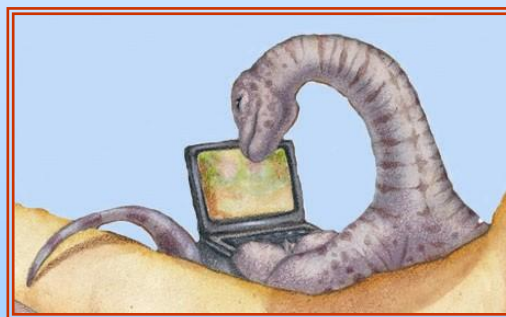


第12章 大容量存储器结构



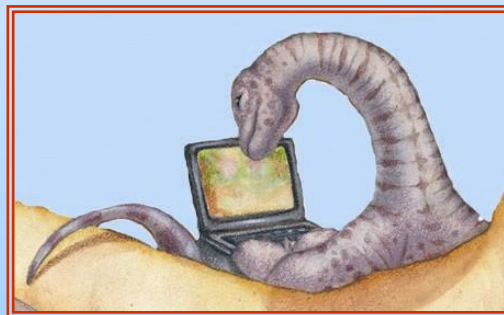


目录

1. 磁盘结构
2. 磁盘调度
3. 磁盘管理
4. RAID结构



1、磁盘结构



磁盘组成

磁盘访问时间

地址映射

磁盘管理



大容量存储设备





磁盘结构

n 盘片

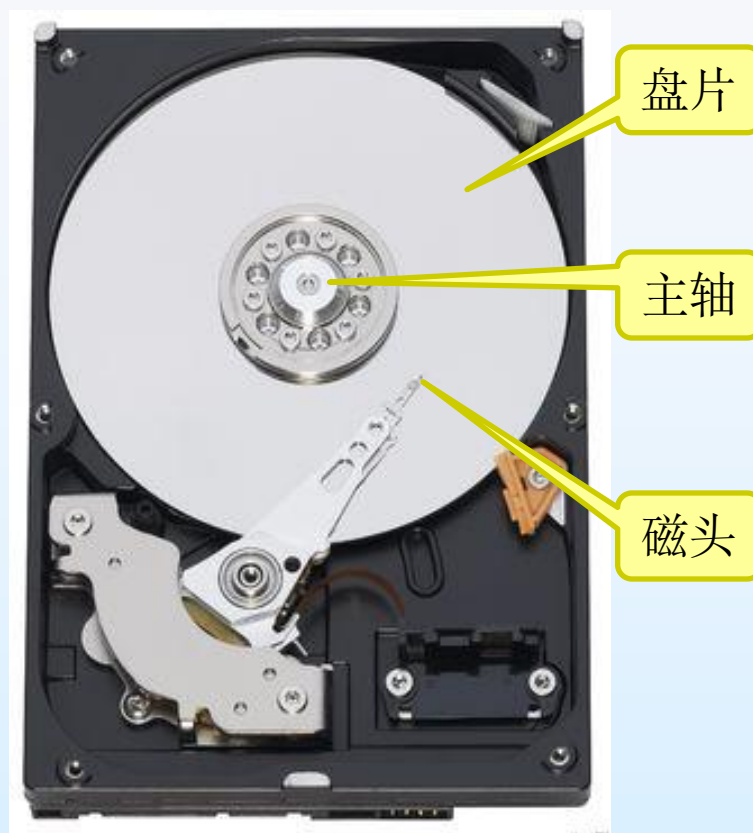
- | 存储数据的介质
- | 正反两面可以存储数据

n 磁头

- | 读写数据，沿磁盘半径移动
- | 有多少盘面就有多少磁头

n 主轴

- | 马达驱动，使盘片旋转
- | 固定速度旋转





磁盘控制器和接口

n 接口

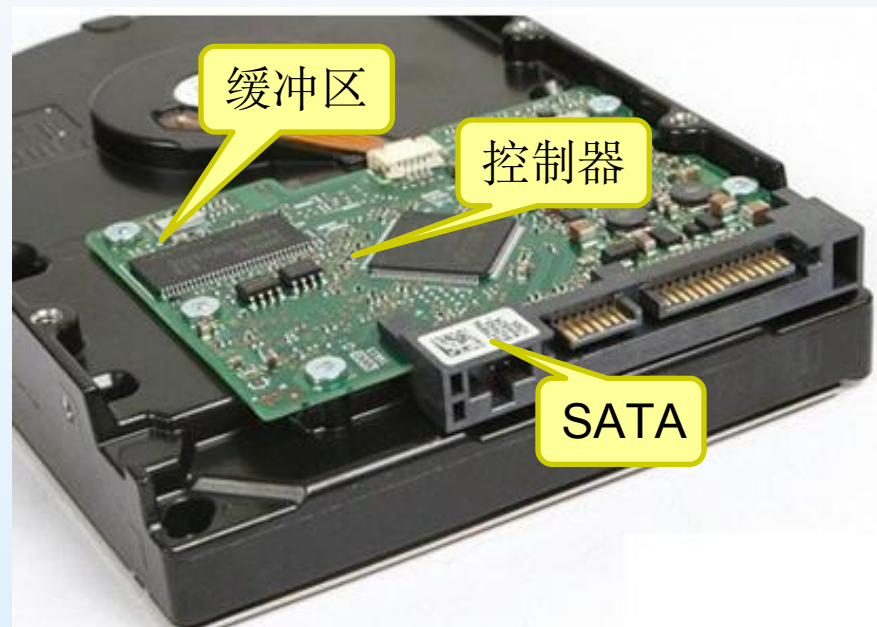
- | EIDE, ATA, SATA, USB, Fibre Channel, SCSI, SAS, Firewire

n 磁盘控制器

- | 控制磁盘的读写等操作

n 缓冲区

- | 利用磁盘缓冲区来暂存数据





盘片结构

n 磁道

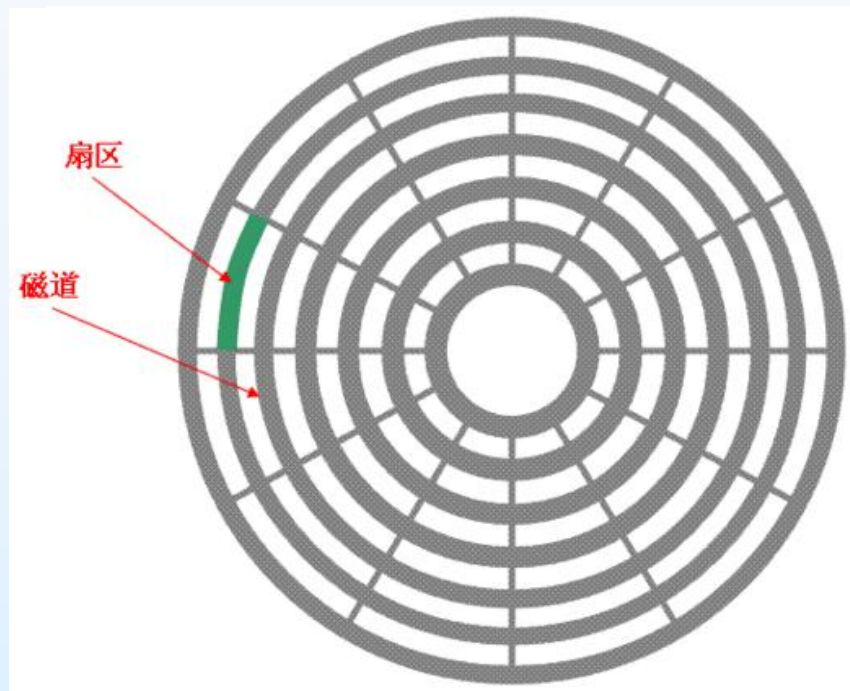
- | 磁头在盘片表面划出的圆形就是一个磁道
- | 每个盘面划分为数目相等的磁道
- | 从盘面外缘“0”开始编号

n 扇区

- | 磁道被等分为若干个弧段，称为扇区
- | 扇区大小：512字节

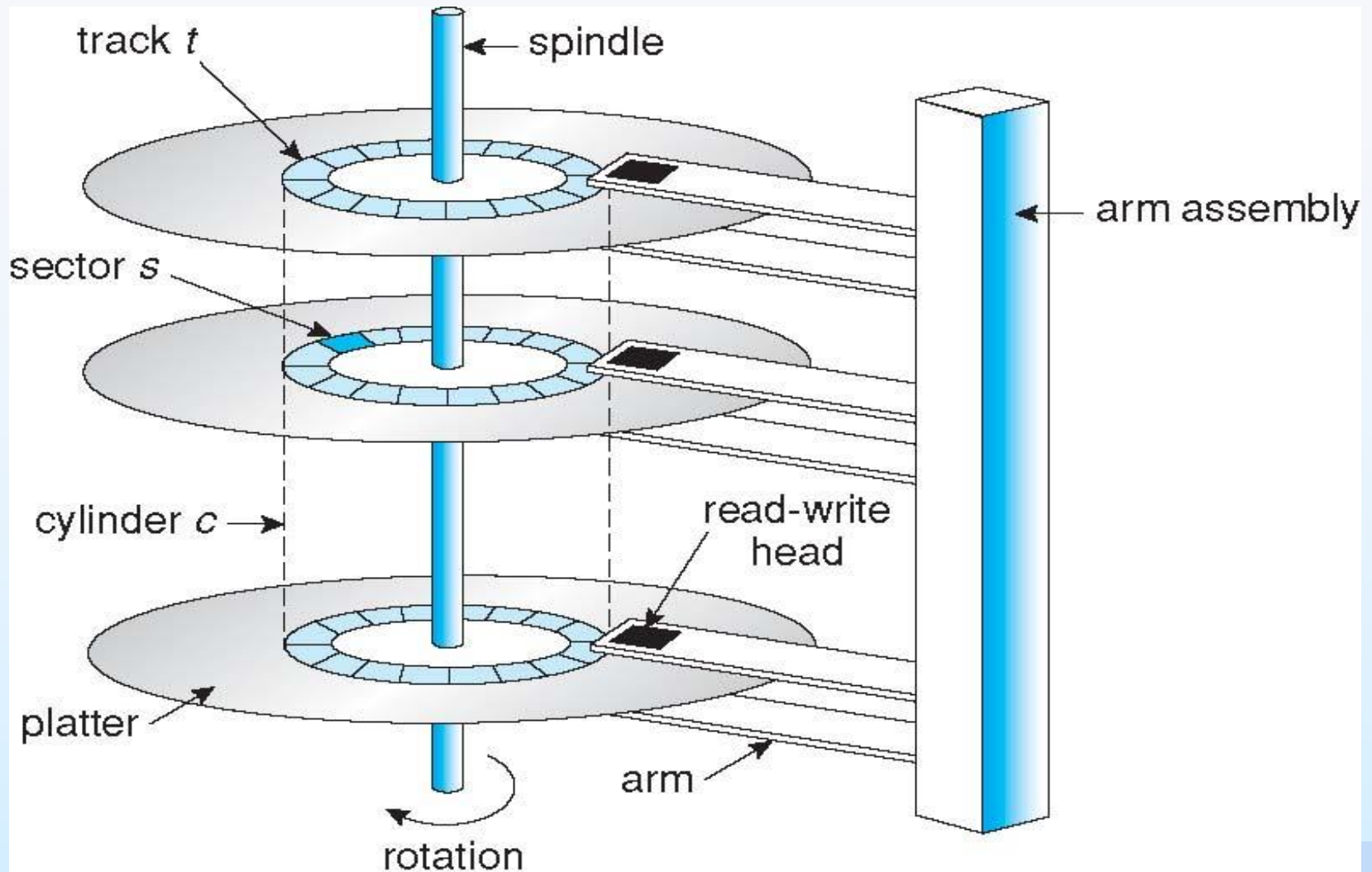
n 柱面

- | 每个盘面上具有相同编号的磁道形成一个圆柱，称为柱面
- | 一个盘面有几个磁道就有几个柱面





磁盘示意图





地址映射关系

n 块号: LBA

n 磁盘地址(CHS): (C,H,S)

| Cylinder (柱面/道C)

| Head (磁头/面H)

| Sector (扇区S)

n SPT:每个磁道最大扇区数

n HPC:最大磁头 (盘面) 数

$$C = LBA \div (HPC \times SPT)$$

$$H = (LBA \div SPT) \bmod HPC$$

$$S = (LBA \bmod SPT) + 1$$





- n 物理扇区号 = ((柱面号 × 磁头数) + 磁头号) × 每磁道扇区数 + 扇区偏移量 - 1
- | 柱面号从0开始编号;
 - | 磁头号从0开始编号;
 - | 扇区号从1开始编号;
 - | 扇区偏移量指的是在某个磁头的某个柱面上的扇区位置偏移量, 从1开始编号。
- n 例如, 假设某个磁盘有16384个柱面, 16个磁头, 每个磁头有63个扇区, 每个扇区大小为512字节, 要计算某个扇区的物理位置, 其CHS地址为 (1234, 5, 6), 则计算过程如下:
- | 物理扇区号 = ((1234 × 16) + 5) × 63 + 6 - 1 = 1244192
 - | 物理位置 = 1244192 × 512 = 637026304





- n 一个磁盘有4个磁片，每个磁片划分为1024个磁道，每个磁道分为256个扇区。每个扇区容量512B，这个磁盘容量为（）





磁盘访问时间

n 定位时间/随机访问时间:

- | **寻道时间**: 移动磁臂到所需磁道时间
 - ▶ 平均寻道时间: 1/3 磁道移动 (1-4ms)
- | **旋转延迟**: 等待扇区移动到磁头下时间
 - ▶ 由磁盘的旋转速度决定
 - ▶ 磁盘旋转速度: 60 – 250转/秒
 - ▶ **RPM (Revolutions Per Minute)**: 每分钟旋转次数, 如: 7200RPM, 即120转/秒
 - ▶ 平均旋转1/2圈时间: $1/(2 \cdot \text{RPM}/60)$
 - ▶ 平均延迟时间: $1/(2 \cdot 7200/60) = 4.17$ 毫秒

Spindle [rpm]	Average latency [ms]
4200	7.14
5400	5.56
7200	4.17
10000	3
15000	2

常用的**RPM**对应的平均延迟时间

❖ 传输时间

- | 传输的数据量乘以传输率
- | 传输率: 传输总字节数除以传输时间
- | 例如: 6 Gb/sec → 1秒可以传输6G位数据

传输1KB数据的传输时间:

$$1K \cdot 8 / 6G = 7.5 \text{ 微秒}$$





磁盘访问时间

基本参数/ESSENTIAL PARAMETER

型号：~~WD60EZRZ~~

容量：6 TB

接口：SATA 6 Gb/s

规格：3.5 英寸

转速：5400 PPM

缓存：64 MB

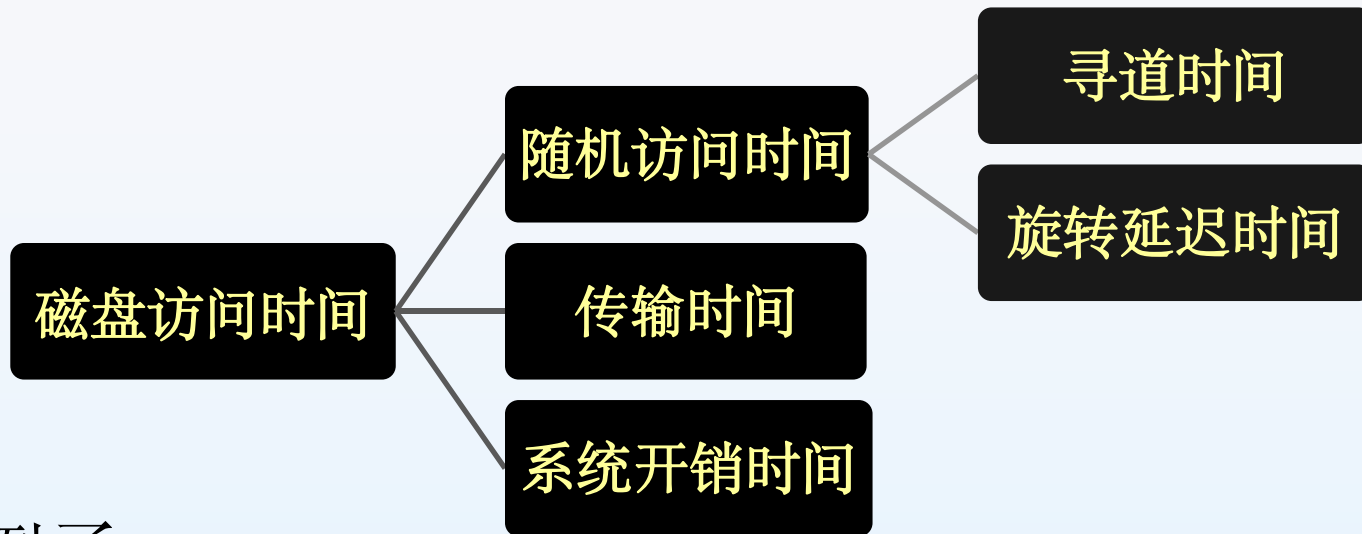
尺寸 (mm)：高*长*宽：26.1*147*101.6

适用系统：台式机 / 一体机电脑





磁盘访问时间



n 例子

- | 7200 RPM 转速， 5ms平均寻道时间， 1Gb/sec传输率， 0.1ms 系统控制开销。那么读取4KB数据块的磁盘访问时间为：

$$\begin{aligned} & 5\text{ms} + \frac{1}{2} \times \frac{1}{(7200/60)\text{sec}} + 4\text{KB} / 1\text{Gb/sec} + 0.1\text{ms} \\ &= 5\text{ms} + 4.17\text{ms} + 0.03\text{ms} + 0.1\text{ms} \\ &= 9.3\text{ms} \end{aligned}$$





n 一个磁盘的传输率为 2Gb/s ，那么传输 1MB 数据需要的传输时间是 ()





磁盘管理

- n 低级格式化（物理格式化）
 - | 将磁盘分成扇区，以便磁盘控制器读写
- n 分区
 - | 将磁盘分成分区
 - | 主分区和扩展分区
- n 高级格式化
 - | 逻辑格式化，创建文件系统
- n 引导块
 - | 自举程序保存在ROM中
 - | 自举程序装载引导块程序



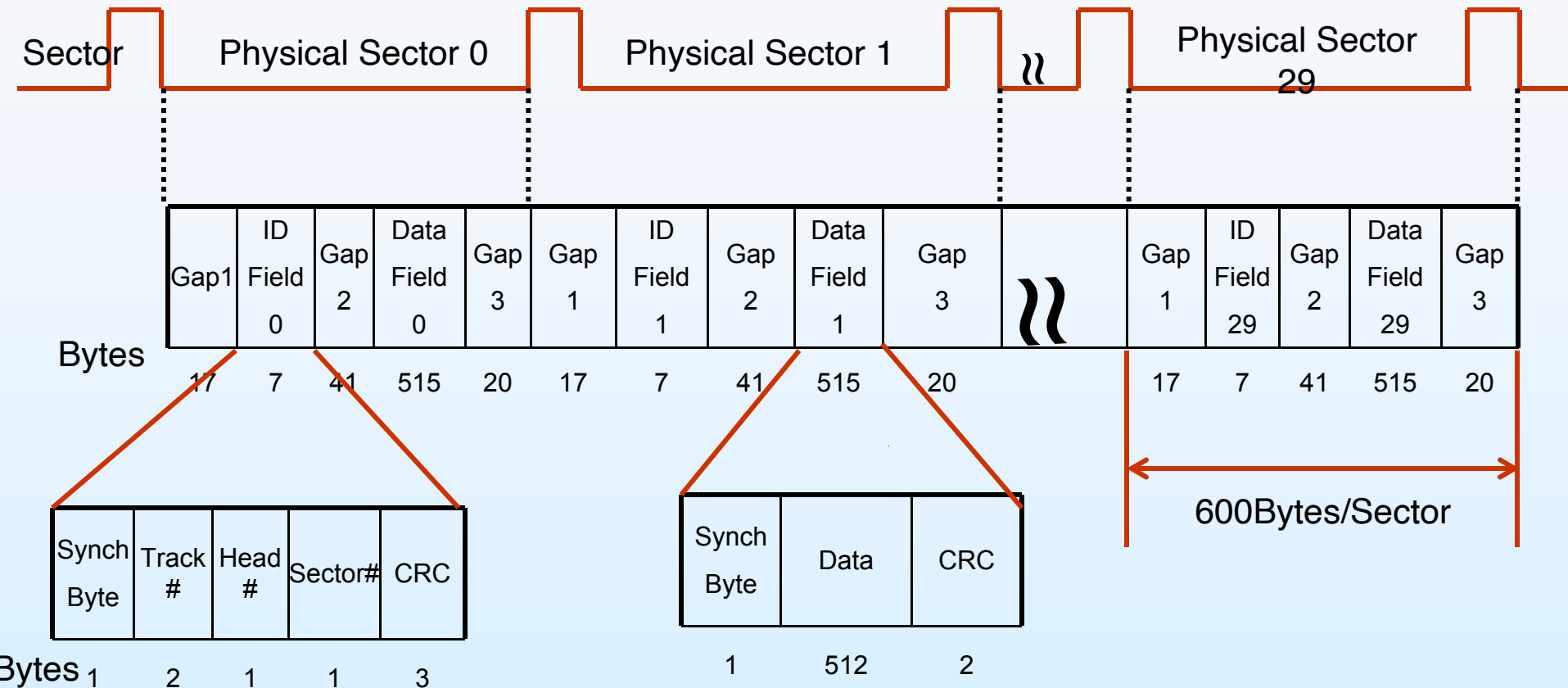


低级格式化



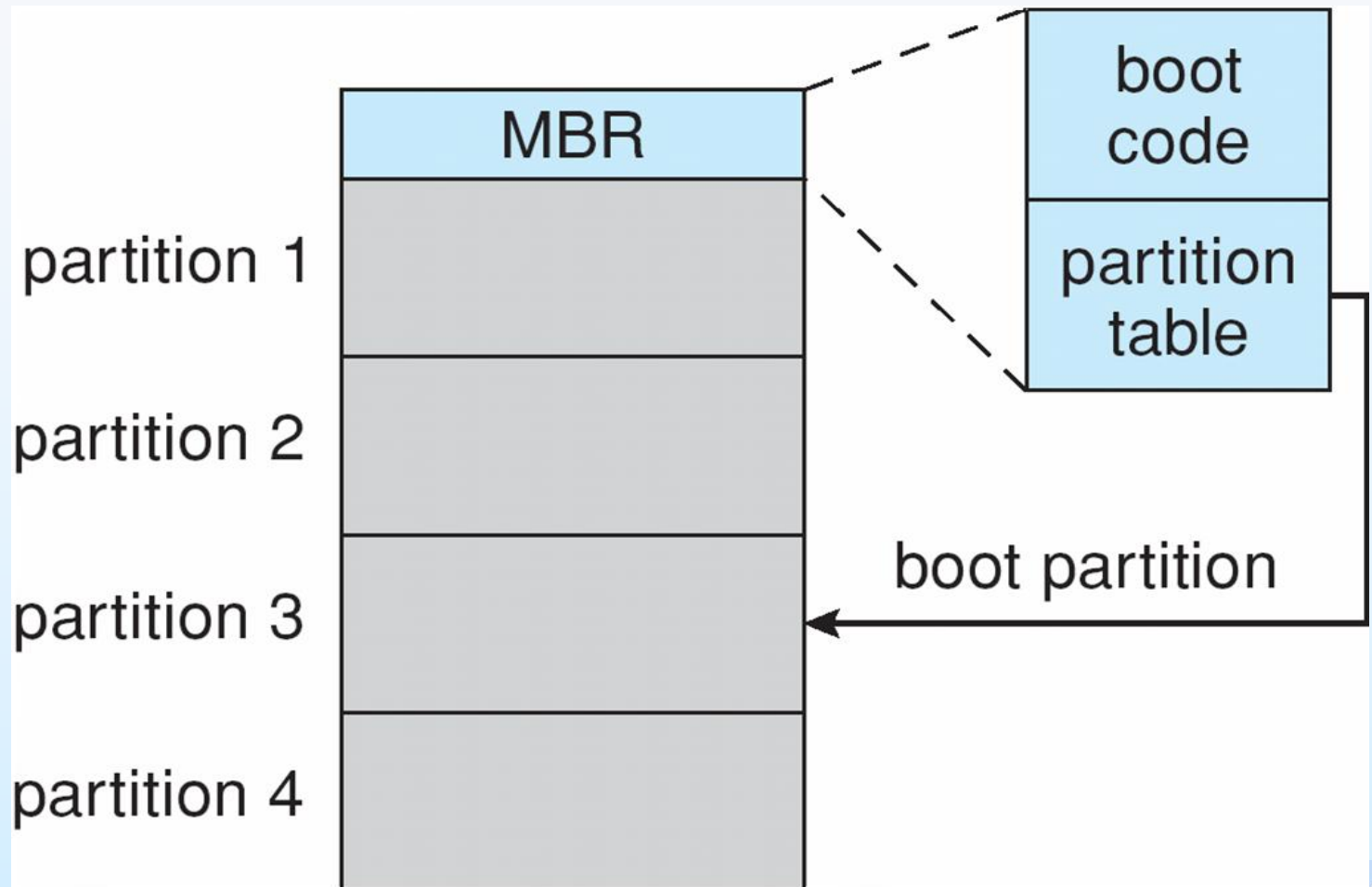


扇区格式





分区 (Windows)





引导区记录 (MBR)

标准 MBR 结构

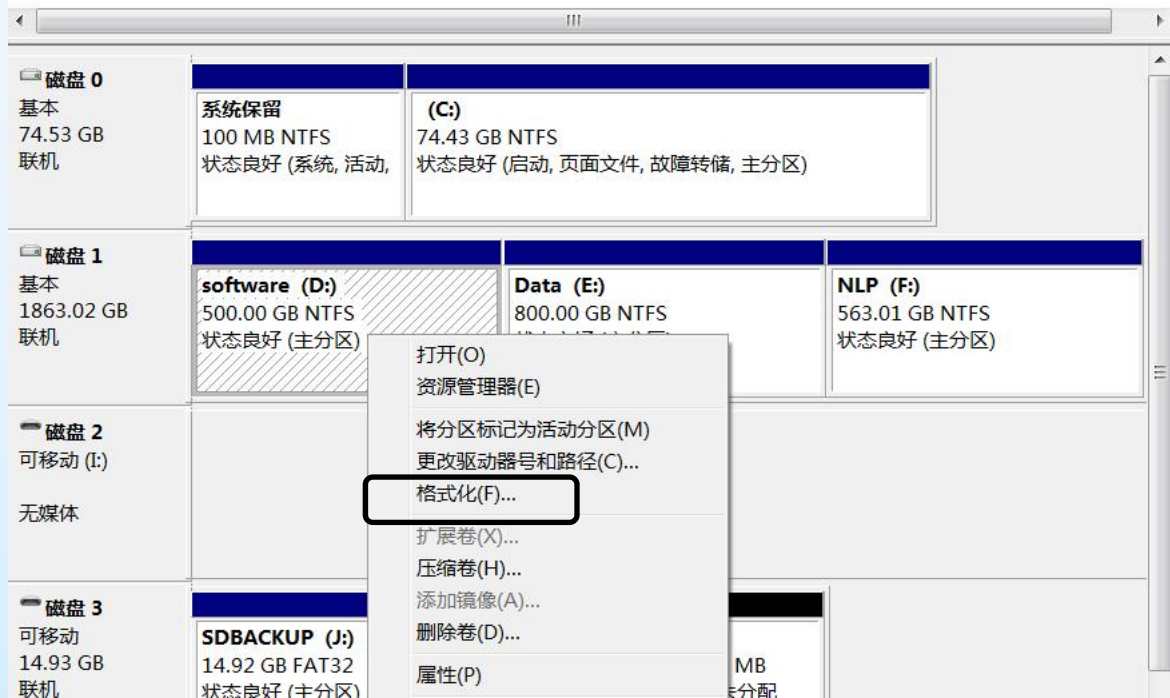
地址			描述	长度
Hex	Oct	Dec		(字节)
0	0	0	代码区	440 (最大 446)
01B8	670	440	选用软盘标志	4
01BC	674	444	一般为空值; 0x0000	2
01BE	676	446	标准 MBR 分区表规划 (四个16 byte的主分区表入口)	64
01FE	776	510	55h	MBR 有效标志
01FF	777	511	AAh	
MBR, 总大小: $446 + 64 + 2 =$				512





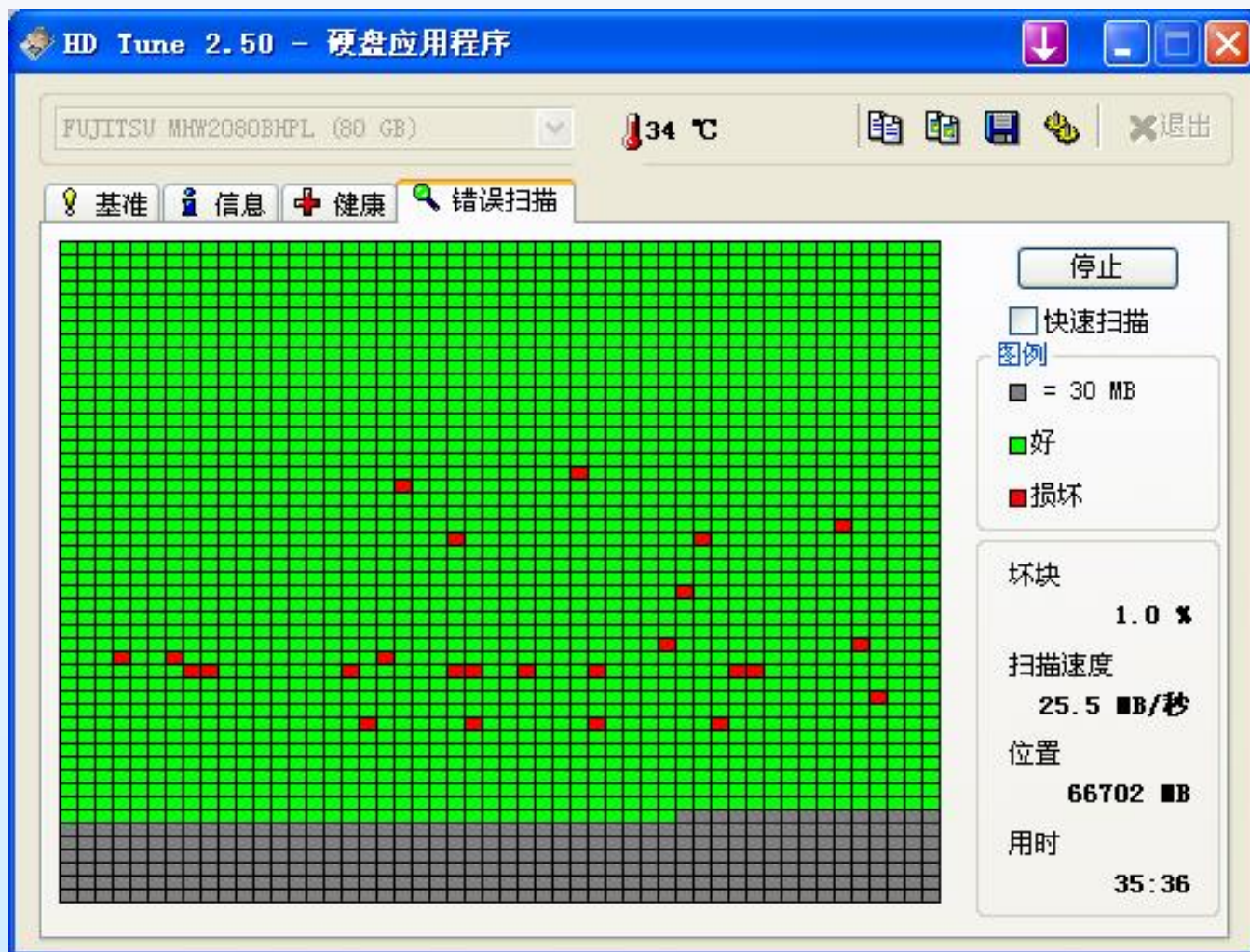
高级格式化

卷	布局	类型	文件系统	状态	容量	可用空间	% 可用	容错	开
(C:)	简单	基本	NTFS	状态良好 (启动, 页面文件, 故障转储, 主分区)	74.43 GB	7.22 GB	10 %	否	0
D3P_SCN (G:)	简单	基本	UDF	状态良好 (主分区)	4.26 GB	0 MB	0 %	否	0
Data (E:)	简单	基本	NTFS	状态良好 (主分区)	800.00 GB	249.75 GB	31 %	否	0
NLP (F:)	简单	基本	NTFS	状态良好 (主分区)	563.01 GB	439.10 GB	78 %	否	0
SDBACKUP (J:)	简单	基本	FAT32	状态良好 (主分区)	14.91 GB	3.40 GB	23 %	否	0
software (D:)	简单	基本	NTFS	状态良好 (主分区)	500.00 GB	456.27 GB	91 %	否	0
系统保留	简单	基本	NTFS	状态良好 (系统, 活动, 主分区)	100 MB	71 MB	71 %	否	0

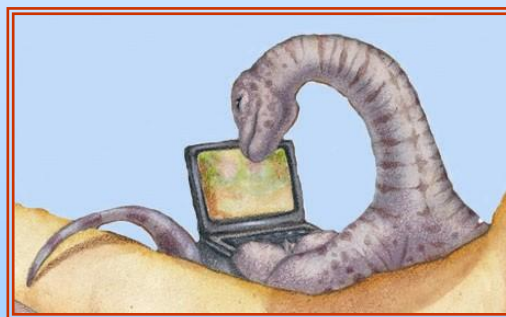




坏块检查



2、磁盘调度





内容

n 磁盘调度

- | 引入磁盘调度的目的是为了降低磁盘访问时间，提高文件系统的效率

n 先来先服务算法

n 最短寻道时间优先算法

n 扫描算法

n RAID

- | 引入**RAID**技术的目的是为了增强数据的可靠性和访问的并行性





磁盘调度

- n 目标：减少磁盘访问时间
- n 访问时间：
 - | 寻道时间：磁头移动到访问扇区所在磁道的时间
 - | 旋转延迟时间：将访问扇区转到磁头下的时间
 - | 传输时间：将数据从磁盘送到内存的时间
- n 寻道时间最小化
- n 寻道时间 \approx 寻道距离





请求系列

n 假定有一个请求序列(0-199道).:

98, 183, 37, 122, 14, 124, 65, 67

磁头当前位置在53

目标：磁头移动距离最小，寻道时间最短





先来先服务算法FCFS

n First Come First Served

n 按照请求提交时间访问

┆ 先提交先访问

┆ 后提交后访问

n 优点

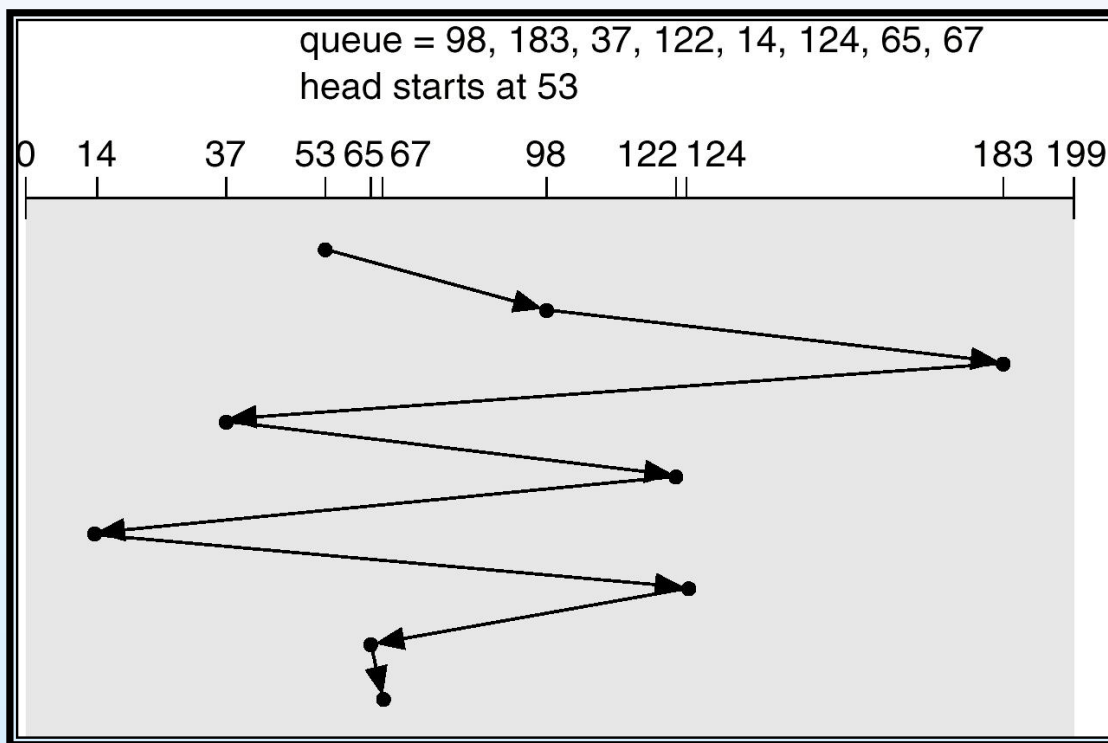
┆ 简单、公平

┆ 易实现

n 缺点

┆ 寻道时间长

┆ 效率低





最短寻道时间优先算法SSTF

n Shortest Seek Time First

n 每次移动到离现在位置最近的磁道

- | 最短寻道时间

- | 最短作业优先 (SJF)

n 优点

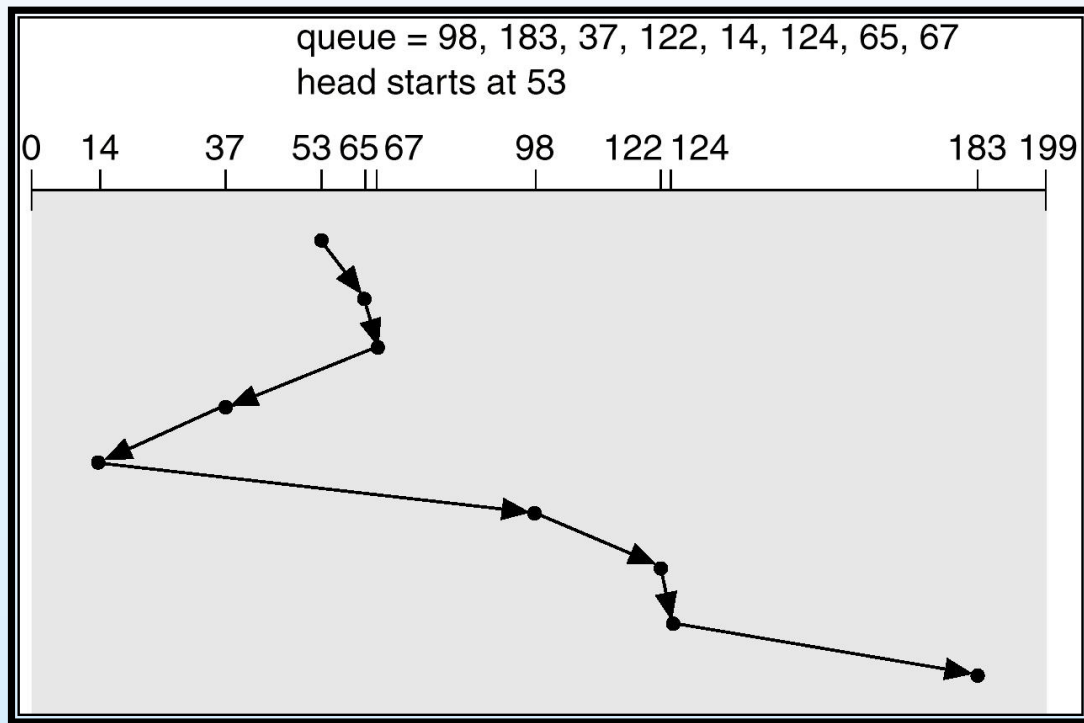
- | 寻道距离短

n 缺点

- | 存在饥饿

- | 磁头频繁变换移动方向

- | 增加寻道时间





扫描算法SCAN

n 磁头从磁盘一端向另一端移动，沿途响应服务请求

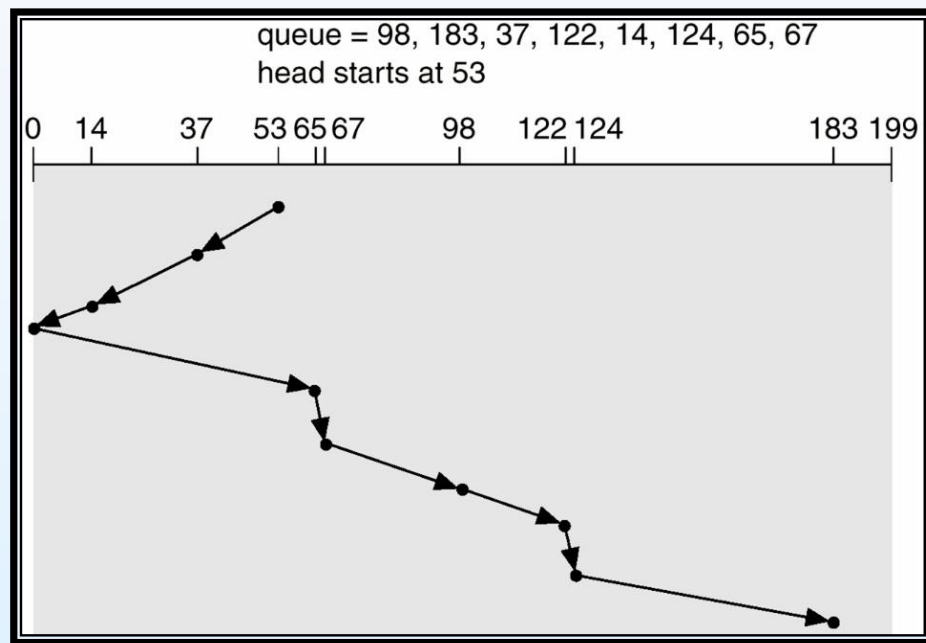
- 到达另一端时，磁头改变移动方向，继续处理
- 磁头在磁盘上来回扫描
- 又称为电梯算法

n 优点

- 同一方向扫描，寻道时间短
- 改变磁头方向少

n 缺点

- 有的请求等待时间长



总的磁头移动为236磁道





循环扫描算法C-SCAN

n Circular Scan

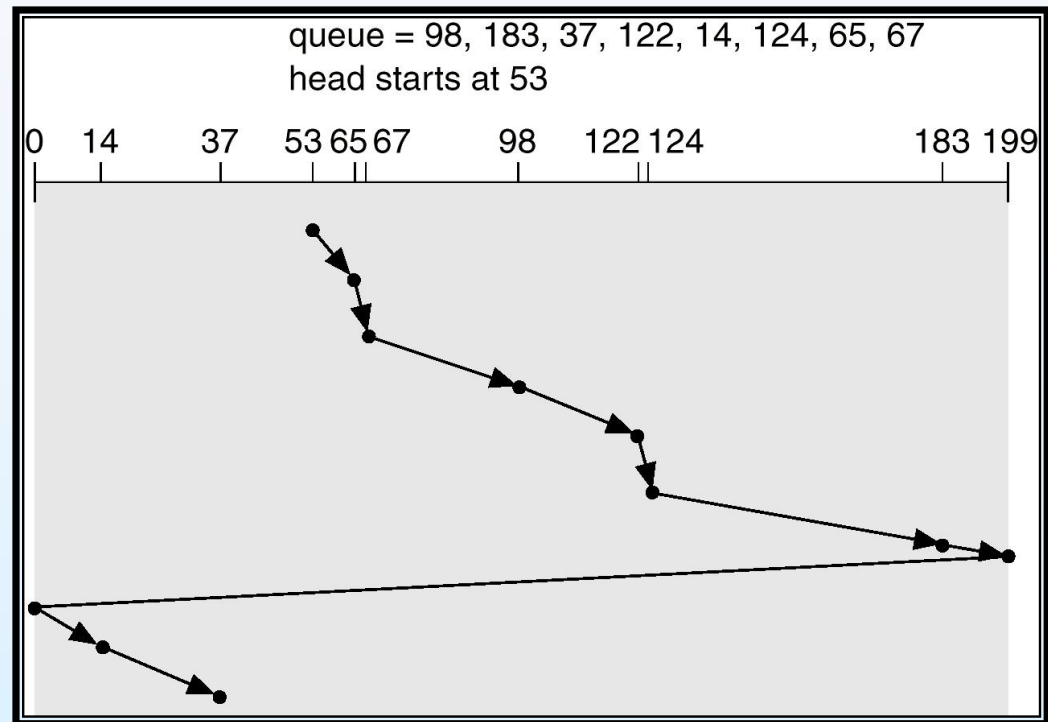
n 单向处理请求

- | 磁头从磁盘外道 (0道) 移到内道过程中处理请求
- | 内道移动到外道的过程中不处理请求

n 优点

- | 更均匀的等待时间

n 从磁道199移动到0的时间很短



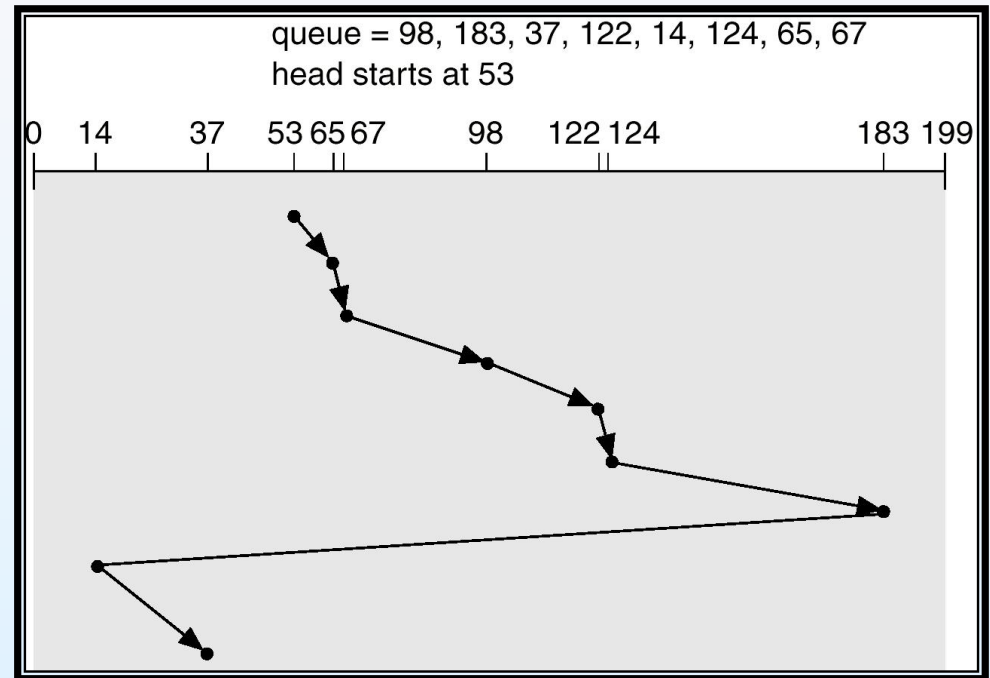
总的磁头移动为382磁道





循环Look算法C-LOOK

- n C-SCAN变形
- n 磁头只移动到一个方向上最远请求为止，而不是继续到磁盘尽头



总的磁头移动为322磁道





磁盘调度算法的选择

- n 磁盘调度性能主要依赖于请求的数量和类型
 - | 磁盘服务请求很大程度上受文件分配方法影响，例如隐式链接的服务请求数就会比较多
 - | SSTF较为普遍且很有吸引力
 - | SCAN和C-SCAN适合磁盘大负荷系统
- n SSTF或LOOK是比较合理的缺省算法





RAID结构

n RAID

- | **Originally** Redundant Arrays of Inexpensive Disks (廉价磁盘冗余阵列)
- | **Now** Redundant Arrays of Independent Disks (独立磁盘冗余阵列)
- | RAID把很多价格较便宜的磁盘组合成一个大容量的磁盘组，利用个别磁盘提供数据所产生加成效果提升整个磁盘系统效能和可靠性。

n RAID卡 (现代CPU集成RAID)

n RAID被分成了多个不同级别

- | RAID0-RAID7
- | RAID01, RAID10, RAID5E, RAID5EE, RAID50





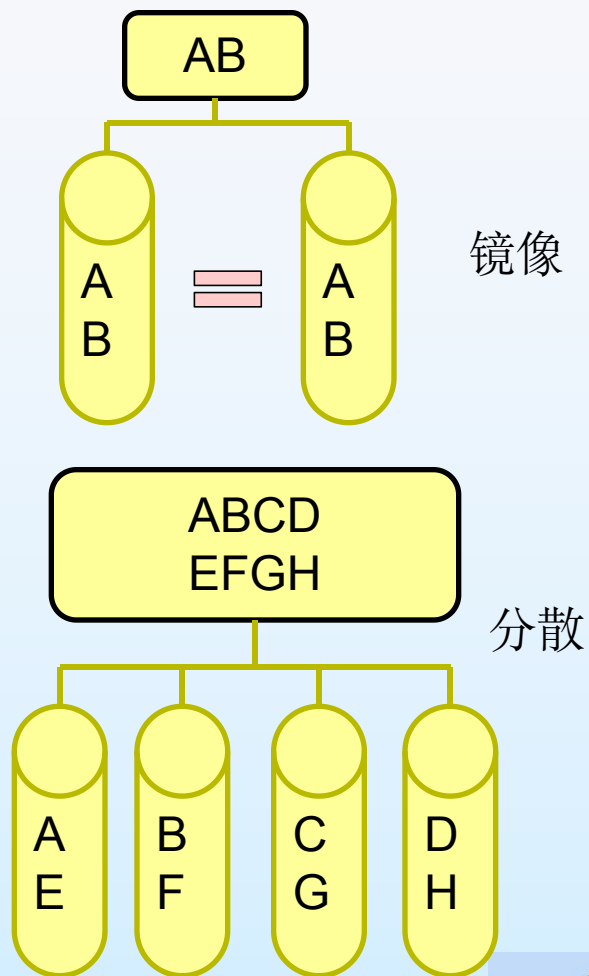
RAID性能

n 可靠性

- | 磁盘可靠性：要求存储在磁盘上的数据不易丢失
- | 引入**冗余**
- | 例如：镜像，把数据在两个磁盘上各存一次

n 性能（**数据分散，并行读写**）

- | 位级分散：数据每个字节的各个位分散在多个磁盘上
- | 块级分散：数据以块为单位分散在多个磁盘上





RAID级别

n RAID 0

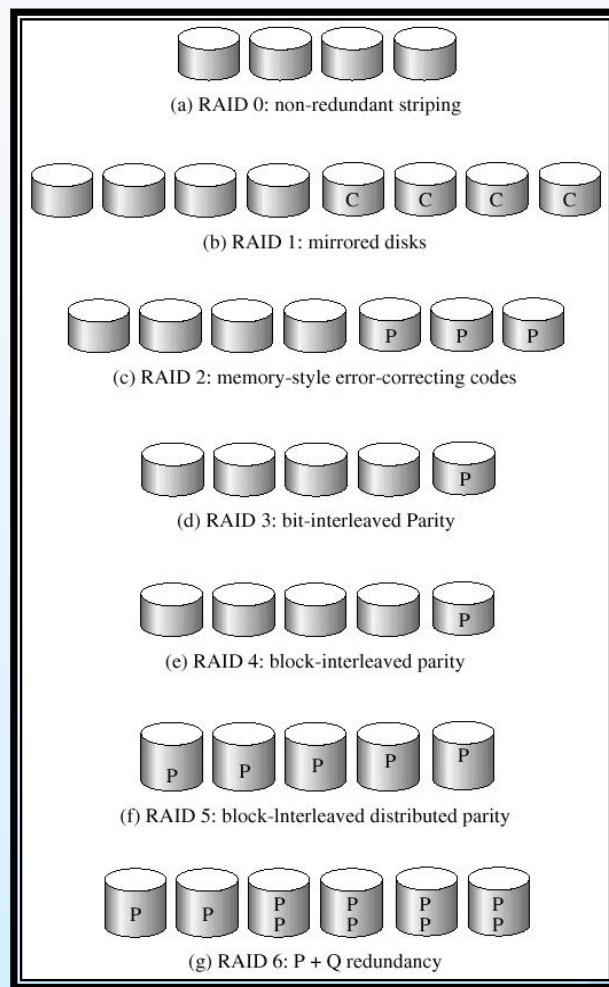
- | 数据分散在多个磁盘上
- | 条状分散技术
- | 提高读写性能

n RAID 1

- | 磁盘镜像
- | 提高可靠性

n RAID 5

- | 分散+校验
- | 校验信息分散在各个磁盘避免对单个校验磁盘的过度使用





RAID级别

n RAID 0

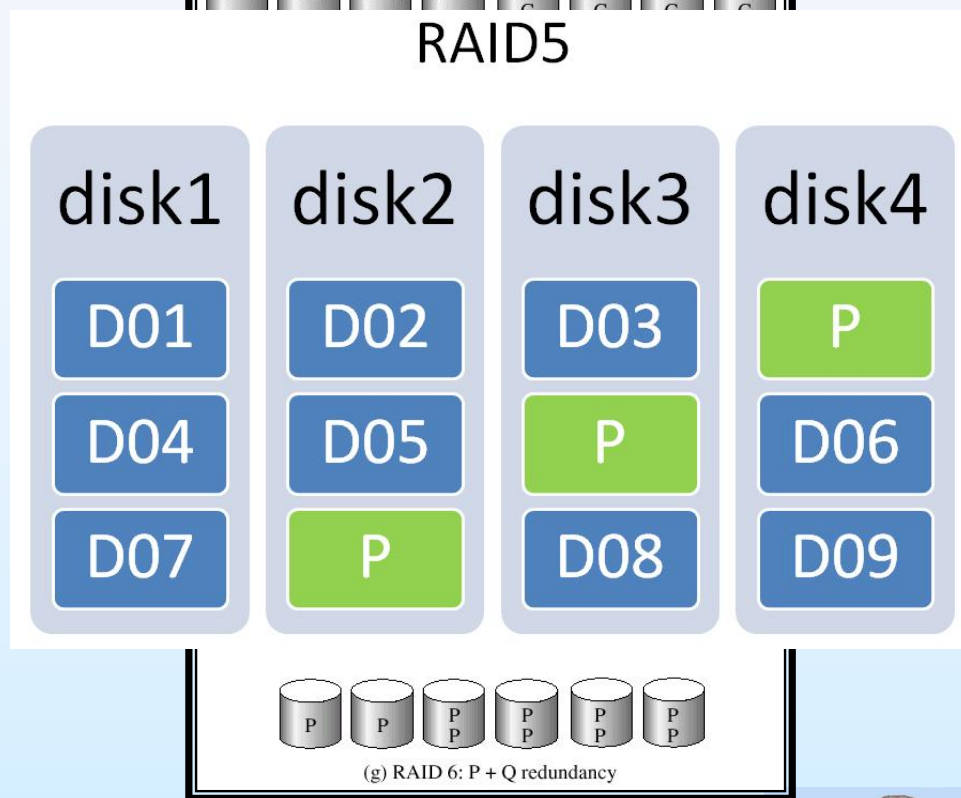
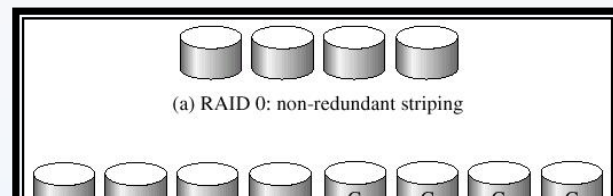
- | 数据分散在多个磁盘上
- | 条状分散技术
- | 提高读写性能

n RAID 1

- | 磁盘镜像
- | 提高可靠性

n RAID 5

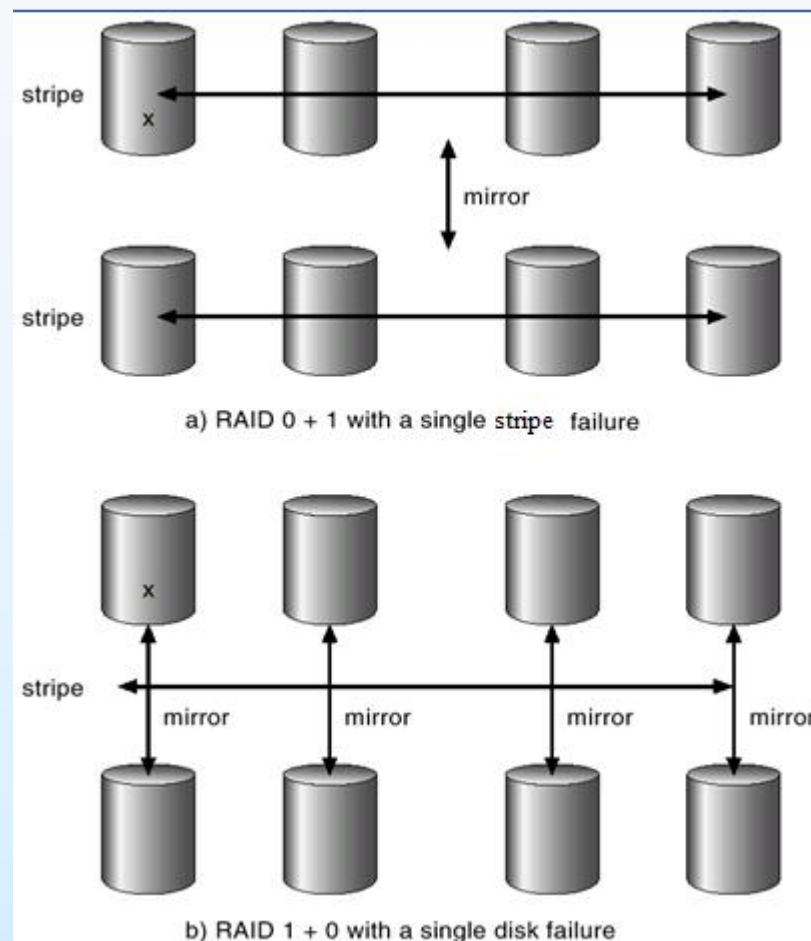
- | 分散+校验
- | 校验信息分散在各个磁盘避免对单个校验磁盘的过度使用





RAID (0 + 1) 和 (1 + 0)

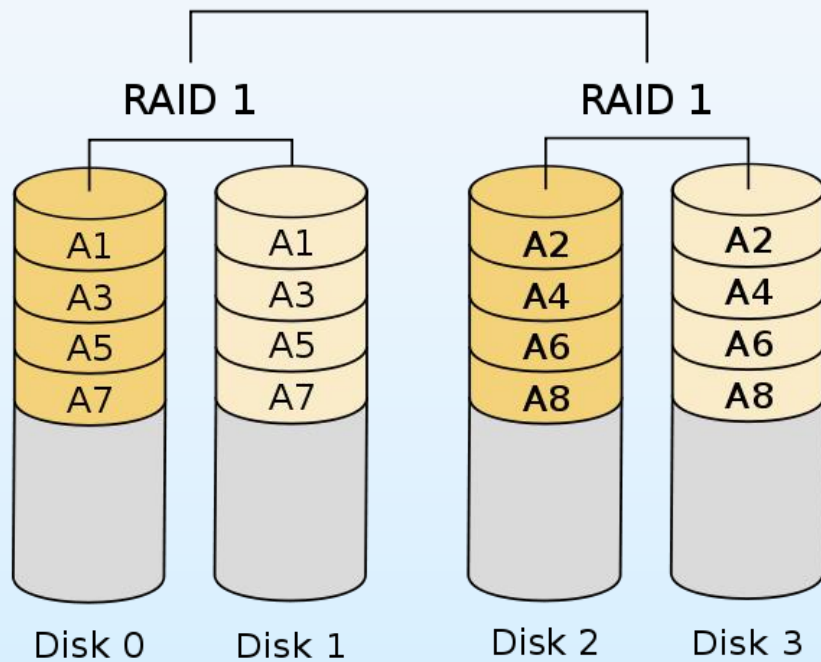
- n RAID0: 性能
- n RAID1: 可靠性
- n 二者兼备
- n RAID01
 - | 先做分散, 再做镜像
 - | 性能好
 - | 但是一个磁盘的故障会导致一条磁盘带不能访问
- n RAID10
 - | 先做镜像, 再做分散
 - | 可靠性好, 一个磁盘的故障不会影响其他磁盘





RAID 1+0

RAID 0



RAID 0+1

RAID 1

