

Unix系统高级编程

系统调用和文件系统

Unit07

系统调用



Unix应用的层次结构



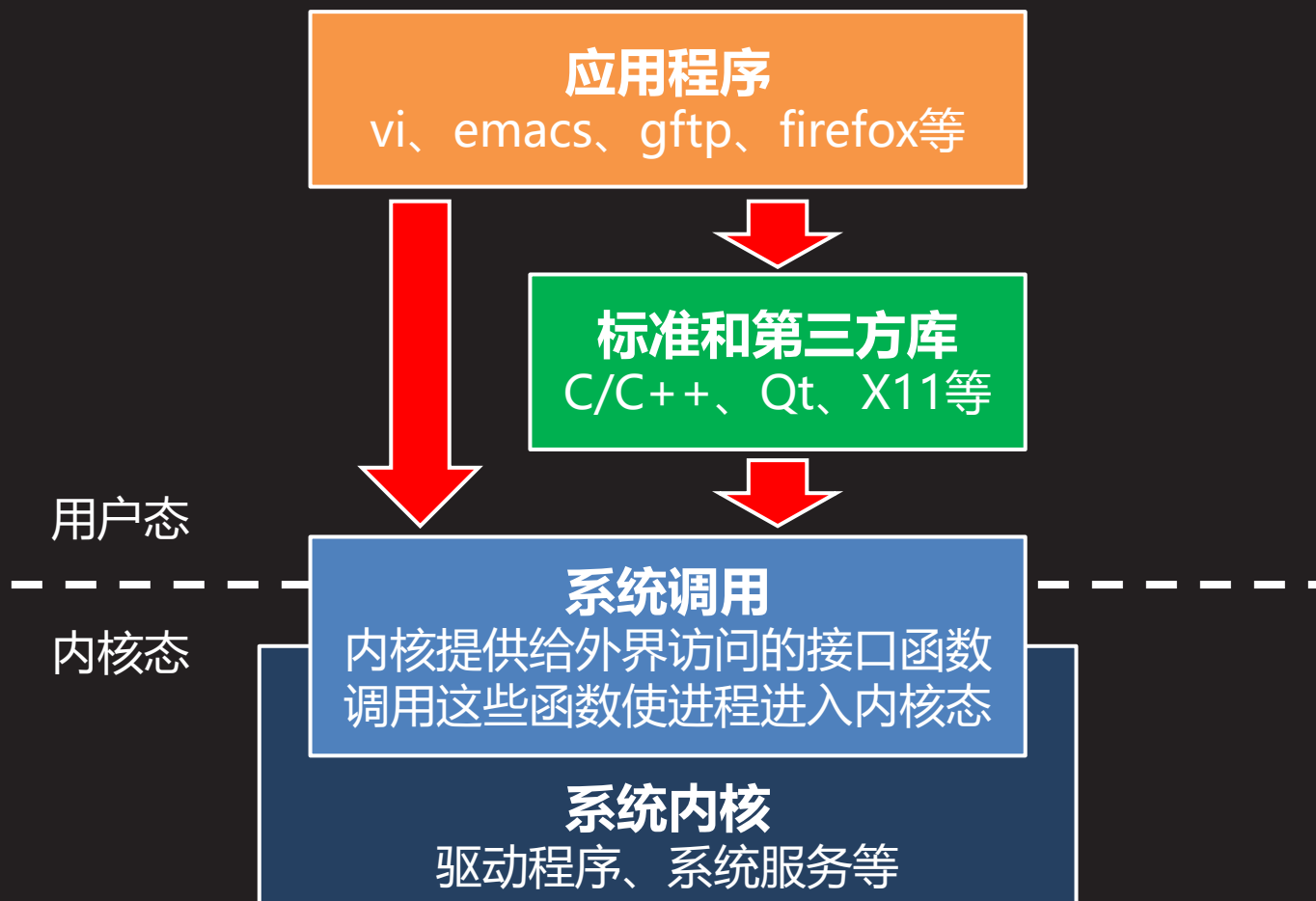
Unix应用的层次结构

- Unix/Linux系统的大部分功能都是通过系统调用实现的，如open、close等
- Unix/Linux的系统调用已被封装成C函数的形式，但它们并不是C语言标准库的一部分
- 标准库函数大部分时间运行在用户态，但部分函数偶尔也会调用系统调用，进入内核态，如malloc、free等
- 程序员自己编写的代码也可以跳过标准库，直接使用系统调用，如brk、sbrk、mmap和munmap等，与操作系统内核交互，进入内核态
- 系统调用在内核中实现，其外部接口定义在C库中，该接口的实现借助软中断进入内核



Unix应用的层次结构（续1）

- 从应用程序到操作系统内核需要经历如下调用链

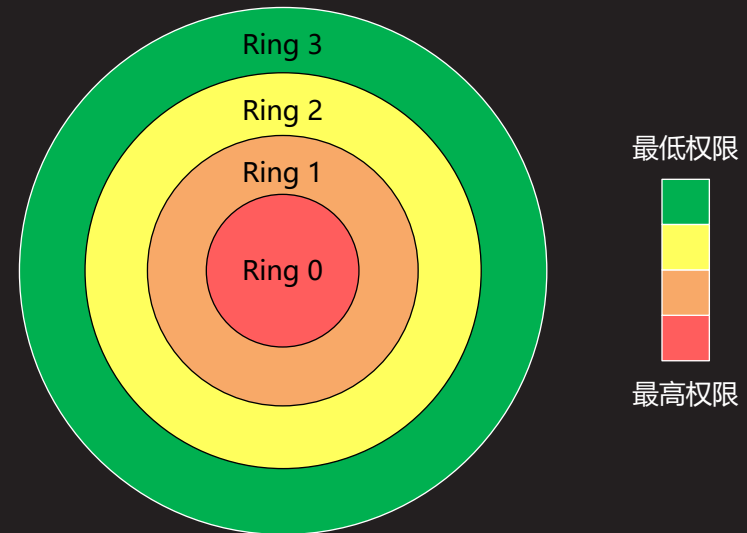


Unix系统调用的流程



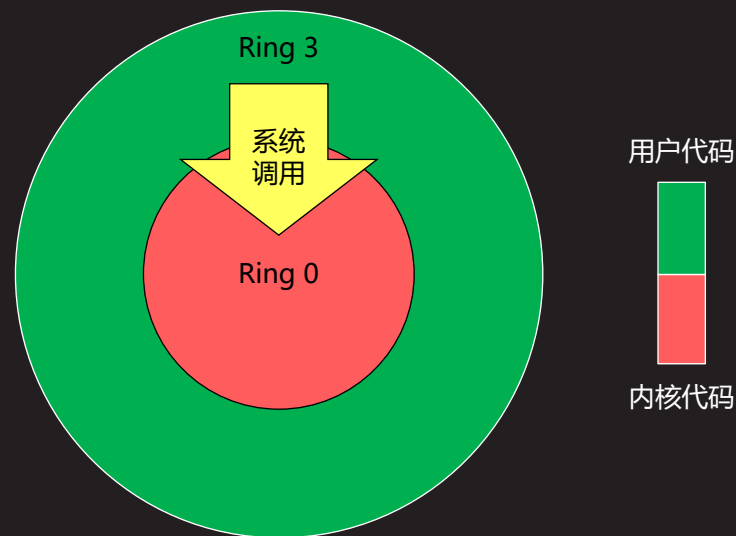
Unix系统调用的流程

- CPU安全保护
 - 在Intel的CPU上运行的代码分为四个安全级别
 - Ring 0 - 设备驱动、内存管理、线程调度等系统核心代码
 - Ring 1 - 设备驱动
 - Ring 2 - 设备驱动
 - Ring 3 - 用户代码
- } Linux未使用
- 工作在Ring 3级的用户代码无法直接访问工作在Ring 0级的核心代码，除非借助于系统调用



Unix系统调用的流程（续1）

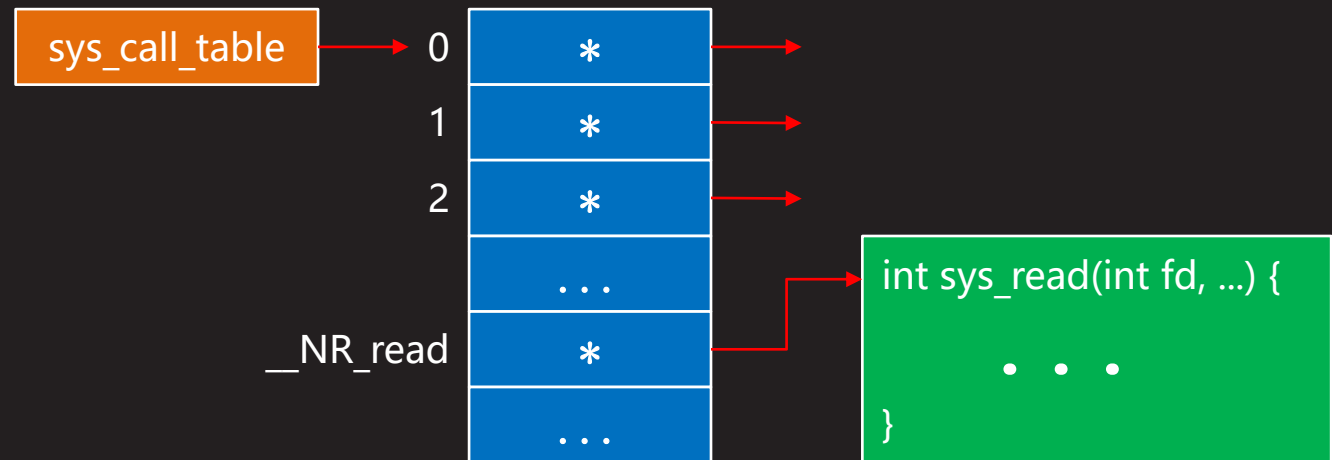
- Linux系统提供了一套实现各种功能的子程序，谓之系统调用
 - 应用程序设计者可以向调用普通C语言函数一样调用系统调用函数
 - 系统调用函数可使其调用者进程从用户态进入内核态
 - 系统调用的具体功能由系统内核实现，以向用户程序提供各种服务
 - 系统调用是用户代码访问内核代码的唯一途径



Unix系统调用的流程 (续2)

- 系统调用表
 - Linux操作系统内核维护着一张全局表sys_call_table
 - 每个条目记录着每个系统调用在内核中实现的入口地址
 - 特定系统调用实现入口的索引即该系统调用的标识号
 - 可将系统调用表视为一个以系统调用标识为下标的函数指针数组，其中每个函数指针指向一个系统调用的实现代码

知识讲解



Unix系统调用的流程 (续3)

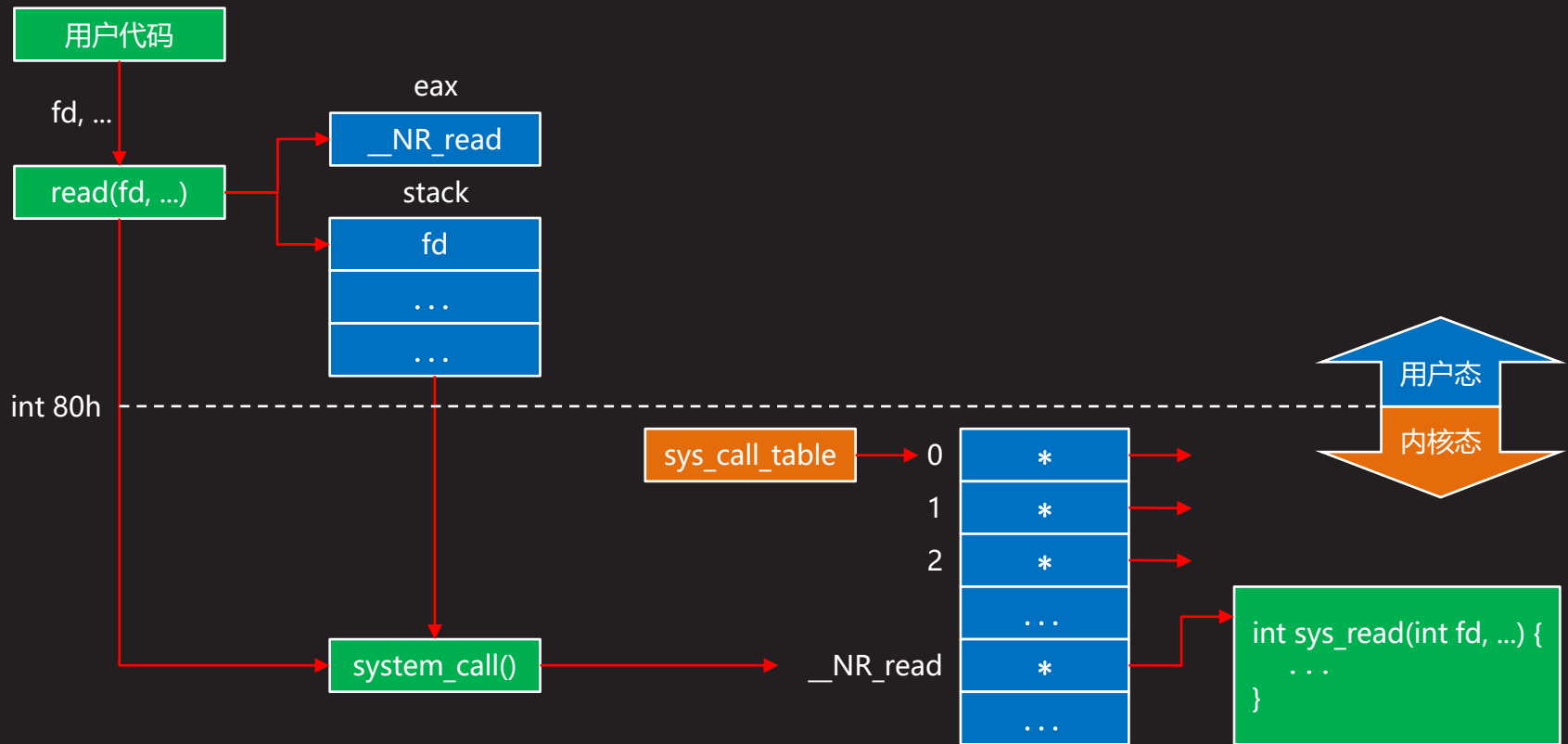
- 系统调用过程
 - 用户代码调用某个系统调用函数，如read(fd, ...):
 - 将系统调用标识__NR_read存入eax寄存器
 - 将文件描述符fd等参数压入栈中
 - 通过int 80h指令触发80h中断
 - 系统内核对80h中断的中断处理函数system_call()被执行
 - 将调用进程从用户态(Ring 3)切换至内核态(Ring 0)
 - 从栈中弹出文件描述符fd等参数
 - 从eax寄存器中获取系统调用标识__NR_read
 - 系统调用表第__NR_read个元素即sys_read()函数的入口地址
 - 通过sys_read()函数的入口地址调用该函数，同时传入参数，并向用户代码返回结果



Unix系统调用的流程 (续4)

- 系统调用过程

知识讲解



测试运行时间



测试运行时间

- 测试应用程序的运行时间分配

```
$ time a.out
```

```
real    0m18.705s  
user    0m0.452s  
sys     0m18.173s
```

- **real** : 总耗时 = 用户空间耗时 + 内核空间耗时 + 等待耗时
- **user** : 用户空间耗时
- **sys** : 内核空间耗时



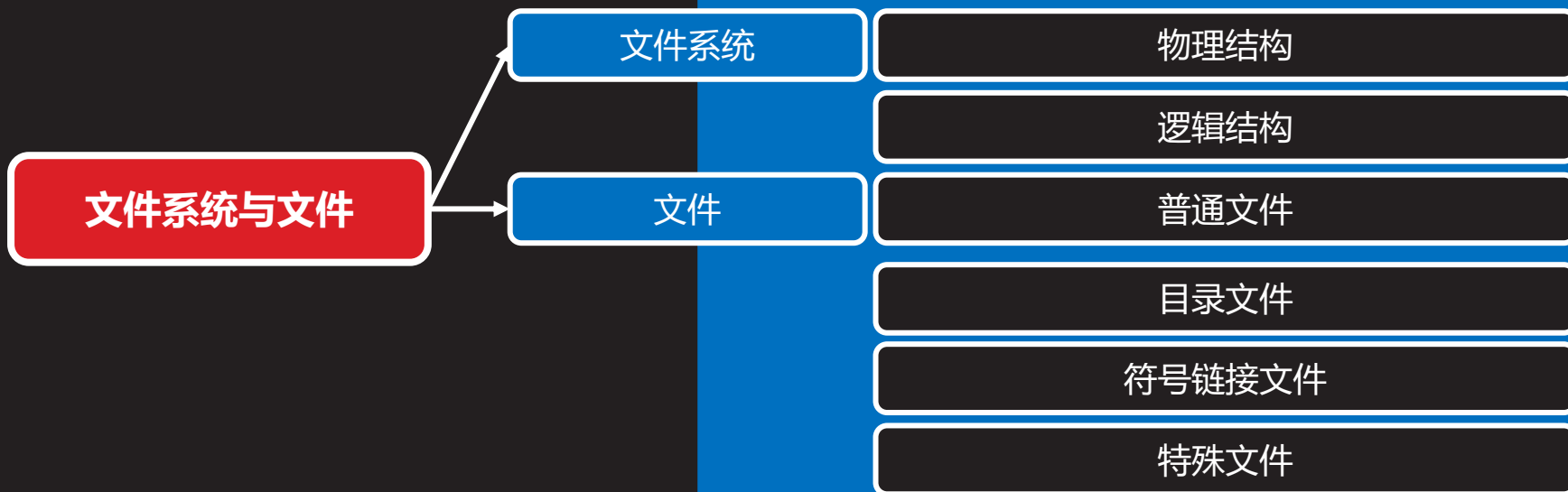
测试运行时间

【参见：elapse.c】

- 测试运行时间



文件系统与文件



文件系统



物理结构

- 硬盘的物理结构

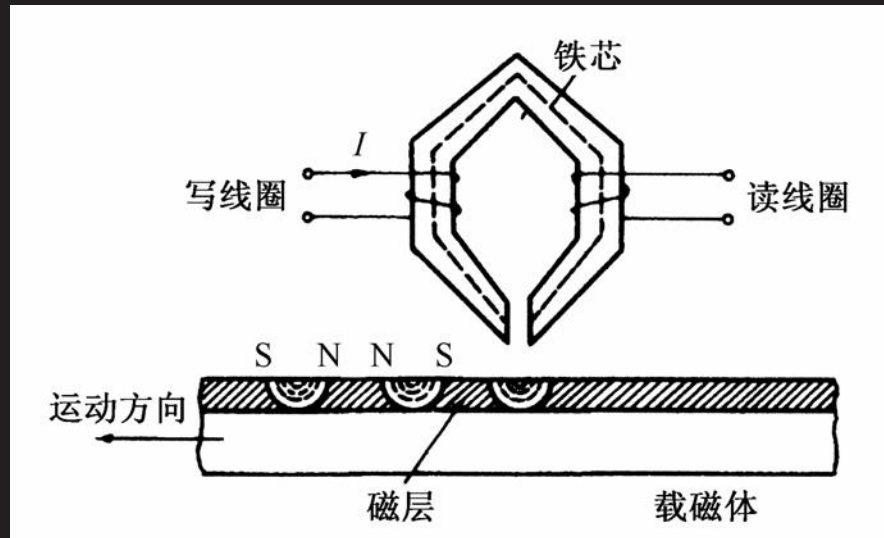
- 驱动臂
- 盘片
- 主轴
- 磁头



物理结构 (续1)

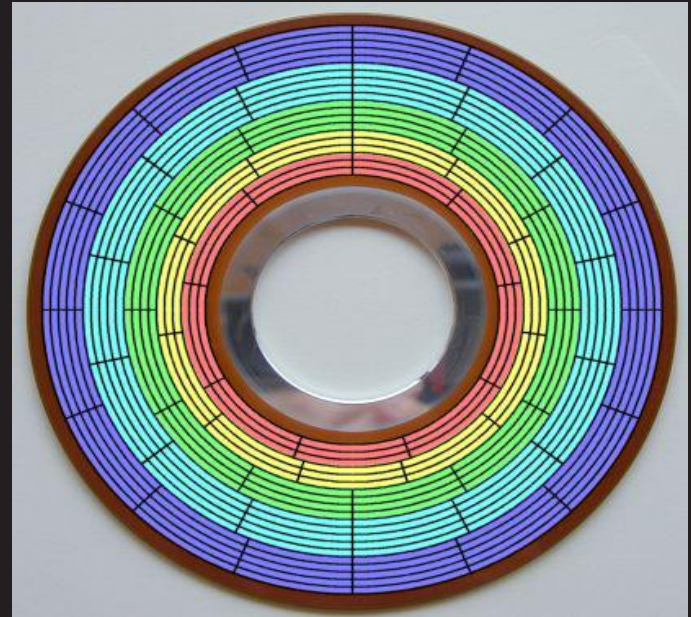
- 磁表面存储器的读写原理
 - 硬盘片的表面覆盖着薄薄的磁性涂层，涂层中含有无数微小的磁性颗粒，谓之磁畴
 - 相邻的若干磁畴组成一个磁性存储元，以其剩磁的极性表示二进制数字0和1
 - 为磁头的写线圈施加脉冲电流，可把一位二进制数字转换为磁性存储元的剩磁极性
 - 利用磁电变换，通过磁头的读线圈，可将磁性存储元的剩磁极性转换为相应的电信号，表示二进制数

知识讲解



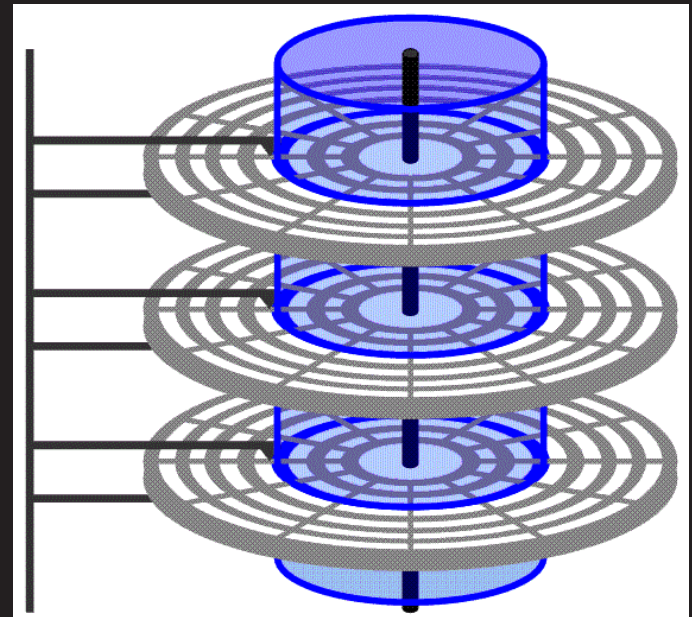
物理结构 (续2)

- 磁道和扇区
 - 磁盘旋转，磁头固定，每个磁头都会在盘片表面划出一个圆形轨迹。改变磁头的位置，可以形成若干大小不等的同心圆，这些同心圆就叫做磁道(Track)
 - 每张盘片的每个面上都有成千上万个磁道
 - 一个磁道，按照512字节等分，其中每个等分叫做扇区(Sector)
 - 扇区是最基本的文件存储单位
 - 每个磁道所包含的扇区数并不相等，越靠近外圈的磁道所包含的扇区越多



物理结构 (续3)

- 柱面、柱面组、分区和磁盘驱动器
 - 硬盘中，不同盘片相同半径的磁道所组成的圆柱称为柱面 (Cylinder)。整个硬盘的柱面数与每张盘片的磁道数相等
 - 硬盘中的每个扇区可由以下三个坐标唯一确定
 - 磁头号：确定哪张盘面(一张盘片有两个盘面各对应一个磁头)
 - 柱面号：确定哪条磁道
 - 扇区号：确定哪个区域
 - 若干连续的柱面构成一个柱面组
 - 若干连续的柱面组构成一个分区
 - 每个分区都建有独立的文件系统
 - 若干个分区构成一个磁盘驱动器



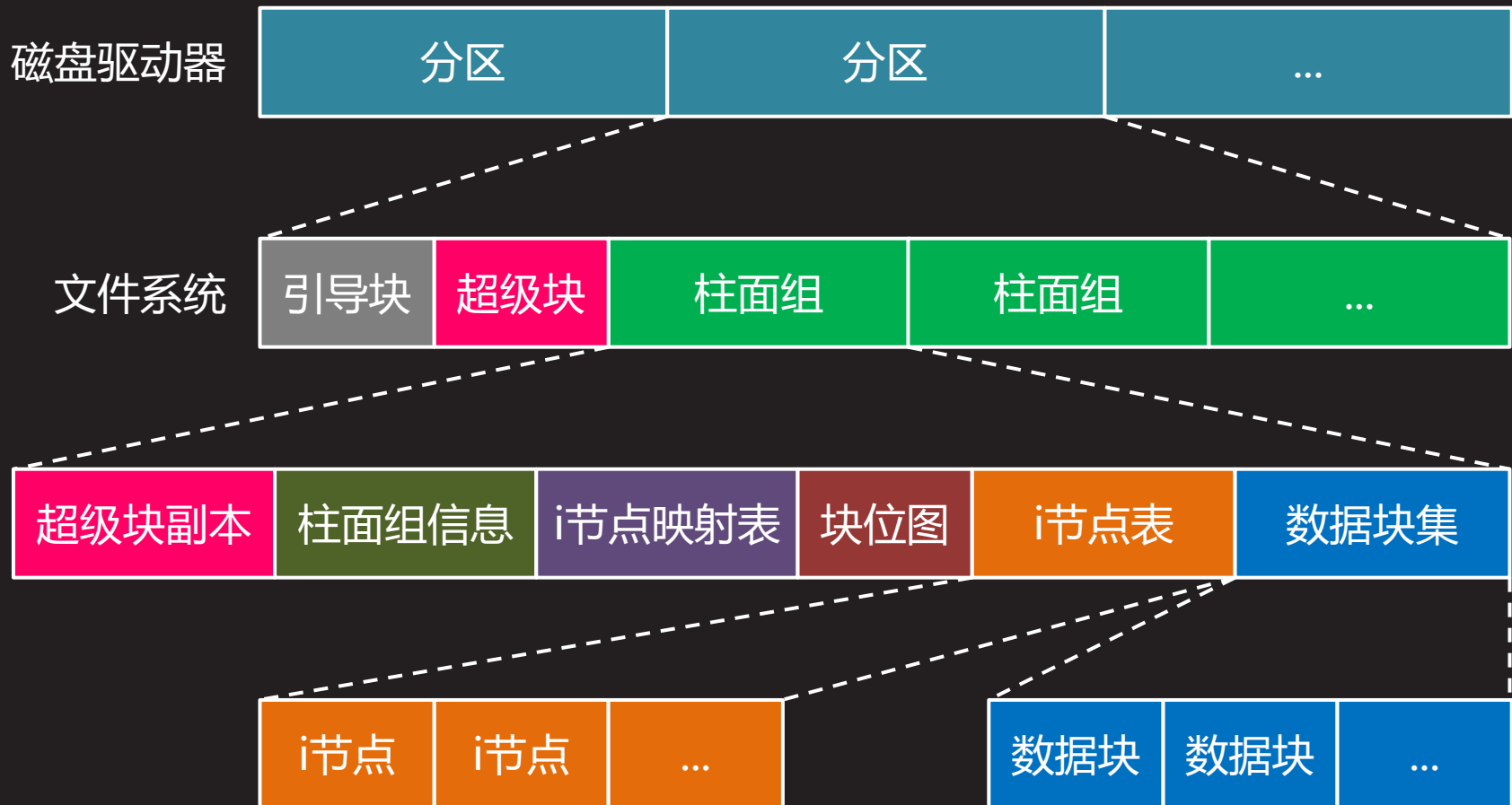
逻辑结构

- 一个磁盘驱动器被划分成一到多个分区，其中每个分区上都建有独立的文件系统，每个文件系统包括
 - 引导块：计算机加电启动时，ROM BIOS从这里读取可执行代码和数据，以完成操作系统自举
 - 超级块：记录文件系统的整体信息，如文件系统的格式和大小，i节点和数据块的总量、使用量和剩余量等等
 - 若干柱面组，其中每个柱面组包括
 - 超级块副本：同上
 - 柱面组信息：柱面组的整体描述
 - i节点映射表：i节点号与i节点磁盘位置的对应表
 - 块位图：位图中的每个二进制位对应一个数据块，用1和0表示该块处于占用或是空闲状态
 - i节点表：包含若干i节点，记录文件的元数据和数据块索引表
 - 数据块集：包含若干数据块，存储文件的内容数据



逻辑结构 (续1)

- UFS (Unix File System)文件系统结构



知识讲解



逻辑结构 (续2)

- i节点, 即索引节点(index node, inode)
 - 磁盘中的每个文件或目录都有唯一的一个的i节点与之对应
 - 每个i节点都有唯一的编号即i节点号, 通过i节点映射表可以查到与每个i节点号相对应的i节点在磁盘上的存储位置
 - 文件名和i节点号的对应关系记录在该文件所在目录的目录文件中, 目录文件中的一条这样的记录就是一个硬链接

```
$ ls -li /etc
```

```
262169 console-setup          262342 os-release
262170 cron.d                     262334 pam.conf
262171 cron.daily                 262225 pam.d
262172 cron.hourly              264643 papersize
262173 cron.monthly            283386 passwd
```



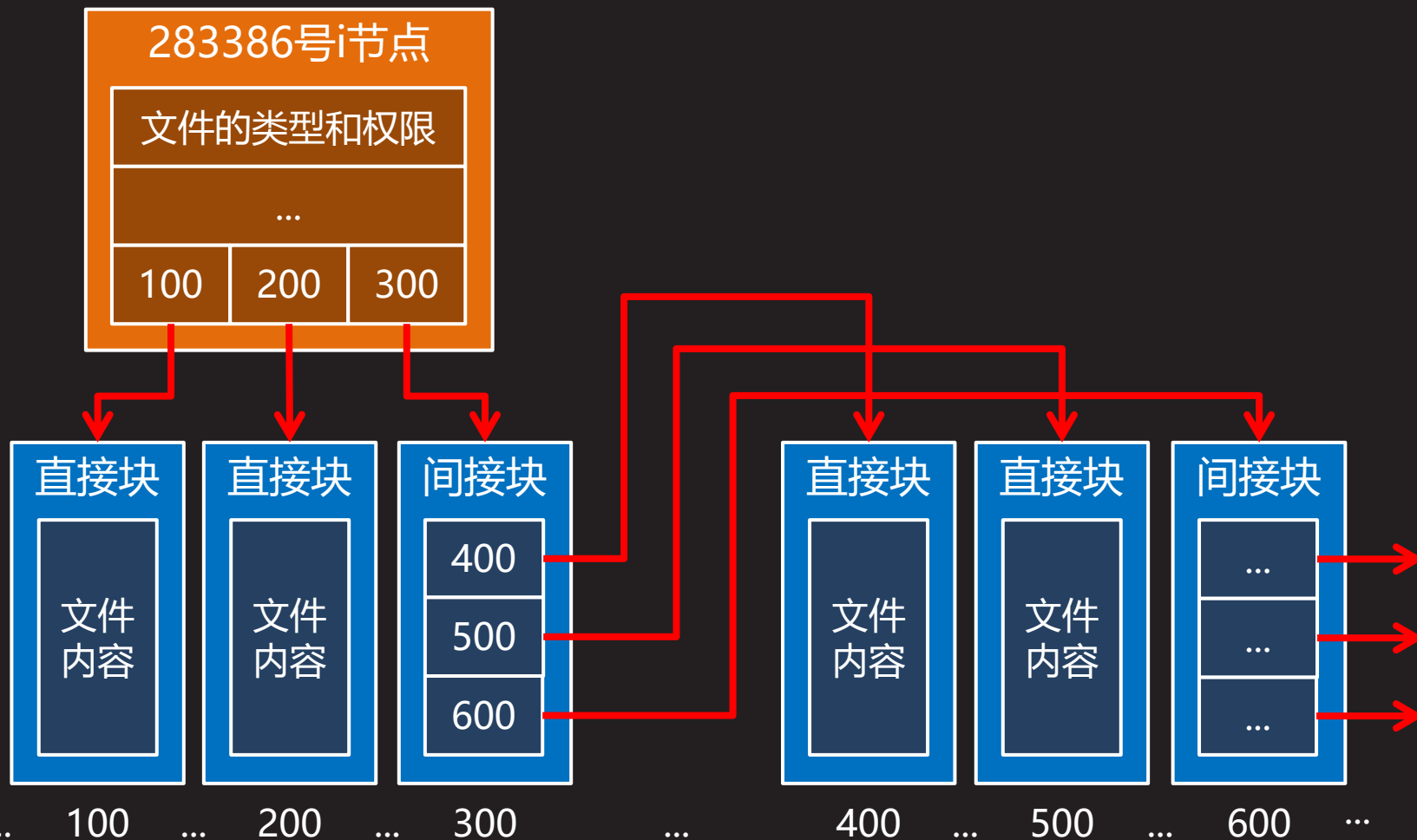
逻辑结构 (续3)

- i节点, 即索引节点(index node, inode)
 - i节点的具体内容包括
 - 文件类型和权限
 - 文件的硬链接数
 - 文件的用户和组
 - 文件的字节大小
 - 文件的最后访问时间、最后修改时间和最后状态改变时间
 - 文件数据块索引表
- 数据块(data block), 512/1024/4096字节
 - 直接块: 存储文件的实际内容数据
 - 文件块: 存储普通文件的内容数据
 - 目录块: 存储目录文件的内容数据
 - 间接块: 存储下级文件数据块索引表



逻辑结构 (续4)

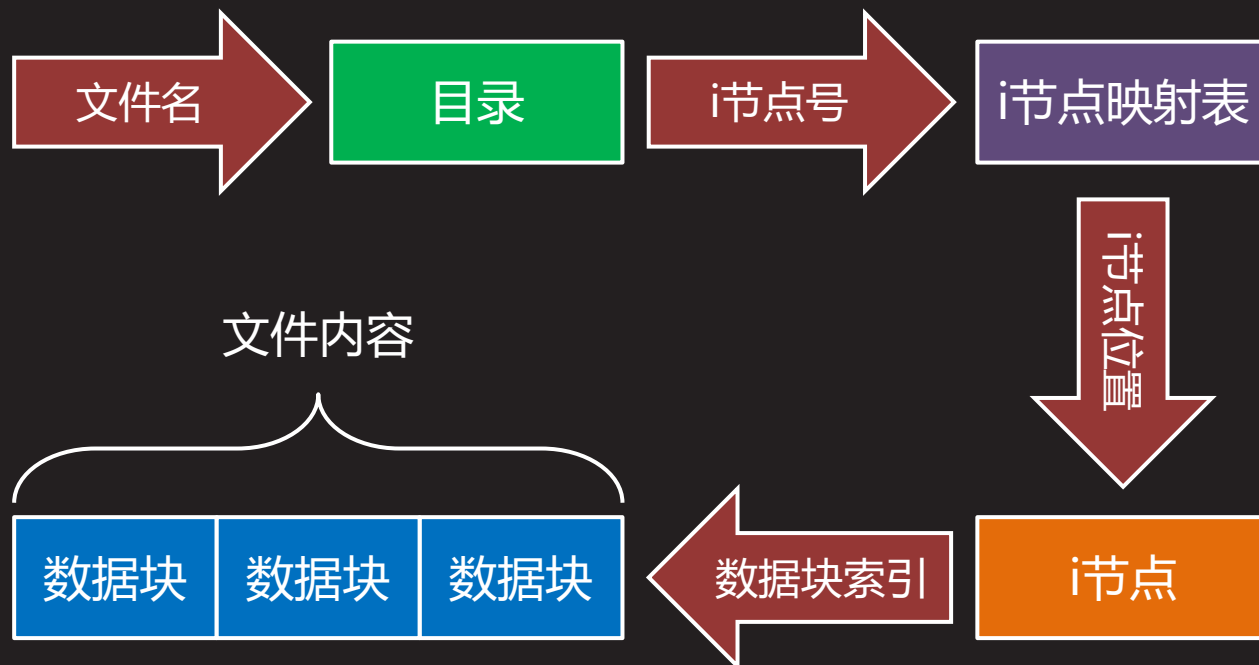
- 通过i节点索引数据块



逻辑结构 (续5)

- 文件访问流程
 - 针对给定的文件名，从其所在目录中可以得到与之对应的i节点号，再通过i节点映射表可以查到该i节点在磁盘上的具体位置，读取i节点信息并从中找到数据块索引，进而找到相应的数据块，最终获得文件的完整内容

知识讲解



文件



普通文件

- 在Unix/Linux系统中通常所见到的文件，如用C/C++语言编写的源代码文件，编译器、汇编器和链接器产生的汇编文件、目标文件和可执行文件，各种系统配置文件，Shell脚本文件，包括音频、视频等数字内容的多媒体文件，乃至包括数据库在内的各种应用程序所维护、管理和使用的数据文件等，都是普通文件
- 一个普通文件包含以线性字节数组方式组织的数据，通常也称为字节流，Unix/Linux文件没有更进一步的组织结构或格式，因此也不存在类似VMS系统中记录的概念
- 文件中的任何字节都可以被读或者写，这些操作皆始于某个处于特定位置的字节，该位置即所谓当前文件偏移



普通文件（续1）

- 当前文件偏移的最大值仅受存储该值的C语言变量的数据类型限制，最新版本的Linux系统已经可以支持到64位
- 组成文件的线性数组里字节的数目，即文件长度或称文件大小，其最大值受限于Linux内核中用于管理文件的C语言代码数据类型的大小，某些文件系统还可能强加自己的限制，将其限定在更小的值
- 操作系统内核并没有对并发文件访问强加任何限制，不同的进程能够同时读写同一个文件，并发访问的结果取决于独立操作的顺序，且通常是不可预测的
- 一个文件包括两部分数据，一部分是元数据，如文件的类型、权限、大小、用户、组、各种时间戳等，存储在i节点中，另一部分是内容数据，存储在数据块中



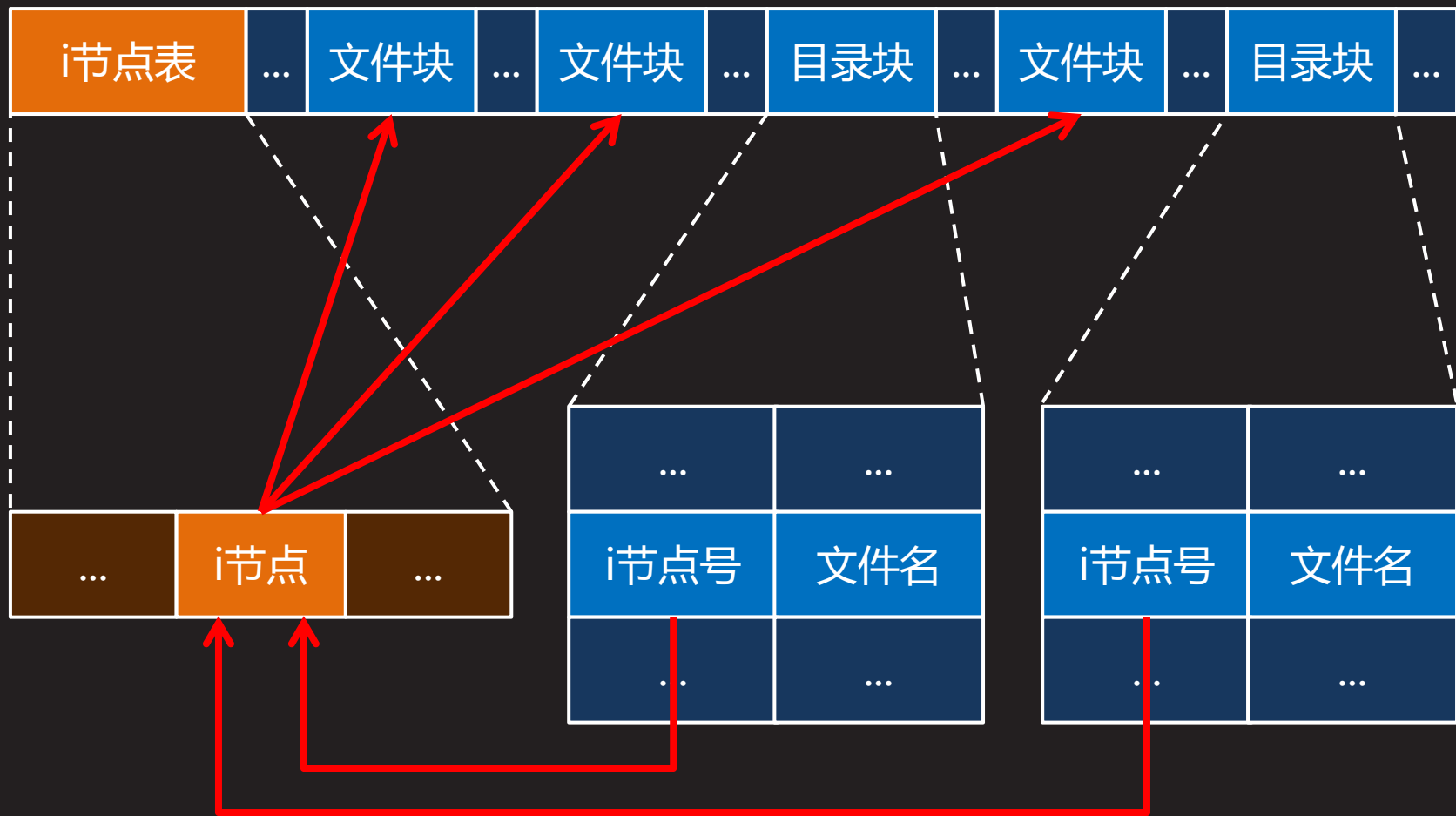
目录文件

- 系统通过i节点号唯一地标识一个文件的存在，但人们更愿意使用有意义的文件名来访问文件，目录就是用来建立文件名和i节点号之间的映射的
- 目录的本质就是一个普通文件，与其它普通文件唯一的区别就是它仅仅存储文件名和i节点号的映射，每一个这样的映射，用目录中的一个条目表示，谓之硬链接
- 既然目录也是文件，那么它同样也有自己的i节点，每个目录的i节点号和它的文件名(目录名)之间的映射，记录在它的父目录中，以此类推，形成了一棵目录树
- 根目录的i节点在i节点表中的存储位置是固定的，因此即使没有父目录，根目录也能被正确检索



目录文件 (续1)

知识讲解



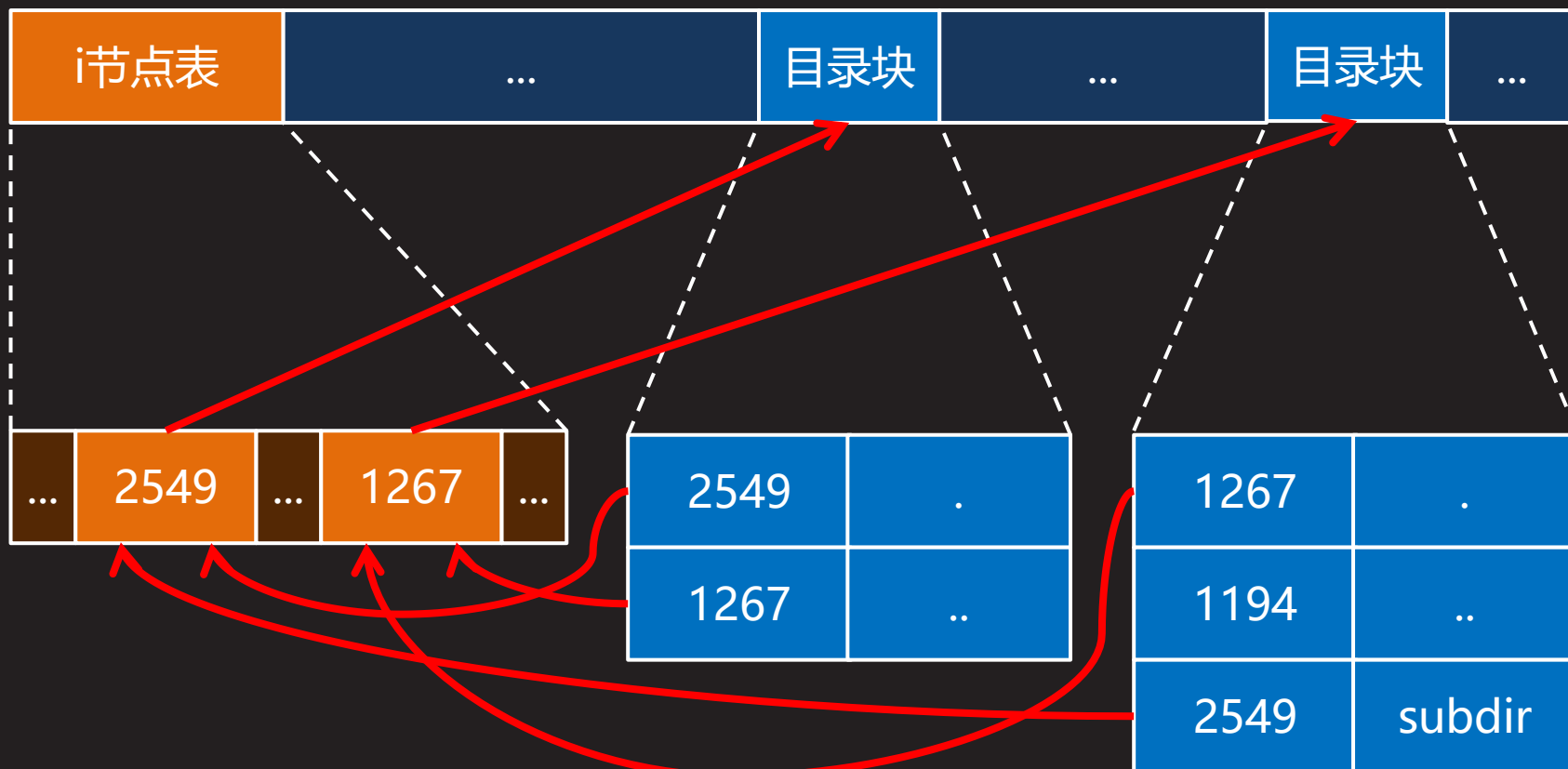
目录文件（续2）

- 当系统内核打开类似 `"/home/tarena/unixc/day01.txt"` 这样的路径时，它会从根目录开始遍历路径中的每一个目录项来查找下一项的i节点号。根目录的i节点可以直接拿到，这样就可以在根目录中找到home目录的i节点号，然后在home目录中找到tarena目录的i节点号，再在tarena目录中找到unixc目录的i节点号，最终在unixc目录中找到day01.txt文件的i节点号并打开该文件
- 如果路径字符串的第一个字符不是 `"/"`，则表示相对路径，基于相对路径的路径解析从当前工作目录开始
- 每个目录中都有两个特殊的条目 `."` 和 `.."` 分别映射该目录本身和其父目录的i节点号，根目录没有父目录，故其 `.."` 和 `."` 一样都映射到根目录本身



目录文件 (续3)

知识讲解



符号链接文件

- 符号链接文件看上去象普通文件，每个符号链接文件都有自己的i节点和包含被链接文件完整路径名的数据块
- 符号链接可以指向任何地方，包括不同文件系统上的文件和目录，甚至不存在的文件和目录(坏链接也是链接)
- 相比于符号链接，硬链接不能跨越文件系统(因为i节点号仅在自己的文件系统内有意义)，且其链接目标必须存在
- 相比于硬链接，符号链接的解析需要更大的系统开销，因为有效地解析符号链接至少需要解析两个文件，符号链接文件本身和它所链接的文件
- 文件系统会为硬链接维护链接计数，但不会为符号链接维护任何东西，因此相较于硬链接符号链接缺乏透明性



特殊文件

- 特殊文件是以文件形式表示的内核对象，具体包括
 - 本地套接字
 - 套接字是网络编程的基础，本地套接字是其面向本机通信的一个变种，它需要依赖文件系统中的一种特殊文件——本地套接字文件
 - 字符设备
 - 设备驱动将字节按顺序写入队列，用户程序从队列中按其被写入的顺序将字节依次读出，如键盘
 - 块设备
 - 设备驱动将字节数组映射到可寻址的设备上，用户程序可以任意顺序访问数组中的任意字节，如硬盘
 - 有名管道
 - 有名管道是一种以文件描述符为信道的进程间通信(IPC)机制，即使是不相关的进程也能通过有名管道文件交换数据



特殊文件 (续1)

- 用ls -l命令查看文件的类型

```
$ ls -l /etc
```

```
-rw-r--r-- 1 root root 1916  9月 14  2013 /etc/passwd
drwxr-xr-x 2 root root 4096  4月 23  2012 opt
```

- -: 普通文件
- d: 目录
- s: 本地套接字
- c: 字符设备
- b: 块设备
- l: 符号链接
- p: 有名管道



总结和答疑

