Btrfs, ZFS, …

Write all the files and inodes sequentially

File A

File B

| Boot Sector | Super Block | Check point #1 | Check point #2 | | | | | inode for A | | | | Inode for B | inodeMap | | | | | | | |

write    sequential
HDD           random access time
.  Write           sequential
HDD              throughput       .

# Log-structured File System (Cont.)

- **Advantages**

  Journaling                , Journaling space    logging
  LFS    data/inode    log              .

  - **(+)** No consistent update problem

  - **(+)** No double writes – an LFS itself is a log!

  - **(+)** Provide excellent write performance – disks are optimized for sequential I/O operations

  - **(+)** Reduce the movements of disk headers further (e.g., inode update and file updates)

- **Disadvantages**

  - **(–)** Expensive garbage collection cost

  - **(–)** Slow read performance

# Disadvantages of LFS

File A

File B

Invalid

Flash        FTL                          .

Write sequentially

update     out of place        .

| Boot Sector | Super Block | Check point #1 | Check point #2 | | | | | | | | | | | inode for A | | Inode for B | inodeMap |

obsolete data    GC                 .
, live data              append        .

■ **Expensive garbage collection cost**: invalid blocks must be reclaimed for future writes; otherwise, free disk space will be exhausted

■ **Slow read performance**: involve more head movements for future reads (e.g., when reading the file A)
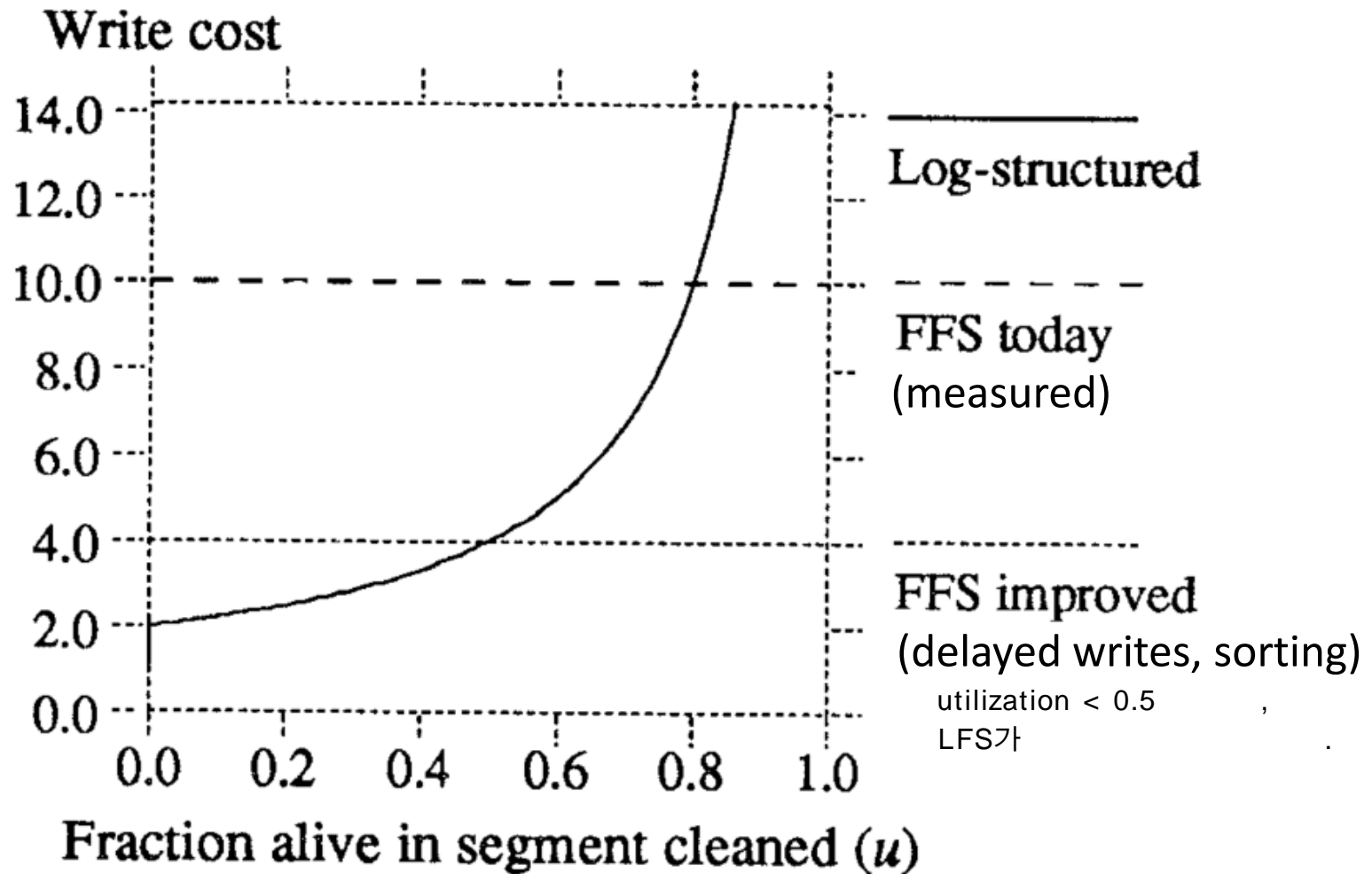
# Write Cost

■ **Write cost with GC is modeled as follows**

■ Note: a segment (seg) is a unit of space allocation and GC

$$\text{write cost} = \frac{\text{total bytes read and written}}{\text{new data written}}$$

GC .
  cost= 1 .
asymptotic. .

$$= \frac{\text{read segs} + \text{write live} + \text{write new}}{\text{new data written}}$$

$$= \frac{N + N*u + N*(1-u)}{N*(1-u)} = \frac{2}{1-u}$$

u= 0      1.
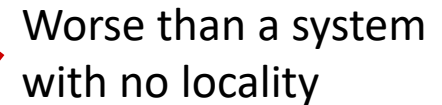u!= 0        .

■ N is the number of segments

■ $\mu$ is the utilization of the segments ($0 \le \mu < 1$)

■ If segments have no live data ($\mu = 0$), write cost becomes 1.0

# Write Cost Comparison



Write cost

Log-structured

FFS today
(measured)

FFS improved
(delayed writes, sorting)
utilization < 0.5            ,
LFS                                  .

Fraction alive in segment cleaned (u)

# Greedy Policy

, victim segments                ?

- The cleaner chooses the *least-utilized segments* and *sorts the live data by age* before writing it out again
- Workloads: 4 KB files with two overwrite patterns
    - (1) **Uniform**: No locality – equal likelihood of being overwritten
    - (2) **Hot-and-cold**: Locality – 10:90



Worse than a system with no locality

No variance:        block    u
      .           Greedy            ,

                              cost      .

The variance in segment utilization

hot- and- cold locality        workload            greedy
      .

# Cost-Benefit Policy

greedy    locality                          , hot                    GC                  invalidated        GC
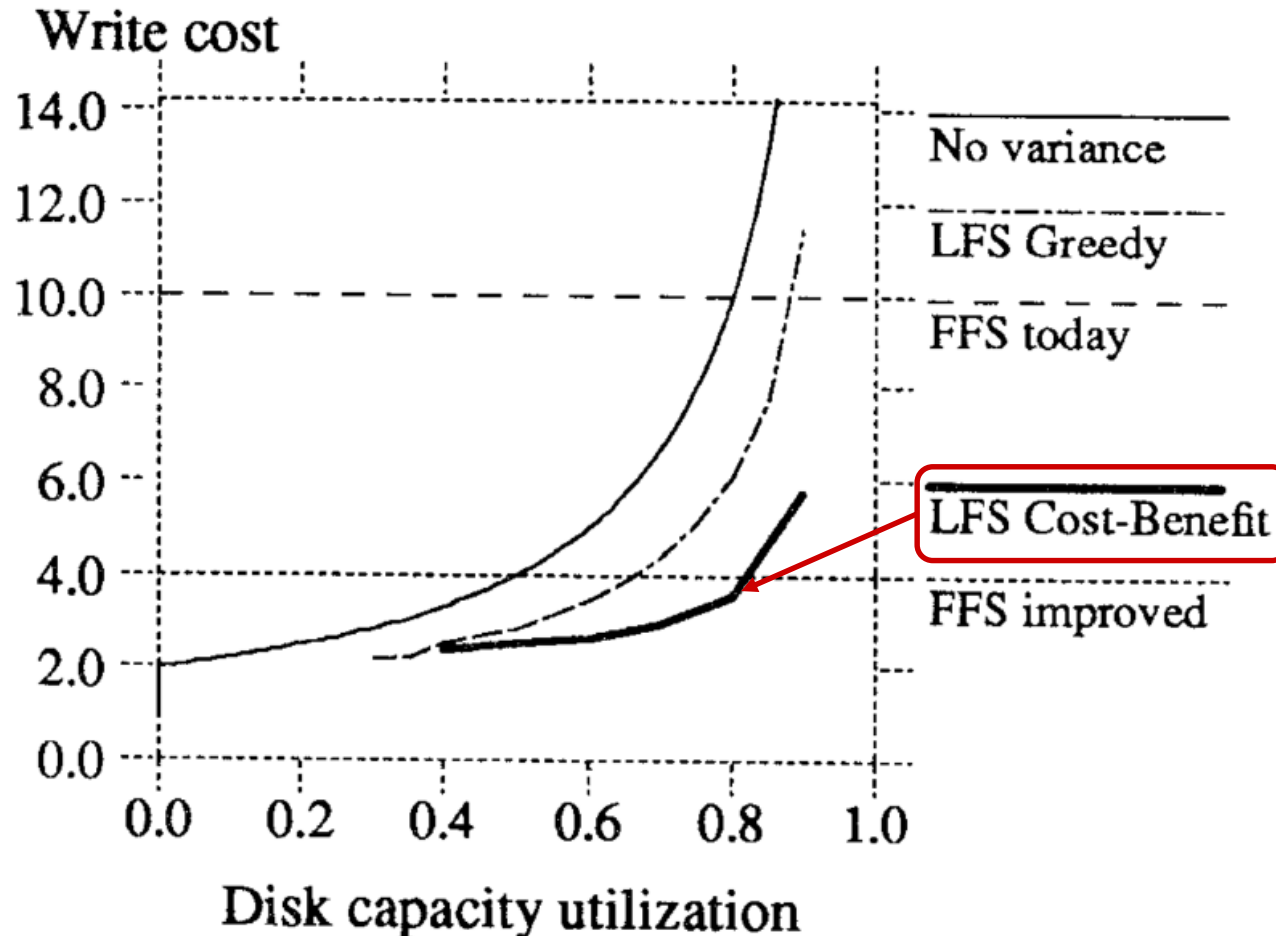              .        ,          cold data    GC                    .            age                      .

- **Hot segments are frequently selected as victims even though their utilizations would drop further**
  - It is necessary to delay cleaning and let more of the blocks die
  - On the other hand, free space in cold segments are valuable
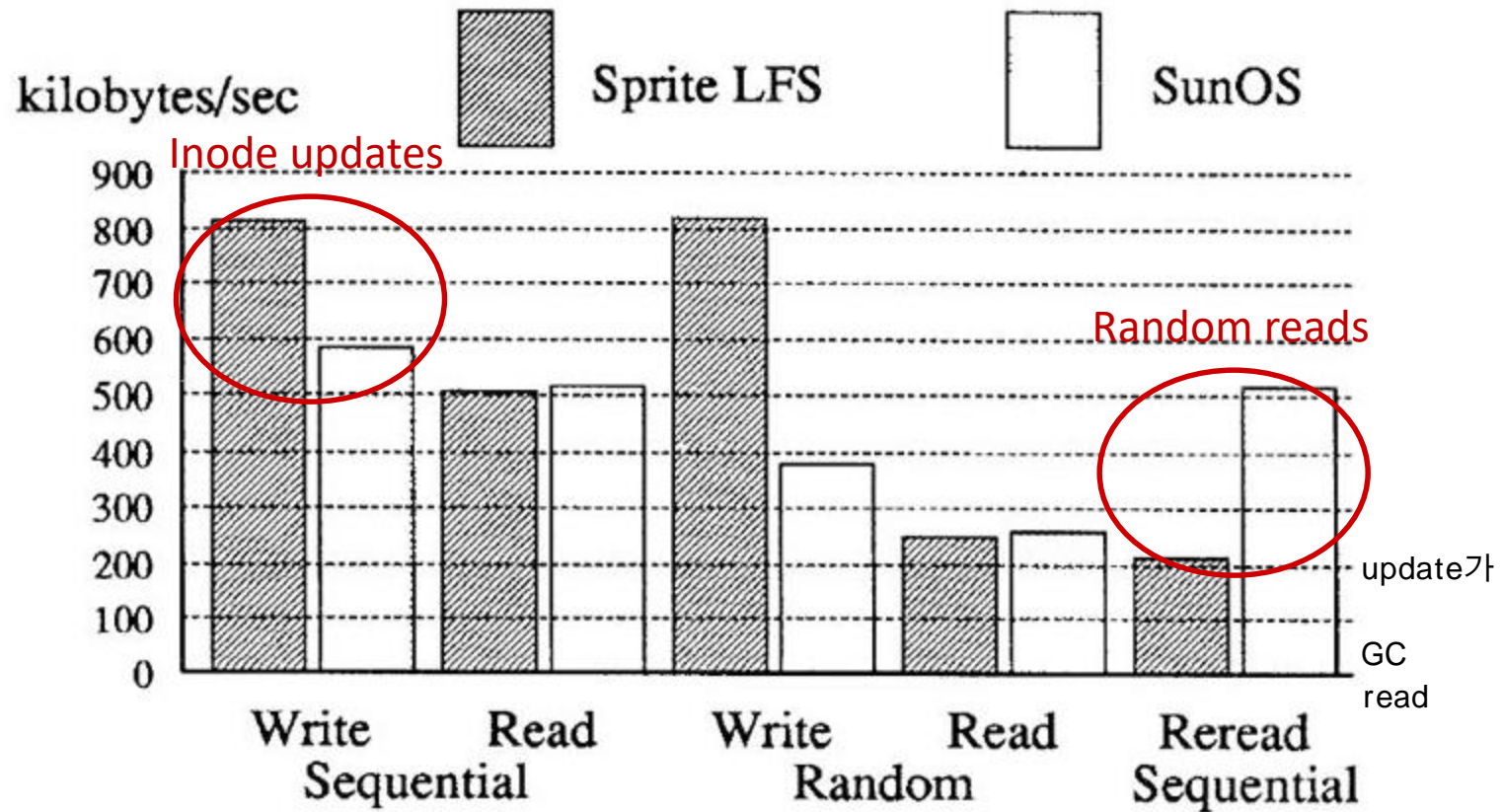

- **Cost-benefit policy:**

$$\frac{\text{benefit}}{\text{cost}} = \frac{\text{free space generated}*\text{age of data}}{\text{cost}} = \frac{(1 - u)*\text{age}}{1 + u}.$$

                                              90                                      ,
                              LFS                      GC victim selection
                                          FTL
                                            .

# Cost-Benefit Policy (Cont.)

*End of Chapter 6*