

# CS 318 Principles of Operating Systems

Fall 2018

## Lecture 16: Advanced File Systems

Ryan Huang



JOHNS HOPKINS

WHITING SCHOOL  
of ENGINEERING

Slides adapted from Andrea Arpaci-Dusseau's lecture

## Automatically Scalable Computation



Distinguished Lecturer  
**MARGO SELTZER**

University of British Columbia

**TUESDAY, NOVEMBER 6, 2018 10:30 AM**  
**Hackerman Hall B-17**

**ABSTRACT:** As our computational infrastructure races gracefully forward into increasingly parallel multi-core and clustered systems, our ability to easily produce software that can successfully exploit such systems continues to stumble. For years, we've fantasized about the world in which we'd write simple, sequential programs, add magic sauce, and suddenly have scalable, parallel executions. We're not there. We're not even close. Professor Seltzer will present a radical, potentially crazy approach to automatic scalability, combining learning, prediction, and speculation. To date, we've achieved shockingly good scalability and reasonable speedup in limited domains, but the potential is tantalizingly enormous.

**BIO:** Margo Seltzer is a Canada 150 Research Chair and Cheriton Family Chair in Computer Systems at the University of British Columbia. Her research interests are in systems, construed quite broadly: systems for capturing and accessing provenance, file systems, databases, transaction processing systems, storage and analysis of graph-structured data, new architectures for parallelizing execution, and systems that apply technology to problems in healthcare. Dr. Seltzer received an A.B. degree in Applied Mathematics from Harvard/Radcliffe College and a Ph. D. in Computer Science from the University of California, Berkeley.

# An Implementation of a Log-Structured File System for UNIX

*Margo Seltzer* – Harvard University  
*Keith Bostic* – University of California, Berkeley  
*Marshall Kirk McKusick* – University of California, Berkeley  
*Carl Staelin* – Hewlett-Packard Laboratories

## ABSTRACT

Research results [ROSE91] demonstrate that a log-structured file system (LFS) offers the potential for dramatically improved write performance, faster recovery time, and faster file creation and deletion than traditional UNIX file systems. This paper presents a redesign and implementation of the Sprite [ROSE91] log-structured file system that is more robust and integrated into the vnode interface [KLEI86]. Measurements show its performance to be superior to the 4BSD Fast File System (FFS) in a variety of benchmarks and not significantly less than FFS in any test. Unfortunately, an enhanced version of FFS (with read and write clustering) [MCVO91] provides comparable and sometimes superior performance to our LFS. However, LFS can be extended to provide additional functionality such as embedded transactions and versioning, not easily implemented in traditional file systems.

## 1. Introduction

Early UNIX file systems used a small, fixed block size and made no attempt to optimize block placement [THOM78]. They assigned disk addresses to new blocks as they were created (pre-allocation) and wrote modified blocks back to their

The log-structured file system, as proposed in [OUST88], attempts to address both of these problems. The fundamental idea of LFS is to improve file system performance by storing all file system data in a single, continuous log. Such a file system is optimized for writing, because no seek is required

# Administrivia

- **Thursday is project hacking day**
  - No class, work on project
  - Office hours still hold
    - Extra office hour today from 3-4pm
- **Lab 3 due Sunday midnight**

# File Systems Examples

- **BSD Fast File System (FFS)**
  - What were the problems with the original Unix FS?
  - How did FFS solve these problems?
- **Log-Structured File system (LFS)**
  - What was the motivation of LFS?
  - How did LFS work?

# Original Unix FS

- From Bell Labs by Ken Thompson

- Simple and elegant:

Unix disk layout



- **Components**

- Data blocks
- Inodes (directories represented as files)
- Free list
- Superblock. (specifies number of blks in FS, counts of max # of files, pointer to head of free list)

- **Problem: slow**

- Only gets **2% of disk maximum** (20Kb/sec) even for sequential disk transfers!

# Why So Slow?

- **Problem 1: blocks too small (512 bytes)**
  - File index too large
  - Require more indirect blocks
  - Transfer rate low (get one block at time)
- **Problem 2: unorganized freelist**
  - Consecutive file blocks not close together
    - Pay seek cost for even sequential access
  - Aging: becomes fragmented over time
- **Problem 3: poor locality**
  - inodes far from data blocks
  - inodes for directory not close together
    - poor enumeration performance: e.g., “ls”, “grep foo \*.c”

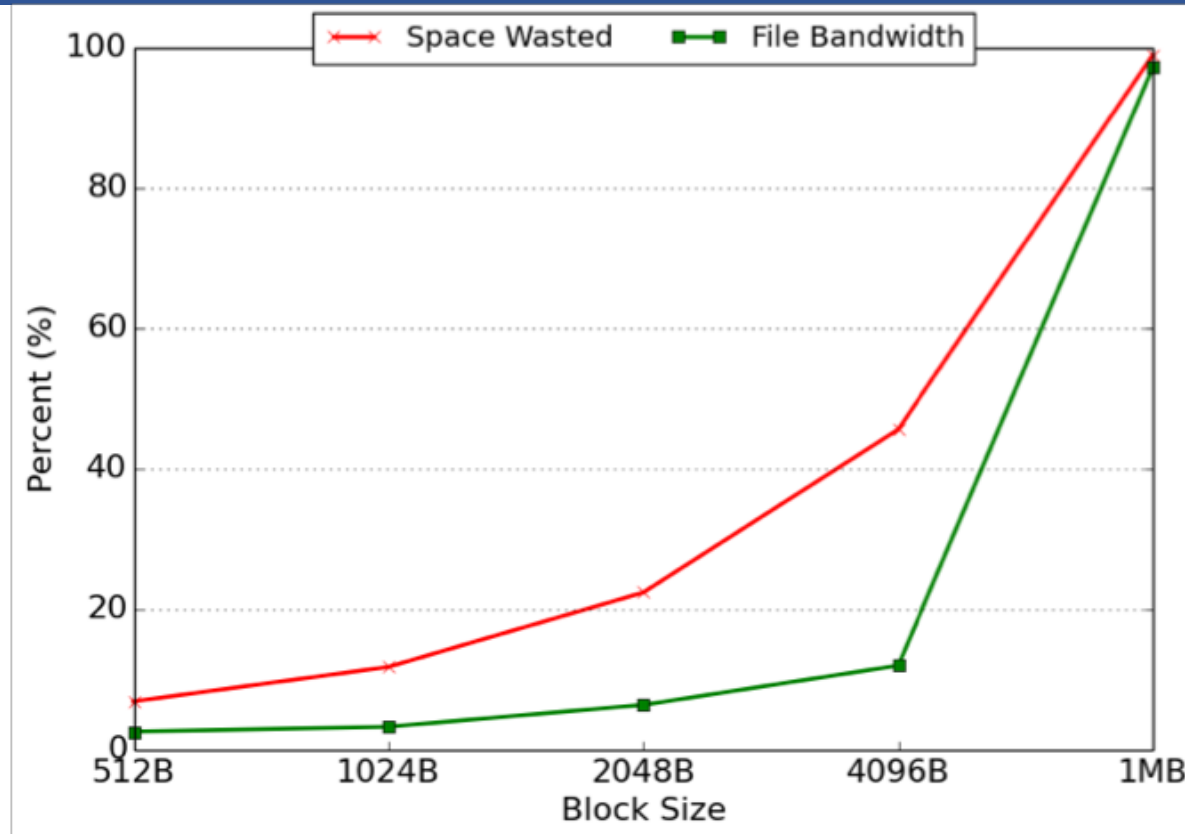
# FFS: Fast File System

- **Designed by a Berkeley research group for the BSD UNIX**
  - A classic file systems paper to read: [[McKusic](#)]
- **Approach:**
  - **measure** an state of the art systems
  - identify and understand the fundamental problems
    - **The original FS treats disks like random-access memory!**
  - get an idea and **build** a better systems
- **Idea: design FS structures and allocation polices to be “disk aware”**
- **Next: how FFS fixes the performance problems (to a degree)**



# Problem 1: Blocks Too Small

Measurement:

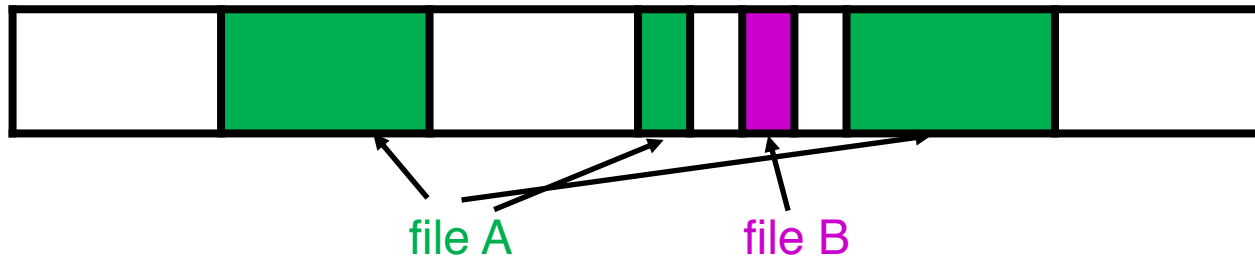


- **Bigger block increases bandwidth, but how to deal with wastage (“**internal fragmentation**”)?**
  - Use idea from malloc: split unused portion

# Solution: Fragments

- **BSD FFS:**

- Has large block size (4096B or 8192B)
- Allow large blocks to be chopped into small ones called “fragments”
- Ensure fragments only used for little files or ends of files



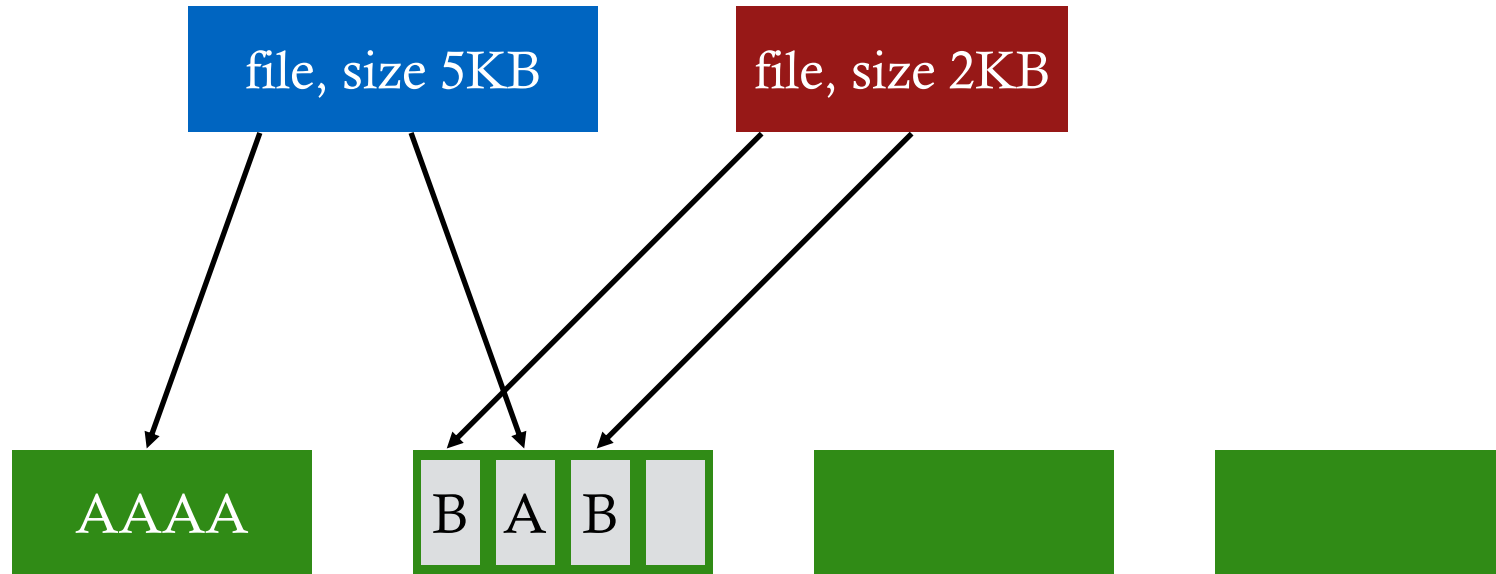
- Fragment size specified at the time that the file system is created
- Limit number of fragments per block to 2, 4, or 8

- **Pros**

- High transfer speed for larger files
- Low wasted space for small files or ends of files

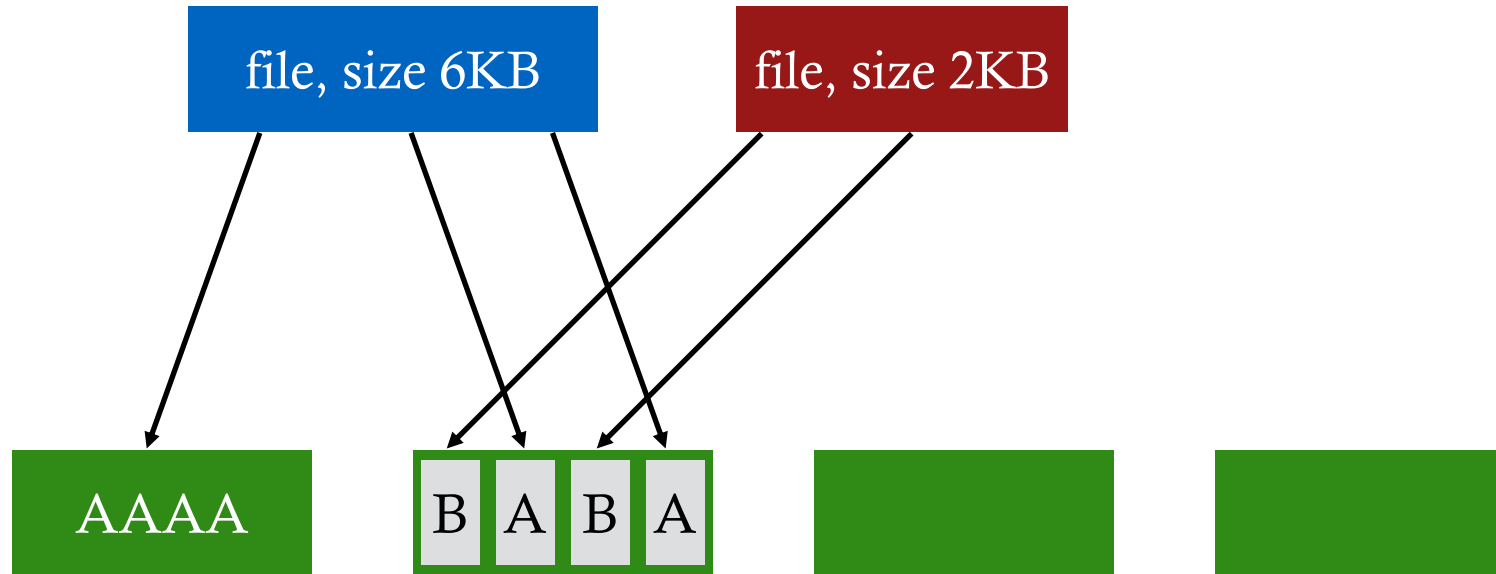
# Fragment Example

Block size: 4096 B  
Fragment size: 1024 B



# Fragment Example

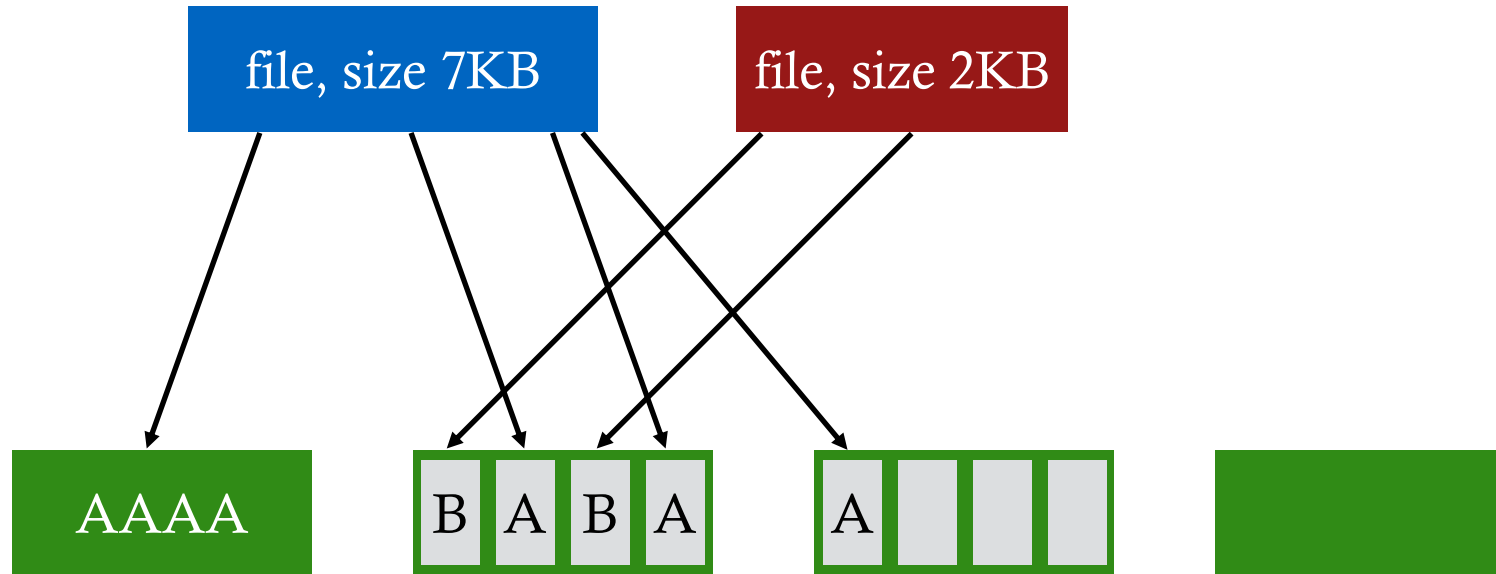
Block size: 4096 B  
Fragment size: 1024 B



append A to first file

# Fragment Example

Block size: 4096 B  
Fragment size: 1024 B



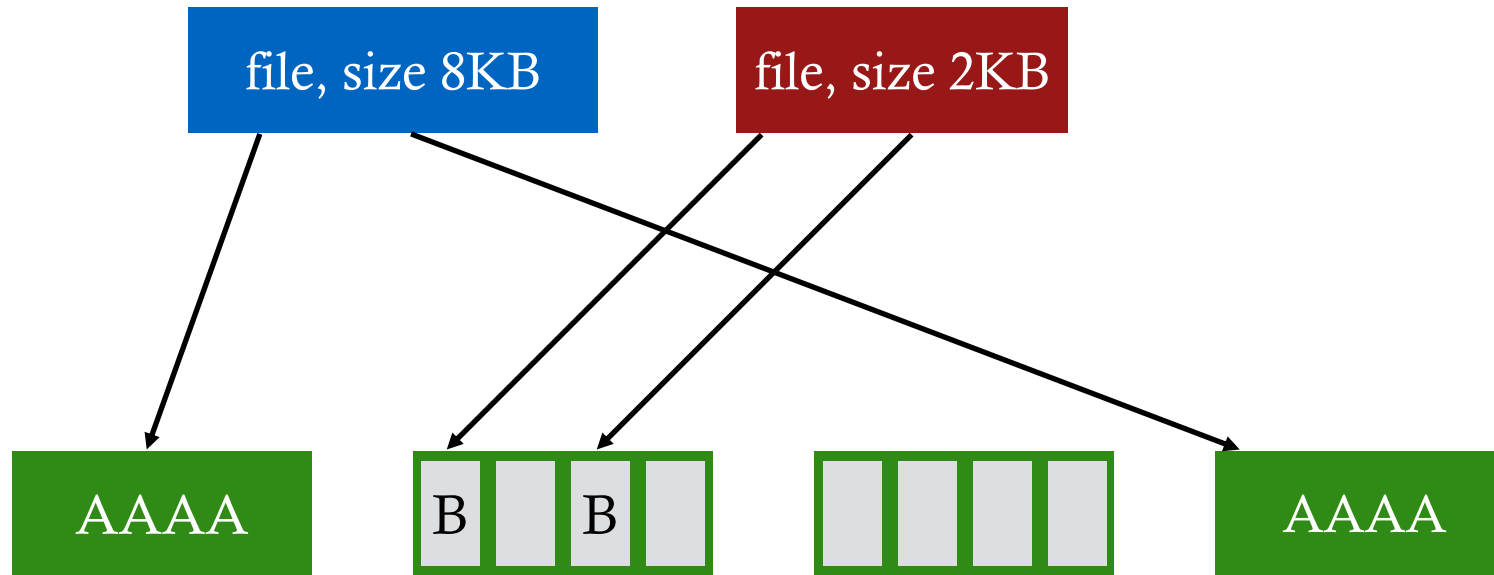
append A to first file

Not allowed to use fragments across multiple blocks!

What to do instead?

# Fragment Example

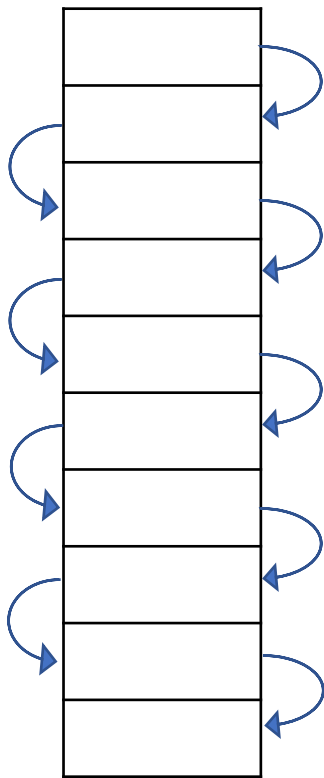
Block size: 4096 B  
Fragment size: 1024 B



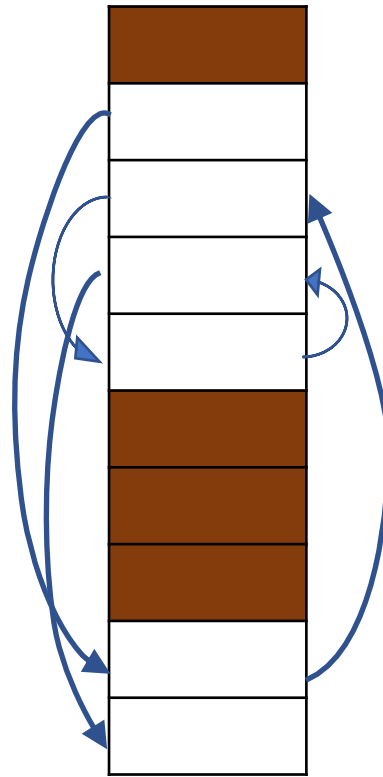
append A to first file,  
copy to fragments to new block

# Problem 2: Unorganized Freelist

- **Leads to random allocation of sequential file blocks overtime**



Initial performance good



Get worse over time

## Measurement:

- New FS: **17.5%** of disk bandwidth
- Few weeks old: **3%** of disk bandwidth

# Fixing the Unorganized Freelist

- **Periodical compact/defragment disk**
  - Cons: locks up disk bandwidth during operation
- **Keep adjacent free blocks together on freelist**
  - Cons: costly to maintain
- **FFS: bitmap of free blocks**
  - Each bit indicates whether block is free
    - E.g., 1010101111111000001111111000101100
  - Easier to find contiguous blocks
  - Small, so usually keep entire thing in memory
  - Time to find free blocks increases if fewer free blocks

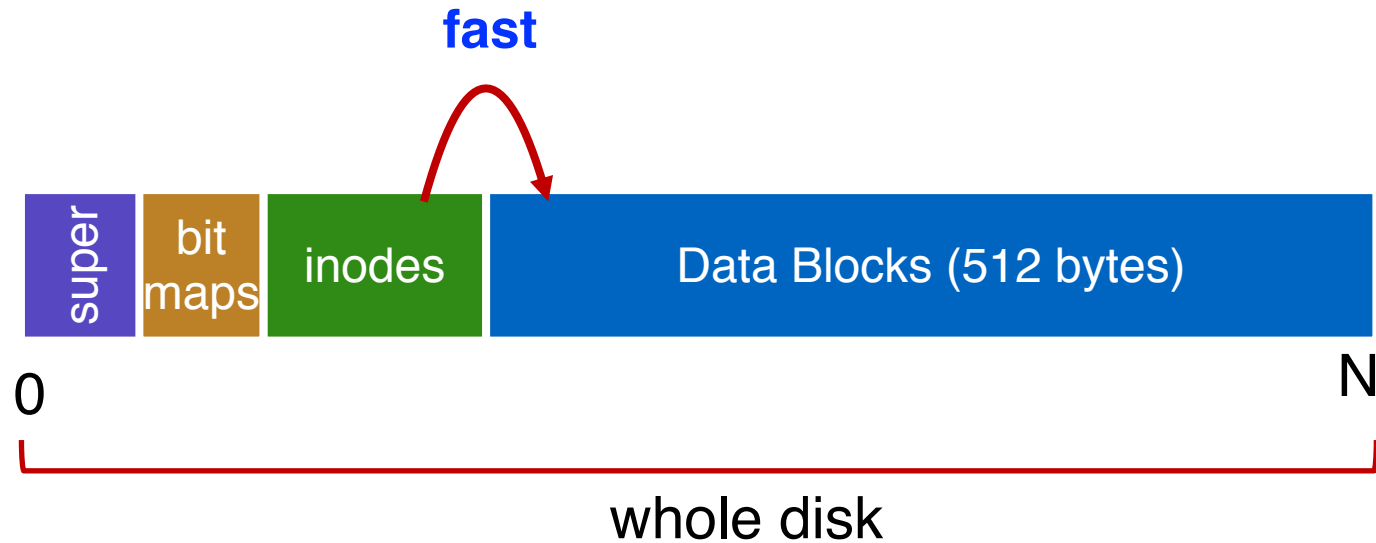


# Using a Bitmap

- **Usually keep entire bitmap in memory:**
  - 4G disk / 4K byte blocks. How big is map?
- **Allocate block close to block x?**
  - Check for blocks near  $\text{bmap}[x/32]$
  - If disk almost empty, will likely find one near
  - As disk becomes full, search becomes more expensive and less effective
- **Trade space for time (search time, file access time)**

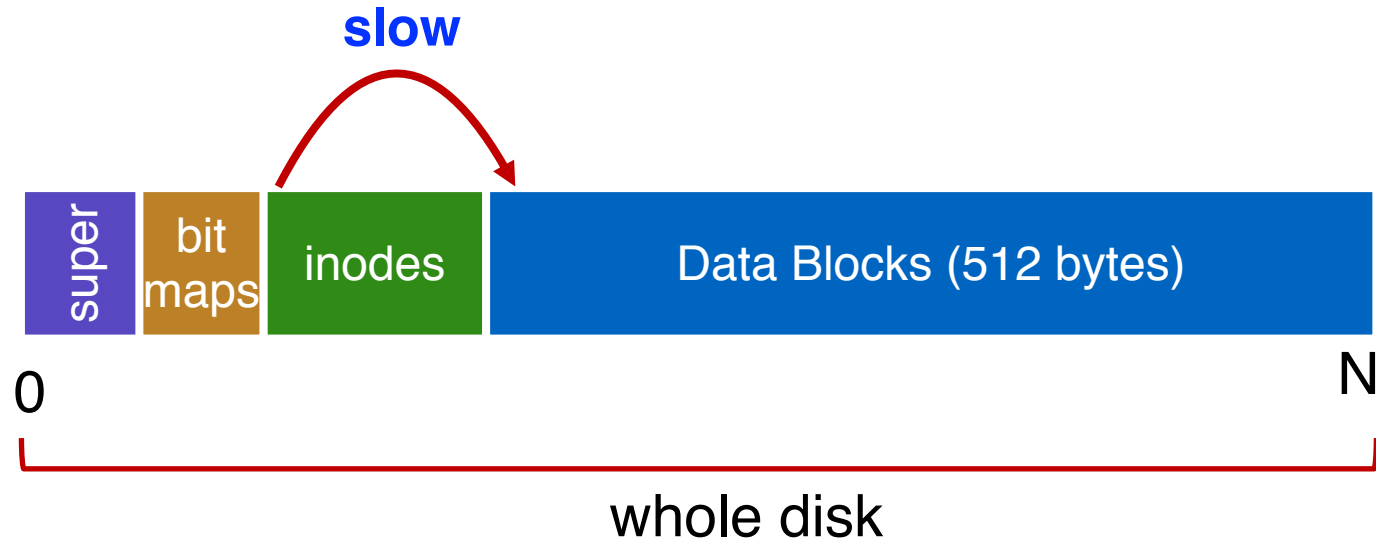


# Problem 3: Poor Locality



- **How to keep inode close to data block?**

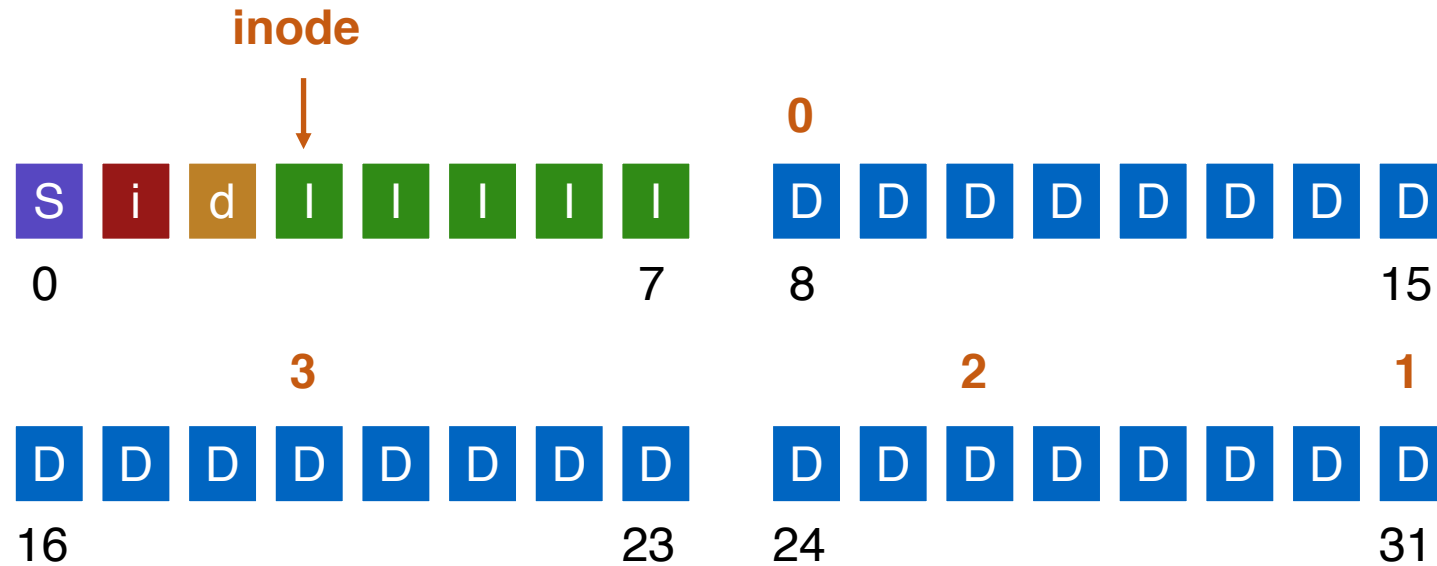
# Problem 3: Poor Locality



- **How to keep inode close to data block?**

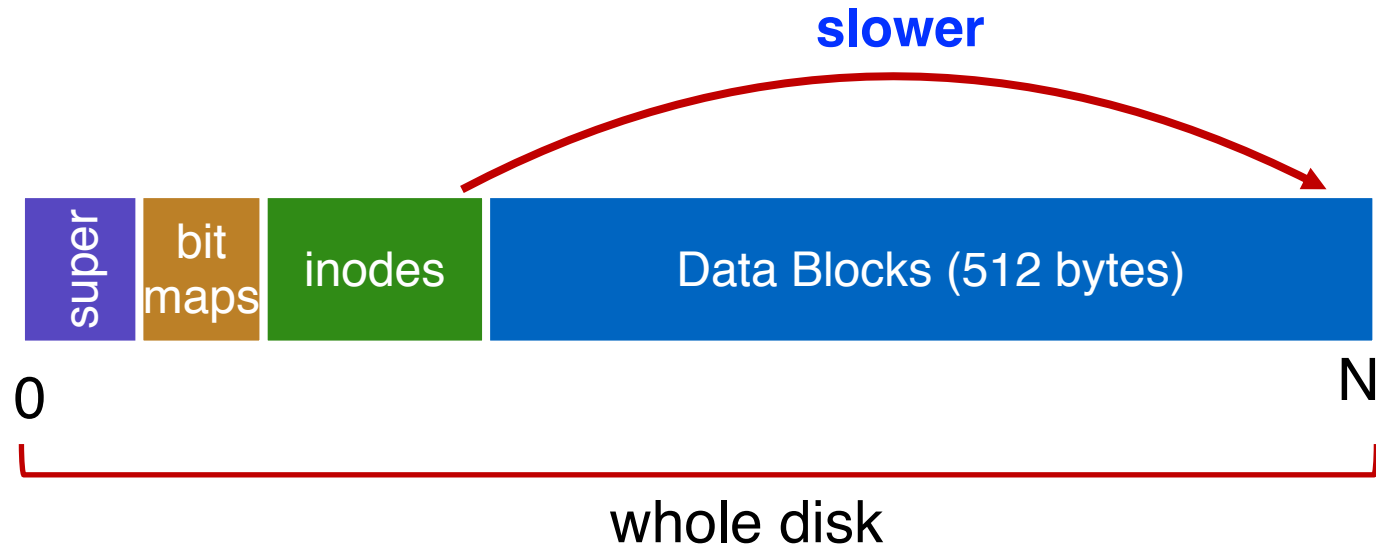
# Problem 3: Poor Locality

- **Example bad layout:**



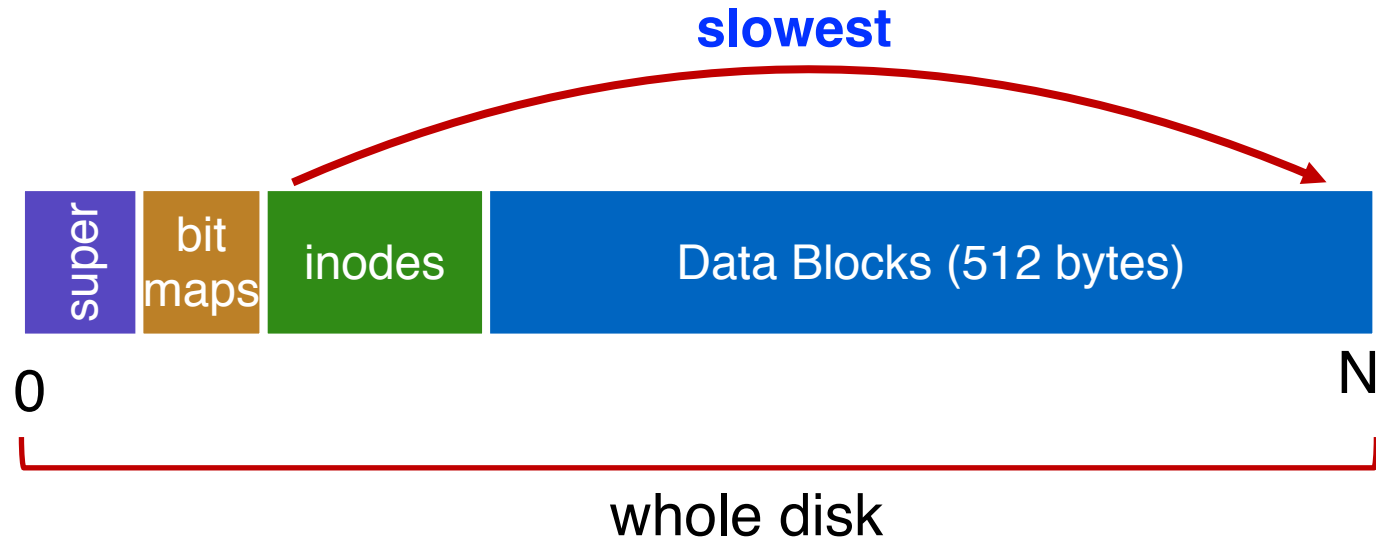
- **How to keep inode close to data block?**

# Problem 3: Poor Locality



- **How to keep inode close to data block?**

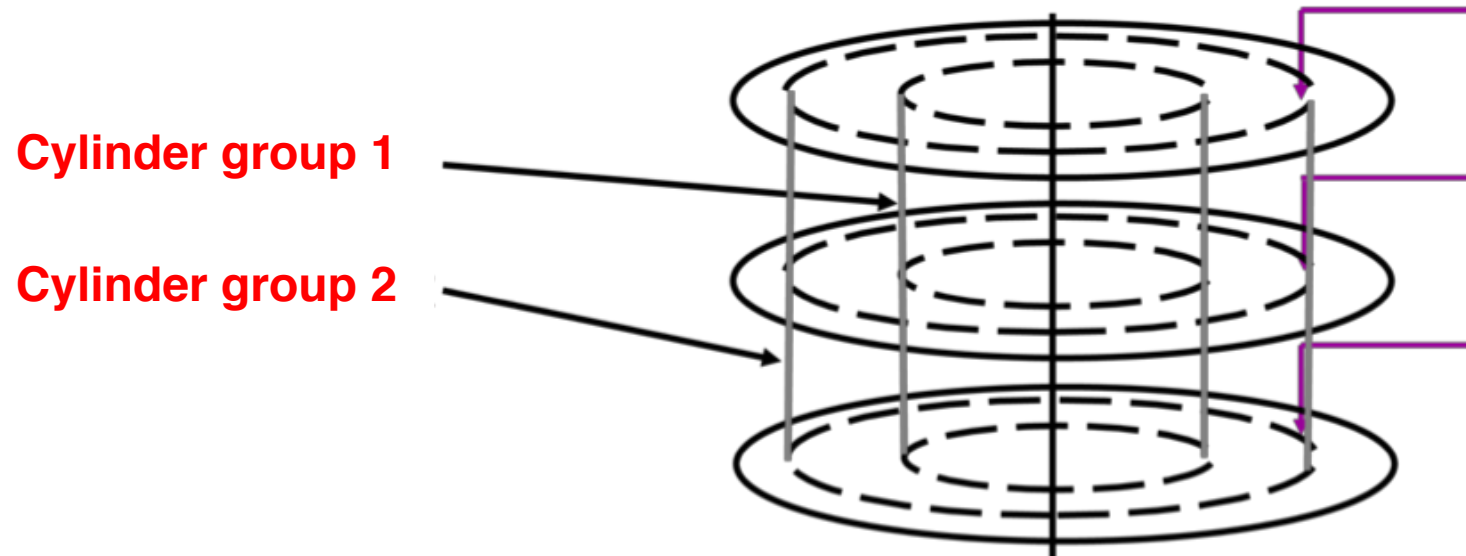
# Problem 3: Poor Locality



- **How to keep inode close to data block?**

# FFS Solution: Cylinder Group

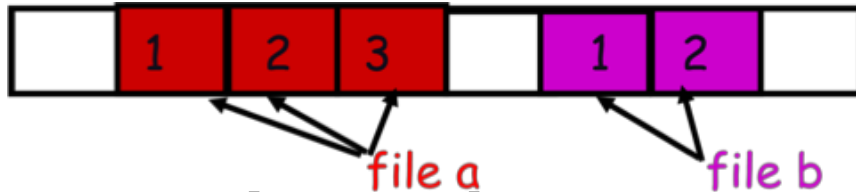
- Group sets of consecutive cylinders into “**cylinder groups**”



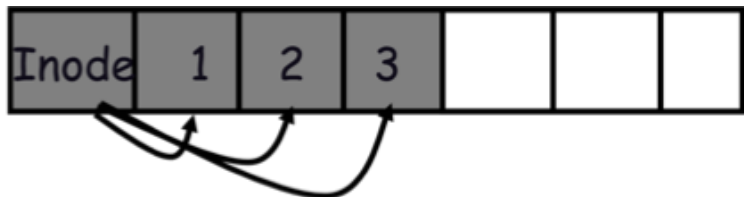
- Key: can access any block in a cylinder without performing a seek. Next fastest place is adjacent cylinder.
- Tries to put everything related in same cylinder group
- Tries to put everything not related in different group

# Clustering in FFS

- **Tries to put sequential blocks in adjacent sectors**
  - (Access one block, probably access next)



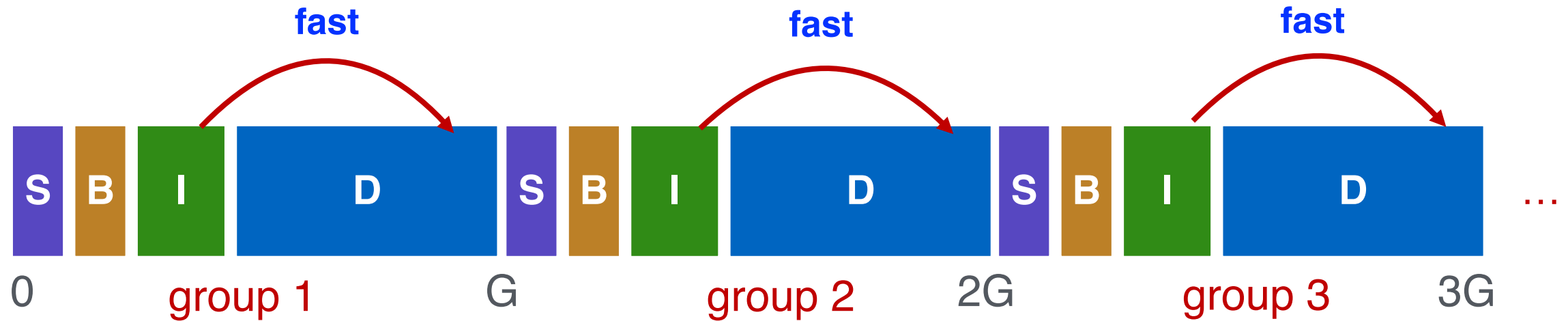
- **Tries to keep inode in same cylinder as file data:**
  - (If you look at inode, most likely will look at data too)



- **Tries to keep all inodes in a dir in same cylinder group**
  - Access one name, frequently access many, e.g., “ls -l”



# What Does Disk Layout Look Like Now?



- **How to keep inode close to data block?**
  - Answer: Use groups across disks
  - Strategy: allocate inodes and data blocks in same group
  - Each cylinder group basically a mini-Unix file system
- **Is it useful to have multiple super blocks?**
  - Yes, if some (but not all) fail

# FFS Results

- **Performance improvements:**
  - Able to get 20-40% of disk bandwidth for large files
  - 10-20x original Unix file system!
  - Stable over FS lifetime
  - Better small file performance (why?)
- **Other enhancements**
  - Long file names
  - Parameterization
  - Free space reserve (10%) that only admin can allocate blocks from

# LFS: Log-structured File System

- **Motivation**

- Faster CPUs: I/O becomes more and more of a bottleneck
- More memory: file cache is effective for reads
- Implication: writes compose most of disk traffic

- **Problems with previous FS**

- Perform many small writes
  - Good performance on large, sequential writes, but many writes are still small, random
- Synchronous operation to avoid data loss
- Depends upon knowledge of disk geometry

- **An influential work designed by Mendel Rosenblum (VMWare co-founder) and John Ousterhout**

# LFS Idea

- **Insight: treat disk like a tape-drive**
  - Best performance from disk for sequential access
- **File system buffers writes in main memory until “enough” data**
  - How much is enough?
  - Enough to get good sequential bandwidth from disk (MB)
  - Unit called a “segment”
- **Write buffered data to new segment on disk in a sequential log**
  - Transfer all updates into a series of sequential writes
  - **Do not overwrite old data on disk: old copies left behind**
  - Write both data and metadata in one operation

# Pros And Cons

- **Pros**

- Always large sequential writes → good performance
- No knowledge of disk geometry
  - Assume sequential better than random

- **Potential problems**

- How do you find data to read?
- What happens to metadata during write?
- What happens when you fill up the disk?

# Read in LFS

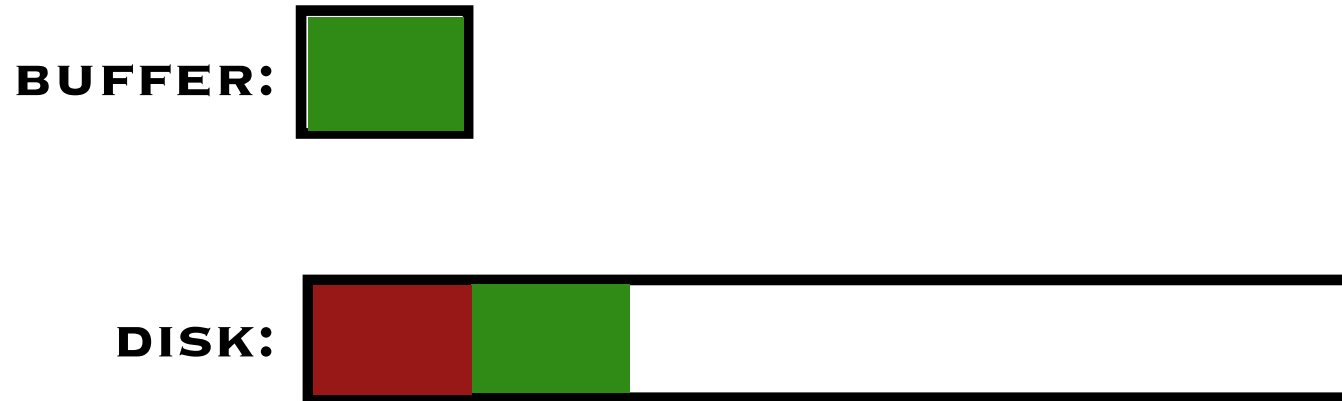
- **Same basic structures as Unix**
  - Directories, inodes, indirect blocks, data blocks
  - Reading data block implies finding the file's inode
    - Unix FS: inodes kept in array
    - LFS: inodes spread around on disk
- **Solution: inode map indicates where each inode is stored**
  - Can keep cached copy in memory
  - inode map written to log with everything else
  - Periodically written to known checkpoint location on disk for crash recovery

# Write in LFS



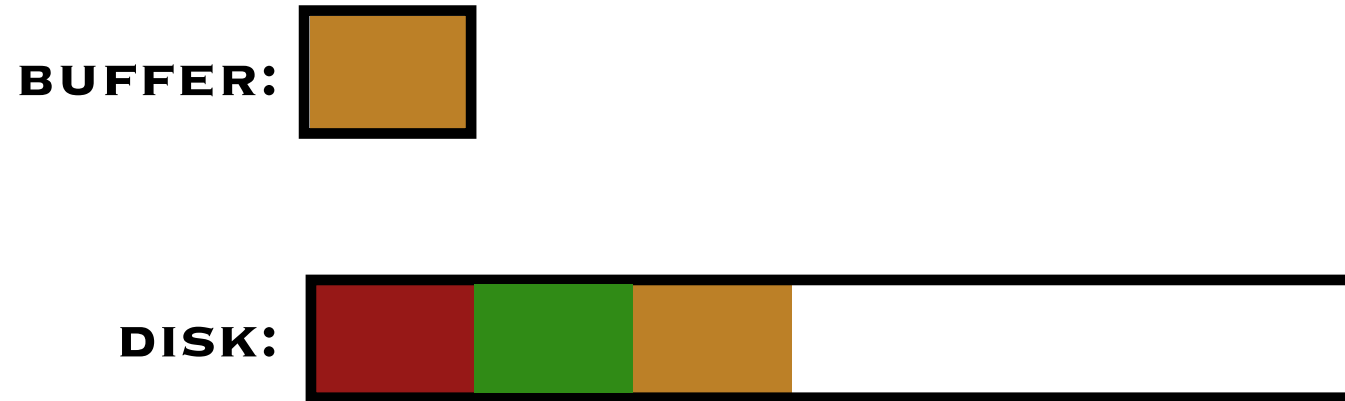
- **Why do we buffer the write?**
  - Sequential write alone is not enough
  - Disk is constantly rotating!
  - Must issue a large number of **contiguous** writes

# Write in LFS

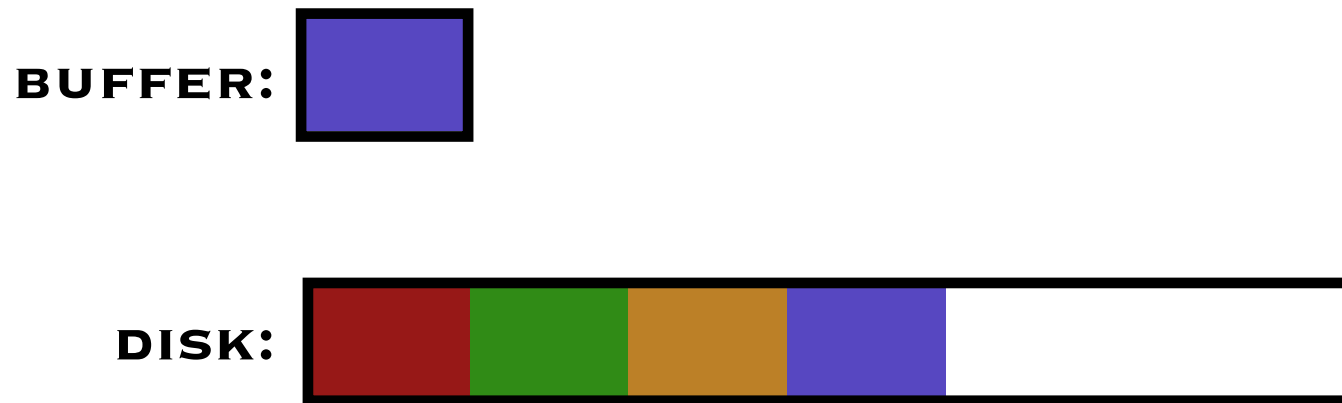




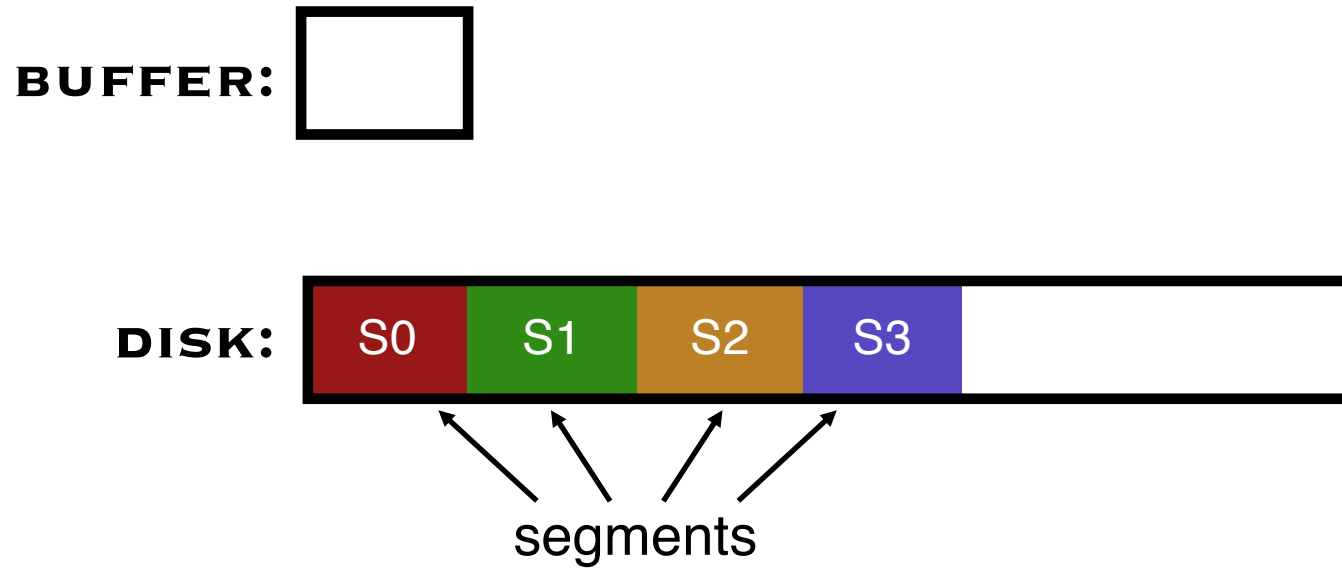
# Write in LFS



# Write in LFS



# Write in LFS



# Data Structures for LFS (attempt 1)



## What data structures from FFS can LFS remove?

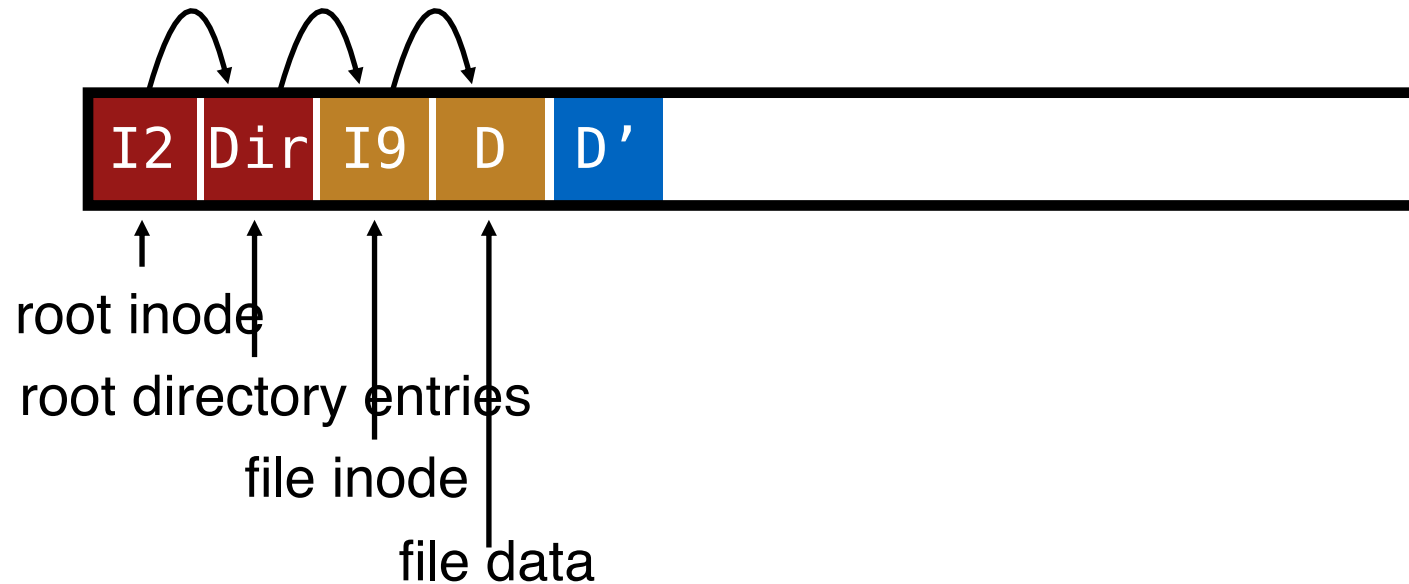
- allocation structs: data + inode bitmaps

## What type of name is much more complicated?

- Inodes are no longer at fixed offset
- Use **current offset on disk** instead of table index for name
- Note: **when update inode, inode number changes!!**

# Overwrite Data in LFS – Attempt 1

- **Overwrite data in /file.txt**



- **How to update Inode 9 to point to new D' ???**

# Overwrite Data in LFS – Attempt 1

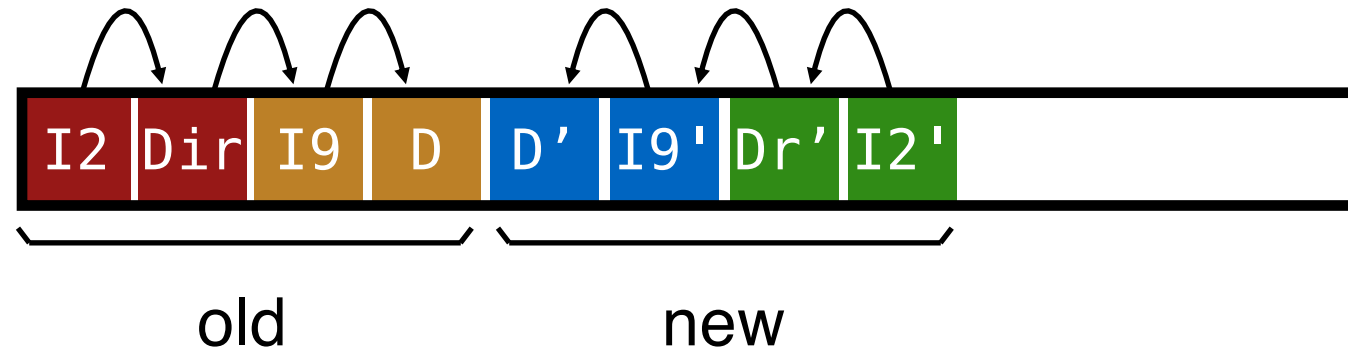
- Overwrite data in /file.txt



- **Can LFS update Inode 9 to point to new D'?**
  - NO! This would be a random write

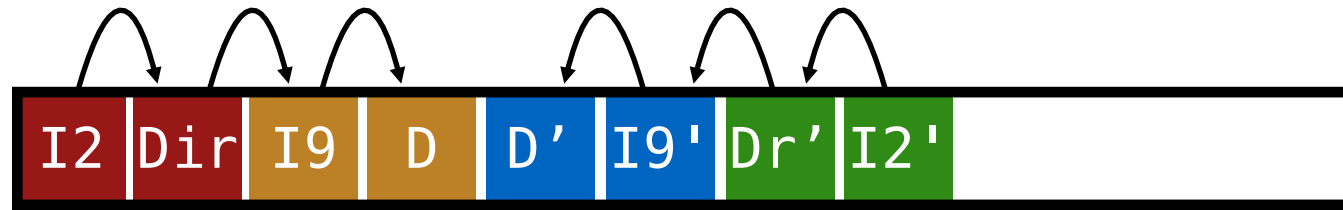
# Overwrite Data in LFS – Attempt 1

- **Overwrite data in /file.txt**



- **Must update all structures in sequential order to log**

# Attempt 1: Problem w/ Inode Numbers



- **Problem:**
  - For every data update, must propagate updates all the way up directory tree to root
- **Why?**
  - When inode copied, its location (inode number) changes
- **Solution:**
  - Keep inode numbers constant; don't base name on offset
- **FFS found inodes with math. How now?**



# Data Structures for LFS (attempt 2)

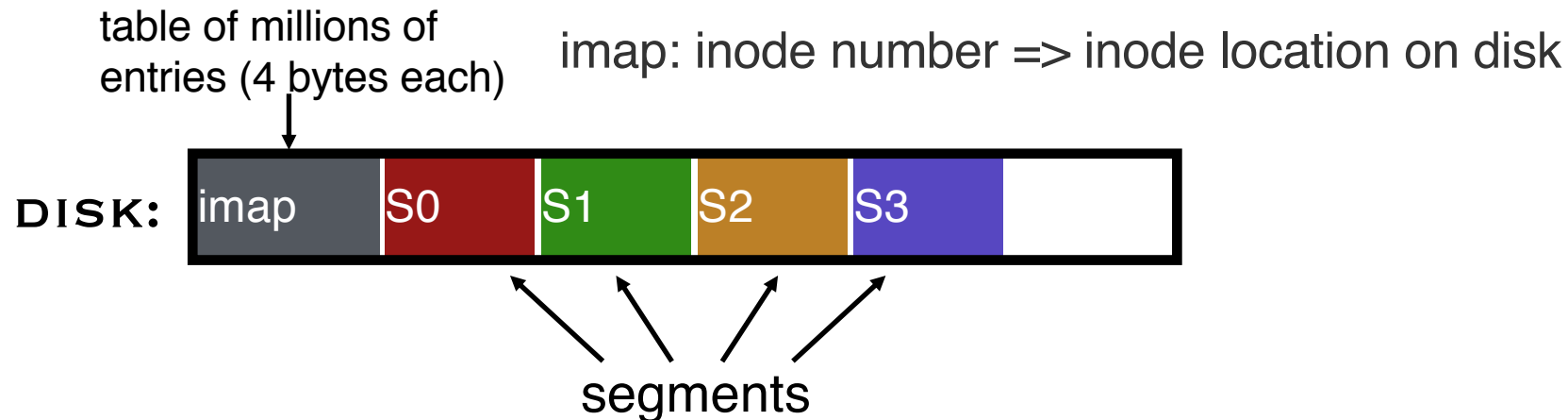
## What data structures from FFS can LFS remove?

- allocation structs: data + inode bitmaps

## What type of name is much more complicated?

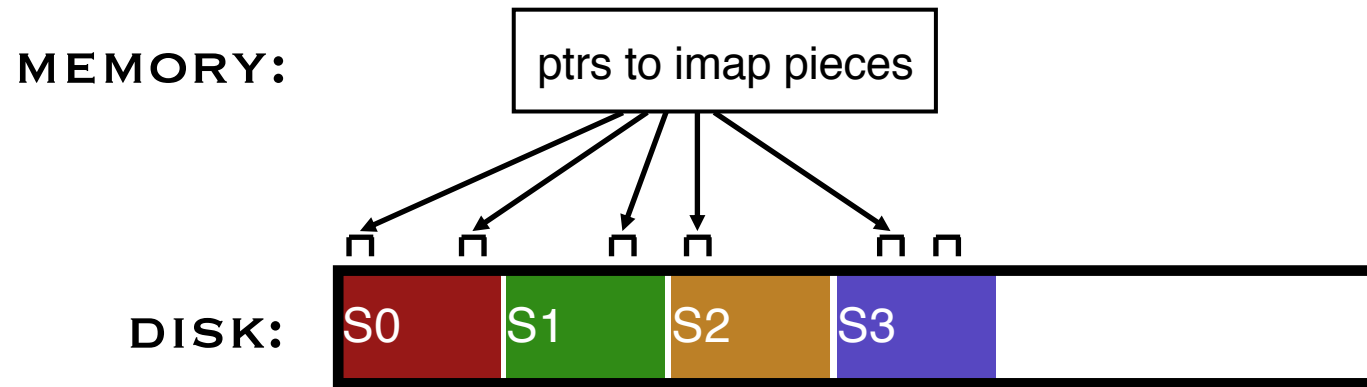
- Inodes are no longer at fixed offset
- Use **imap** structure to map:
  - inode number => **most recent** inode location on disk

# Where to keep Imap?



- **Where can imap be stored? Dilemma:**
  1. imap too large to keep in memory
  2. don't want to perform random writes for imap
- **Solution: Write imap in segments**
  - Keep pointers to pieces of imap in memory

# Solution: Imap in Segments



- **Solution:**

- Write imap in segments
- Keep pointers to pieces of imap in memory
- Keep recent accesses to imap cached in memory

# Disk Cleaning

- **When disk runs low on free space**
  - Run a disk cleaning process
  - Compacts live information to contiguous blocks of disk
- **Problem: long-lived data repeatedly copied over time**
  - Solution: partition disk in to segments
  - Group older files into same segment
    - Do not clean segments with old files
- **Try to run cleaner when disk is not being used**

# Next Time...

- **Read Chapter 42**