

Persistence: File System Implementation

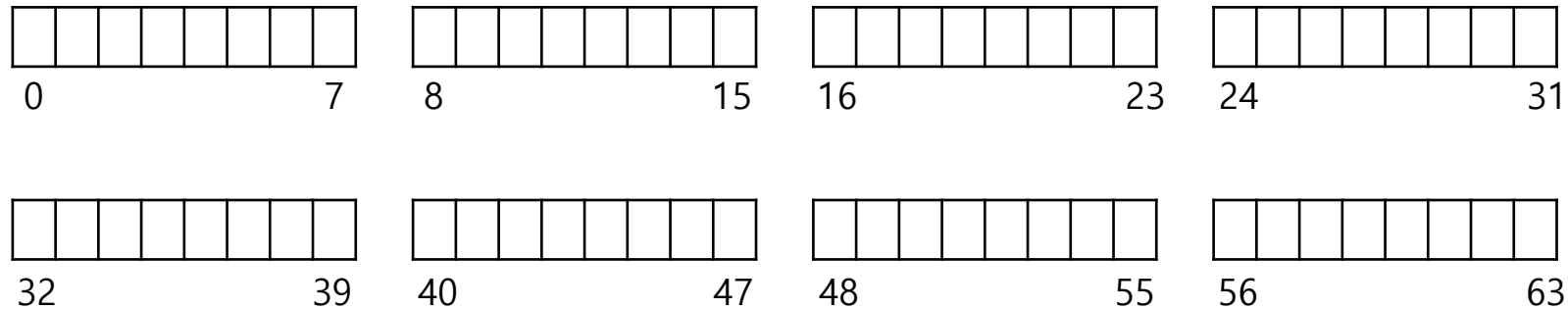
Sridhar Alagar

Review: Files

- Files: sequence of bytes with logical addresses
- Inode: contains meta data; one inode per file
- Directory: maps file name to inode number
- Per process open file table: file descriptor for every open file

File System Organization

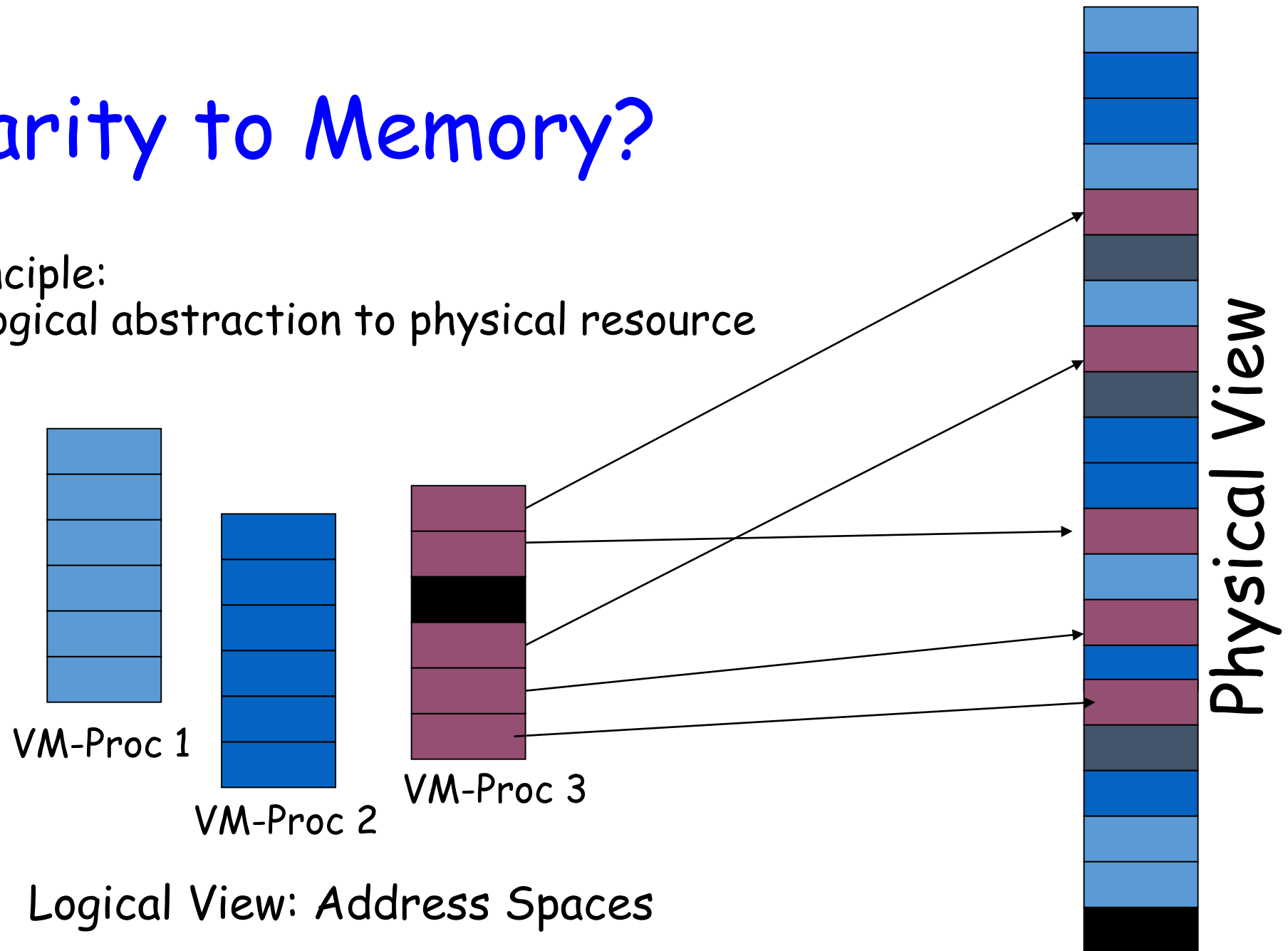
Given: disk is divided into fixed size **blocks**



Want: some structure to map files to disk blocks

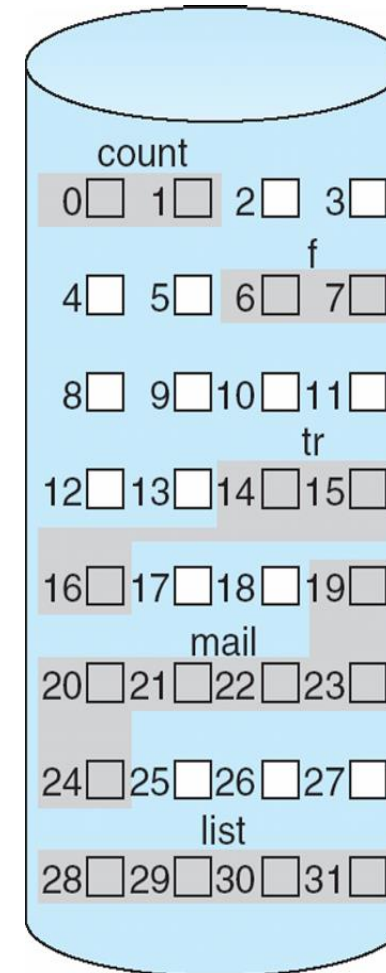
Similarity to Memory?

Same principle:
map logical abstraction to physical resource



Contiguous Allocation

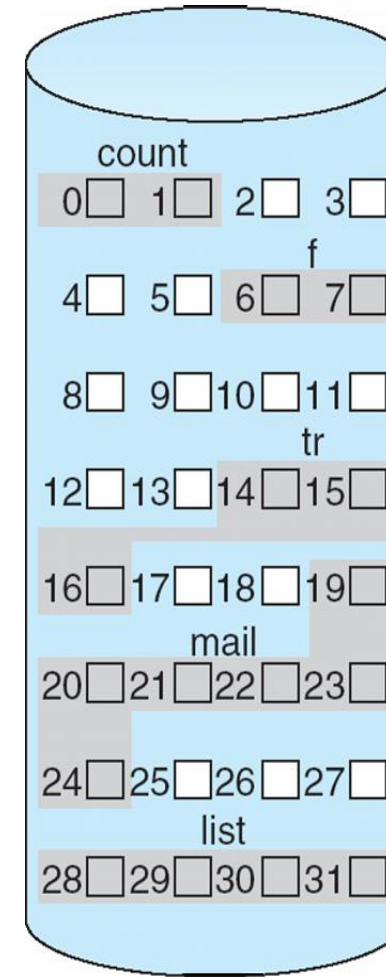
- Allocate required number of blocks(sectors) contiguously
- Simple; little overhead
- Good performance for sequential access
- Easy calculation for random access



file	start	length
count	0	2
tr	14	3
mail	19	6
list	28	4
f	6	2

Contiguous Allocation - Drawbacks

- External fragmentation
- May not be able to grow file
- Need compaction



file	start	length
count	0	2
tr	14	3
mail	19	6
list	28	4
f	6	2

Allocation - Indexed

- Allocate blocks (non-contiguously)
- Keep index table to locate the blocks of a file
- No external fragmentation
- Where is the index table stored?
 - **inode**

Index table for D

4
8
12
15

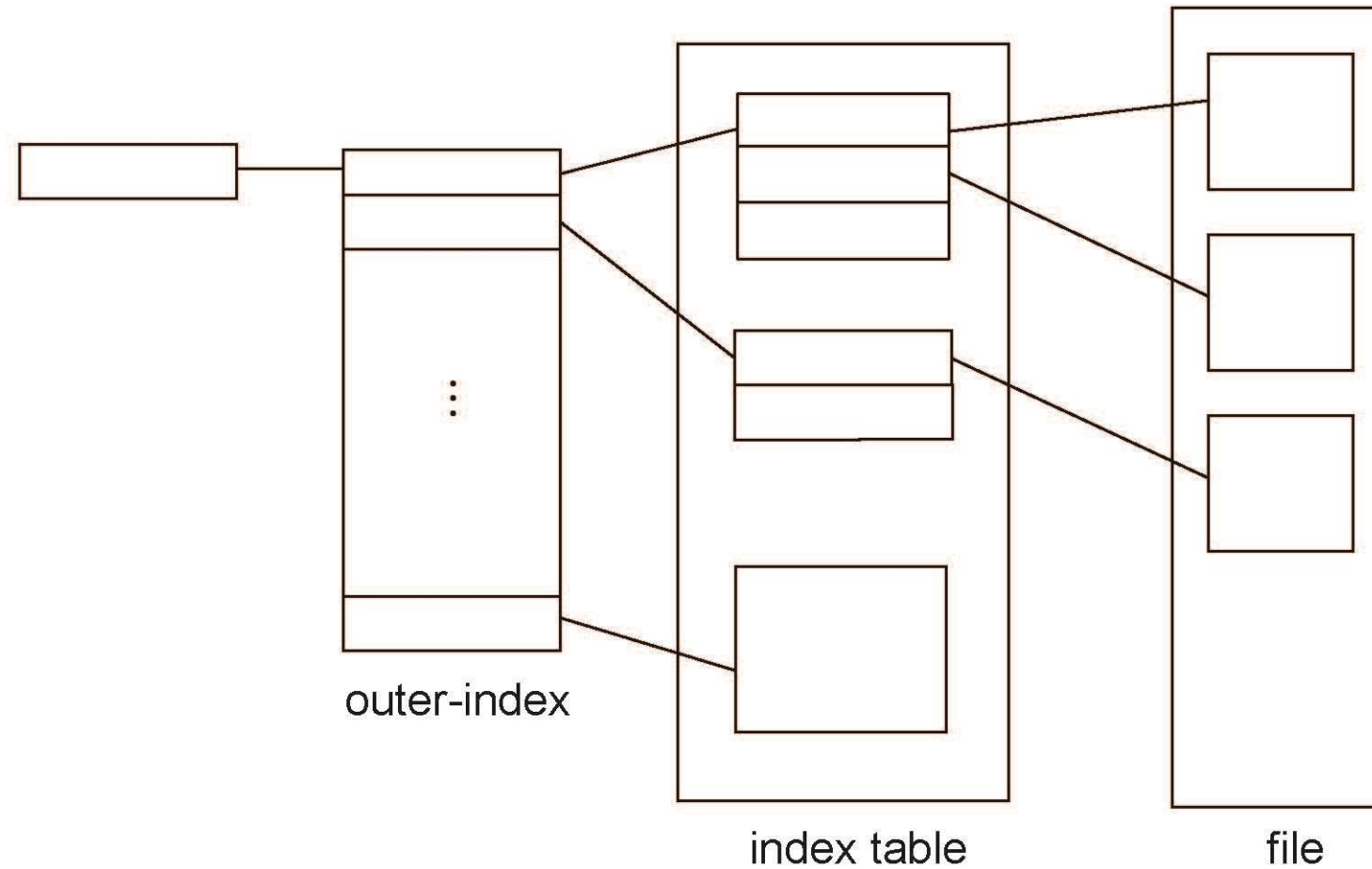
Blocks

	A
	B
	C
4	D
	C
8	D
12	D
15	D

How big is the index table?

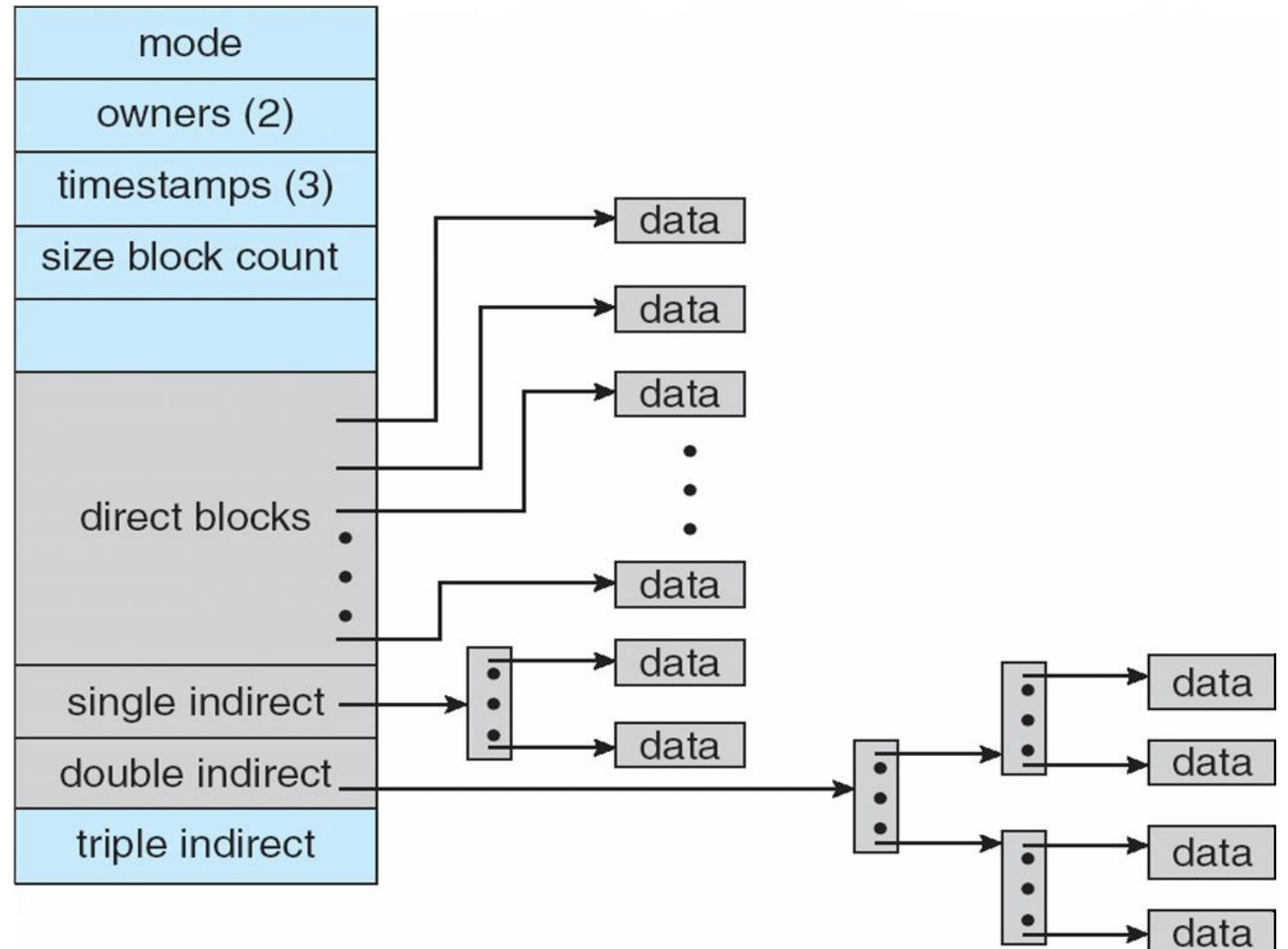
- Assume max file size 1 MB and block size 512 bytes
 - What is the size of index table?
-
- Assume max file size 4 GB and block size 4K bytes
 - What is the size of index table?

Multi-level Index Table

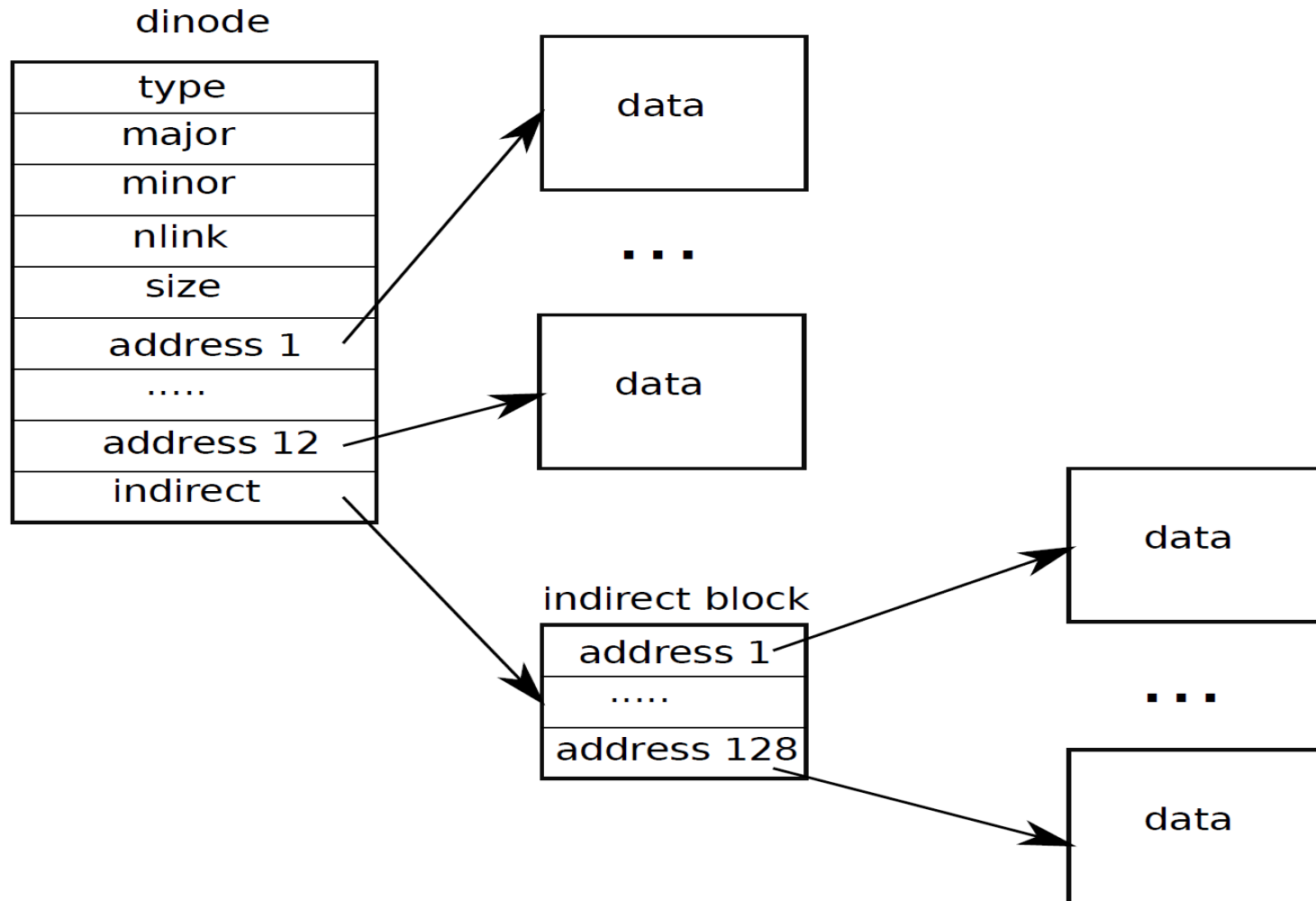


Dynamic Allocation

- Overhead (blocks for indirect pointers) allocated only when the file size grows
- For small file size direct blocks would suffice

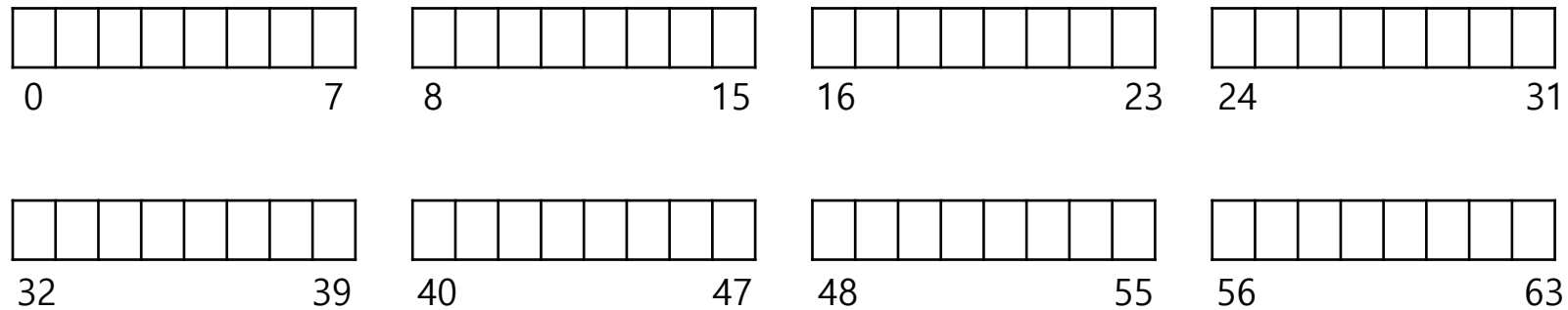


XV6 - Index table



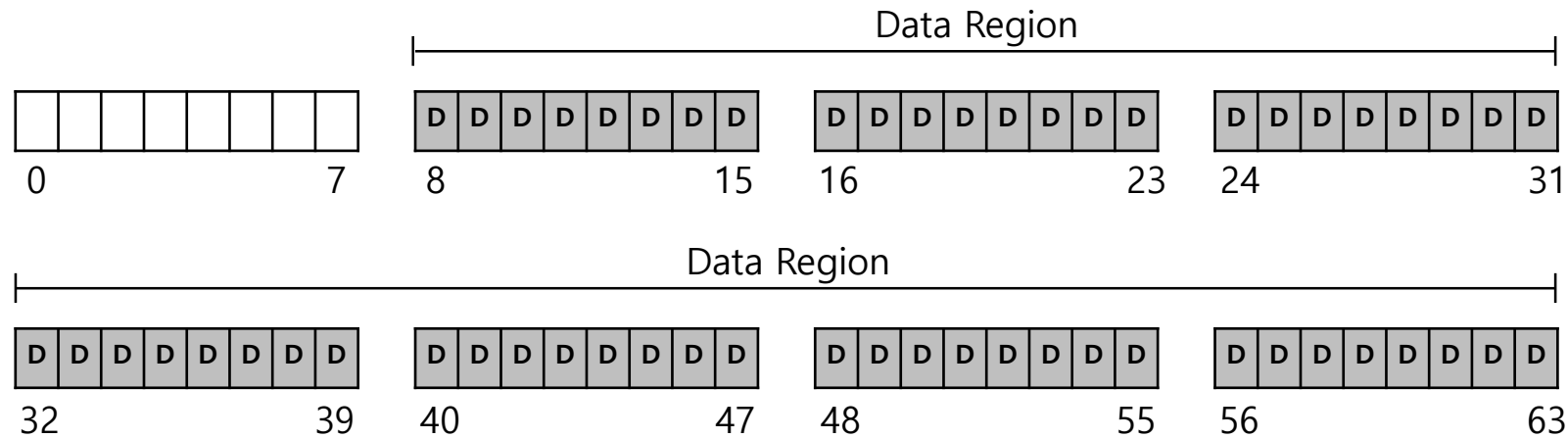
A simple File System

- 64 blocks of size 4 KB each
- Indexed allocation is used



Data region in file system

- Reserve blocks (**data blocks**) to store user data



Where to store **inodes** in the file system?

How many inodes in one block?

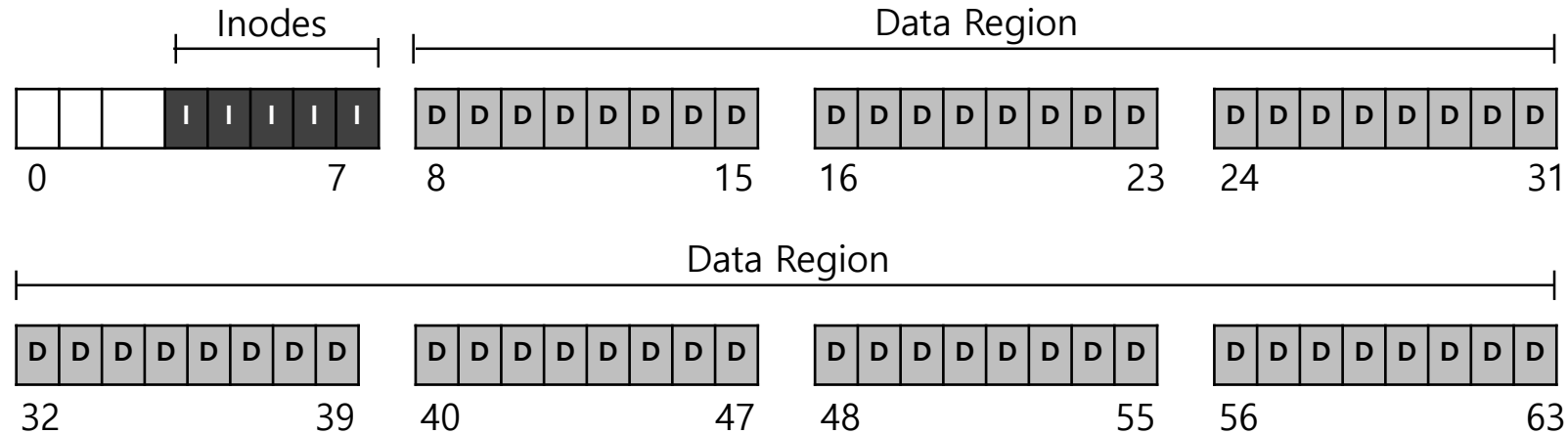
- Each inode is typically 256 bytes (depends on the FS)
- 4KB disk block
- 16 inodes per inode block.

Block 1

inode 0	inode 1	inode 2	inode 3
inode 4	inode 5	inode 6	inode 7
inode 8	inode 9	inode 10	inode 11
inode 12	inode 13	inode 14	inode 15

Inode table in file system

- Reserve some blocks (3 - 7) for inodes



FS Organization: The inode

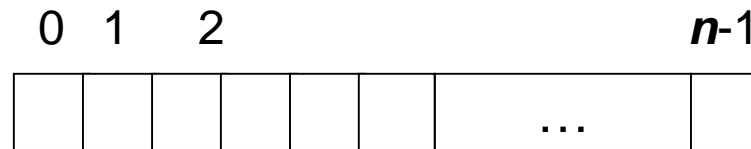
- How to find an inode, say, 32?

The Inode table

				iblock 0				iblock 1				iblock 2				iblock 3				iblock 4			
				0	1	2	3	16	17	18	19	32	33	34	35	48	49	50	51	64	65	66	67
				4	5	6	7	20	21	22	23	36	37	38	39	52	53	54	55	68	69	70	71
				8	9	10	11	24	25	26	27	40	41	42	43	56	57	58	59	72	73	74	75
				12	13	14	15	28	29	30	31	44	45	46	47	60	61	62	63	76	77	78	79
0KB	4KB	8KB	12KB	16KB				20KB				24KB				28KB				32KB			

How to find free data blocks/inodes?

- Use bit vector or bit map. 1 bit per block



$$\text{bit}[i] = \begin{cases} 1 \Rightarrow \text{block } i \text{ free} \\ 0 \Rightarrow \text{block } i \text{ occupied} \end{cases}$$

Block number = (number of bits per word) * (number of 0-value words) + offset of first 1 bit

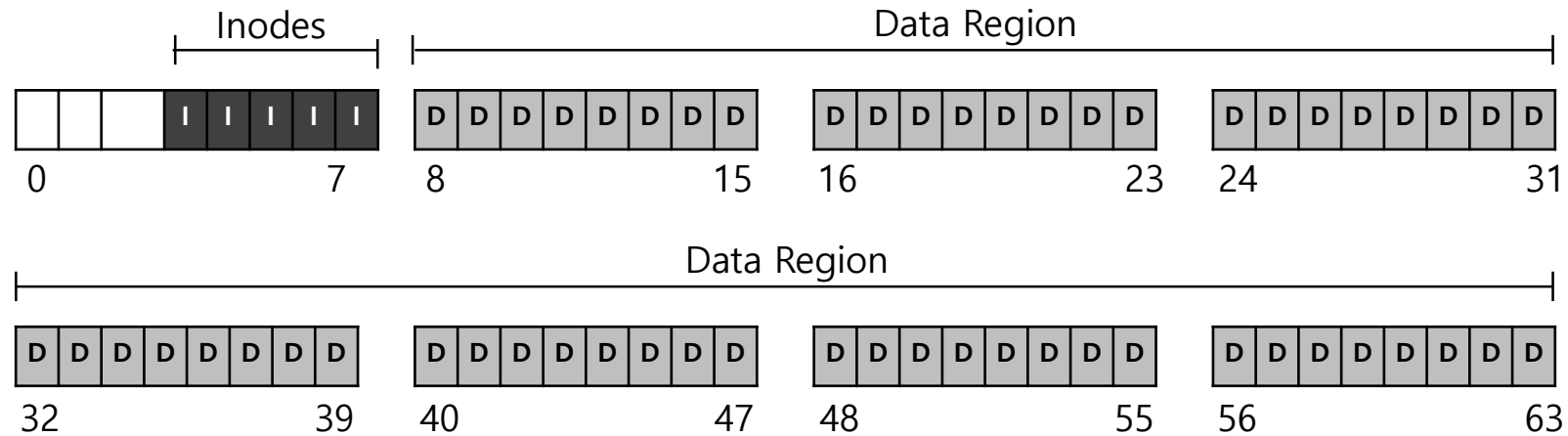
CPUs have instructions to return offset of first “1” bit in a word

How big is bitmap?

- In our example: 56 bits data blocks, and 80 for inodes
- Disk size = 1 TB, block size = 4KB. Bitmap size?
 - block size = 4KB = 2^{12} bytes
 - disk size = 2^{40} bytes (1 terabyte)
 - $2^{40}/2^{12} = 2^{28}$ bits (or 32MB)
 - $2^{13} = 8192$ blocks

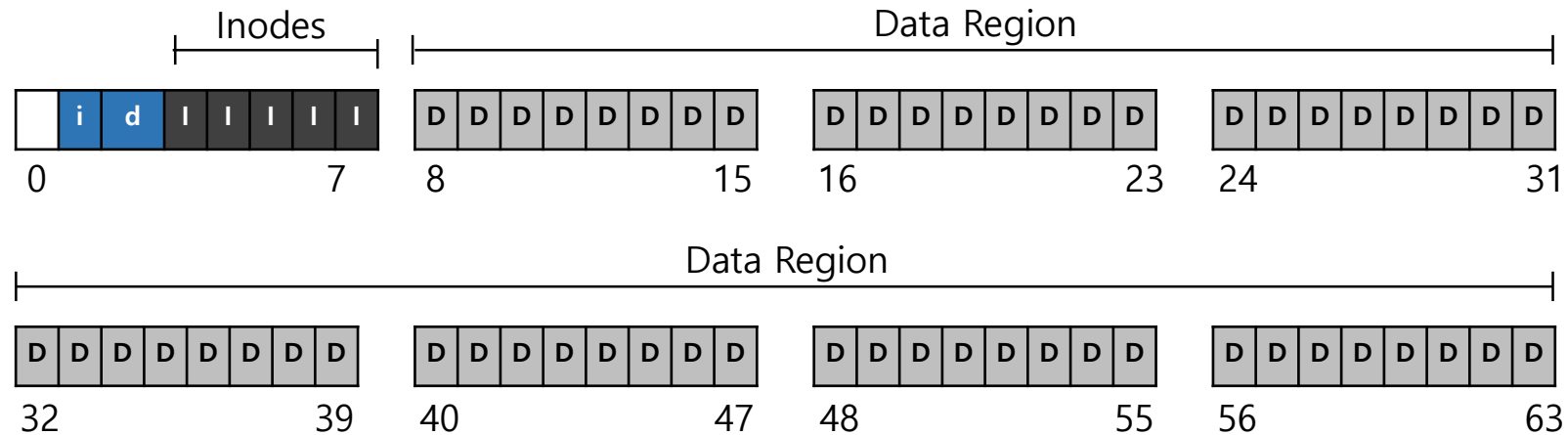
Where are bitmaps stored?

- Bitmap for data blocks
- Bitmap for inodes



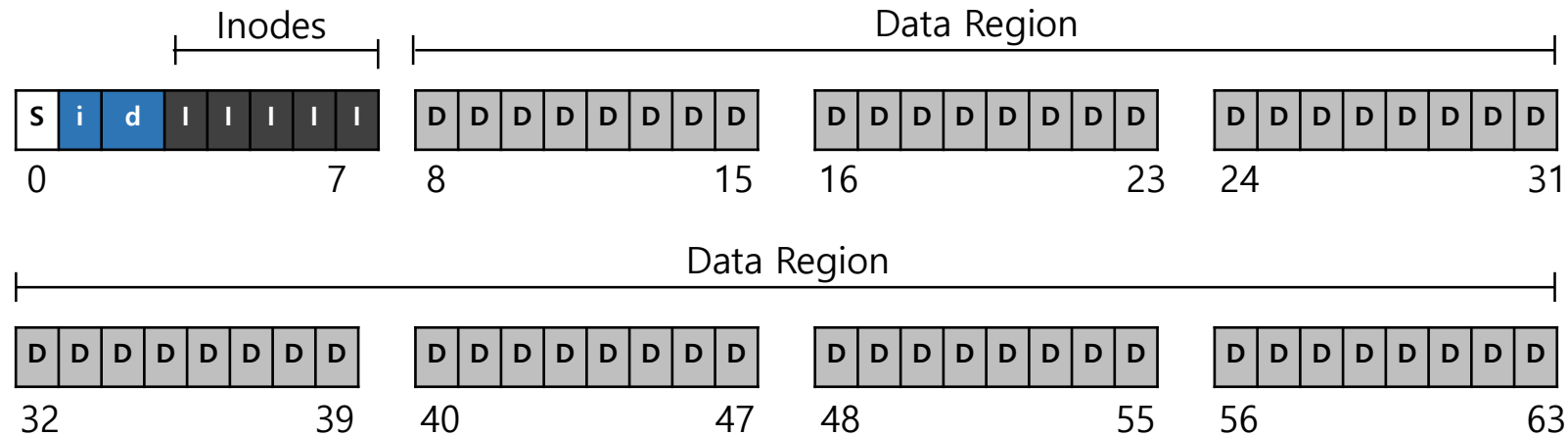
Where are bitmaps stored?

- Bitmap for data blocks is stored in block 1
- Bitmap for inodes is stored in block 2



Superblock

- Superblock contains information about the file system
 - disk size, block size, # inodes, beginning of inode table, etc.
- While mounting, superblock is read first



On-Disk Structure

Super Block

Data Bitmap

Inode Bitmap

Inode Table

Data Block
directories indirects

Operations

- What on-disk structures these operations change?
 - Create()
 - Open()
 - Write()
 - Read()
 - Close()

create /foo/bar: What needs to be read and written?

data bitmap	inode bitmap	root inode	foo inode	bar inode	root data	foo data
		read			read	
			read			read
	read write					write
			write	read write		

open /foo/bar

data bitmap	inode bitmap	root inode	foo inode	bar inode	root data	foo data	bar data
		read			read		
			read				
				read		read	

write to /foo/bar (assume file exists and has been opened)

data bitmap	inode bitmap	root inode	foo inode	bar inode	root data	foo data	bar data
read write				read			write
				write			

read /foo/bar - assume opened

data bitmap	inode bitmap	root inode	foo inode	bar inode	root data	foo data	bar data
				read			
				write			read

close /foo/bar

data bitmap	inode bitmap	root inode	foo inode	bar inode	root data	foo data	bar data

nothing to do on disk!

Disclaimer

- Some of the materials in this lecture slides are from the lecture slides by Prof. Arpaci, Prof. Youjip, and other educators. Thanks to all of them.