

并处理机分配给新高优先级进程。
短作业优先，周转时间较短，基本没有饥饿

Process synchronization

临界资源：一次仅允许一个进程使用资源
Entry section：检查是否可进入临界区；
Critical：访问临界资源的代码
Exit：访问临界区标志的清除；
Remainder：代码中剩余部分；
同步：直接制约关系
互斥：异步制约关系
临界区问题的解决必须满足三个要求：**忙则等待 (mutual exclusion)**；**空闲让进 (progress)**；**有限等待 (bounded waiting)**。
让权等待非必须 (不要忙等)。

Peterson 算法 (不饥饿)

只用于两个进程的情况，并且假设 load 和 store 是原子操作，是一种软件解决方法。

```
The two processes share two variables:
• int turn;
• boolean flag[2];
The variable turn indicates whose turn it is to enter the critical section.
The flag array is used to indicate if a process is ready to enter the critical section.
flag[i] = true implies that process P_i is ready!
Process P_i
P_i:
do {
  flag[i] = true;
  turn = i;
  while (flag[i] and turn == i);
  critical section
  flag[i] = false;
  remainder section
} while (1);
P_j:
do {
  flag[j] = true;
  turn = j;
  while (flag[j] and turn == j);
  critical section
  flag[j] = false;
  remainder section
} while (1);
```

Meets all three requirements; solves the critical-section problem for two processes!
现代 OS 不适用(编译器优化可能交换指令)

Bakery Algorithm (面包房算法)

```
choosing[i]: 进程 i 是否正在获取排队号
number[i]: 进程 i 的排队号, 0 表示不需要
do {
  choosing[i] = true;
  number[i] = max(number[0], number[1], ..., number [n - 1]) + 1;
  choosing[i] = false;
  for (j = 0; j < n; j++)
  { while (choosing[j]);
    while (number[j] != 0 && (number[j][i] < (number[j][j])));
  }
  critical section
  number[i] = 0;
  remainder section
} while (1);
```

硬件同步方法

1. 中断屏蔽：在临界区前关闭开启中断。
多处理器：**Memory barriers**(an instruction forcing any change in memory to be propagated (made visible) to all other processors)

2. 硬件指令方法

Test_and_set 存在 busy waiting, lock 变量共享 (swap)，存在忙等问题

```
boolean TestAndSet(boolean &target) {
  while (1) {
    boolean rv = target;
    target = true;
    return rv;
  }
};
while (1) {
  key = TRUE;
  while (key == TRUE)
    Swap(&lock, &key);
  critical section
  lock = false;
  remainder section
}
```

硬件方法优点：进程数任意，支持多个临界区，简单易验证正确性；缺点：无法让权等待，可能饥饿，可能死锁。
解决饥饿：

```
Boolean waiting(n); lock; // to be initialize to false, waiting[]: 排队队列
while(1) {
  waiting[i]=true;
  key= true;
  while (waiting[i] && key) key=TestAndSet(lock);
  waiting[i]=false;
  critical section;
  j=(i+1)%n;
  while (j != i && !waiting[j]) j=(j+1)%n;
  if (j == i) lock = false;
  else waiting[j] = false;
  remainder section;
}
```

自旋锁 spinlock：适用于多处理器，当一个进程欲访问已被其他进程锁定的资源时，进程循环检测该锁是否被释放
软件 (bad)
单标志法：不同进程只能交替进入
双标志法先检查：不同进程可能同时进入

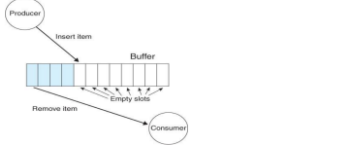
```
P_i进程:
while(flag[i]){
  flag[i]=TRUE;
  critical section;
  flag[i]=FALSE;
  remainder section;
}
P_j进程:
while(flag[j]){
  flag[j]=TRUE;
  critical section;
  flag[j]=FALSE;
  remainder section;
}
```

semaphores

```
两个操作 wait/P/down 申请和 signal/V/up 释放。
信号量分为计数信号量，整型信号量(有忙等的)，二值信号量 (互斥锁(mutex locks))，记录型信号量；
wait(semaphore *S){
  S->value--;
  if (S->value <= 0) {
    add this process to S->list;
    block();
  }
}
signal(semaphore *S){
  S->value++;
  if (S->value <= 0) {
    remove a process P from S->list;
    wakeup(P);
  }
}
```

具有忙等的信号量值非负，**解决忙等**如上信号量可以为负，负数的绝对值代表等待该信号量的进程数，0 代表无资源可用。Block()用于自我阻塞，wake_up()用于进程唤醒，**Wait 和 signal 成对出现**，互斥要求在同一进程出现，同步要求在不同进程。**连续的 wait 需要注意顺序**。同步 wait 和互斥 wait 相邻时，要先同步 wait，不然可能会死锁
前驱关系：利用同步的信号量思想，控制语句间的前驱顺序
优先级倒置(priority inversion)：当优先级较低的进程持有较高优先级进程所需的锁时的调度问题。**优先级继承**：将等待进程中的最高优先级临时分配给持有锁的进程。

Bounded-Buffer Problem

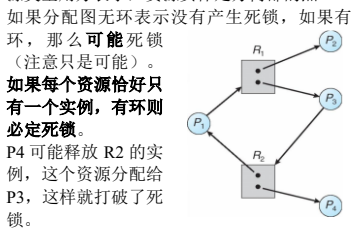


```
do {
  produce an item in nextp
  wait(empty);
  wait(mutex);
  add nextp to buffer
  signal(mutex);
  signal(nextp);
} while (1);
do {
  wait(full);
  wait(mutex);
  remove an item from buffer to nextc
  signal(mutex);
  consume the item in nextc
} while (1);
do {
  wait(mutex);
  readcount++;
  if (readcount == 1)
    writing is performed
  wait(mutex);
  while (flag[i]){
    flag[i]=TRUE;
    critical section;
    flag[i]=FALSE;
    remainder section;
  }
} while (TRUE);
```

Dining-Philosophers Problem 哲学家进餐
典型的同步问题，问题描述：N 个哲学家坐在圆桌，每个哲学家和邻居共享一根筷子；哲学家吃饭要用身边的两只筷子一起吃；邻居不允许同时吃饭；哲学家只会思考或者吃饭。可能解决 (死锁/饥饿) 方案：最多只允许 4 个哲学家坐在桌上/临界区内必须同时拿起两根筷子/使用非对称的解决方法：奇数先拿左手，偶数先拿右手；加入吃饭时碗的限制。

Deadlock 死锁

四个必要条件：
Mutual exclusion (互斥条件)
hold and wait (请求并保持条件)
No preemption (不剥夺条件)
circular wait (循环等待)
资源分配图，由点 V 和边 E 组成，V 被分为两部分：系统活动进程的集合/系统所有资源类型的集合。进程 P_i 到资源 R_j 的有向边记为 P_i->R_j，表示进程 P_i 正在申请资源类型 R_j 的一个实例，表示请求边；资源 R_j 到进程 P_i 的有向边表示资源类型 R_j 的一个实例已经分配给进程 P_i，表示分配边。进程用圆表示，资源类型用方表示，资源实体是方内部的点。如果分配图无环表示没有产生死锁，如果有环，那么可能死锁 (注意只是可能)。
如果一个资源恰好只有一个实例，有环则必定死锁。
P4 可能释放 R2 的实例，这个资源分配给 P3，这样就打破了死锁。



死锁处理

保证系统不进入死锁：prevention/avoidance 允许进入死锁但可恢复：detection/recovery。Unix/Linux/Windows 三个系统都不考虑死锁，

产生死锁时再进行处理。
死锁预防 Prevention
通过破坏死锁产生的必要条件来预防死锁。Mutual exclusion：非共享资源必须互斥，而共享资源不需要互斥，也不导致死锁。Hold and wait：**采用预先静态分配**，进程在执行前就要申请并获得所有需要的资源。缺点：低资源利用率、可能饥饿。不破坏循环等待) No preemption：如果一个进程占有资源并且申请了另一个不能立即分配的资源，那么它现已分配的资源都可以被抢占，即被隐式释放了。
Circular wait：**采用顺序资源分配法**，给资源设置显式序号，请求必须按照资源序号递增的方式进行。

死锁避免

要求每个进程声明它**可能需要的每种类型的资源的最大数量**。死锁避免算法动态检查资源分配状态，确保不会出现循环等待。资源分配状态由可用和已分配资源的数量以及进程的最大需求定义。
安全状态：对于所有进程，如果存在一个安全序列，那么系统就处于安全状态。对于进程序列 P1,P2,...,Pn，如果对于每一个 Pi,Pi 仍然可以申请的资源数小于当前可用的资源加上所有进程 Pj(j>i)所占有的资源，那么这一序列是安全序列。
安全状态->没有死锁：不安全状态->可能有死锁：避免->保证系统永远不进入非安全状态。
资源分配图, single instance 死锁避免算法：
Single instance: 每种资源只有一个引入一种新的 claim edge 需求边, Pi->Rj 表示进程 Pi 在未来可能请求资源 Rj, 用虚线表示。当进程真正请求资源时，用请求边覆盖掉需求边。当资源被分配给进程后，用分配边来覆盖掉请求边，当资源被释放后，分配边恢复为需求边。

假设进程 Pi 申请资源 Rj。只有在需求边 Pi ->Rj 变成分配边 Rj->Pi 而不会导致资源分配图形成环时，才允许申请。
Banker, 银行家算法：
每个进程实现说明最大需求：进程请求资源时可能会等待；进程拿到资源后必须在有限时间内释放它们。
N 进程数，m 资源类型的种类数；
Available 现有的 Max: 最大需求 Allocation: 当前分配 Need: 当前需求
没有操作系统使用 banker 死锁避免。
安全状态检测算法：

1. 设 work 和 finish 分别是长度为 m 和 n 的向量，初始化：work=available, finish[i]=false; 算法需要 m*n*n 的操作数量级确定系统状态
资源请求算法：*
死锁检测
允许系统进入死锁状态的话，那么系统就需要提供检测算法和恢复算法。
死锁定理：S 为死锁状态的充分条件是，尚且仅当 S 状态的资源分配图是不可完全简化的
单实体资源检测算法还是使用**资源分配图**有环来判断是否死锁
多实体资源类型检测算法 (类似银行家算法)：
数据结构：Available, allocation 是一样的，request: n*m 的矩阵，表示当前各进程的资源请求状况，request[i][j]=k 表示 P_i 正在请求 k 个资源 R_j。

1. 设 work 和 finish 分别是长度为 m 和 n 的向量，初始化：work=available, 如果 allocation[i] 非 0, finish[i]=false 否则初始化为 true;
2. 寻找 i 满足 finish[i]=false 且 request[i]<=work, 如果 i 不存在跳到第四步;
3. work=work+allocation[i], finish[i]=true, 返回第二步;
4. 如果某个 finish 是 false, 那么系统处于死锁状态，且对应下标的进程 P_i 死锁。
算法需要 m*n*n 的操作数量级确定系统状态
死锁检测算法的应用
检测算法的调用时刻及频率取决于：死锁发生频率及以思索发生时受影响的进程数。如果经常发生死锁，那么就要经常调用检测。如果在不确定的时间调用检测算法，资源图可能有很多环，通常不能确定哪些造成了死锁

死锁恢复 (解除)

进程终止 (撤销进程)
两种方法来恢复死锁：终止所有死锁进程(一次终止一个进程直到不死锁。许多因素影响终止进程的选择：优先级/进程已经计算了多久/还要多久完成/进程使用了哪些类型的资源等等)
抢占资源 (资源剥夺法)
选一个进程挂起：代价尽量最小化；
回滚：回退到安全状态；
饥饿：确保被挂起的进程不会长时间得不到资源而产生饥饿
进程回退法
让一个或多个进程回退到足以回避死锁的状态，回退时进程自愿释放资源

Main Memory 内存

编译：将用户源代码编译成若干目标模块链接：将编译后的一组目标模块以及需要的库函数链接在一起 (**逻辑地址的产生**)
装入：程序装入内存运行
静态链接：编译后所有目标模块都是从 0 开始编址，需要修改相对地址使多个模块的地址具有一致把外部调用符号转变为相对地址
装入时动态链接：边装入边链接，便于实现目标模块的共享、修改与更新
运行时动态链接：未被用到的目标模块，都不会被调入内存和链接到模块上
绝对装入：只适合单道程序环境，编译程序将产生绝对地址的代码
可重定位装入 (静态重定位)：地址信息可以产生一定的偏移；作业装入内存时，要求分配全部完整的内存空间，程序在内存中不能移动
动态运行装入 (动态重定位)：支持程序在内存中移动，将地址转换推迟到程序真正需要运行时进行；需要借助重定位寄存器实现功能
层次存储中主存 cache 寄存器为 volatic 易失的。逻辑地址/虚地址/相对地址：由 CPU 生成，首地址为 0，逻辑地址无法在内存中读取信息。物理地址/实地址/绝对地址：内存中存储单元的地址，可以直接寻址。物理地址中的逻辑地址空间是通过一对**基址寄存器 (base register)**和**界限地址寄存器 (limit register)**控制。
Memory-Management Unit (MMU)
就是将虚拟地址映射到物理地址的硬件设备。在 MMU 中，base 寄存器叫做重定位寄存器 (relocation register)，用户进程送到内存前，

都要加上重定位寄存器的值。PA = relocation register + LA。用户程序只能处理 LA，永远看不到真的 PA。
Dynamic Loading (动态加载)
进程大小会受到物理内存大小的限制，为了更好的空间使用率，采用动态加载。一个子程序只有在调用时才被加载，所有子程序都可有重定位的形式存储在磁盘上，需要的时候装入内存中。
Swapping (交换技术)
进程可以暂时从内存中交换到备份存储 backing store 上，当需要再次执行时再调回。需要动态重定位 dynamic relocate 备份内存：是快速硬盘，可以容纳所有用户的所有内存映像，并为这些内存映像提供直接访问，如 Linux 交换区 windows 的交换文件 pagefile.sys
Roll out roll in: 如果有一个更高优先级的进程需要服务，内存交换出低优先级的进程以便装入和执行高优先级进程；介导程序在辅存与内存之间的交换。
交换时间的主要部分是转移时间 transfer time。总转移时间与所交换的内存大小成正比。系统维护一个就绪的可立即运行的进程队列，并在磁盘上有内存映像。
正常情况下，禁止交换；当空闲内存低于某个阈值时，启用交换块出；当空闲内存增加至一定数量时，停止换出

Contiguous Allocation (连续分配)

单一连续分配
内存通常分为两个区域：一个驻留 resident 操作系统，一个用于用户进程 (仅有一道用户程序)，由于中断向量一般位于低内存，所以 **OS 也放在低内存**。
重定位寄存器用于保护各个用户进程以及 OS 的代码和数据不被修改。Base 是 PA 的最小值；limit 包含了 LA 的范围，每个 LA 不能超过 Limit。MMU 地址映射是**动态的**。
Multiple-partition allocation (可能外部碎片也可能内部)：分区式管理将内存划分为多个连续区域叫做分区，每个分区放一个进程。
固定分区：又分为分区大小相等、分区大小不等，产生内部碎片
动态分区 (可变分区分配)：
动态划分内存，在程序装入内存时切出一个连续的区域 hole 分配给进程，分区大小恰好符合需要。往往产生外部碎片。
操作系统需要维护一个表，记录哪些内存可用哪些已用。从一组可用的 hole 选择一个空闲 hole 的常用算法 first/best/worst/next-fit 四种：First: 分配到第一个足够大的；best (最容易产生碎片)：分配到最小的足够大的；worst: 分配到最大的；next: 从上次分配的地址开始查找的 first-fit。Best 和 worst 都要搜索整个 list，除非按 size 排序好了。First 和 best 在时间和空间利用率都比 worst 好。

紧凑技术：对系统中的磁盘碎片进行整理。
碎片 fragmentation
first 和 best 都存在外部碎片的问题。外碎片指所有的总可用内存可以满足请求，但是并不连续。内碎片是进程内部无法使用的内存，这是由于零头和块大小造成的。
分页存储管理
分页允许进程的 PA 空间非连续；将物理内存分为固定大小的块，叫做**帧 frame/物理块/页框**，将逻辑内存也分为同样大小的块叫做**页**

块)，不断套娃。
空闲盘块分配时，从第一个成组链块开始分配，指针逐一下移直至第一个成组链块的空闲盘块全部分配完毕，再分配第二个成组链块。
盘块的回收，记录回收盘号，不断上移指针，若当前成组链块已满，则回收盘块充当前一个成组链块。

页面缓冲 page buffer

将文件数据作为页而不是磁盘块缓冲起来到虚存。

Unified Buffer Cache

A unified buffer cache uses the same page cache to cache both memory-mapped pages and ordinary file system I/O

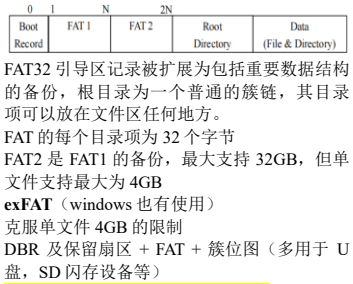
恢复 Recovery

一致性检查，将目录结构数据与磁盘数据块比较，发现并纠正不一致。
用系统程序将磁盘数据备份到另一个设备。然后从该设备恢复。

日志结构的文件系统

日志文件系统记录文件系统的更新为事务。事务会被写到日志里。事务一旦写入日志就是已经 commit 了，否则文件系统则还未更新。

FAT32 磁盘结构



FAT32 引导区记录被扩展为包括重要数据结构的备份，根目录为一个普通的簇链，其目录项可以放在文件区任何地方。
FAT的每个目录项为 32 个字节
FAT2 是 FAT1 的备份，最大支持 32GB，但单文件支持最大为 4GB
exFAT (windows 也有使用) 克服单文件 4GB 的限制
DBR 及保留扇区 = FAT + 簇位图 (多用于 U 盘, SD 闪存设备等)

NTFS file system (windows 使用)

引导扇区	主文件表(MFT)	文件区	MFT 前 16 个无数据文件备份	文件区

ReFS 文件系统

旨在最大限度地提高数据可用性、跨不同的工作负荷高效地扩展到大型数据集，并提供数据完整性，使其能够恢复损坏。旨在解决存储方案的扩展问题。

功能	ReFS	NTFS
最大文件名称长度	255 个 Unicode 字符	255 个 Unicode 字符
最大路径名称长度	32K Unicode 字符	32K Unicode 字符
文件大小上限	35 PB (pb)	256 TB
最大大小	35 PB	256 TB

文件系统集合

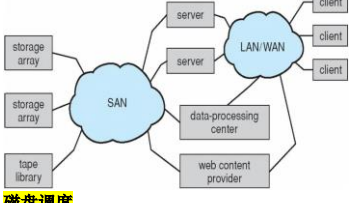
- FAT (MS-DOS 文件系统)
- FAT32 (VFAT) , exFAT (64位)
- NTFS (NT 文件系统)
- ReFS (Resilient File System) 复原文件系统 (Windows 8/6, server) S51K/S52K (AT&T UNIX 系统)
- ext (minix 文件系统)
- ext2, ext3, ext4 (linux 文件系统, Android)
- proc. sysfs (linux 虚拟文件系统)
- yaofs (Yet Another Flash File System)
- ReiserFS (Linux 一种日志文件系统)
- HPFS (OS/2 高性能文件系统)
- HFS (BSD UNIX 的一种文件系统)
- HFS+ (Mac OS, iOS 文件系统)
- iso9660 (通用的光盘文件系统)
- NFS (网络文件系统)
- VFS (Linux 虚拟文件系统)
- VFS 是物理文件系统与服务之间的一个接口，它屏蔽各类文件系统的差异，给用户和程序提供一个统一的接口
- ZFS (Open Solaris 文件系统)
- LTFS (Linear Tape File System, 线性磁带文件系统) 为磁带提供了一种通用、开放的文件系统。
- APFS (Apple File System) 苹果新一代的文件系统，MAC OS 10.13 以后，iOS 10.3 以后

Mass storage system 大容量存储

Host controller in computer uses bus to talk to disk controller built into drive or storage array
磁盘的 0 扇区是最外面的磁道的第一个扇区。逻辑块时最小传出单位 512B
Platter 表面被分成了很多环形 track (磁道)，再细分为 sector (扇区)。在一个 arm position 下的 track 组成一个 cylinder (柱面)。
Disk Attachment
三种方式：DAS(Direct/Host Attached Storage) /NAS(Network Attached Storage 网络附加存储)/SAN(Storage-Area Network 存储区域网)
Host-attached storage accessed through I/O ports talking to I/O busses
Network-attached storage (NAS) is storage made available over a network rather than over a local connection (such as a bus)



SAN: Multiple hosts attached to multiple storage arrays – flexible



磁盘调度

Position time = accesstime= seek time + rotational latency
IO time = position time + transfer time+control delay
Seek time 寻道时间，磁头移动到包含目标扇区的柱面的时间。
旋转延迟 rotational latency: 旋转到目标扇区的时间 (average latency=1/2*latency=1/2*1/rpm *60s, 即平均情况下，旋转半圈，磁盘的转速一般以 round per minute(rpm) 表示)。
传输时间 transfer time: 数据传输时间。
FCFS 先来先服务：算法公平，不会存在饥饿。
SSTF (shortest seek time first) 最短寻道时间优先：处理靠近当前磁头位置的请求，本质上和 SJF 一样，有可能一些请求会永远无服务，

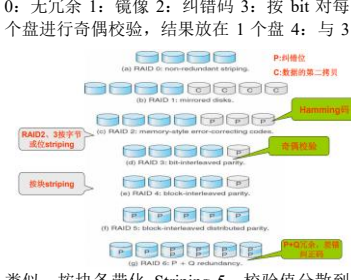
时间非最优。
SCAN (电梯调度算法)：从磁盘一端到另一端服务请求，对所有路上经过的柱面进行处理。到达另一端时改变移动方向，继续处理。
C-SCAN 磁头从一端移动到另一端，到了另一端就马上返回到磁盘开始，返回路径中不服务。
LOOK: 磁头从一端到另一端，到达另一端最远的服务就不继续向前，开始折返服务。
C-LOOK: 磁头从一端到另一端，到达另一端最远服务就立即返回到磁盘开始的第一个服务，返回路径不服务。
调度算法选择：
SSTF 一般来说比较好，通用而自然。
SCAN C-SCAN 对于高负荷 IO 磁盘表现更好。
SSTF 和 LOOK 都是默认算法的合理选择
减少延迟时间: 给用户提供一种通号, 错位命名

磁盘管理

低级格式化 (物理格式化、磁盘初始化)：将磁盘划分为扇区 (使用特殊的物理结构)，头部尾部添加一些磁盘控制器使用信息。会保留一些块作为备用 (扇区备用)。
分区：将磁盘的若干柱面分区，每个分区的起始扇区和大小记录于主引导记录的分区表中。

逻辑 (高级) 格式化: 创建文件系统。引导扇区产生在这里。
启动块 Boot Block: 启动块初始化系统，引导 (Bootstrap 自举)程序存储在 ROM 中，引导程序装载程序。
坏块 Bad Block: 处理方法:format, chkdsk 指令
交换空间管理
Windows: 存在 pagefile.sys 文件中，简单但低效 (外部碎片)
Linux: 独立的磁盘分区 SWAP 分区 (内部碎片)

RAID



0: 无冗余 1: 镜像 2: 纠错码 3: 按 bit 对每个盘进行奇偶校验, 结果放在 1 个盘 4: 与 3 类似，按块条带化 Striping 5: 校验值分散到各个盘 6: P+Q 冗余，差纠错码

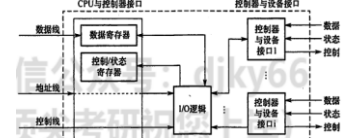
三级存储 Tertiary storage device

Low cost is the defining characteristic of tertiary storage. Generally, tertiary storage is built using removable media Common examples of removable media are floppy disks and CD-ROMs; other types are available
Swap space 三级存储 Tertiary storage device 虚存使用硬盘空间作为主存。两种形式：普通文件系统：win 都是 pagefile.sys 独立硬盘分区 linux solaris 都是 swap 分区。还有一种方法：建立在 raw 的磁盘分区上，这种速度最快。
性能: Sustained bandwidth 大传输的平均速率 字节/时间。Effective bandwidth IO 时间下的平均速

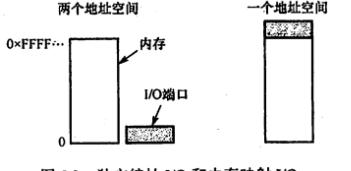
率。前者是数据真正流动时的速率，后者是驱动器能够提供的能力，一般驱动器带宽指前者。
NVM Devices 包括固态硬盘 (Solid-state disks, SSD), USB drives (thumb drive, flash drive), DRAM disk replacements 等。
NVM devices 比 HDD 更可靠，也更贵，可能寿命更短，容量更小，但是快很多。标准总线可能会太慢，所以有的 SSD 设计成直接连接到系统总线，如 PCI。
没有需要「移动」的部分，因此没有 seek time or rotational latency。

IO system

IO 接口：串行、并行、USB、键盘、硬盘
IO 端口：设备与计算机通信的连接点。端口地址。I/O 端口通常有四种寄存器，即状态、控制、数据输入与数据输出寄存器。
I/O 方式：程序 IO，中断 IO (包括同步/异步)，DMA 方式，通道方式。
I/O 设备设备分类
块设备：信息交换以数据块为单位 (结构设备，磁盘 (使用 DMA I/O) 等)；可寻址，可随机读写。
字符设备：信息交换以字节为单位 (无结构设备，鼠标键盘打印机等)，不可寻址常使用中断 I/O。(多为独占设备，互斥共享)
I/O 接口 (设备控制器) 位于设备与 CPU 之间
设备控制器与 CPU 接口：数据线 (数据寄存器/状态寄存器)、地址线、状态线。
设备控制器与设备的接口：一个设备控制器可以连接多个设备，也有数据、地址、状态三种类型信号。
I/O 逻辑：用于实现对设备的控制，对 CPU 的 I/O 命令进行译码并选择相应设备进行控制。



I/O 端口：设备控制器中可被 CPU 直接访问的寄存器。
数据寄存器：CPU 与设备间的数据缓冲。
状态寄存器：执行结果/设备状态寄存器。
控制寄存器：由 CPU 写入控制命令。
编址方式



独立编址 I/O 和内存映射 I/O
独立编址：为每个端口分配 I/O 端口号，所有端口形成 I/O 端口空间。
统一编址 (内存映射 I/O)：每个端口分配唯一的内存地址，通常靠近地址空间顶端。
IO 控制方式
Polling 程序直接控制方式 (字为单位)：CPU 对外设状态进行轮询 polling，直到确定

该字已经在 I/O 控制器的数据寄存器中。CPU 利用率相当低。
中断驱动方式 (字为单位)：用户程序→系统调用处理程序→设备驱动程序→中断处理程序，中断处理后结束后存放内容在内存缓冲区
DMA (direct memory access, 块为单位) 用于避免程序 I/O for large data movement, 需要 DMA controller, 绕过 CPU 来直接在 IO 设备和内存之间传输数据
传送的数据，直接从设备进入内存，只需在传送一个/多个数据块的开始/结束时，CPU 进行干预处理即可。
任务分配给 DMA 控制器，在 DMA 开始传输时，主机向内存中写入 DMA 命令块，然后 DMA 自己操作内存总线，直接向内存进行传输。传输结束后，DMA 发送中断信号给处理器。
其中 DMA 的 4 类寄存器
命令状态寄存器 (CR)：接受 CPU 的 I/O 指令以及控制信息与设备状态。
内存地址寄存器 (MAR)：输入时存放内存起始目标地址；输出时存放设备的内存源地址。
数据寄存器 (DR)：暂存内存与设备间的数据。
数据计数器 (DC)：存放本次传送的字节数。
IO 分类: block I/O(read, write, seek); character I/O (stream, keyboard, clock); memory-mapped file access; network sockets
I/O 软件层级结构
分为四层：中断处理程序，设备驱动程序，与设备无关的操作系统软件，以及用户软件 (指用户空间的 IO 软件)
用户层 I/O 软件：与用户交互的接口 (I/O 库函数)

设备独立性 (设备无关性) 软件：用于实现用户程序与设备驱动器的统一接口、设备命令、设备的保护及设备的分配与释放等，提供必要的存储空间
可以实现逻辑设备名到物理设备名的映射，便于实现重用，增加设备分配的灵活性
设备驱动程序：每类设备配置一个设备驱动程序，通常以**进程**的形式存在；设备驱动程序向上层用户程序提供一组标准接口，设备具体的差别被设备驱动程序所封装 (上层同一操作可操纵不同的硬件设备)

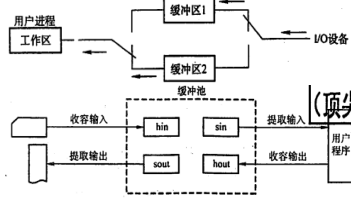
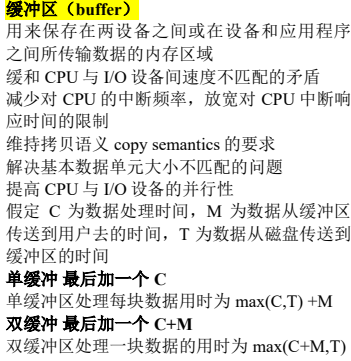
中断处理程序：用于保存被中断进程的 CPU 环境，转入相应的中断处理程序进行处理，处理完毕再恢复被中断进程的现场后，返回到被中断进程。
最后是硬件设备进行相应的操作。
时钟和定时器 Clocks and Timers 提供获取当前时间/获取已经逝去的时间/设置定时器以在 T 时触发操作 X 三种函数。
测量逝去时间与触发器操作的硬件称为可编程间隔定时器 (programmable interval timer)

内核 IO 子系统

有关服务：IO 调度、Buffering、caching、spooling 虚拟化、device reservation、error handling
负责：文件和设备命名空间的管理，文件和设备访问控制，操作控制 (for example,a modern cannot seek(), 文件系统空间的分配, 设备分配, 缓冲、高速缓存、假脱机, I/O 调

度, 设备状态监控、错误处理、失败恢复, 设备驱动程序的配置和初始化
设备独立性软件
用户程序的设备独立性是：用户程序不直接使用物理设备名 (或设备的物理地址)，而只能使用逻辑设备名；而系统在实际执行时，将逻辑设备名转换为某个具体的物理设备名，实施 I/O 操作。
I/O 软件的设备独立性是：除了直接及设备打交道的低层软件之外，其他部分的软件并不依赖于硬件。I/O 软件独立于设备，就可以提高设备管理软件的设计效率。
逻辑设备是实际物理设备属性的抽象，它并不限于某个具体设备。

Disk cache (磁盘高速缓存)：利用内存中的存储空间来暂存从磁盘中读到的一系列盘块中的信息。逻辑上属于磁盘，物理上是驻留在内存中的盘块。
两种形式：在内存中开辟一个单独的空间作为磁盘高速缓存；把未利用的内存空间作为一个缓冲区。
缓冲区 (buffer) 用来保存在两设备之间或在设备和应用程序之间所传输数据的内存区域
缓冲 CPU 与 I/O 设备间速度不匹配的矛盾 减少对 CPU 的中断频率，放宽对 CPU 中断响应时间的限制
维持拷贝语义 copy semantics 的要求 解决基本数据单元大小不匹配的问题 提高 CPU 与 I/O 设备的并行性 假定 C 为数据缓冲区，M 为数据从缓冲区传送到用户去的时间，T 为数据从磁盘传送到缓冲区的时间



高速缓存与缓冲区的相同点与对比

高 速 缓 存		
相 同 点		都介于高
区别	存放的数据	存放的是低速设备上的某些数据的复制数据，即高速缓存上有的，低速设备上未必有
	目的	高速缓存存放的是高速设备经常访问的数据，若高速设备要访问的数据不在高速缓存中，则高速设备就要访问低速设备

缓冲区

高速设备和低速设备之间
存放的是低速设备传向高速设备的数据 (或相反)，而这些数据在低速设备 (或高速设备) 上却不一定有备份，这些数据再从缓冲区传送到高速设备 (或低速设备)

设备的分配与回收
每当进程向系统提出 I/O 请求时，设备驱动程序按照一定的策略，把其所需的设备及其有关资源 (如缓冲区、控制器和通道) 分配给该进程。

