

文件系统

Section 5

文件

支持数据的持久化存储是现代操作系统的重要功能之一

- 应用程序数据 (可执行文件、动态链接库 ...)
- 用户数据 (文档、图片、视频 ...)
- 系统数据 (配置文件 ...)



Folder Name



FileName.pdf



FileName.docx



FileName.xls



FileName.ppt



FileName.mp3



FileName.mp4



FileName.jpg



FileName.zip



FileName.java

文件

文件 (File) 是操作系统为存储信息而创建的一个抽象概念

- 一个文件就是一个按名存取的字符序列 (an array of bytes)
 - 命名 (naming) 是文件抽象概念的一个重要特征
 - 用一个具有意义的名称来引用系统中的特定数据
- 提供了一种将信息持久化存储并在随后进行访问的机制

文件

文件由两部分组成

- **文件数据 (File Data):** 字符序列
 - 用户可以对其进行创建、读取、写入和删除等操作
 - 具体内容由应用程序负责解释 (for flexibility)
- **文件属性 (File Metadata):** 用于支撑文件功能的其他信息
 - **Size:** 文件大小 (字符序列的长度)
 - **Owner:** 文件的所有者
 - **Protection:** 不同用户对文件的访问权限
 - **Time:** 文件的创建时间、最近访问时间、最近修改时间
 - ...

文件

围绕文件这一概念提供一组特定的文件操作

- `open()`: 打开一个文件
- `read()`: 顺序读取一个打开文件中的若干字节
- `write()`: 顺序向一个打开文件中写入若干字节
- `lseek()`: 移动当前的偏移量
- `fsync()`: 将当前对文件的修改立即写回磁盘
- `close()`: 关闭一个已打开的文件
- `stat()`: 获取文件的属性信息
- ...

文件

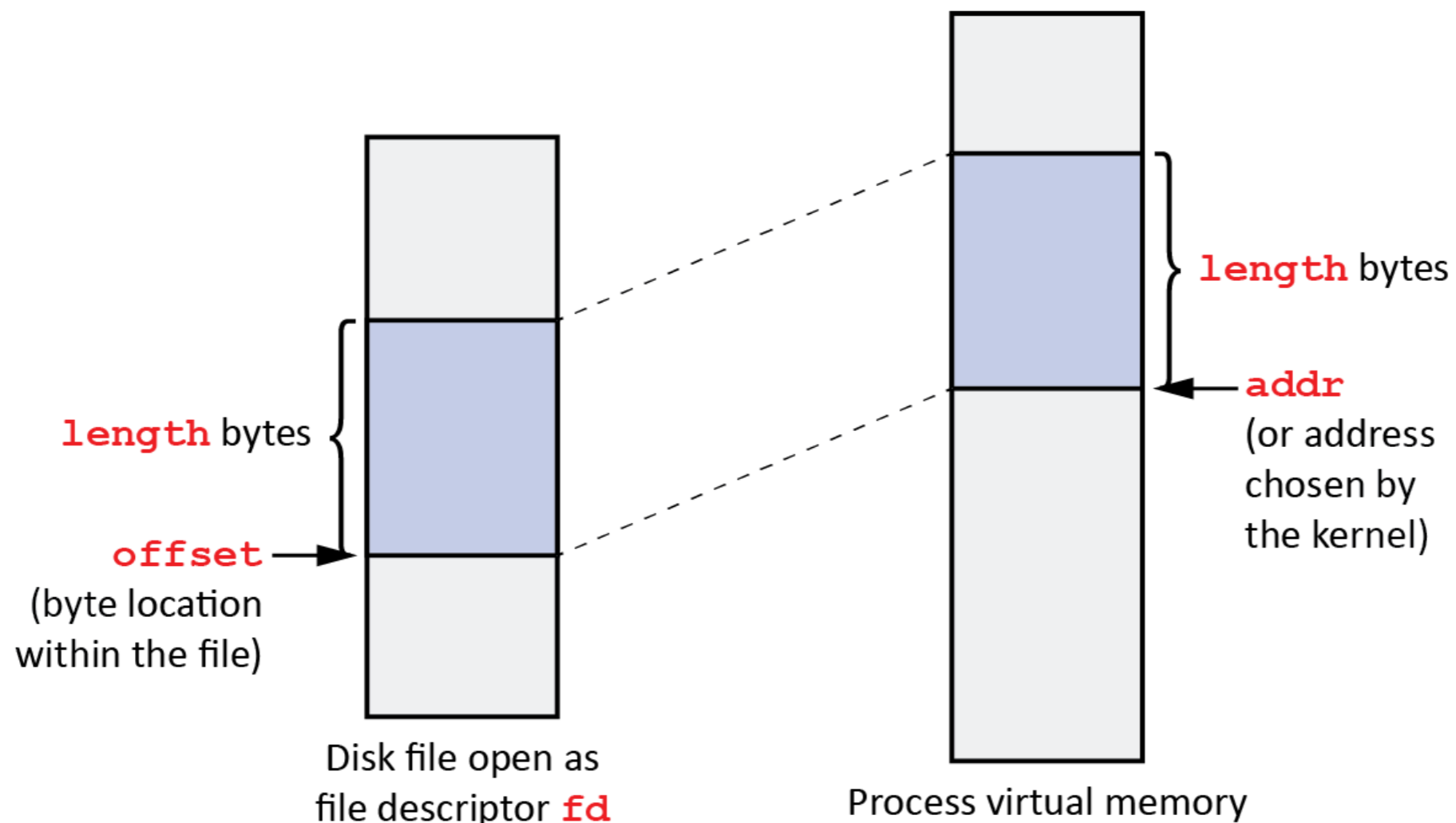
Unix 的 Everything is a File (an array of bytes / an object / a file descriptor)

- 使用同一组 APIs (with a file descriptor) 来操作不同的资源
 - - : 普通文件、d : 目录文件、b : 块设备文件、c : 字符设备文件、p : 命名管道文件、l : 符号链接文件
 - Shell 中的 `stdin (fd = 0)`、`stdout (fd = 1)`、`stderr (fd = 2)`
 - `/proc` : 内核内部数据结构的接口
e.g., `/proc/{PID}/maps`、`/proc/{PID}/mem`
 - `/dev` : 设备驱动的接口
e.g., `/dev/random`、`/dev/null`

文件

内存映射文件 (Memory-Mapped File): 将文件字节序列映射到进程虚拟地址空间，随后进程可以使用内存访问指令来直接访问文件内容

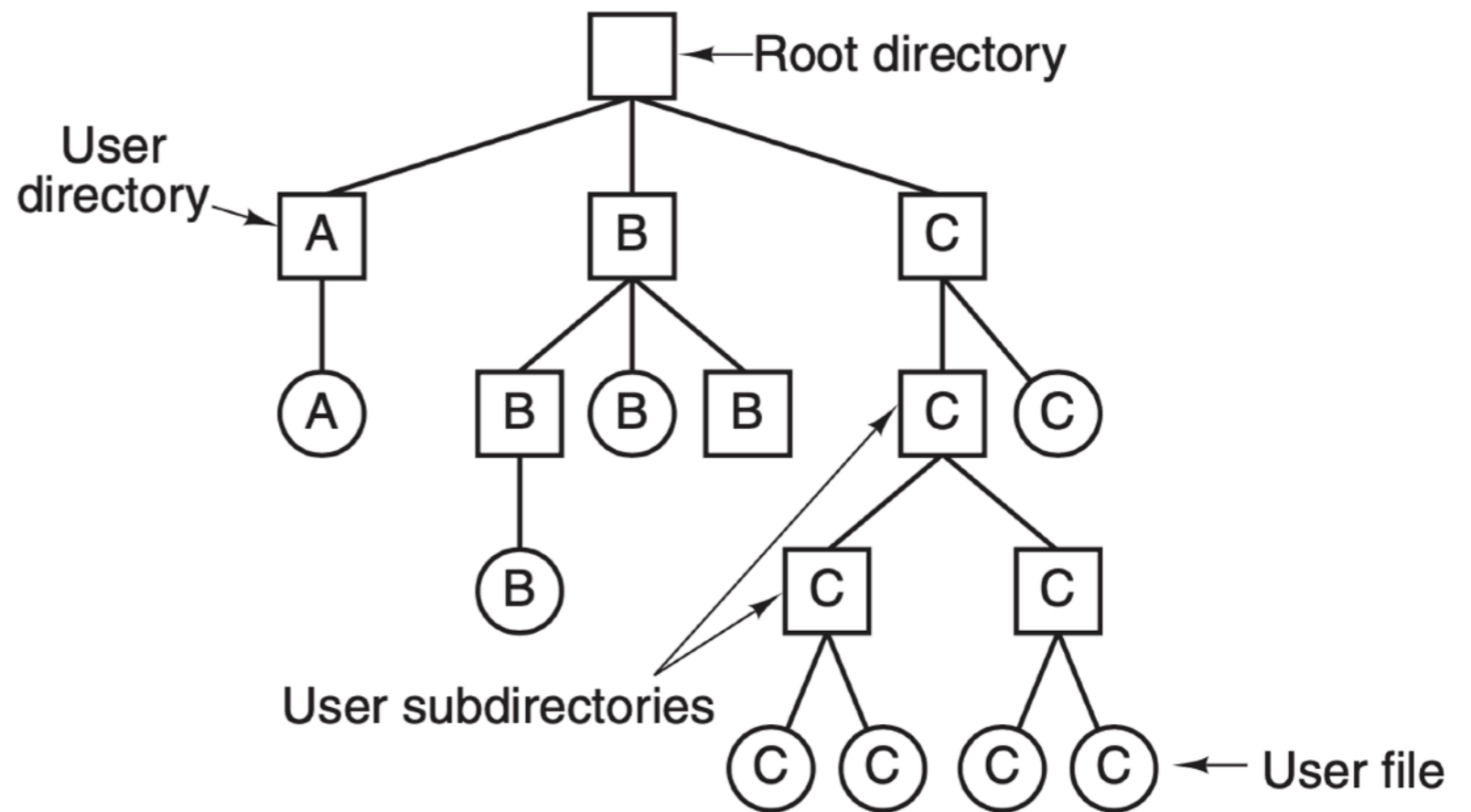
```
void *mmap(void *addr, size_t length, int prot,  
           int flags, int fd, off_t offset)
```



目录

在文件系统中通常使用目录 (directory) 来管理文件

- 从用户视角，可以将逻辑相关的数据放在同一个目录
- 通过在目录中创建子目录，形成一种目录树 (directory tree) 层次结构
 - 绝对路径和相对路径
 - Unix 中的 . 和 ..



目录

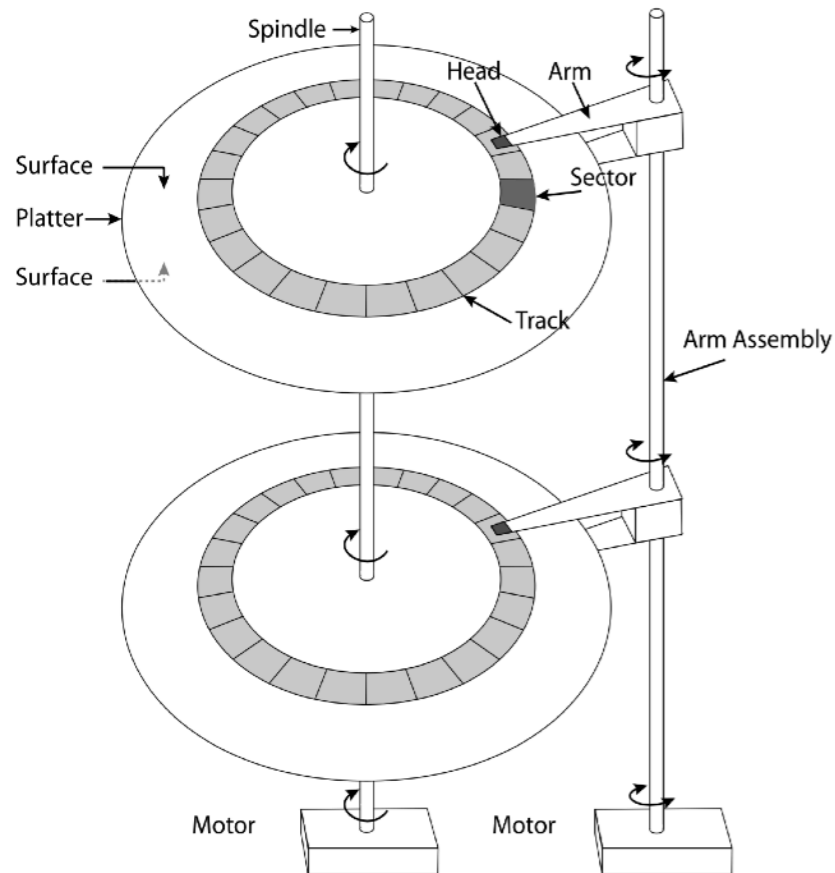
围绕目录也有一组特定的目录操作

- `mkdir()`: 创建一个目录
- `rmdir()`: 删除一个目录
- `opendir()`: 打开一个目录
- `readdir()`: 返回已打开目录中的下一个目录项
- `link()`: 将一个已有文件链接到一个目录项 (创建文件)
- `unlink()`: 删除目录中的一个目录项 (删除文件)
- ...

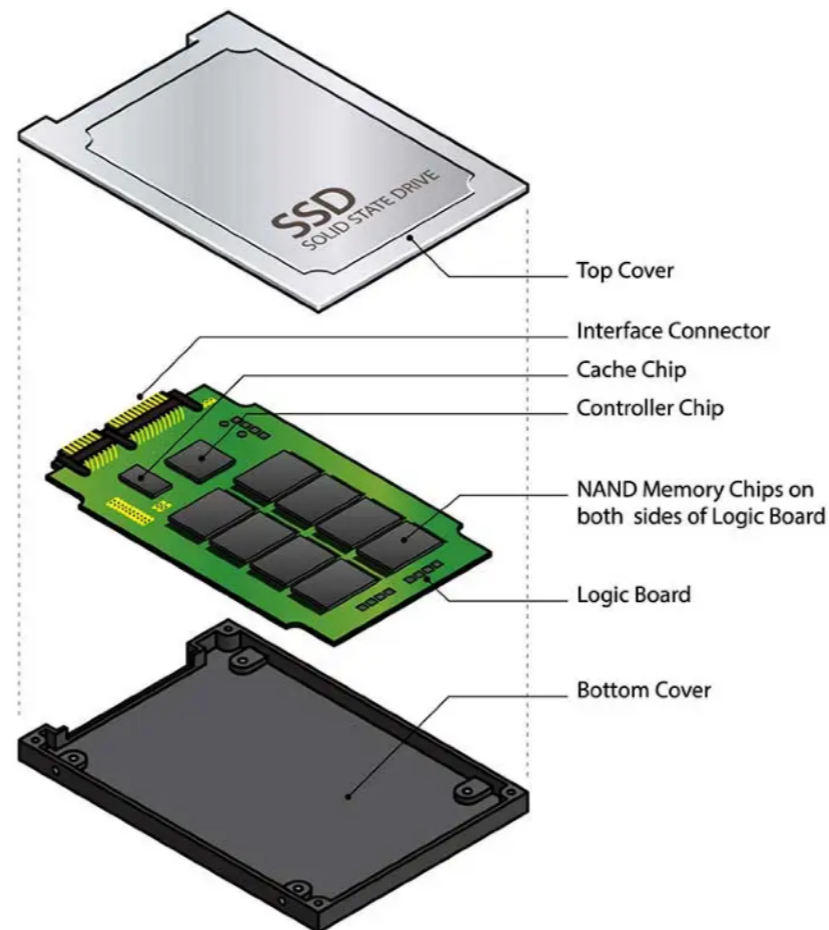
但在持久化存储的最底层



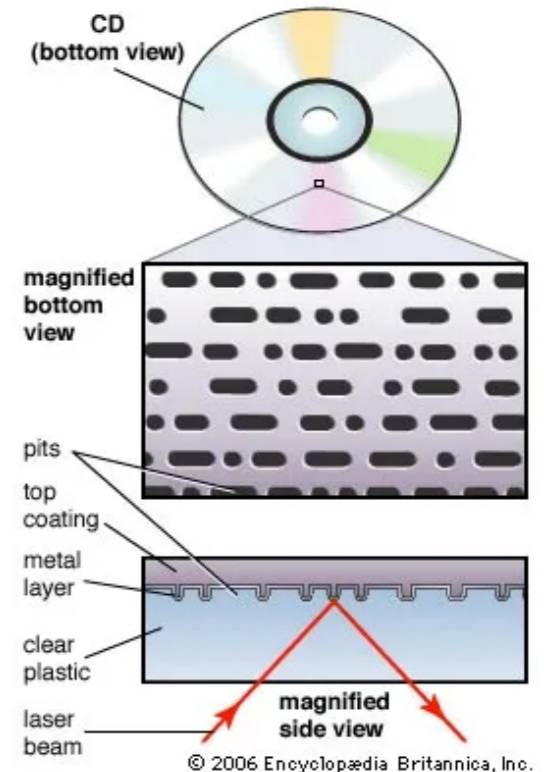
磁盘
(Hard Disk)



固态硬盘
(Solid State Drive)



光盘
(Compact Disk)

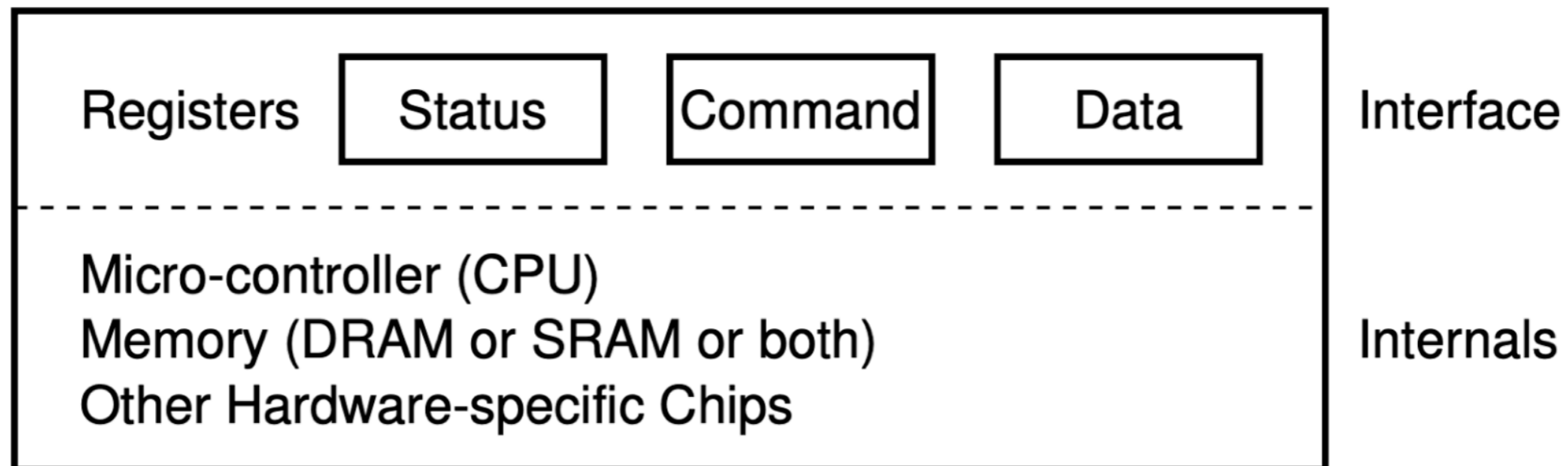


© 2006 Encyclopædia Britannica, Inc.

设备接口

设备提供由若干寄存器组成的设备接口 (interface)

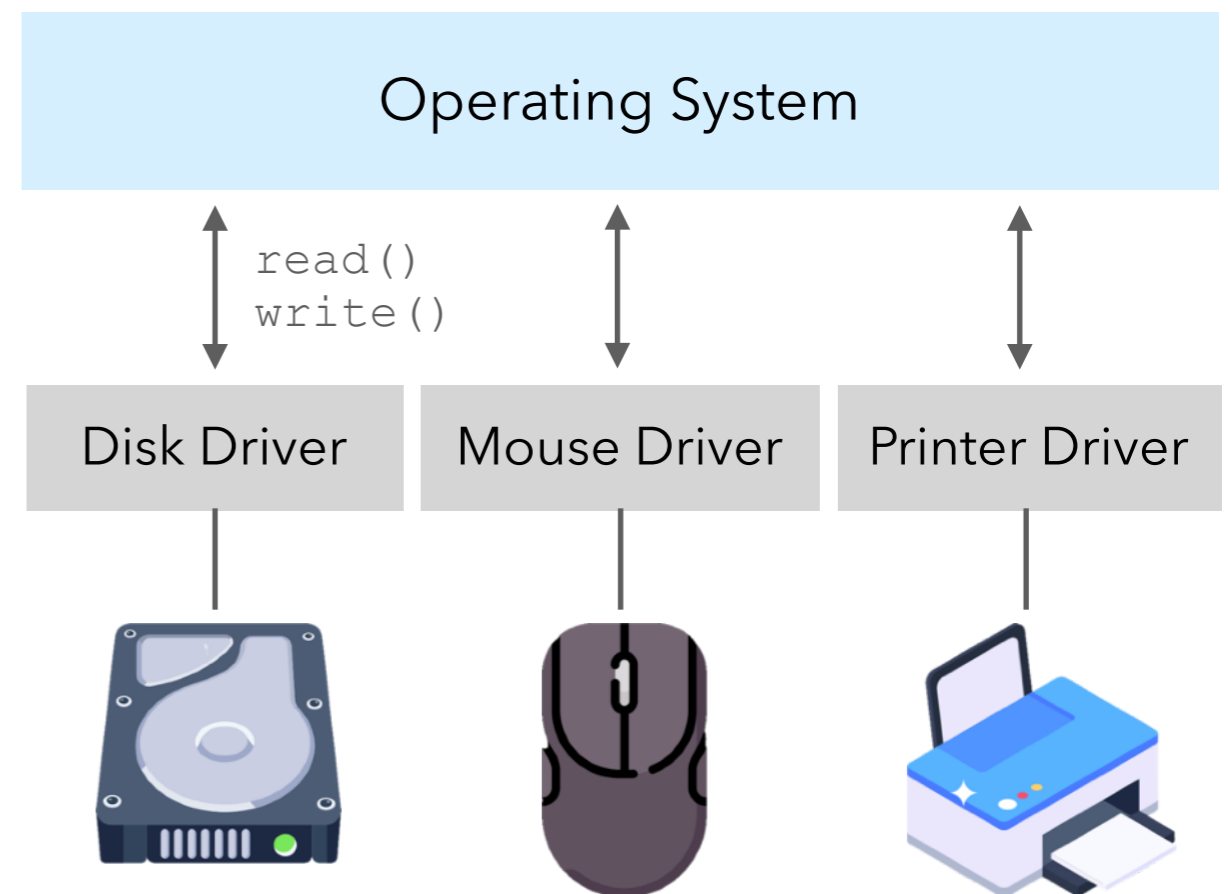
- 状态寄存器 (status register): 设备当前的状态
- 命令寄存器 (command register): 指示设备执行特定的任务
- 数据寄存器 (data register): 将数据传递给设备或从设备获取数据



驱动程序

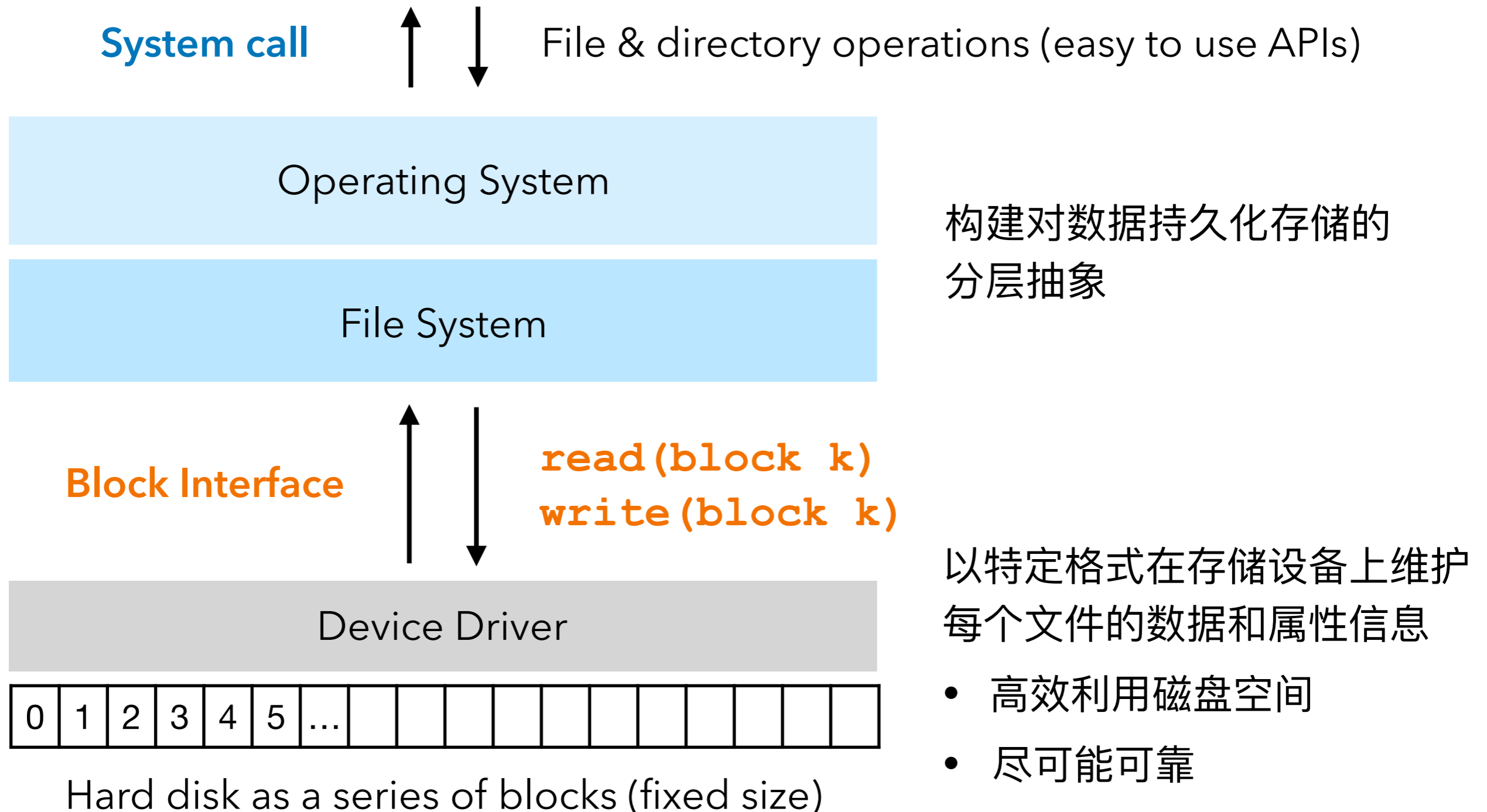
操作系统依靠**设备驱动 (Device Driver)** 实现对底层硬件的抽象

- 通过设备硬件接口与特定设备进行交互
- 向外暴露一组通用的 APIs
 - `read()`, `write()` 和 `ioctl()`
 - 只要实现了这些接口, 就可以看作是一个“设备”
(`/dev/zero`, `/dev/null`, ...)
- 现代操作系统中 70%+ 的代码都是驱动程序代码



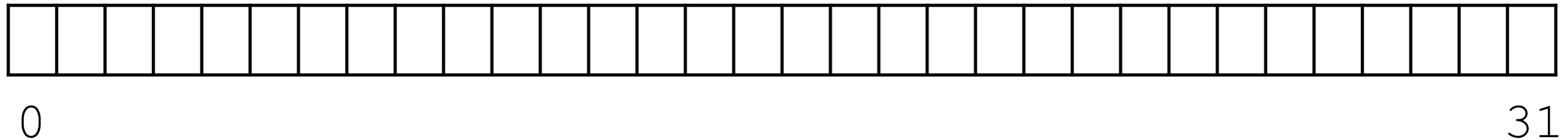
文件系统

实现文件接口并负责管理文件数据和属性的系统即为文件系统



文件系统的实现

File System Layout



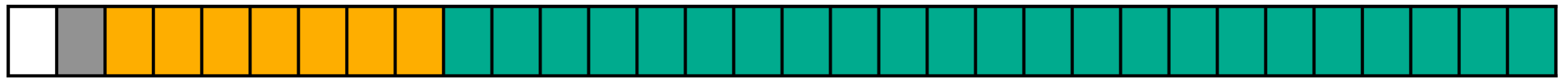
将整个磁盘划分为若干固定大小的磁盘块 (Disk Blocks)

文件系统的实现

File System Layout

Data Blocks

分配一定数量的 blocks 用于存储文件数据 (File Data)



0

File Attributes and Allocation Structure

分配一定数量的 blocks 用于存储文件属性 (File Metadata)、以及空闲空间管理相关的信息

31

Super Block

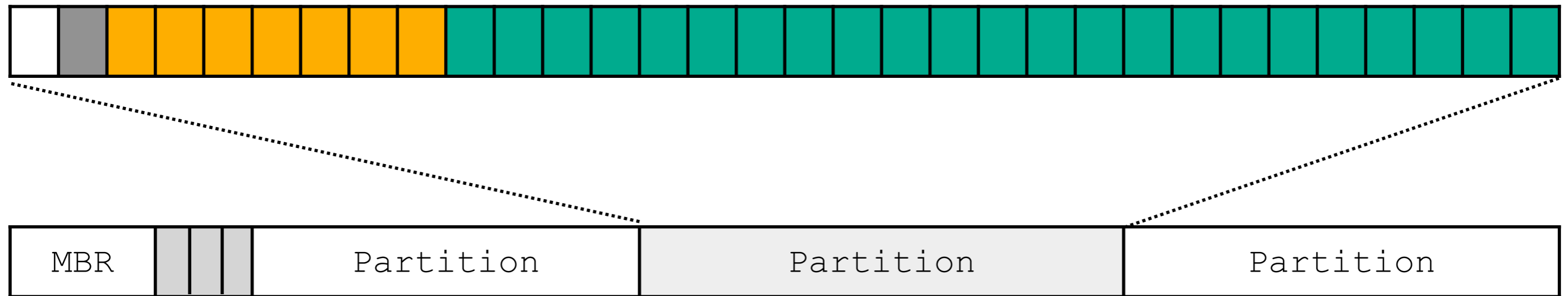
存储文件系统的 metadata 信息

Boot Block

存储 Boot OS 相关的信息 (empty if no OS)

文件系统的实现

File System Layout



Disk Partitions

不同 partition 可格式化为不同文件系统

Partition Table

存储每个 partition 的属性信息 (one is marked as active)

Master Boot Record (MBR)

存储 Boot Computer 相关的信息

文件系统的实现

File Control Block

文件系统需要设计和维护特定的**数据结构**来实现对文件的存储和管理

- **文件控制块 (File Control Block)**
 - 记录文件属性 (File Metadata) 信息
 - 记录文件数据 (File Data) 在磁盘中的存储位置
 - 该结构同样需要在磁盘中持久化存储 (分配某些磁盘块)
- 在 Unix 文件系统中, 该结构被称为 **inode (index node)**

文件系统的实现

File Control Block

Size	Name	What is this inode field for?
2	mode	can this file be read/written/executed?
2	uid	who owns this file?
4	size	how many bytes are in this file?
4	time	what time was this file last accessed?
4	ctime	what time was this file created?
4	mtime	what time was this file last modified?
4	dtime	what time was this inode deleted?
2	gid	which group does this file belong to?
2	links_count	how many hard links are there to this file?
4	blocks	how many blocks have been allocated to this file?
4	flags	how should ext2 use this inode?
4	osd1	an OS-dependent field
60	block	a set of disk pointers (15 total)
4	generation	file version (used by NFS)
4	file_acl	a new permissions model beyond mode bits
4	dir_acl	called access control lists

A simplified **ext2** inode structure

文件的实现方式

记录 File Data 在磁盘中的存储位置：即文件的第 N 个 bytes (第 M 个 logical block) 存储在 File System Layout 中的哪个 Data Block

- 连续方式 (Contiguous Allocation)
- 链表方式 (Linked List Allocation)
 - FAT 文件系统的 File Allocation Table
- 索引方式 (Index Allocation)
 - Unix 文件系统的 inode Structure

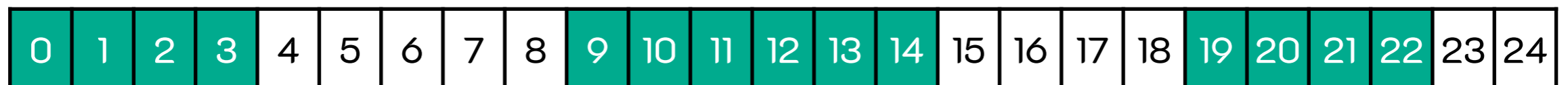
文件的实现方式

Contiguous Allocation

使用磁盘上的若干连续 Data Blocks 来存储 File Data

- ✓ 仅需记录 number of first block + total number of blocks
- ✓ 有较好的顺序读写性能、且能快速定位随机 Data Block

Data Blocks



`a.txt`

`list`

`main.c`

`start = 0`

`length = 4`

文件的实现方式

Contiguous Allocation

使用磁盘上的若干连续 Data Blocks 来存储 File Data

- ✓ 仅需记录 number of first block + total number of blocks
- ✓ 有较好的顺序读写性能、且能快速定位随机 Data Block
- ✗ 但显然不够灵活
 - 在创建文件时需确定文件大小
 - 还会存在严重的外部碎片问题
 - 在特定场景下仍是很好的选择 (e.g., CD-ROM)

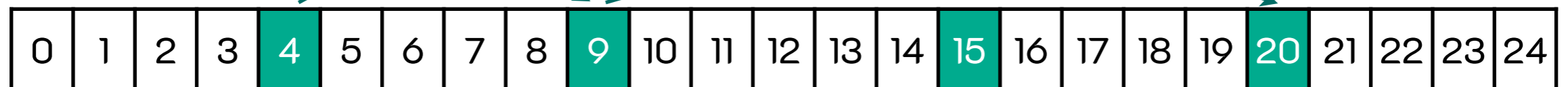
文件的实现方式

Linked List Allocation

使用 **Data Blocks** 的一个链表来存储 File Data

- ✓ 仅需记录 number of first block
- ✓ 可以把所有磁盘 Data Blocks 都利用起来 (没有外部碎片)

Data Blocks



`a.txt`

`start = 4`

文件的实现方式

Linked List Allocation

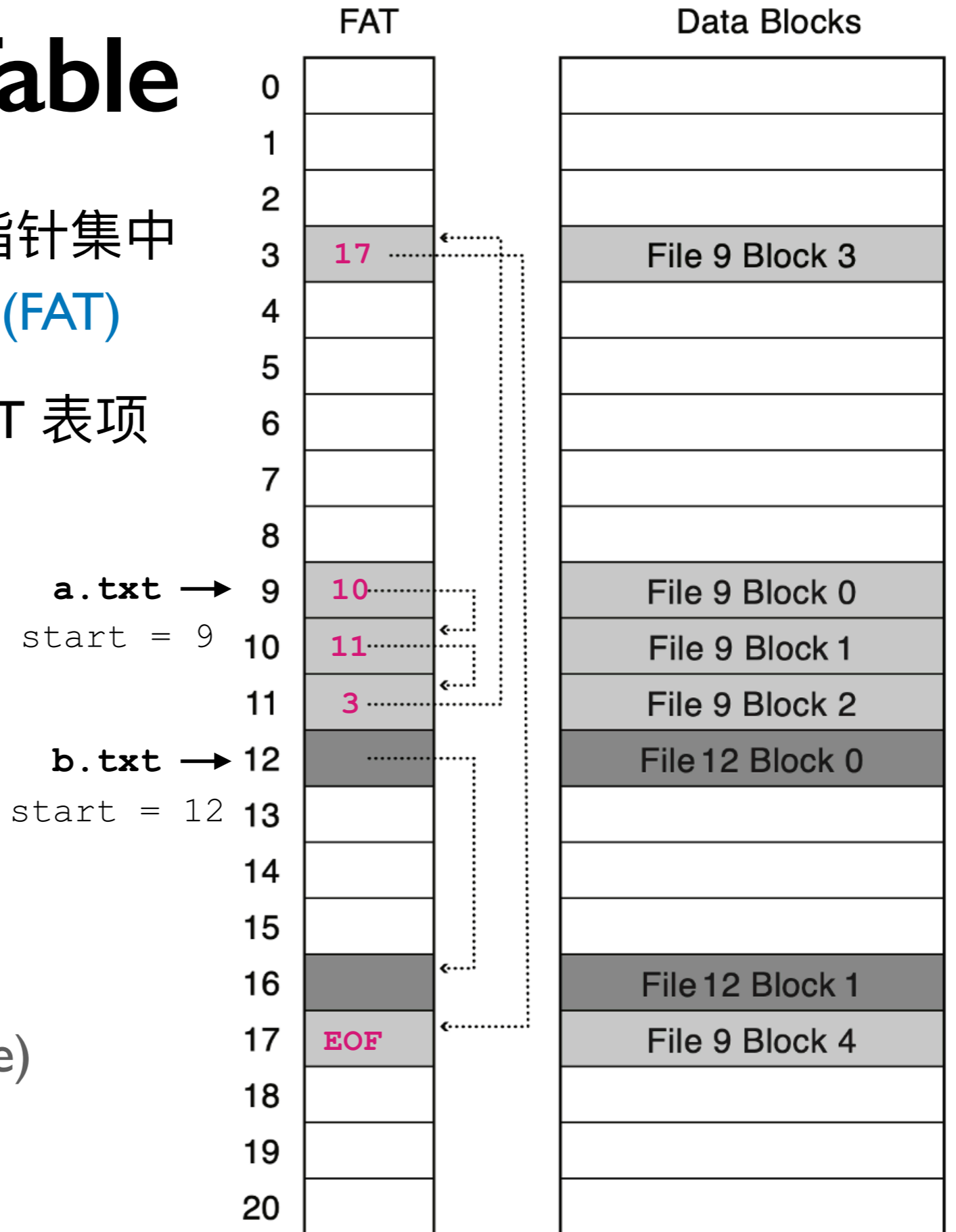
使用 **Data Blocks** 的一个链表来存储 File Data

- ✓ 仅需记录 number of first block
- ✓ 可以把所有磁盘 Data Blocks 都利用起来 (没有外部碎片)
- ✗ 比较适合小文件，但是
 - 定位一个大文件的随机 Data Block 需要遍历链表
 - 每个 Data Block 中需要额外存储一个 next 指针 (Data Block 的可用大小不再是 2^k)
 - 指针损坏导致文件后续数据丢失

File Allocation Table

把 Data Blocks 中的所有 next 指针集中存放到一个 **File Allocation Table (FAT)**

- 每个 Data Block 对应一个 FAT 表项
- 每个 FAT 表项存储
 - 文件下一个 Data Block 编号 (pointer to the next FAT entry of the file)
 - 文件终止标记 (end of file)
`0xFFF`
 - 或空闲空间标记 (free space)
`0x000`



File Allocation Table

- 有效改进了简单的 Linked List 实现
 - 遍历链表 ➡ 只需访问 FAT (缓存在内存) 而不需频繁读磁盘
 - 可用空间 ➡ Data Block 中不需要预留额外空间存储 next 指针
 - 指针损坏 ➡ 在磁盘中存储多份 FAT
- 但是，随着磁盘存储空间的变大，在内存中缓存整个 FAT 会导致较高的存储开销
 - 对于 1TB (2^{40}) 大小的磁盘和 4KB (2^{12}) 大小的 Disk Block
 - 总共需要 2^{28} 个 FAT 表项
 - 如果每个表项占用 4 bytes (FAT 32)，则整个 FAT 大小为 1GB

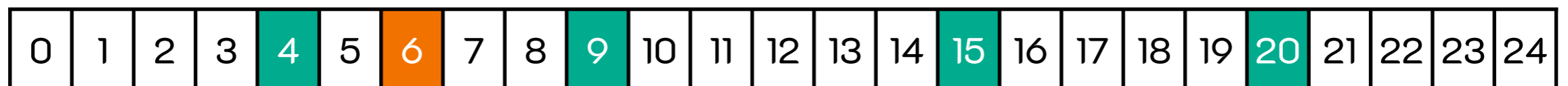
文件的实现方式

Indexed Allocation

使用一个记录 Data Blocks 编号的数组来存储 File Data 的索引信息

- ✓ 仅需记录存储该数组结构的 Data Block 的编号
- ✓ 快速定位任意 Data Block (在内存中缓存索引结构)
- ✓ 需要缓存的大小只和当前已打开文件数相关 (而不是整个磁盘大小)

Data Blocks



a.txt

index = 6

4
9
20
15
END

第 i 项对应文件第 i 个
Data Block 的编号

文件的实现方式

Indexed Allocation

使用一个记录 **Data Blocks** 编号的数组来存储 File Data 的索引信息

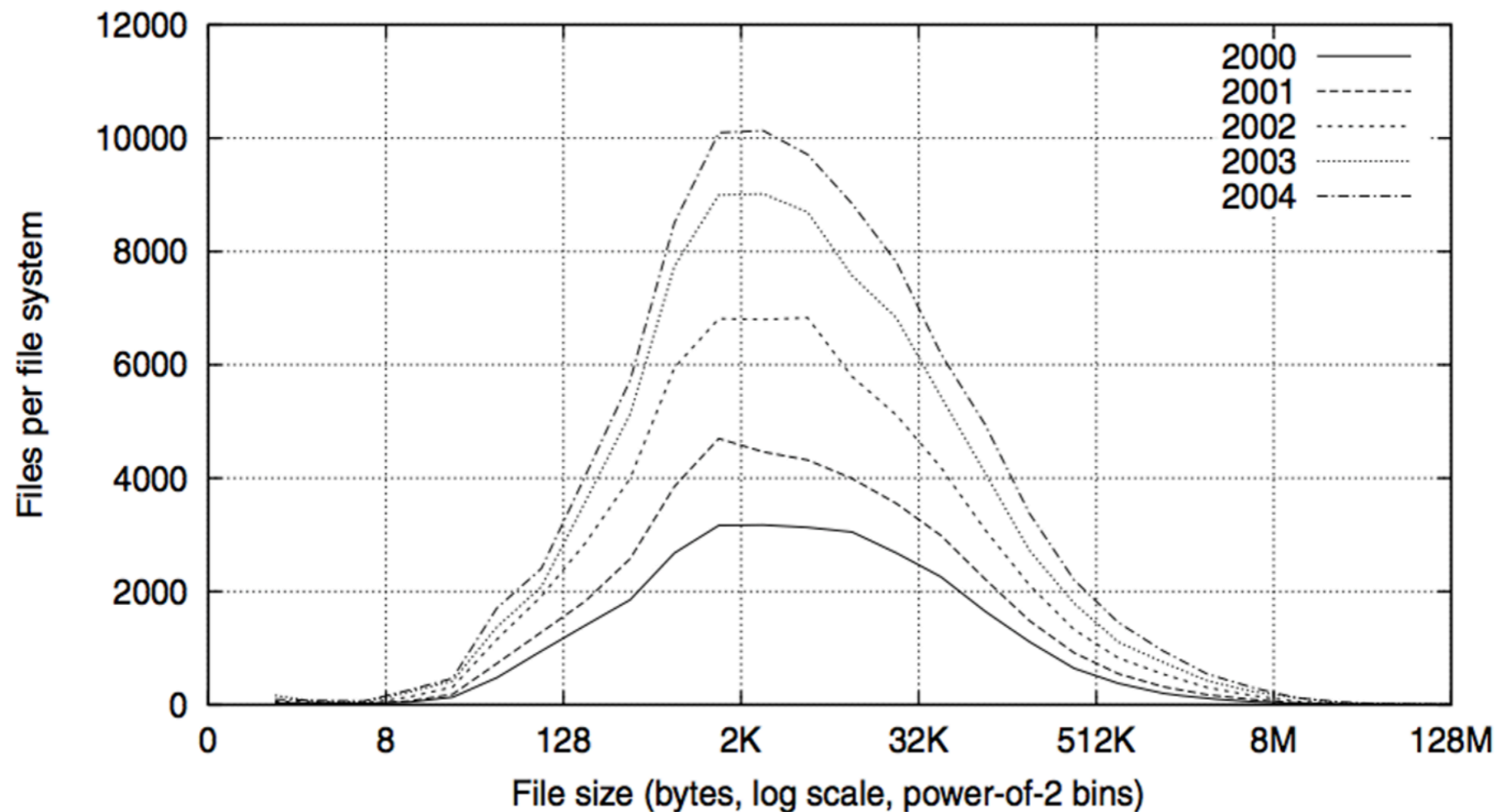
? 但是，如何决定所需数组结构的长度

- 文件系统中的文件大小不定
 - 既有小文件 (few bytes ~ KB)
 - 又有大文件 (several GB / TB)
- 通常希望 File Metadata 是固定长度 (便于管理)

文件的实现方式

Characteristics of Real-World File Systems

- 大部分文件都是小文件 (~2K is the most common size)



Agrawal N, Bolosky W J, Douceur J R, et al. A five-year study of file-system metadata. ACM Transactions on Storage (TOS), 2007, 3(3): 9-es.

文件的实现方式

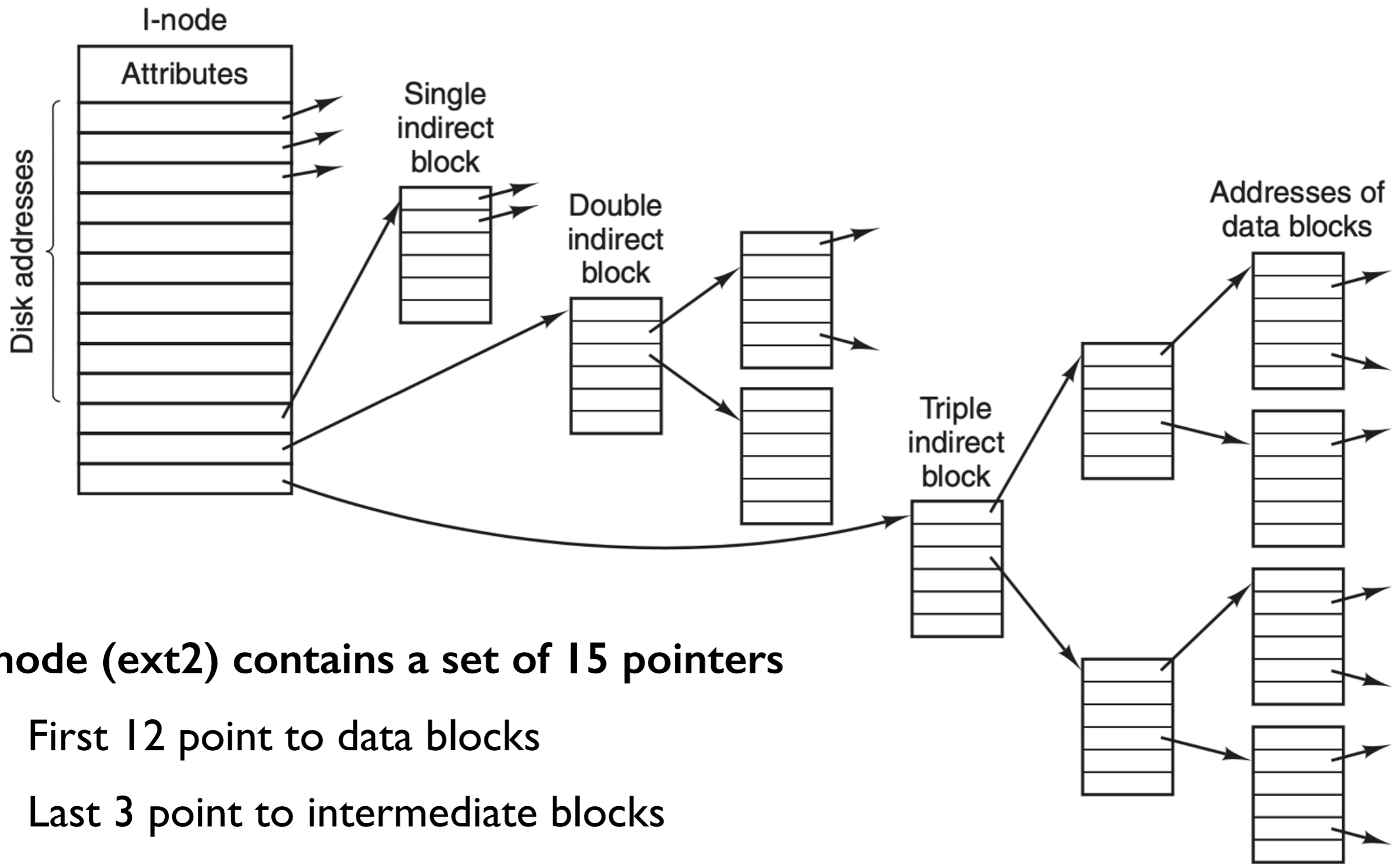
Characteristics of Real-World File Systems

- 大部分文件都是小文件 (~2K is the most common size)
- 少量大文件占用了大部分存储空间
- 文件数量通常很多 (almost 100K on average)
- 目录中通常仅包含少量文件 (most have 20 or fewer)
- 平均文件大小在不断增长
- 文件系统通常会占据硬盘一半的空间

inode

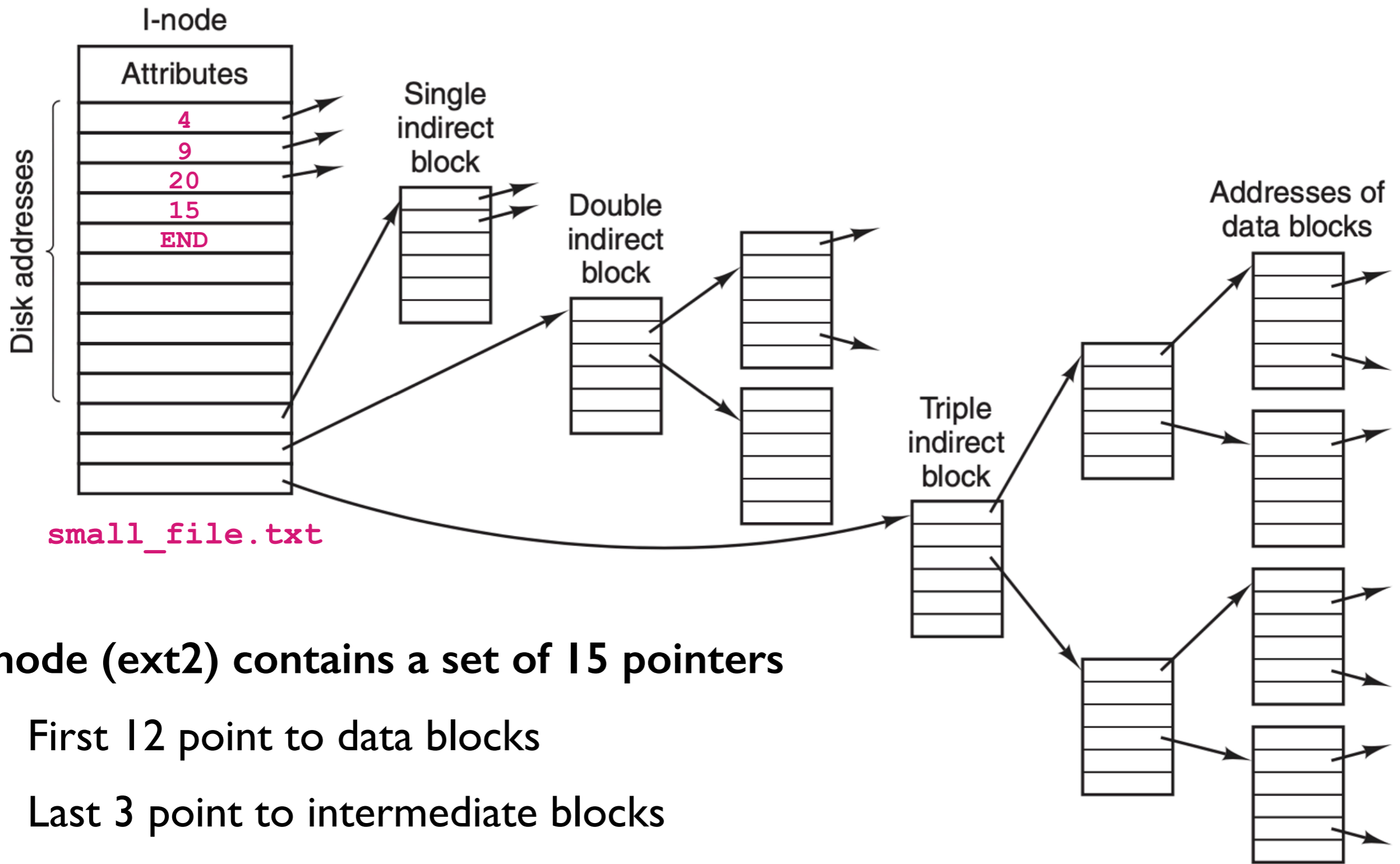
对小文件 (fast path) 和大文件 (slow path) 都提供良好支持

- 在 inode 结构中维护一组固定数目的 pointers
 - **Direct pointers**: 直接指向包含文件数据的 Data Block
 - **Indirect pointers**: 指向一个包含 direct pointers 的 Data Block
- 支持不同文件的快速随机访问
 - 小文件直接使用 direct pointers → an array structure
 - 大文件额外使用 indirect pointers → a tree structure



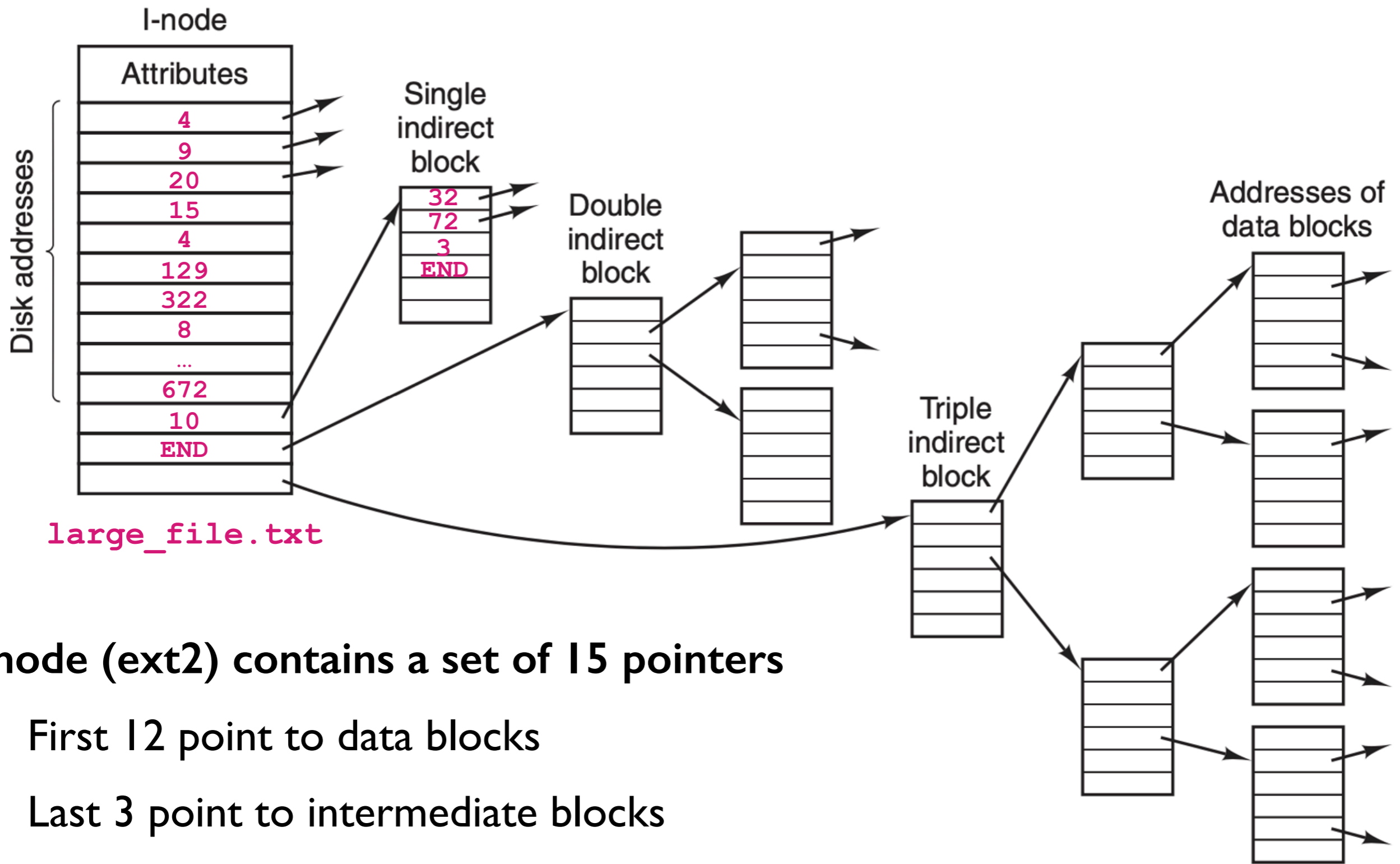
inode (ext2) contains a set of 15 pointers

- First 12 point to data blocks
- Last 3 point to intermediate blocks
 - #13: single indirect pointer
 - #14: double indirect pointer
 - #15: triple indirect pointer



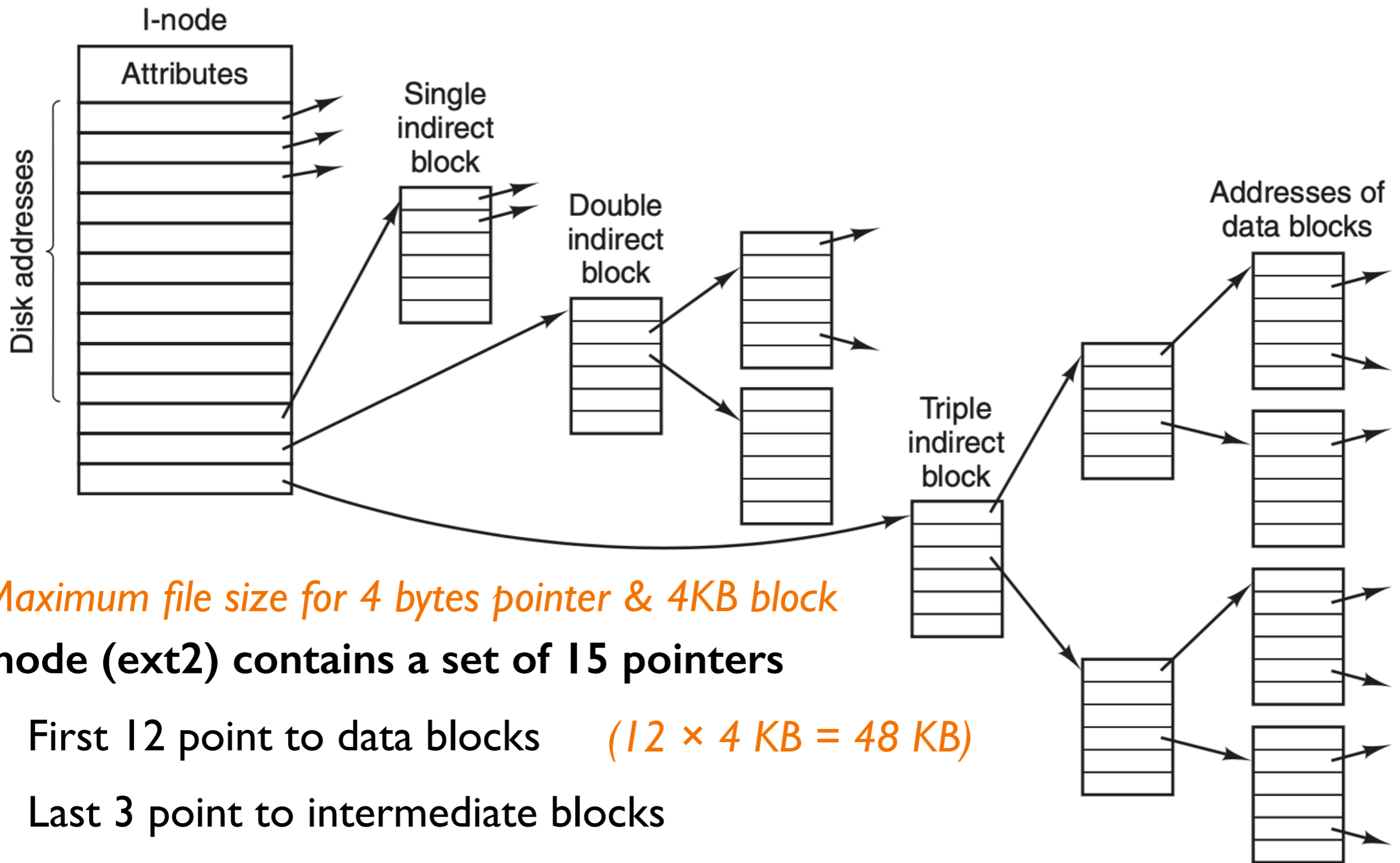
inode (ext2) contains a set of 15 pointers

- First 12 point to data blocks
- Last 3 point to intermediate blocks
 - #13: single indirect pointer
 - #14: double indirect pointer
 - #15: triple indirect pointer



inode (ext2) contains a set of 15 pointers

- First 12 point to data blocks
- Last 3 point to intermediate blocks
 - #13: single indirect pointer
 - #14: double indirect pointer
 - #15: triple indirect pointer



Maximum file size for 4 bytes pointer & 4KB block

inode (ext2) contains a set of 15 pointers

- First 12 point to data blocks $(12 \times 4 \text{ KB} = 48 \text{ KB})$
- Last 3 point to intermediate blocks
 - #13: single indirect pointer $(2^{10} \times 4 \text{ KB} = 4 \text{ MB})$
 - #14: double indirect pointer $(2^{10} \times 2^{10} \times 4 \text{ KB} = 4 \text{ GB})$
 - #15: triple indirect pointer $(2^{10} \times 2^{10} \times 2^{10} \times 4 \text{ KB} = 4 \text{ TB})$

inode

File Size & Size on Disk

文件大小 (File Size) 和 占用磁盘空间 (Size on Disk)

- 文件大小 (`ls -lh`) ➡ 文件字符序列的长度
- 占用磁盘空间 (`du -h`) ➡ 存储该文件使用了多少 Data Blocks
 - Data Block 是文件存储空间分配的基本单位
 - 对于大文件，还需要分配 Data Block 来存储 Indirect Pointers
 - 此时 Size on Disk 大于 File Size
 - 在利用 extent 后大文件也能做到 Size on Disk 等于 File Size (让每个 pointer 指向若干连续的 Data Blocks)

inode

File Size & Size on Disk

甚至还可以 Size on Disk 小于 File Size

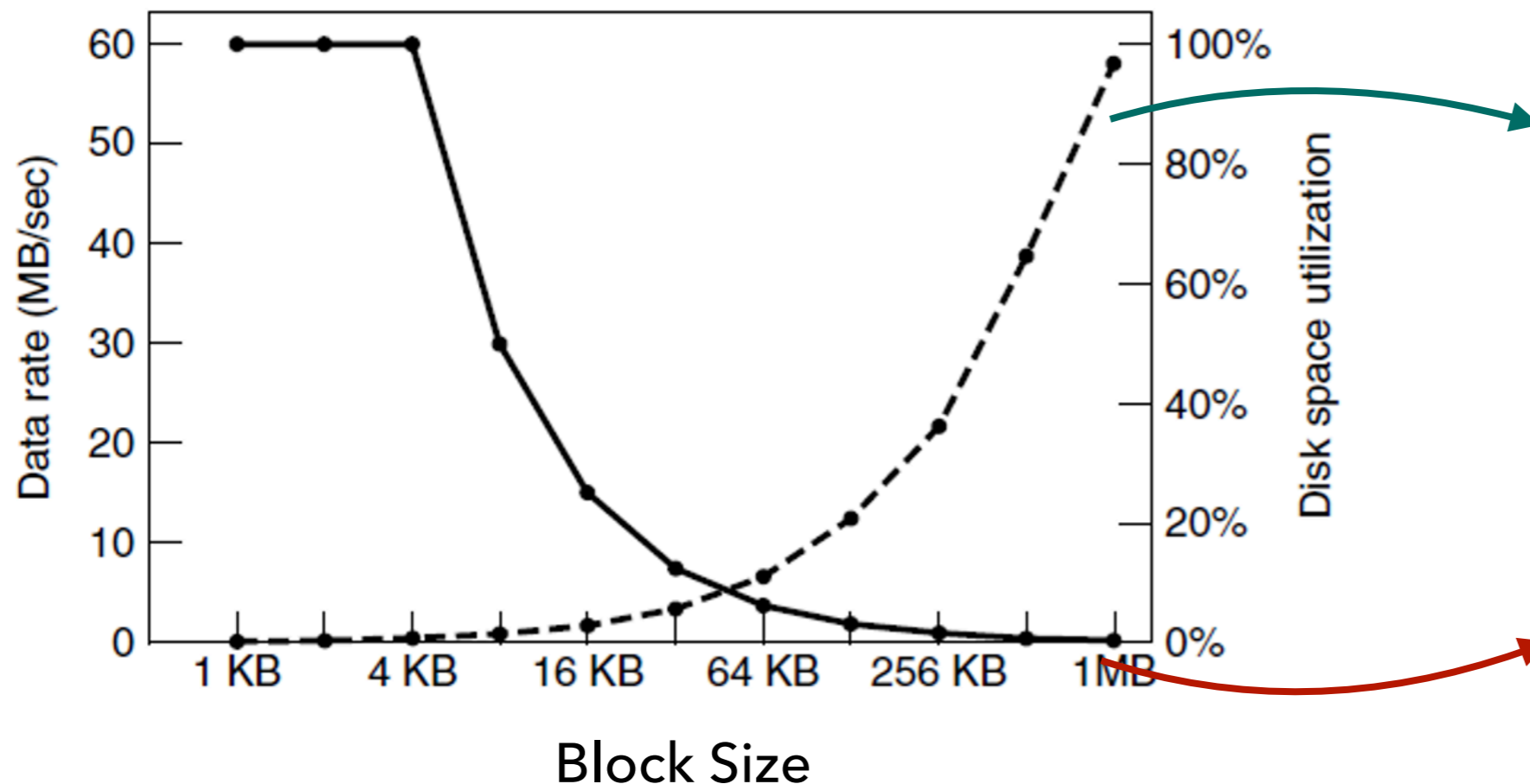
- 不占用磁盘空间存储文件?



Block Size

文件系统可按不同的磁盘块大小 (Block Size) 来进行格式化

- 设置太大会浪费磁盘空间 (存储小文件时的内部碎片)
- 设置太小会导致较高的磁盘读写开销 (大文件需随机读写多个磁盘块)



To store a 4KB file

Data rate goes up almost linearly with block size (the more data fetched, the better)

More spaces are wasted for larger blocks

文件系统的实现

File and Directory

- 在已知 File Metadata 的情况下 (low-level name), 我们可以利用相关结构信息找到 File Data 存储在磁盘上的哪些 Data Block
 - FAT 文件系统 → number of first block + FAT
 - Unix 文件系统 → inode (direct & indirect pointers)
- 如何根据文件名 (human-friendly name) 来找到 File Metadata ?

文件系统的实现

File and Directory

在文件系统中，**目录 (Directory)** 是一种用于存储文件名和文件底层结构间映射关系的特殊文件

- 目录文件也是一个文件
 - 同样具有 File Data 和 File Metadata
- 但目录文件的 File Data 在文件系统中具有特殊含义
 - 记录该目录下的每个文件对应的属性信息
 - 该信息属于文件系统内部维护的数据，通常不允许用户直接使用文件操作进行访问

FAT 文件系统

Directory Implementation

文件的属性 (File Metadata) 信息直接存储在目录文件中

- 目录文件的内容就是一个目录项 (directory entry) 的列表, 其中每个目录项对应一个文件

dir

```
| - a.txt  
| - list/  
| - main.c
```

File Name	Attributes	First Data Block
a.txt	...	9
list	...	102
main.c	...	82

File data of the `dir` file

FAT 文件系统

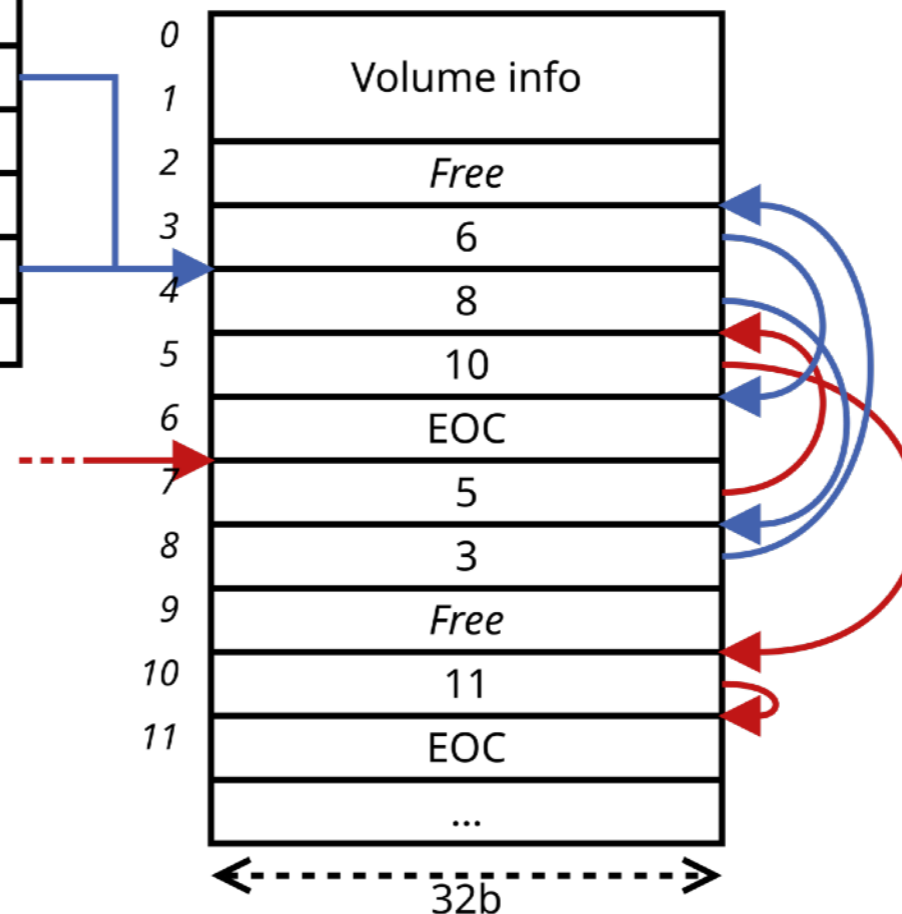
Directory Implementation

- 对于 FAT 32 文件系统，每个目录项大小为 32 bytes

Directory table entry (32B)

Filename (8B)
Extension (3B)
Attributes (1B)
Reserved (1B)
Create time (3B)
Create date (2B)
Last access date (2B)
First cluster # (MSB, 2B)
Last mod. time (2B)
Last mod. date (2B)
First cluster # (LSB, 2B)
File size (4B)

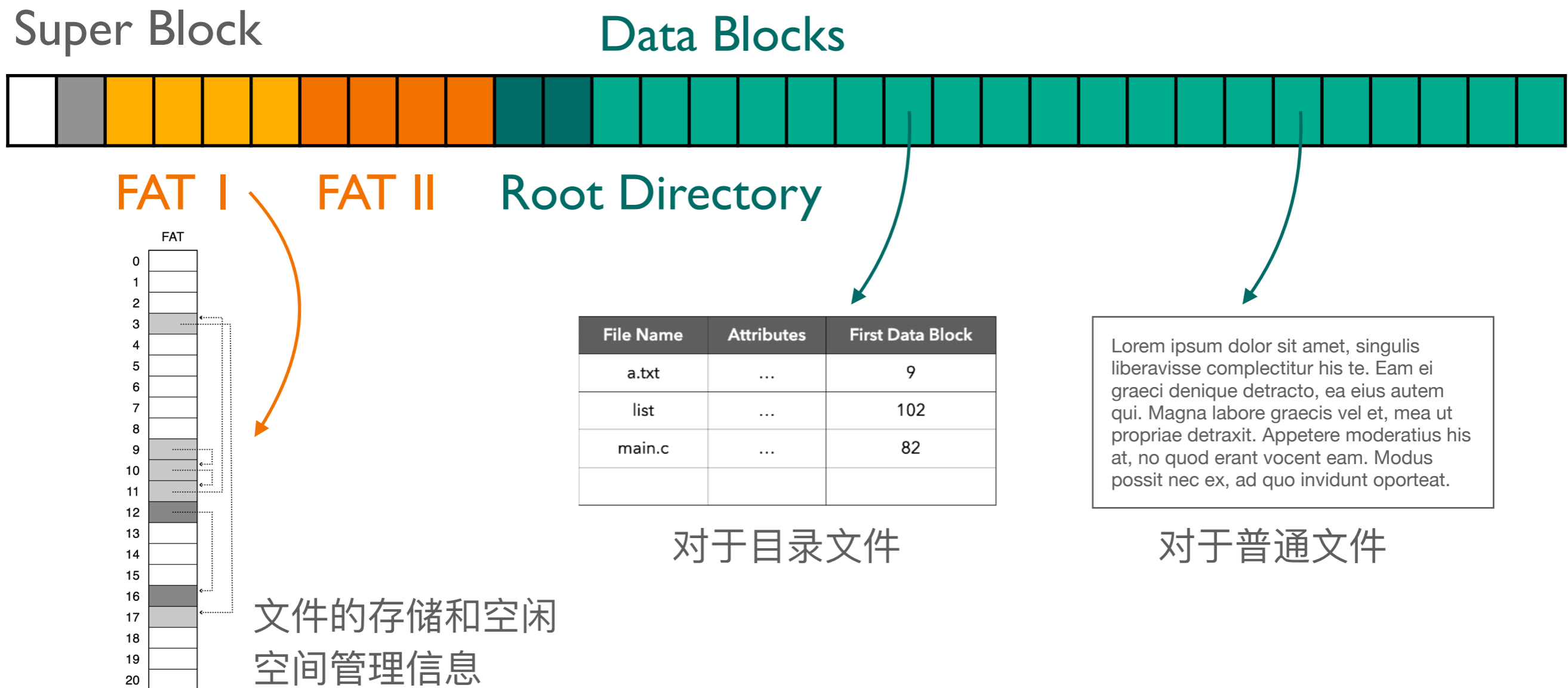
File allocation table



FAT 文件系统

File System Layout

FAT 文件系统的磁盘布局：结合其文件 (by File Allocation Table) 和目录 (Save File Metadata in Directory) 的实现方式



Unix 文件系统

Directory Implementation

文件的属性信息存储在一个专门的结构中 (inode)

- 每个目录项记录 $\langle \text{file name, inode number} \rangle$ 的对应关系
- 为根目录分配一个特殊的 inode

dir

```
| - a.txt  
| - list/  
| - main.c
```

File Name	inode Number
a.txt	6782
list	32
main.c	1982

File data of the `dir` file

inode
6782

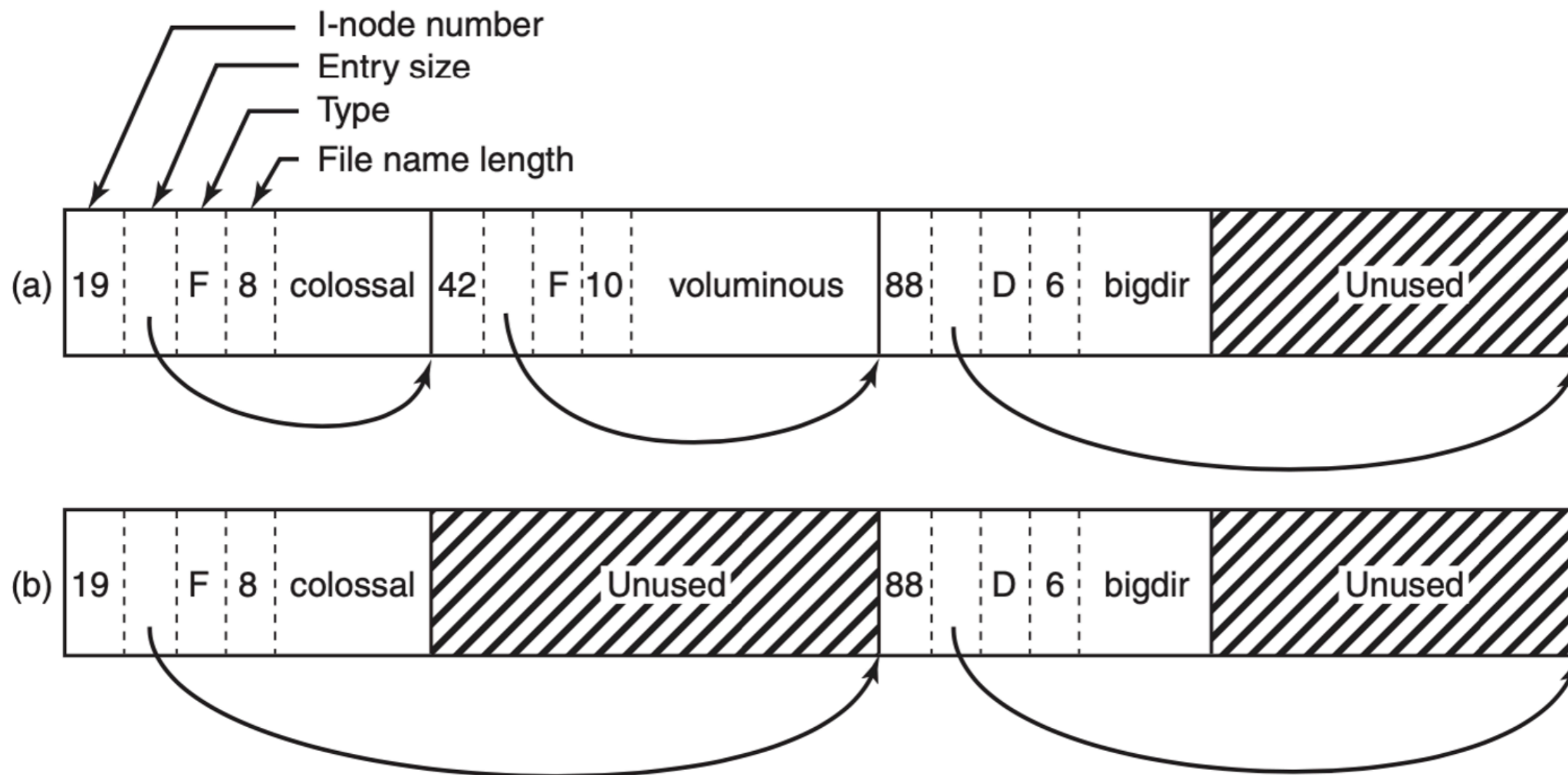
inode
32

inode
1982

Unix 文件系统

Directory Implementation

- 可在目录项中额外记录文件名长度来支持不定长文件名

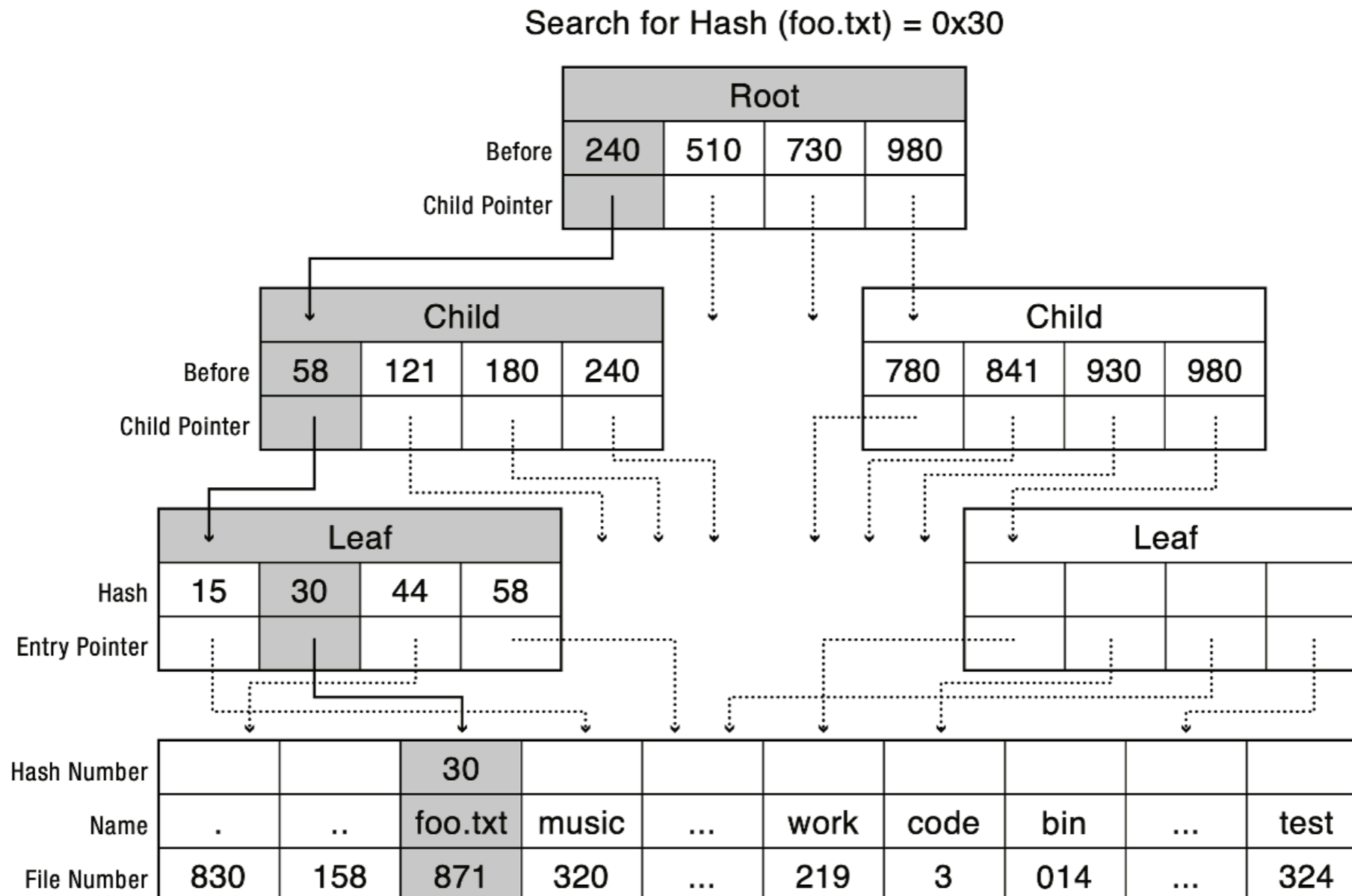


ext2 example

Unix 文件系统

Directory Implementation

- 为了加速大目录文件的目录项查找速度，可进一步建立〈file name, inode number〉的某种索引结构 (例如使用 B+ tree)



Unix 文件系统

File System Layout

Unix 文件系统的磁盘布局：结合其文件 (by Direct/Indirect Pointers in inode structure) 和目录 (Save File Metadata in inode) 的实现方式

Super Block

Data Blocks



inode Table (file metadata)

所有文件的 inode 统一集中存储 (每个 Block 可存储多个 inode),
为每个 inode 分配一个编号 (可根据编号直接定位 inode 存储位置)

Allocation Structure (e.g., one bitmap for inode, one for data)

需要额外的数据结构来进行空闲空间管理
(哪些 inode / Data Block 空闲 / 已分配)

Unix 文件系统

File System Layout

在 Unix 文件系统中查找 `/usr/ast/mbox` 的过程 (Resolve File Name)

Root directory

1	.
1	..
4	bin
7	dev
14	lib
9	etc
6	usr
8	tmp

Looking up
usr yields
i-node 6

I-node 6
is for /usr

Mode size times
132

I-node 6
says that
/usr is in
block 132

Block 132
is /usr
directory

6	.
1	..
19	dick
30	erik
51	jim
26	ast
45	bal

/usr/ast
is i-node
26

I-node 26
is for
/usr/ast

Mode size times
406

I-node 26
says that
/usr/ast is in
block 406

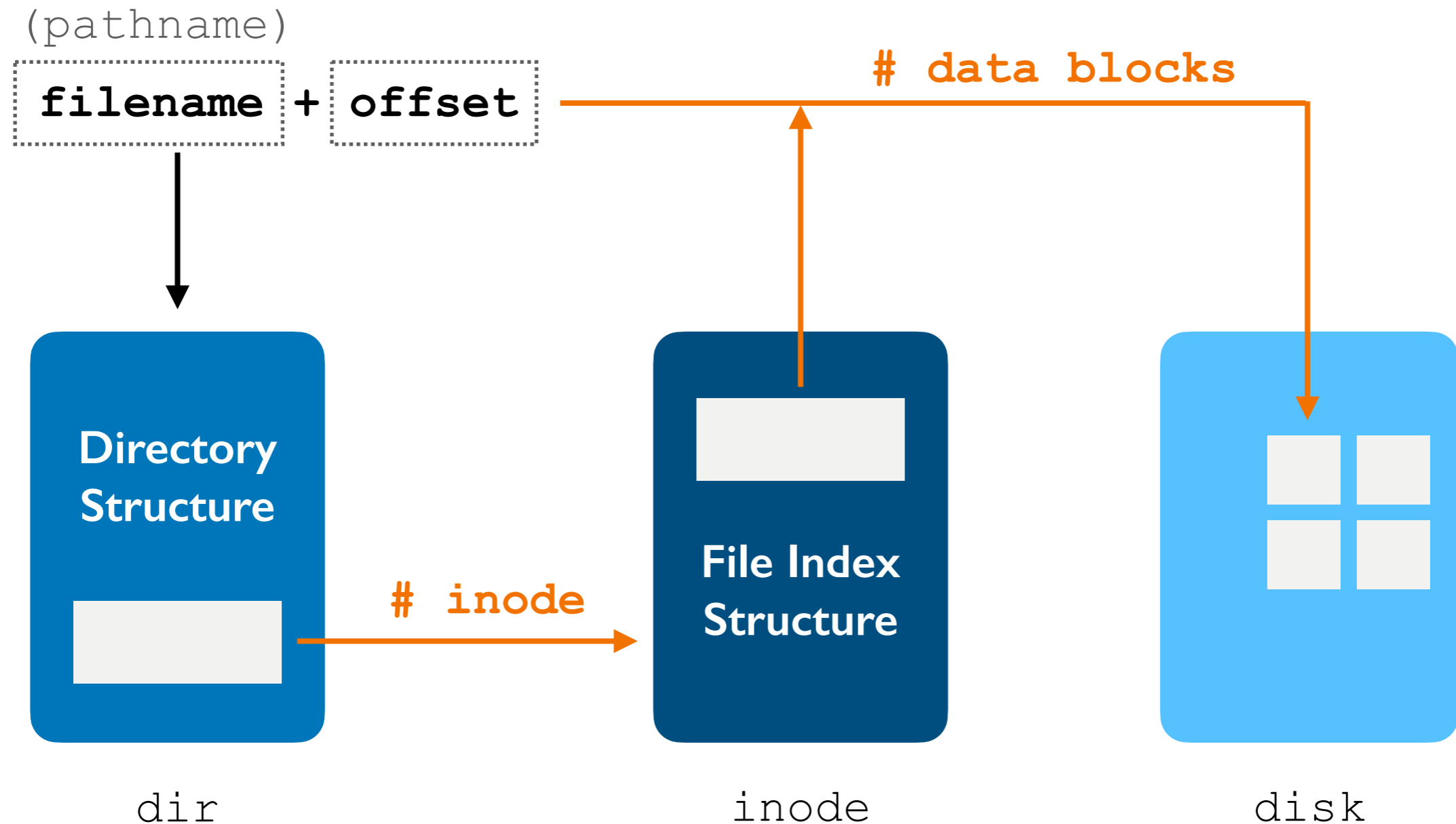
Block 406
is /usr/ast
directory

26	.
6	..
64	grants
92	books
60	mbox
81	minix
17	src

/usr/ast/mbox
is i-node
60

Unix 文件系统

Two Key Abstractions: File & Directory



Quiz

对于如下 Unix 文件系统，其所能支持的最大文件大小是多少？单个目录中最多能包含多少个文件？

- 磁盘 Block Size 为 1 KB，其中划分 4096 个 Blocks 用于存储 inodes
- inode 结构中存储 User ID (2 bytes)、时间戳 (12 bytes)、文件类型和保护信息 (4 bytes)、Link Count (2 bytes)、文件大小 (4 bytes)、以及文件 Data Blocks 索引结构
 - 索引结构中包含 8 个 direct pointer、1 个 indirect pointer、以及 1 个 double indirect pointer (4 bytes for each pointer)
- 目录项 (Directory Entry) 仅存储文件名和 inode 编号，每个目录项占 16 bytes