

Today

Distributed Subsystems (Lesson 7)

- ✓ * Global memory system
- ✓ * Distributed Shared Memory

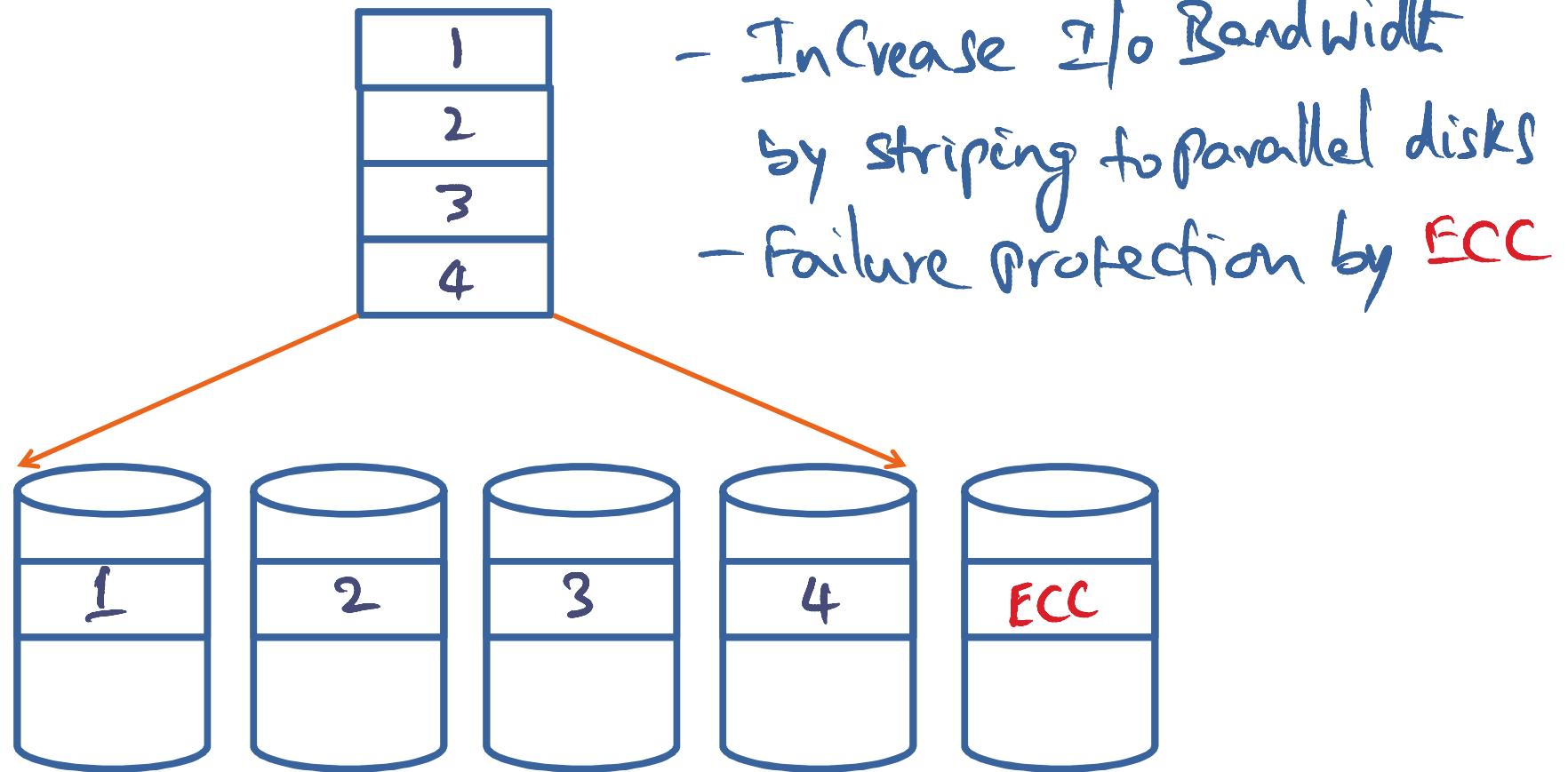
⇒ * Distributed file Systems

Wednesday

Failures + recovery (Lesson 8)

* Please watch LRVM lecture *

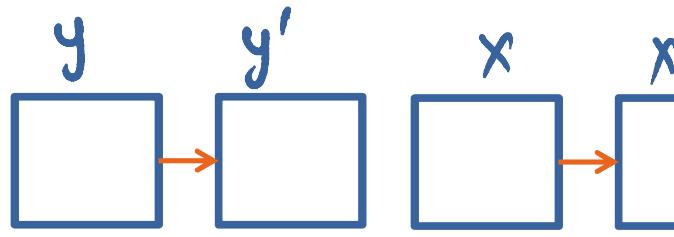
Preliminaries: Striping a file to multiple Disks



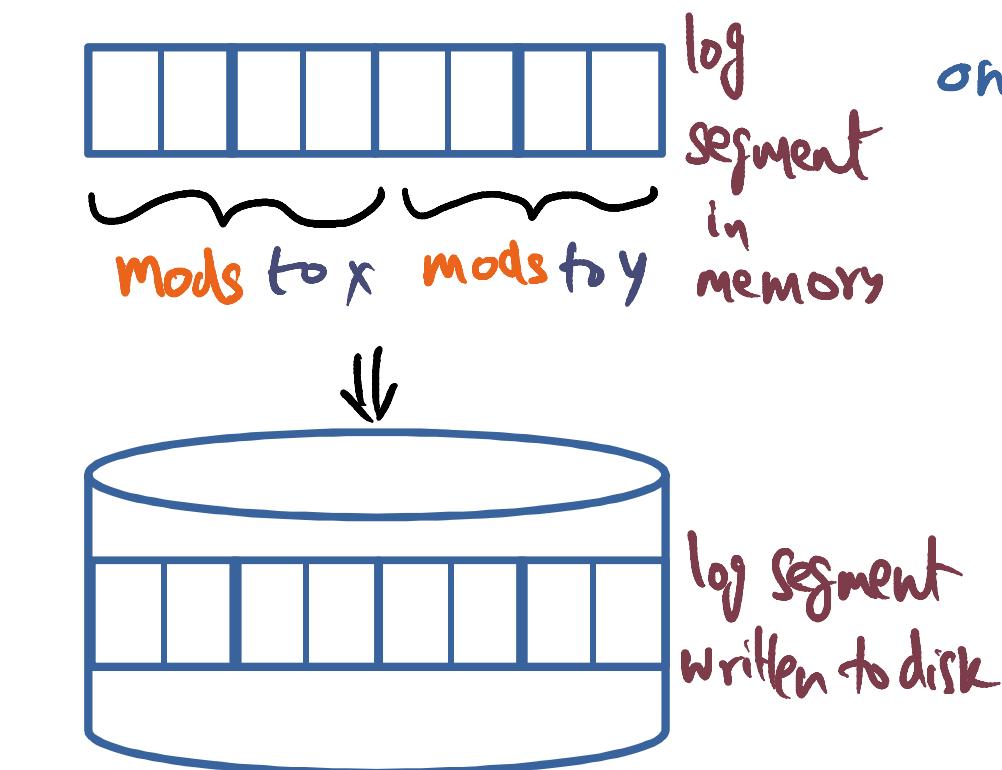
Drawbacks

- Cost
- Small write Problem

Preliminaries: Log Structured File System



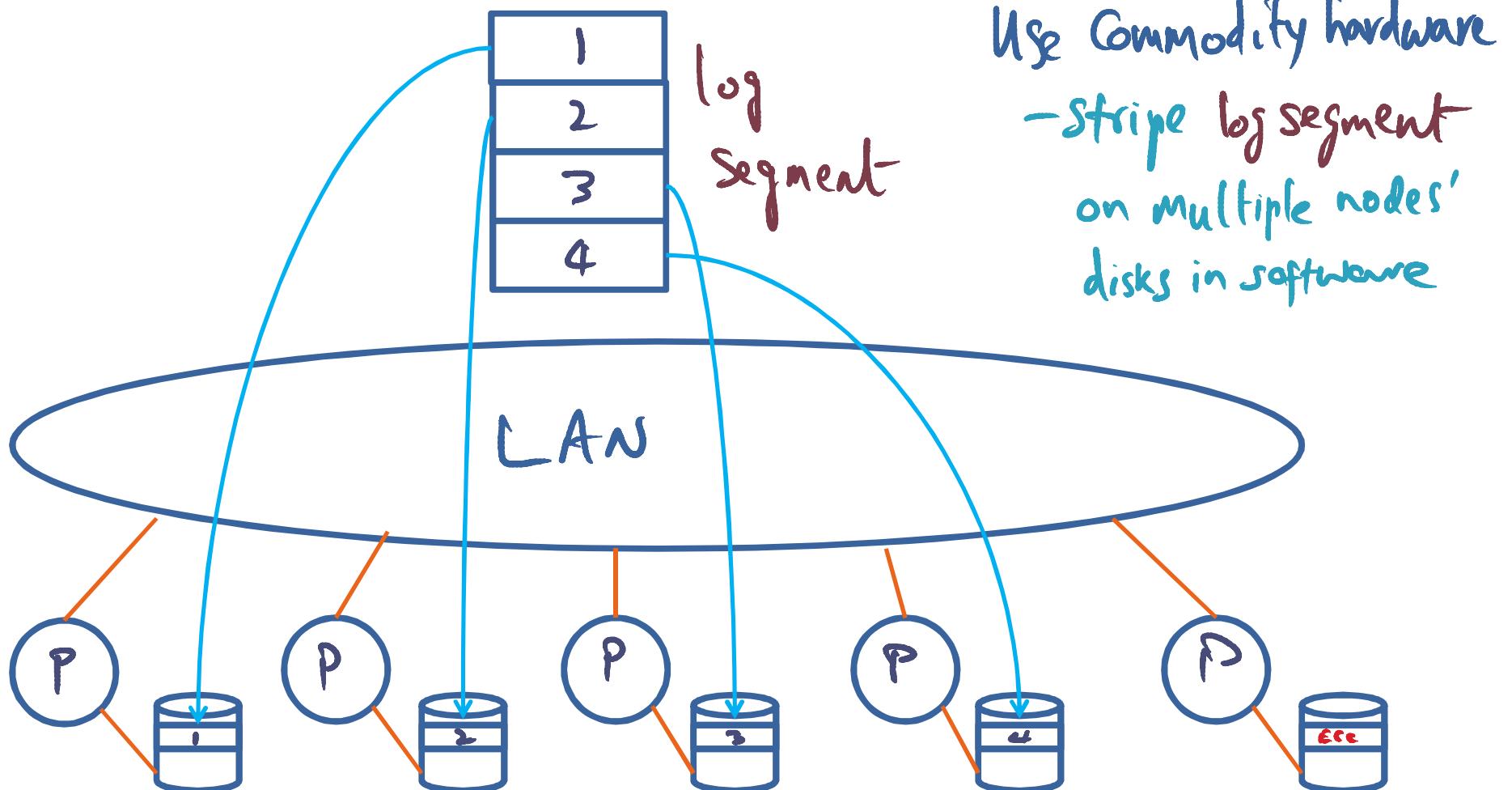
- Buffer changes to multiple files in one contiguous log Segment data structure
- Flush log segment to disk once it fills up or periodically
solves small write problem



Preliminaries: Software RAID

Febra File System (UC-Berkeley)

- Combines LFS + RAID

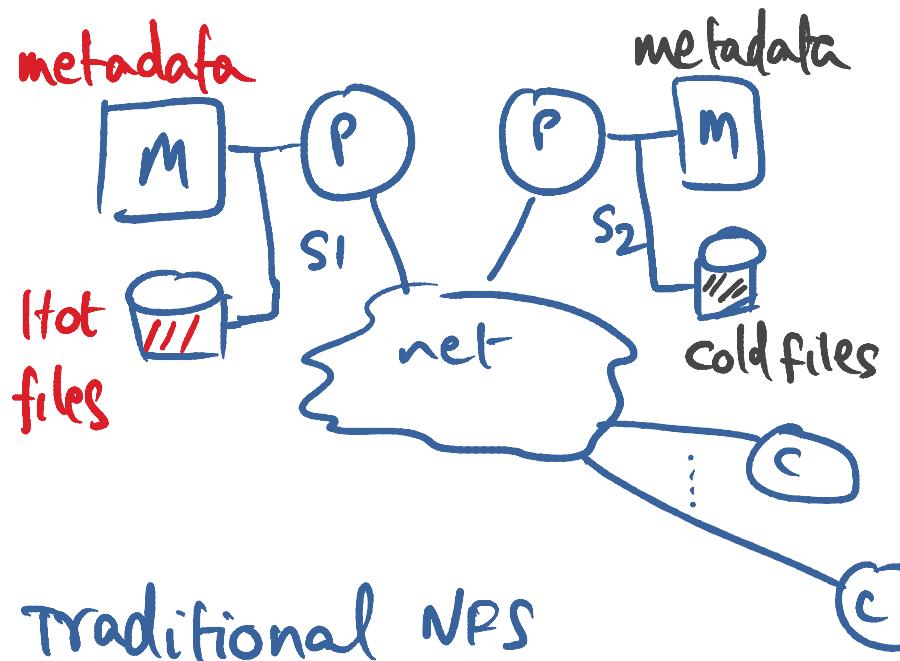


Putting them all together Plus More

xFS - a DFS

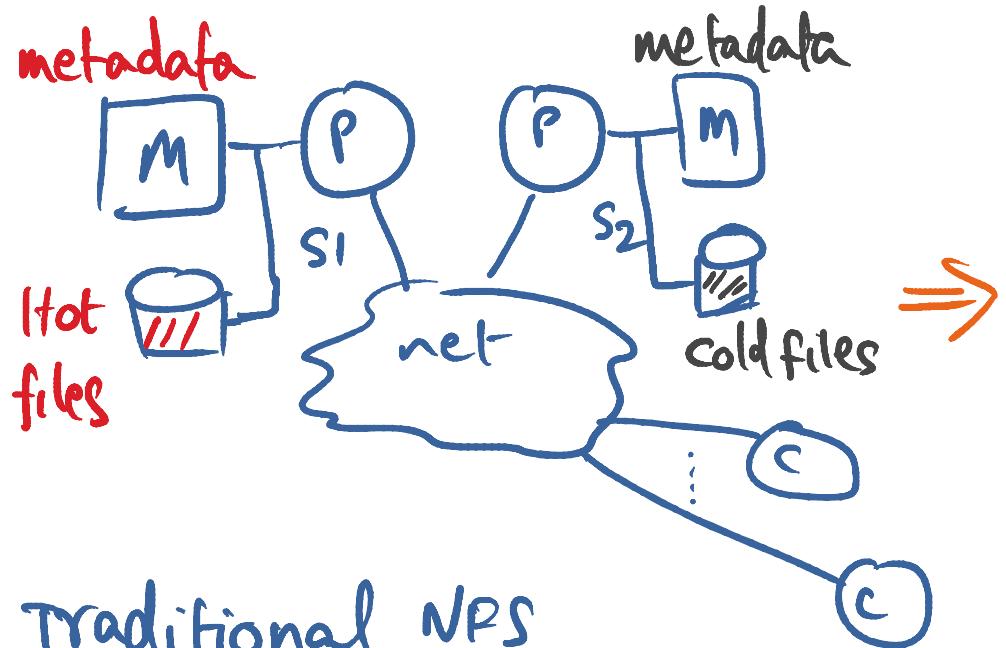
- log based striping (from Zesra)
- cooperative caching (from prior UCB work)
- dynamic management of data + metadata
- Subsetting storage servers
- distributed log cleaning

Dynamic Management



Traditional NPS
with centralized server(s)
- memory contents
+ metadata, file cache,
client caching directory

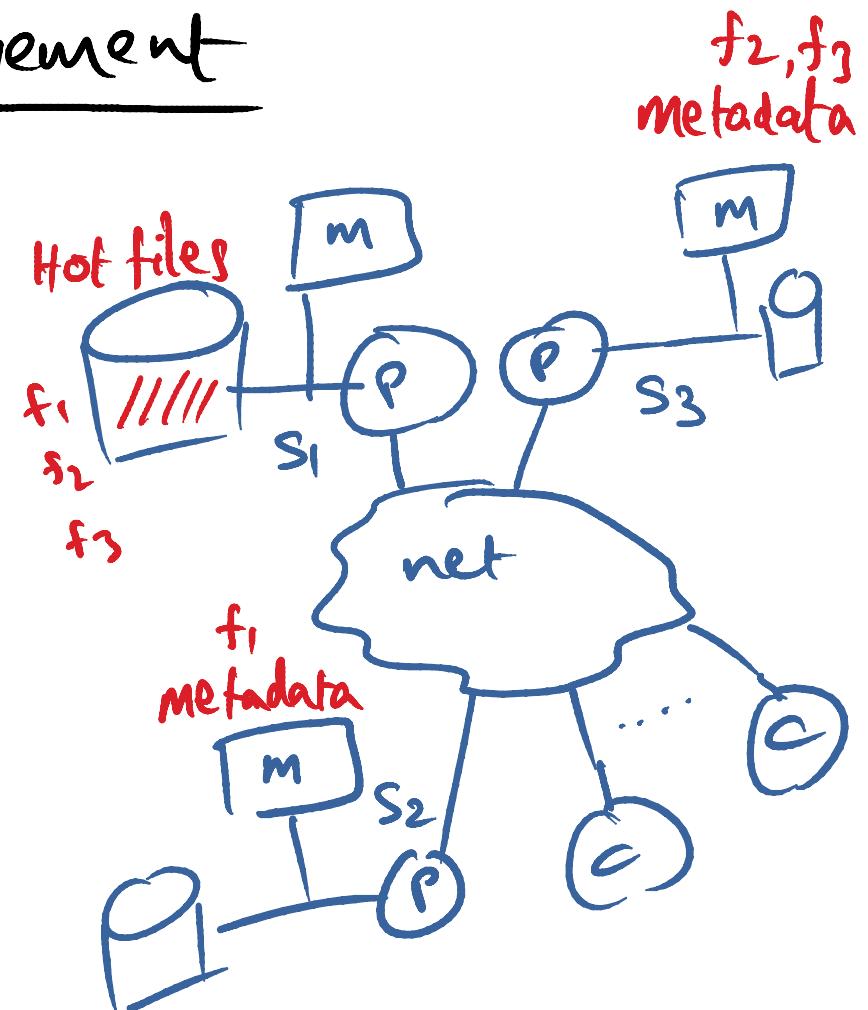
Dynamic Management



Traditional NPS

with centralized server(s)

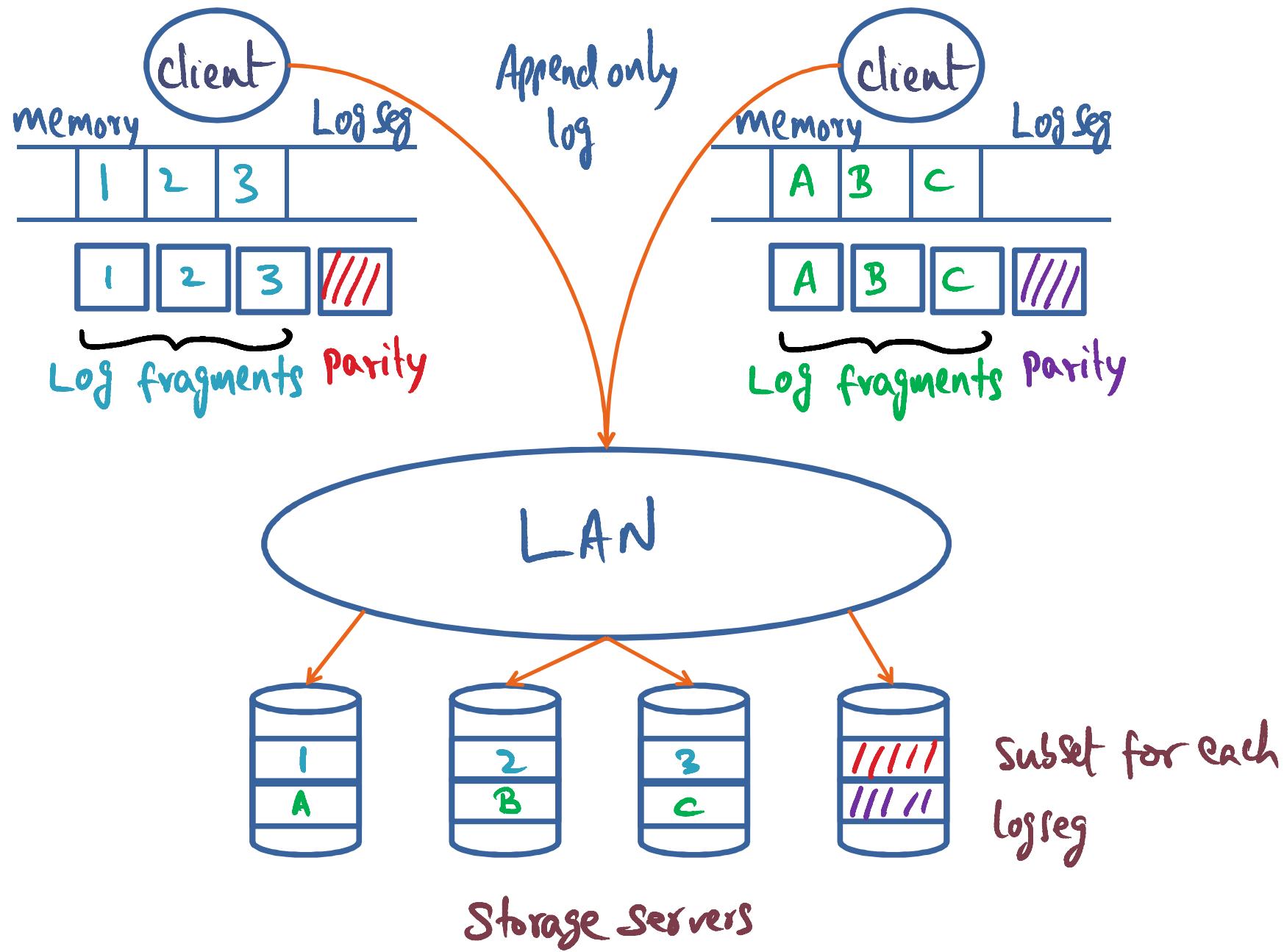
- memory contents
 - * metadata, file cache,
client caching directory



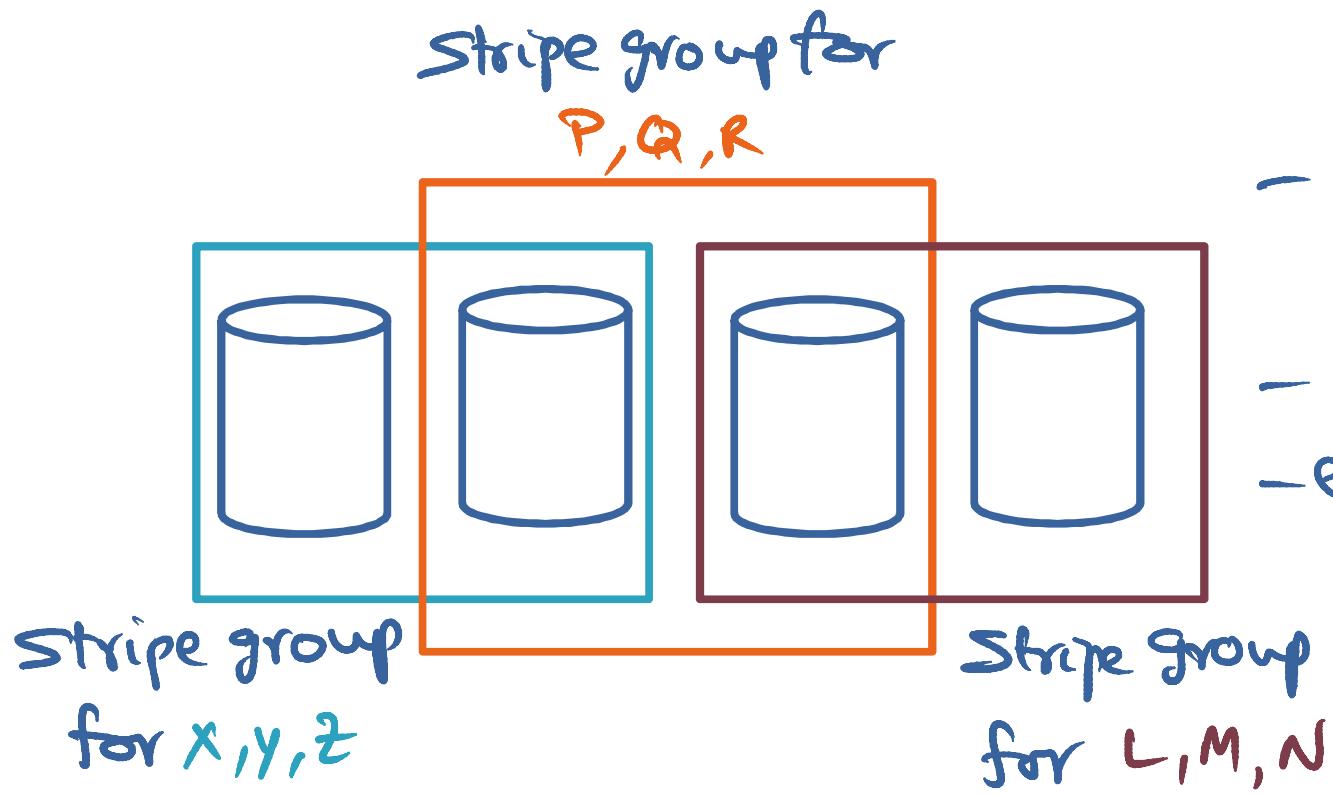
XFS

- metadata management
 - dynamically distributed
- cooperative client file caching

Log based Striping and stripe groups



Stripe Group

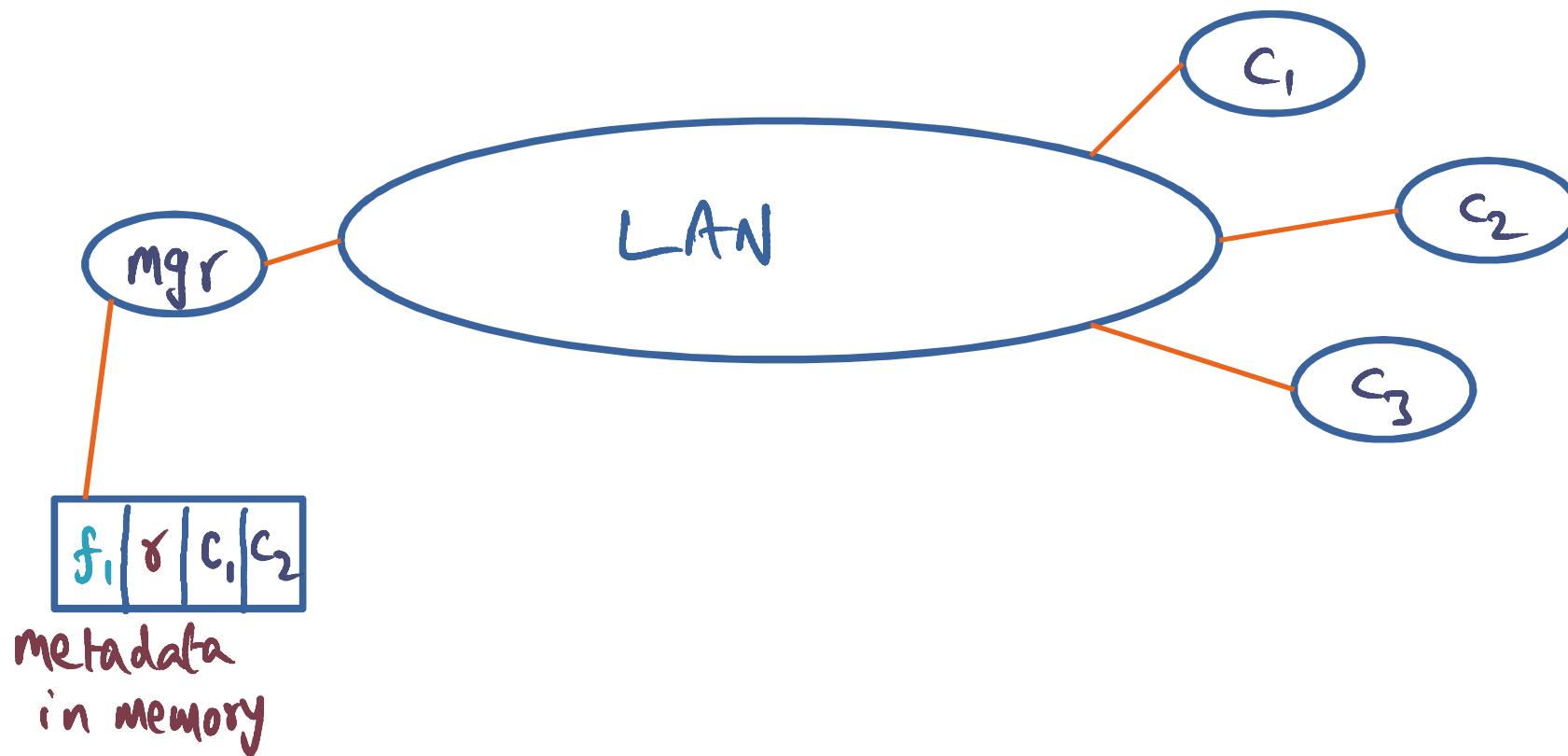


- subset servers into stripe groups
- parallel client activities
- increased availability
- efficient log cleaning

cooperative Caching

Cache Coherence

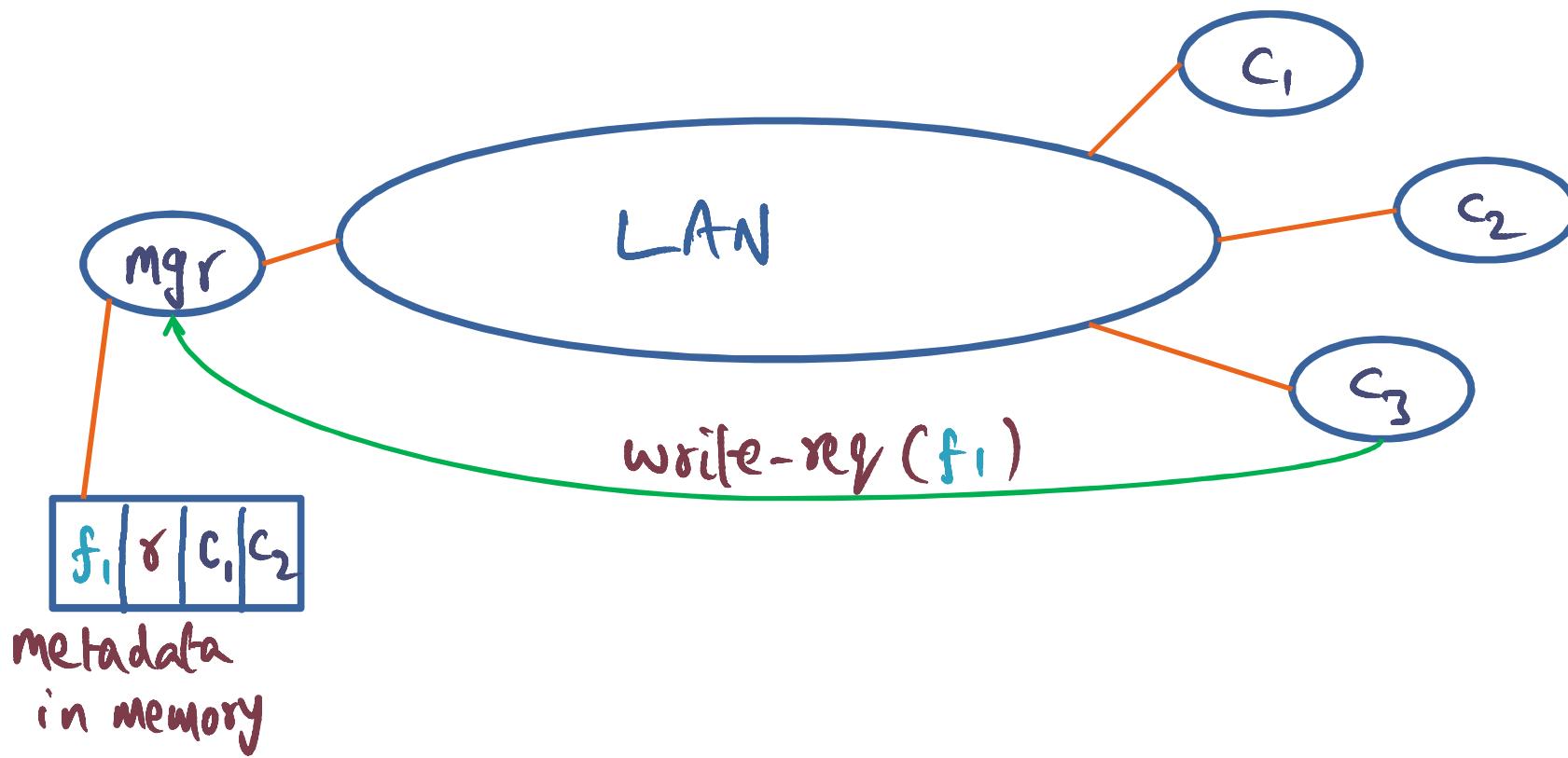
- single writer, multiple readers
- file block unit of coherence



cooperative caching

Cache Coherence

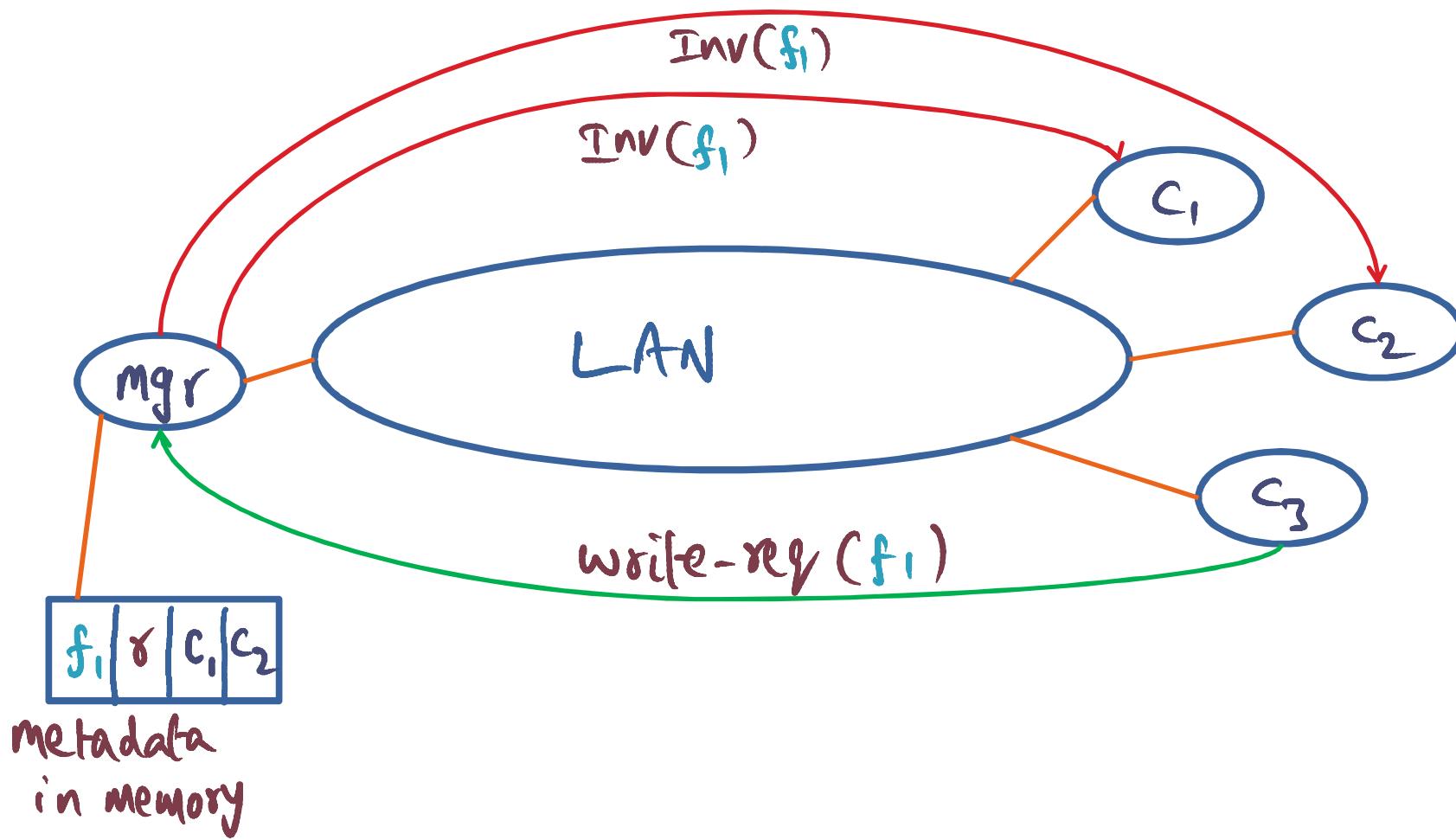
- single writer, multiple readers
- file block unit of coherence



cooperative caching

Cache Coherence

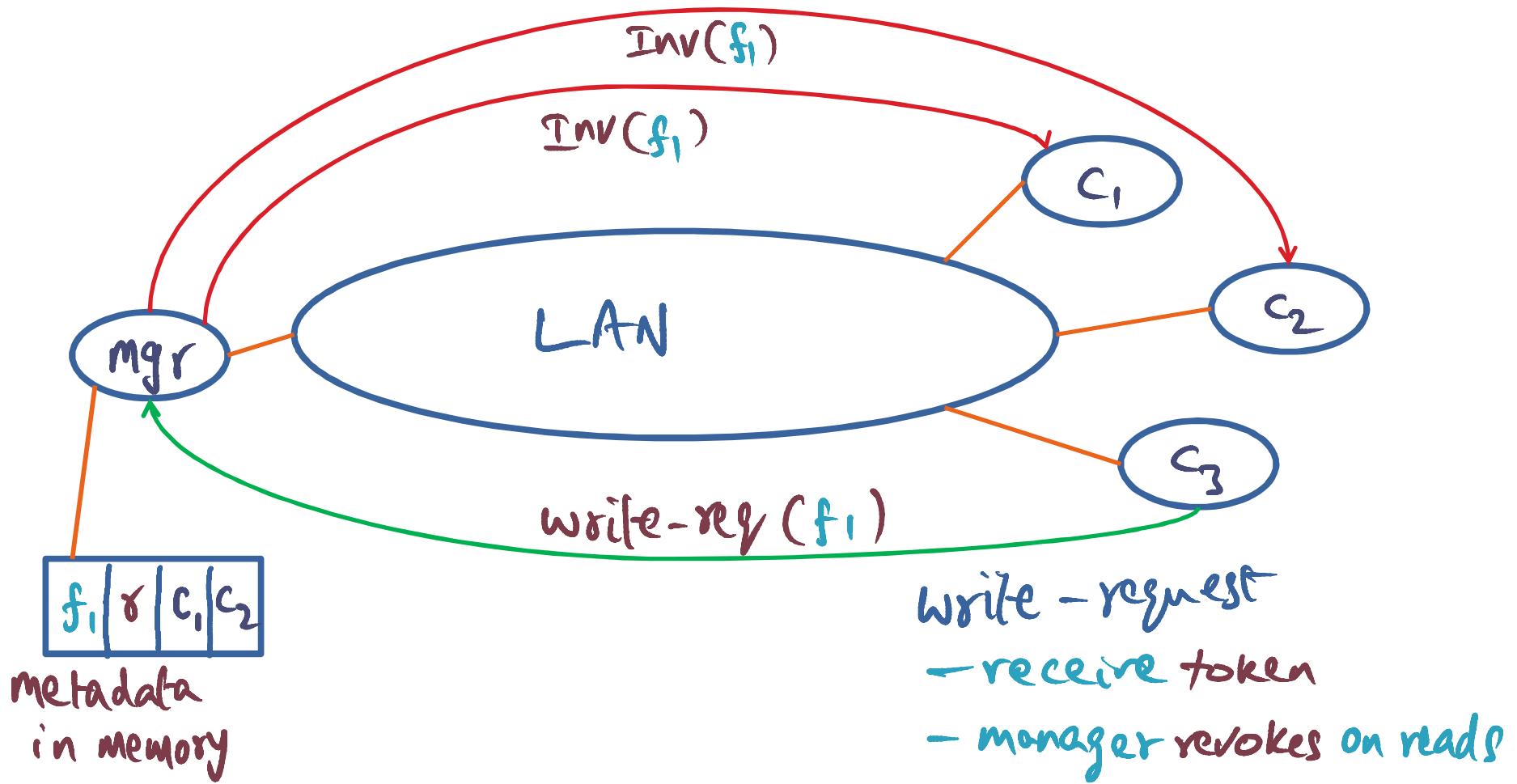
- single writer, multiple readers
- file block unit of coherence



cooperative caching

Cache Coherence

- single writer, multiple readers
- file block unit of coherence



Log cleaning

Log Segment evolution

1' 2' 5'

seg1

writes to file blocks 1, 2, 5
—may belong to different files

Log cleaning

Log Segment evolution

X' 2' 5'

seg1

writes to file blocks 1, 2, 5
— may belong to different files

1" 3' 4'

seg2

block 1 overwritten \Rightarrow Kill old block
writes to blocks 3 + 4

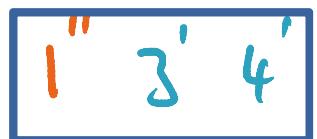
Log cleaning

Log Segment evolution



seg1

writes to file blocks 1, 2, 5
— may belong to different files



seg2

block 1 overwritten \Rightarrow Kill old block
writes to blocks 3 + 4



seg3

block 2 overwritten \Rightarrow Kill old block

Log cleaning

Log Segment evolution



seg1

writes to file blocks 1, 2, 5
—may belong to different files



seg2

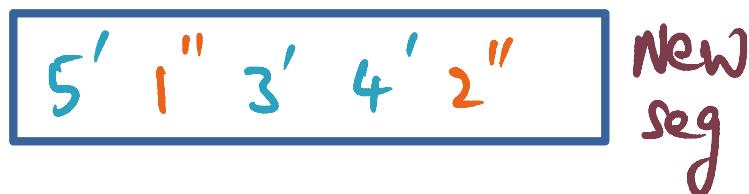
block 1 overwritten \Rightarrow Kill old block
writes to blocks 3 + 4



seg3

block 2 overwritten \Rightarrow Kill old block

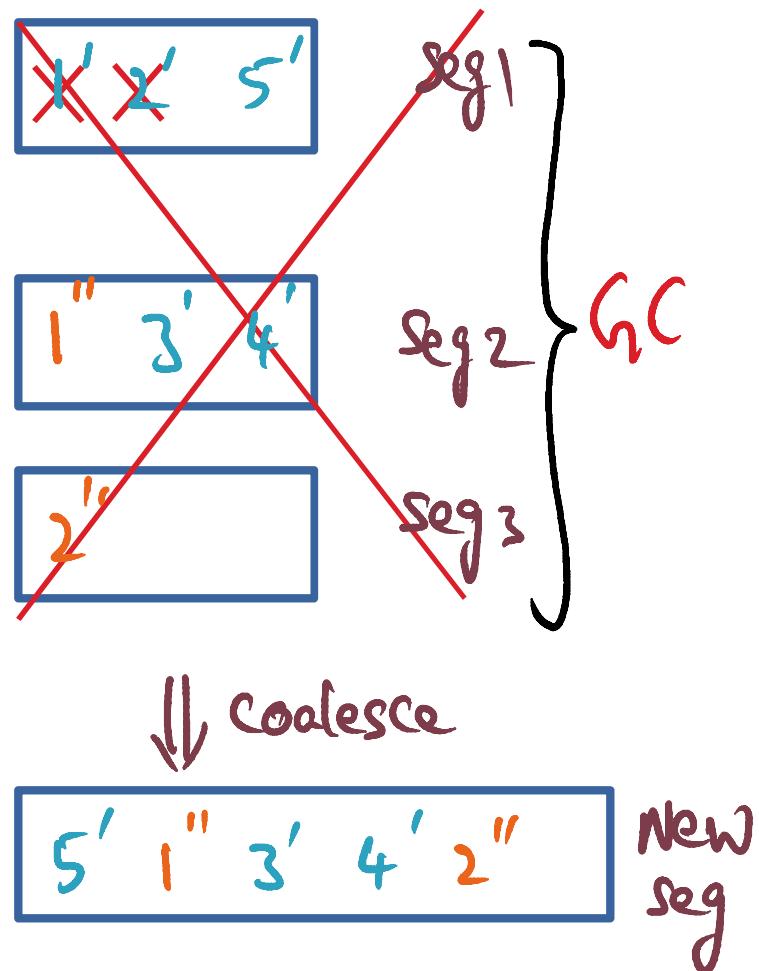
\Downarrow coalesce



Aggregate all "live" blocks into
new segment

Log cleaning

Log Segment evolution



writes to file blocks 1, 2, 5
—may belong to different files

block 1 overwritten \Rightarrow Kill old block
writes to blocks 3 + 4

block 2 overwritten \Rightarrow Kill old block

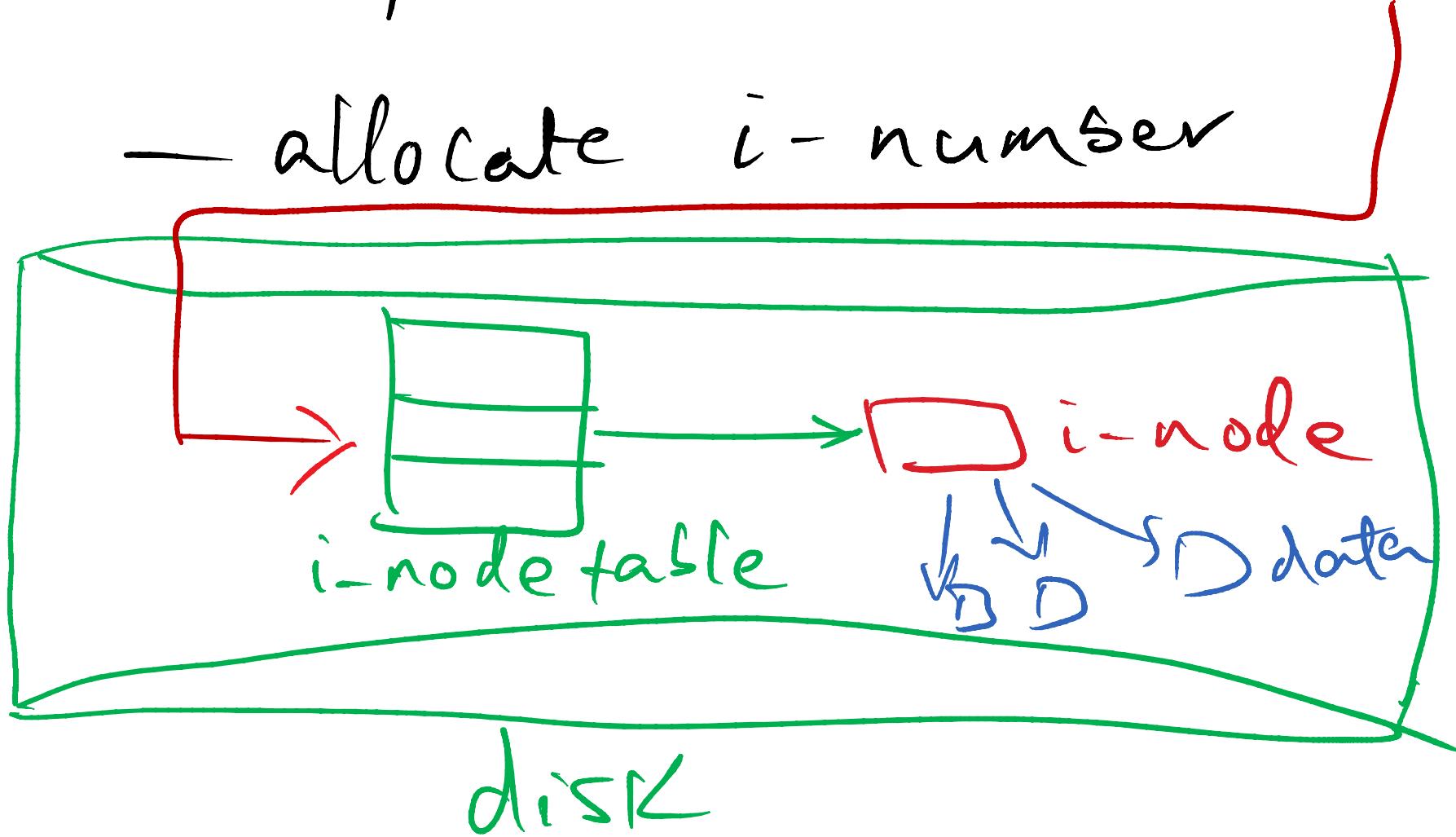
Aggregate all "live" blocks into
new segment

In Unix

on file creation

file table :
name → i-number

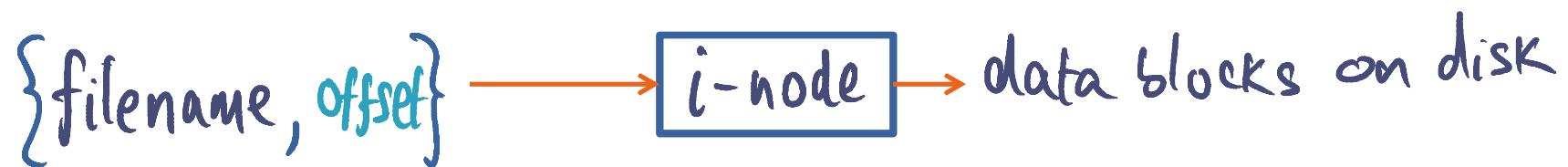
— allocate i-number



File System

- data structures persistant
- But bring them into memory for manipulation
- Applies to data +
metadata (file table
i-node)

Unix File System

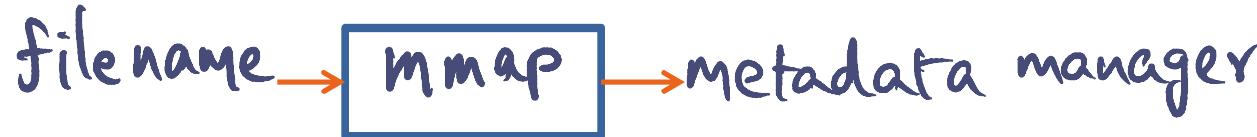


Group activity

- XFS data structures (where are they? Replicated?)
 - File dir
 - Mmap
 - Imap
 - i-node
 - Stripe group
 - logseg
 - Data blocks

- File Dir (on the server that hosts the file) (i.e., on the physical machine where a logical volume resides)
 - Human readable name to i-number mapping
 - File name to index number; xFS stores directories in regular files, allowing a client to learn an index number by reading a directory. In xFS, the index number listed in a directory determines a file's manager. When a file is created, its index number is chosen such that the file's manager is on the same machine as the client that creates the file.
- Manager map (replicated globally)
 - Index number to manager mapping
- Imap (partitioned among managers)
 - Used by manager for file to on-disk log location mapping (split among managers)
- Stripe group map (replicated globally)
 - Log segment id to set of servers mapping

Client node action

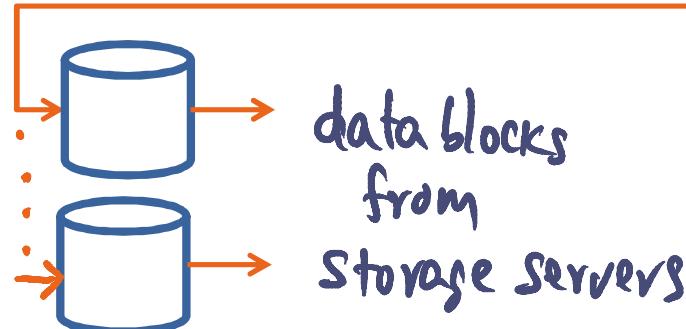
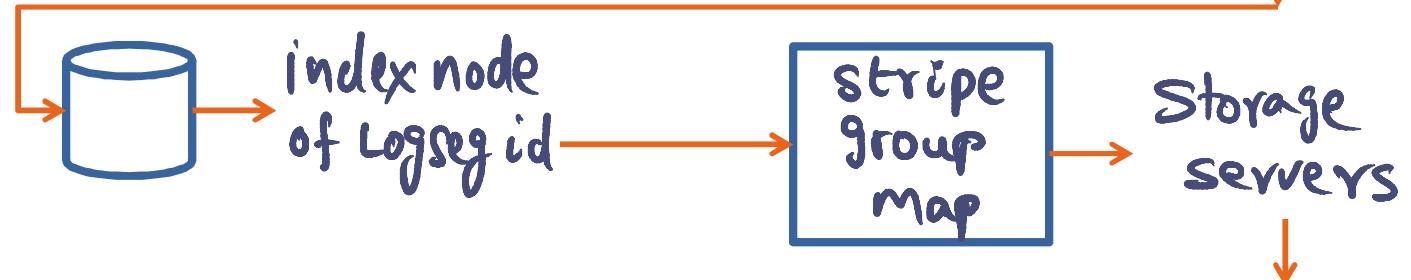
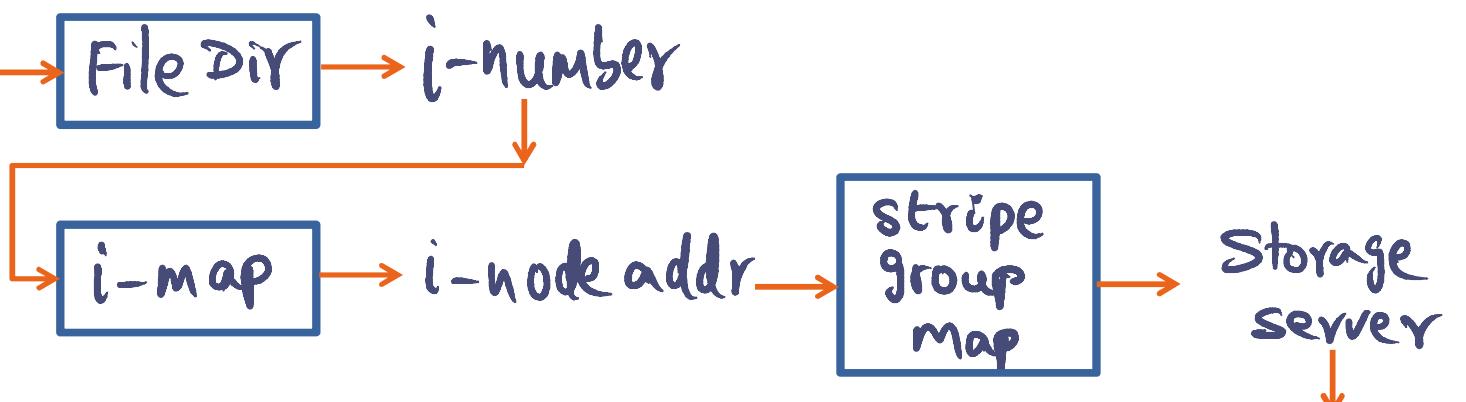


xFS

data

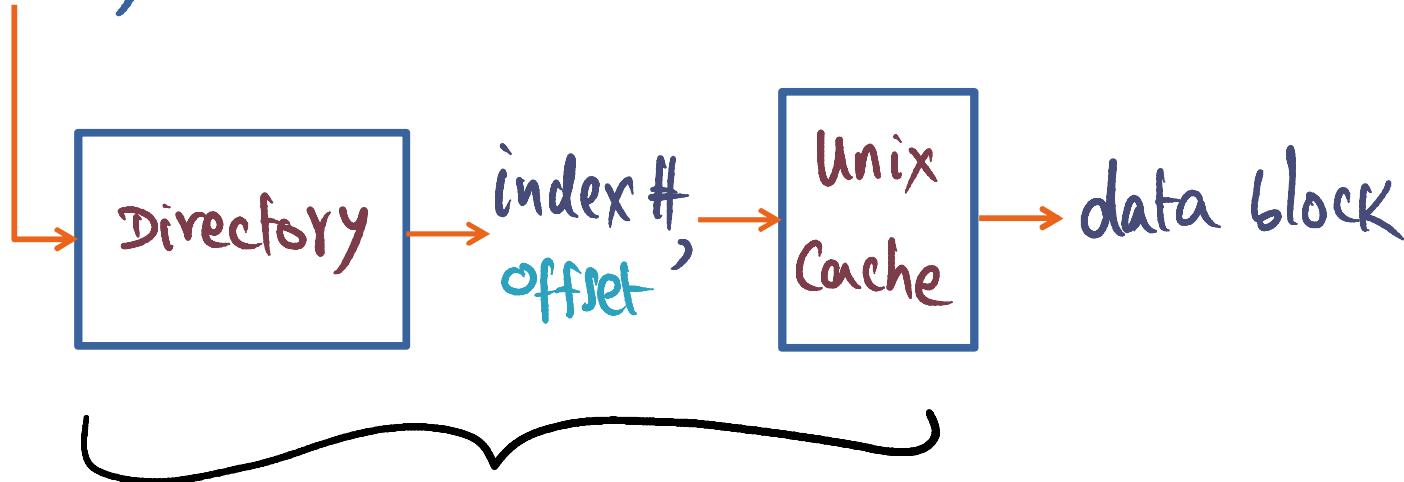
Structures

Manager node actions



Client Reading a file – Own Cache

{filename, offset}



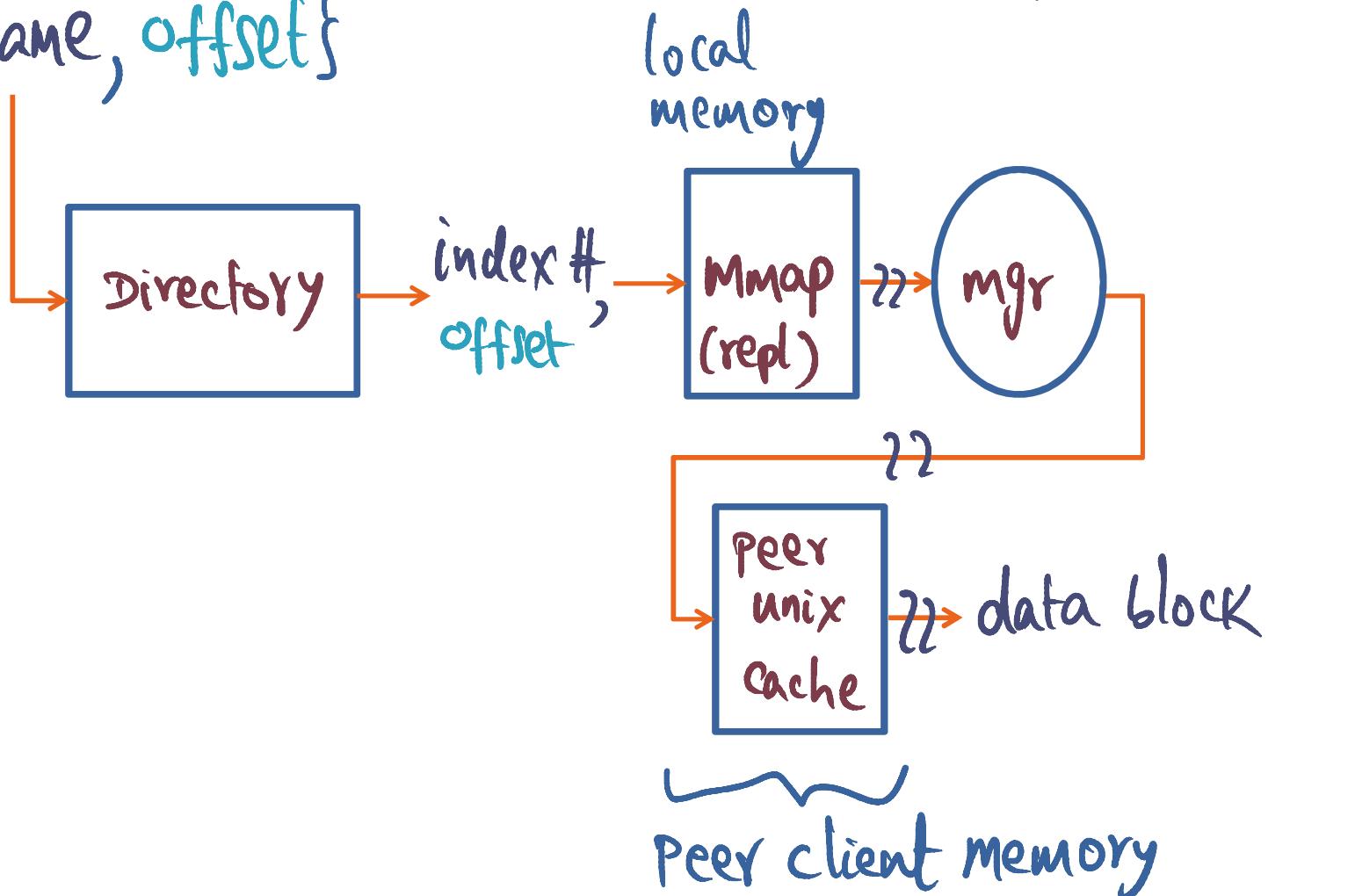
Fastest path for file access

All in client memory

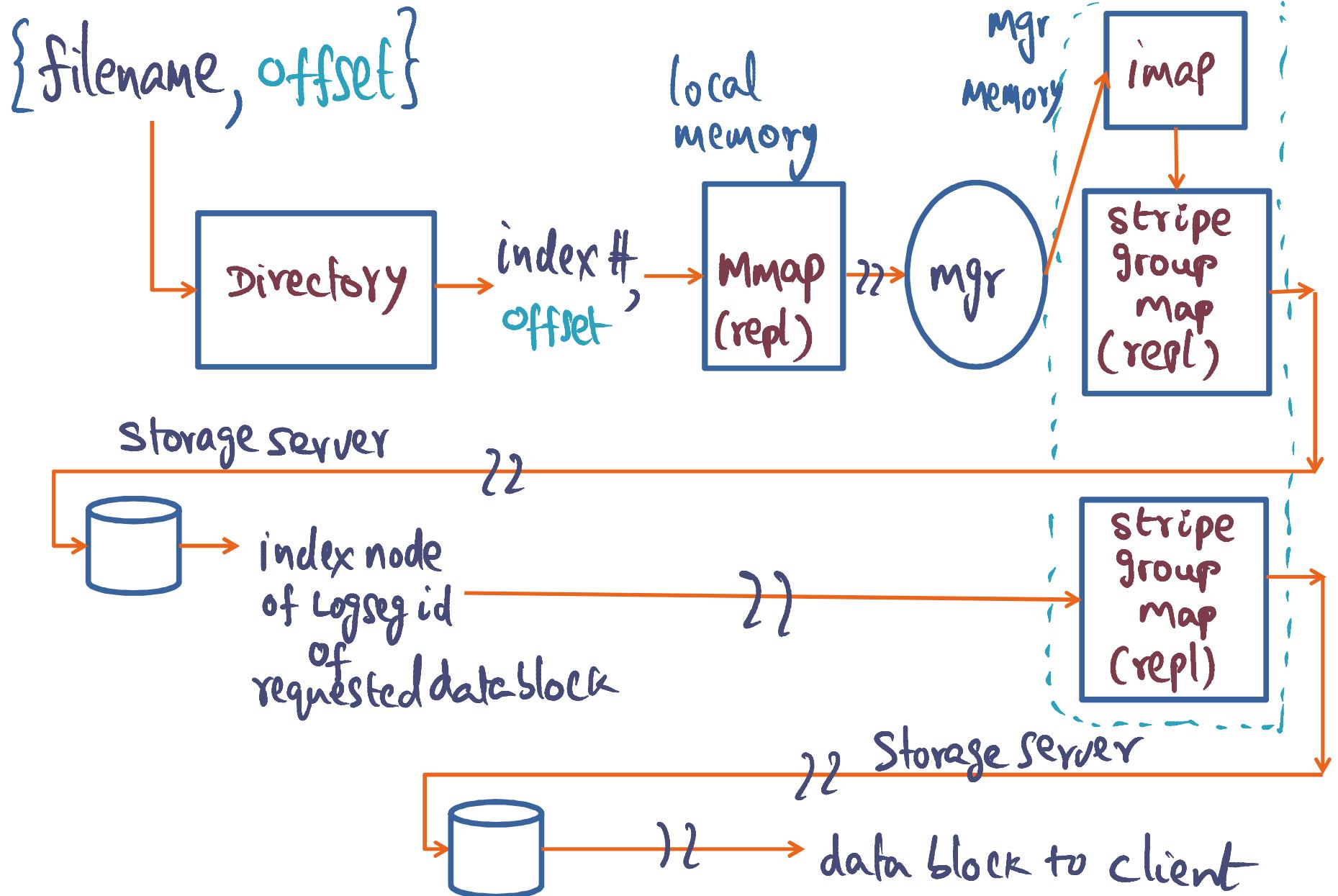
locally cached

Client Reading a file - get from peer cache

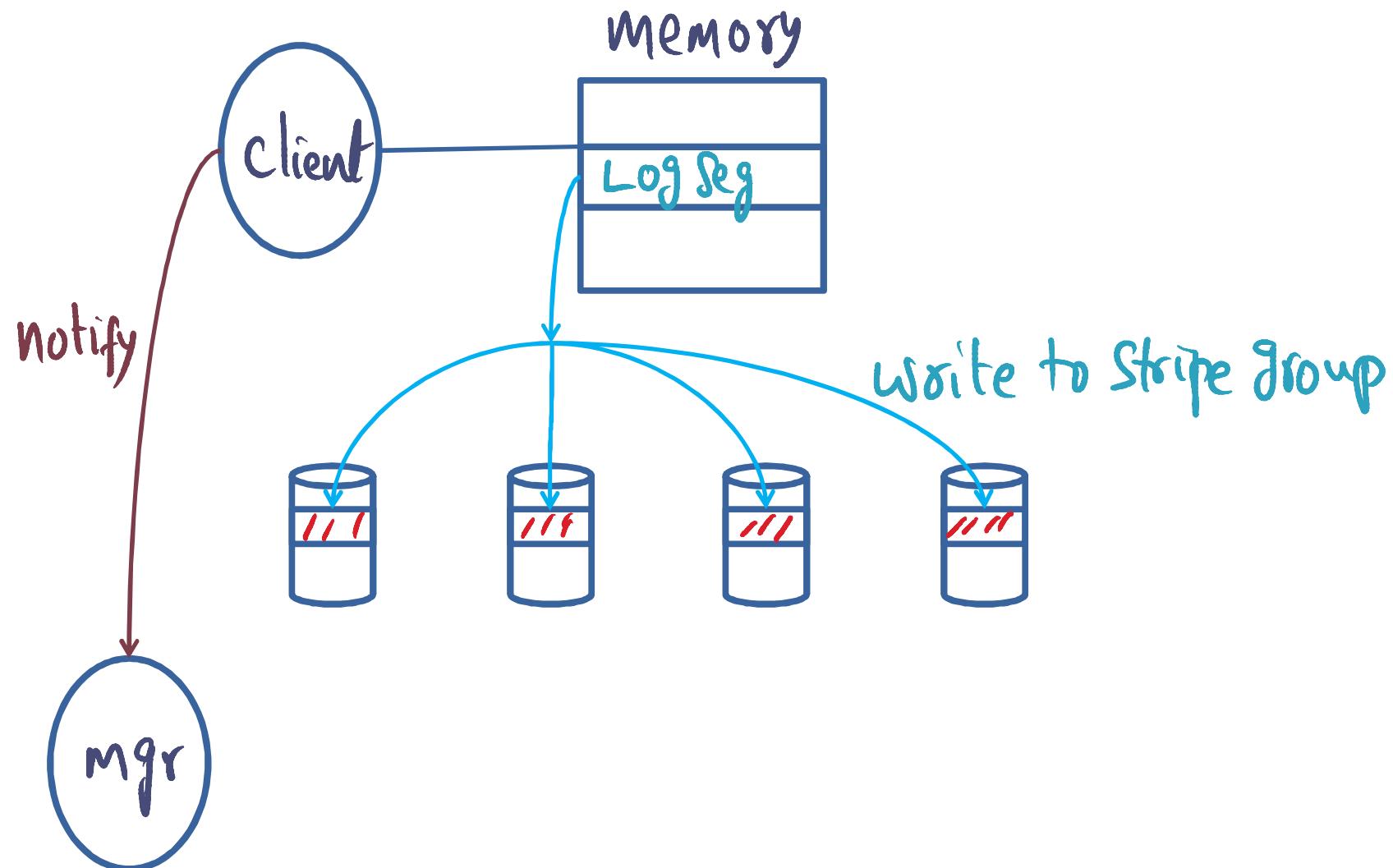
{filename, offset}



Client Reading a file – the real long way!



Client writing a file



Key Takeaways

- Network file systems are important components of any computing environment – corporate or university. There are companies that have sprung up (such as **NetApp**) solely to peddle scalable NFS products.
- Design and implementation of distributed file systems, in particular how to make the implementation **scalable** by removing centralization and utilizing the memory available in the nodes of a LAN intelligently.
- Such **techniques** for identifying and removing bottlenecks are the **re-usable nuggets** that can be taken and applied to the design and implementation of other distributed subsystems in addition to file systems themselves.

Overall, creative ways to fully utilize idle memory available in the nodes of a LAN.