

《大话存储》

Notes

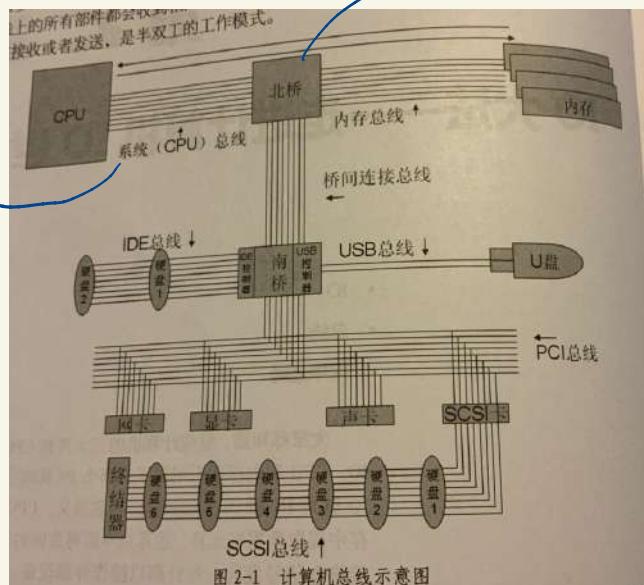
P12

总线

高速 PCI 2.0 直接连北桥



前端总线
(快)



I/O 总线
(慢)

总线 (bus) = { 基本单向
每根导线 半双工 }

{ 数据总线 : 决定位宽 }

控制总线

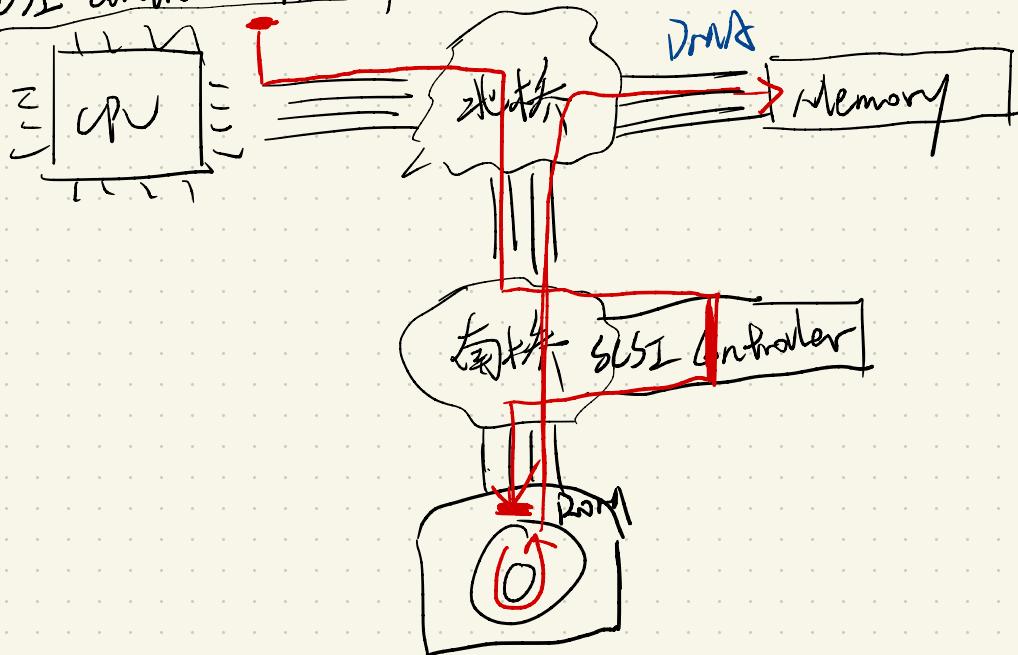
{ 仲裁总线 (中断) : 决定某时刻某条线
谁来使用 }

地址总线

PCI Multiplex + 中断共享

网络 = 连 + 找 + 发

SUSI Controller Driver



向 SCSI 硬盘 某地址读取的数据流。

P21 磁头

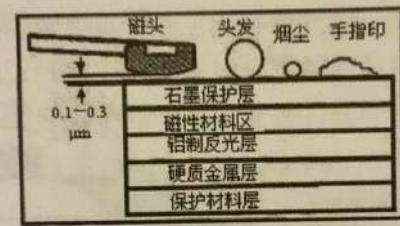


图 3-3 磁头厚度示意图

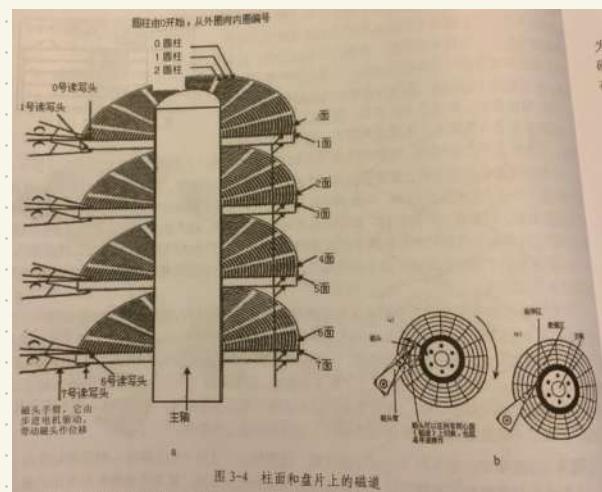


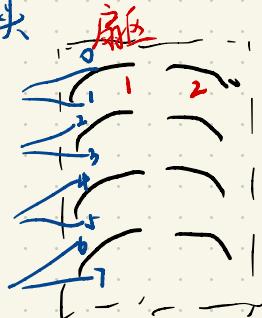
图 3-4 柱面和盘片上的磁道

② 磁道(扇区)存放的代码

读写最小单位为扇区 512B / 4KB

- 低级格式化：盘面上划分磁道-扇区
- (高级)格式化：进行文件系统的标记

磁头



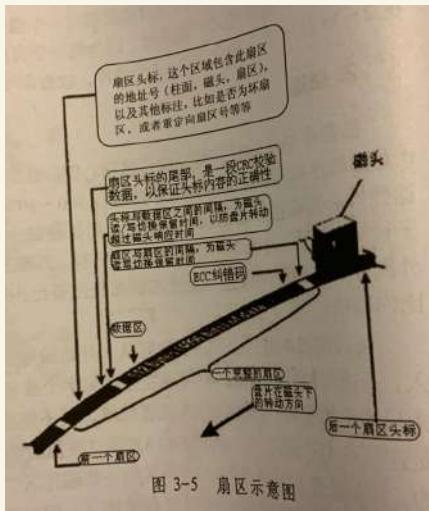
RW) 次序：

柱面(磁道) \Rightarrow 磁头(盘面) \Rightarrow 扇区

柱面

$$\text{Delay} = T_{\text{寻道}} + T_{\text{旋转}} + T_{\text{读取}}$$

*



Cyl 地址 (内部)

LBA 地址 (对外线性)

Virtualization!

交叉因子编号：

防止扇区间隔时间控制

器来不及处理余留数据

P32 磁头 Scheduling

FCFS (First Come First Serve)

SSTF (Shortest Seek Time First)

SCAN / C-SCAN (回旋扫描)

LOOK / -LOOK (回旋不到头)

R51 Drivers

Bios 中内置简化版可用驱动，启动完之后再 load OS 中的复杂完整的驱动

$$IOPS = \frac{\text{Irene Depth}}{\text{IO Latency}}$$

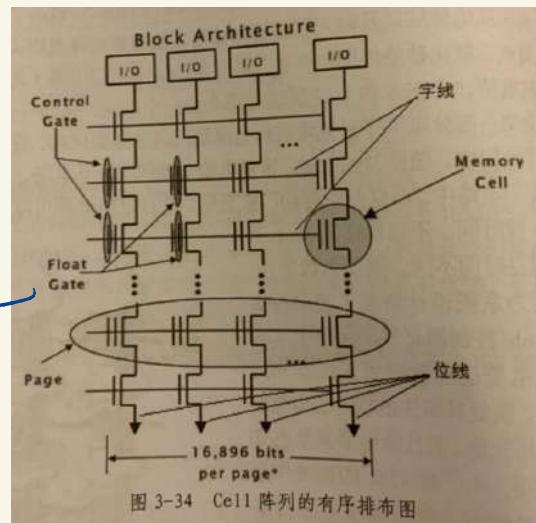
$$n \times \text{latency} \approx n \cdot 2.5 \text{Gb/s}$$

{ 高 IOPS 在随机大量小 RW 时才优势明显
高 Bandwidth 在传输大块连续数据时优势
二者不对立，可以兼有。

P57

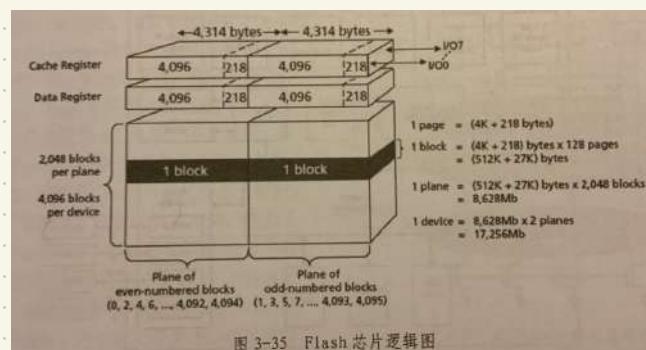
SSD 固态硬盘

擦写寿命为 10¹²
平均读写 10¹¹



* 最小 R 单位
一个 Page
最小 W 单位
一个 Block

MLC Multilevel 和 SLC single-level 带成本低但更不 robust



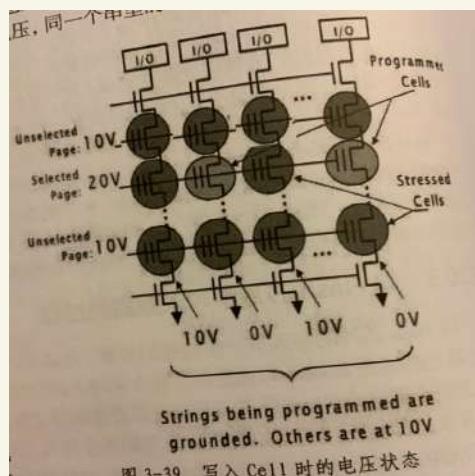
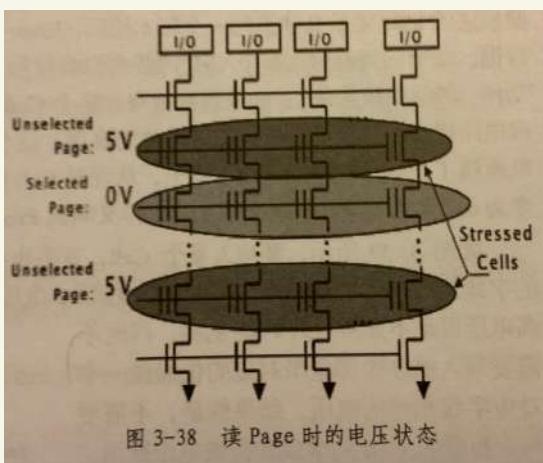
示例规格：

一行为一个 page = 4KB + 218B cells

一个 Block = 128 Pages

奇/偶 Plane = 2,048 Blocks (各自)

一个 SSD Device = 2 Planes + --



R

给位线 pre-charge \Rightarrow 让其自然漏电
 \Rightarrow 充了负电荷的 cell 漏电电流越弱
 \Rightarrow 一段时间后电压较高 $\Rightarrow V_{ret} < V_{bit}$
 \Rightarrow 遗留下“0”

W

Erase 整个 Block \Rightarrow Program
 (全部放电为 1)

原因：寿命，1 个 cell 改写次数有限，

所以用 Append 写的方式平衡了，Row

写放大

Write Amplification 使编写 (update) \Rightarrow Append + 查定位句。
 垃圾回收点：随机大量小 W 时性能变差
 GC

都是钱的问题。做到 SRAM 那样的 byte addressable 会单位成本很贵，”（面积也大了）

NOR Flash 即独立位线，可独立寻址但面积更大，写效率也低了。

Free space 保留是缓解的方法之一 (5%~20%)

Delayed Write 等等 --- Flash-aware FS ---

P68 处理 bad 损坏

参数 8b@±12b \Rightarrow 每±12b 需配合纠正
8bits 错误的纠正码

P78

RAID

Redundant Array of Inexpensive Disks
Independent

∇ <https://zh.m.wikipedia.org/zh/RAID>

软体 RAID vs. 硬件 RAID 卡

(集成于 SCSI 卡 / 南桥
或独立成卡)

校验型 RAID 需要初始化

P140 | 卷管理 — LVM 方案

物理盘（或 RAID Controller 暴露给 OS 的逻辑盘，OS 亦认为其为物理盘），称“物理卷”



多个物理卷组成一段逻辑上连续的
Volume Group



VG 又被分割为数个“物理分区”
Physical Partitions



PP 被对应到“逻辑区块” Logical Partition

{ 1 PP → 1 LP
多个 PP → 1 LP

mb/s 的单位



多个 LP 组成一个逻辑卷 Logical Volume

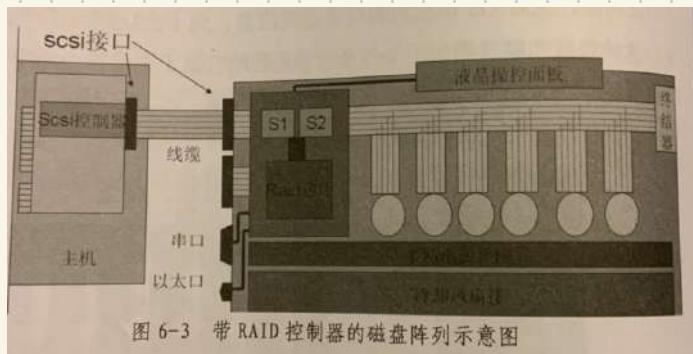
好处：可以随时在线地增减 LV 大小或创建新的，不受限于物理磁盘或 RAID 配置

问：1 boot 挂载到 sda 的前 20 个磁道
成为设备 sda1，不参与 LVM
而后 sda2, sdb, sdc ... 可拼在一起被分割
映射层级越来越多，更方便但也更复杂以致
Tradeoff

P15 文件系统

在逻辑卷上又提供一层「文件」的抽象能力
比较熟悉，skip之。

P16 磁盘阵列 - storage node 独立



前一台端；内外接口灵活性

人们把不论红绿的硬件层面生成的虚拟磁盘
均称为“LUN”；软件层面生成的则称“Volume”(卷)

· 双控制器 ⇒ “脑裂”问题 (split brain)

P170

SAN (Storage Area Network)

独立的磁盘阵列配备了自己复杂的控制器，
成为了一个半独立于主机 cluster 之外的存储
cluster，即称 SAN。

P186

Fibre channel，略

以太网交换网络 \Rightarrow 扩展性 \uparrow

面向无连接 \rightsquigarrow 面向有连接

+ 资源高效复用 - 资源浪费，维护成本

- 数据丢失，鲁棒性 - 可用性低

+ 保障到达

串行在高速底层连接上 $>$ 并行传输

P228

SAS (Serial-Attached SCSI)

Speed Contest

(平均速度)

• 原始 HBA

: SCSI 1.2Gbps

• FC ~ SAS

: SAS 1.5Gbps

• 固态 SSD

: 6G ~ 8Gbps

: SAS 3Gbps

ref

• IB 已上 40Gbps

• PCIe 4.0 16Gbps

P284

NFS 协议 \Rightarrow NAS (Network Attached Storage)

SNMP over FC $\xrightarrow[\text{cost}]{\text{performance}}$ NAS over Ethernet

DAS 即 原始的 Directly Attached Storage

P310

阶段总结，值得一看

P384

Virtualization 简谈

带内 (In-band) v.s. 带外 (Out-band)

控制指令与数据
共享线路

e.g., IP headers, NFS

控制指令通过独立

e.g., Bus, SANFS

P414

集群目标

高可用性 HAC

负载均衡 LBC

高性能 HPL

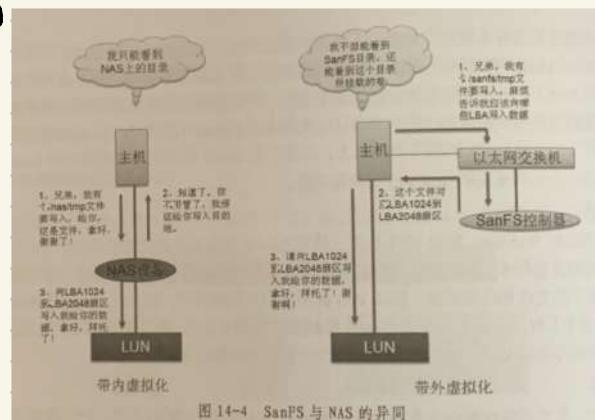


图 14-4 SanFS 与 NAS 的异同

P445

Scale-up vs. Scale-out



单机内堆砌
更多更高的
算力

集群化，互取更多
较弱的单机
※

P446

DSS 对象存储：数据本地化的外置存储

P486

数据保护

· 块级保护 — 磁盘镜像

· 文件级保护 — 文件复制备份

快照 (Snapshot) 技术
+ clone

① 复制 metadata

(保证备份过程短，可用性高)

Row +
/ Cw

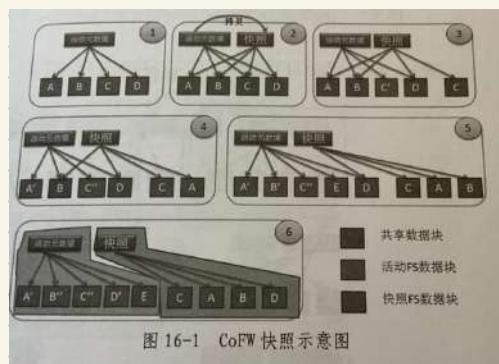


图 16-1 CoFS 快照示意图

P490

FS concepts:

Bitmap, inode, ---

P522

连续数据保护 CDP

P576

备份系统总结

备份目标

备份通路

备份引擎

备份策略

P610

“容灾”不仅包括安全的备份，亦包括将备份的数据尽快恢复生产状态。（同步？）

同步 vs 异步复制

P678

数据管理 | 数据存储 的软硬分离

Data Cooker

包括 snapshot, clone

包括

分层

Recovery (sync)

硬件产品

Virtualization

--

Over subscription

De provision

De duplication (去重)

--

Tiering

应用级 - 文件级 - 块级，无Tradeoff

P748

IO 路径：

Bypass FS

APP

初始化逻辑，Networking, --

文件操作

OS/FS

同异步，缓存，Pagecache, TLB, --

Lock (Synchronization) 问题

Block I/O

LUN 管理

Memory map, RAID, --

mapping

Drivers

In-device queuing, --

设备操作

设备

Hardware device optimizations

(可看出各级均有缓存问题 tradeoff)

P806

IO 性能诊断！值得一看并借鉴，定位 Bottleneck 在哪里。

P88H

对“云”的定义显得过于宽泛了，
个人认为缺溢些反而好，不炒作概念。

胜在灵活、节约、共用

弱在安全、性能(网络带宽之限)

P91H

分、合 的系统很到位

后记

一些补充和展望，值得一看。(2019)