





该图展示了 Linux 进程控制块 (PCB) 的状态转换逻辑。中心是一个名为“拥有CPU”的方框。从该框出发，有两条向上指向“等待资源”方框的箭头，分别标注为“等待资源 并且不想被打断”和“等待资源”。从“拥有CPU”框出发，有三条向下的箭头：一条指向“继续执行”方框，一条指向“TASK\_STOPPED TASK\_TRACED”方框（标注为“进程被调试跟踪”），以及一条指向“进程退出”方框（标注为“进程退出”）。“继续执行”方框有一个指向“拥有CPU”框的反馈箭头。从“进程退出”框，有两条指向“EXIT\_ZOMBIE”和“EXIT\_DEAD”框的箭头。在图的左下角，有文字“Linux 进程控制块 PCB=task\_struct”。







明这个 inode 对应的对象可以被删除。删除磁盘的数据块，磁盘的 inode 以及 VFS 的 inode。  
put\_super(): 由于当前的文件系统的卸载而释放当前的超级块对象。  
write\_super(): 更新当前的超级块对象的内容。  
stats(): 返回当前 mount 的文件系统的一些统计信息  
remount\_fs(): 按照一定的选项重新 mount 文件系统  
clear\_inode() 和 put\_inode 类似，但是也删除包含数据在内的内存对 inode 中的结构。  
umount\_begin(): 开始 umount 操作，并中断它的 mount 操作，用于网络文件系统  
物理文件系统的 inode 在外存中并且是长期存在的， VFS 的 inode 对象在内存中，它仅在需要时才建立，不再需要时撤消。物理文件系统的 inode 是静态的，而 VFS 的 inode 是动态结构

```
struct inode {
    struct list_head i_hash; /* inode hash 链表指针 */
    struct list_head i_list; /* inode 链表指针 */
    struct list_head i_dentry; /* dentry 链表 */
    kdev_t i_dev; /* 主设备号 */
    unsigned long i_ino; /* 外存的 inode 号 */
    umode_t i_mode; /* 文件类型和访问权限 */
    nlink_t i_nlink; /* 该文件的链接数 */
    uid_t i_uid; /* 文件所有者的用户标识 */
    gid_t i_gid; /* 文件的所有用户组标识 */
    kdev_t i_dev; /* 次设备号 */
    off_t i_size; /* 文件长度，以字节为单位 */
    time_t i_atime; /* 文件最后一次访问时间 */
    time_t i_mtime; /* 文件最后一次修改时间 */
    time_t i_ctime; /* 文件创建时间 */
    unsigned long i_blksize; /* 块尺寸，以字节为单位 */
    unsigned long i_blocks; /* 文件的块数 */
    unsigned long i_version; /* 文件版本号 */
    unsigned long i_npages; /* 文件在内存中占用的页面数 */

    struct semaphore i_sem; /* 文件同步操作作用的信号量 */
    struct inode_operations *i_op; /* 指向 inode 操作函数入口表的指针 */
    struct super_block *i_sb; /* 指向该文件系统的 VFS 超级块 */
    struct wait_queue *i_wait; /* 文件同步操作等待队列 */
    struct file_lock *i_flock; /* 指向文件锁定链表的指针 */
    struct vm_area_struct *i_mmap; /* 文件使用的虚存区域 */
    struct page **i_pages; /* 指向文件占用内存页面结构体链表 */

    unsigned long i_count; /* 使用该 inode 的进程计数 */
    unsigned short i_flags; /* 该文件系统的超级块标志 */
    unsigned short i_writect; /* 写计数 */
    unsigned char i_lock; /* 对该 inode 的锁定标志 */
    unsigned char i_dirt; /* 该 inode 的修改标志 */
    unsigned char i_pipe; /* 该 inode 表示管道文件 */
    unsigned char i_sock; /* 该 inode 表示套接字 */
    unsigned char i_seek; /* 未使用 */
    unsigned char i_update; /* inode 更新标志 */
    unsigned char i_conدمanded;

}
```

VFS 的 inode 与某个文件的对应关系是通过设备号 i\_dev 与 inode 号 i\_ino 建立的，它们唯一地指定了某个设备上的一个文件或目录。

VFS 的 inode 是物理设备上的文件或目录的 inode 在内存中的映像。这些特有信息是各种文件系统的 inode 在内存中的映像。如 EXT2 的 ext2\_inode\_info 结构。

i\_lock 表示该 inode 被锁定，禁止对它访问。i\_flock 表示该 inode 对应的文件被锁定。i\_flock 是指向 file\_lock 结构链表的指针，该链表指出了一系列被锁定的文件。

VFS 的 inode 组成一个双向链表，全局变量 first\_inode 指向链表的表头。在这个链表中，空闲的 inode 总是从表头加入，而占用的 inode 总是从表尾加入。

系统还设置了一些管理 inode 对象的全局变量，如：  
max\_inodes 给定 inode 的最大数量，  
nr\_inodes 表示当前使用的 inode 数量，  
nr\_free\_inodes 表示空闲的 inode 数量。

### inode 操作

create: 只适用于目录 inode，当 VFS 需要在“inode”里面创建一个文件(文件名在 dentry 里面给出)的时候被调用。VFS 必须已经检查过文件名在这个目录里不存在。

lookup: 用于检查一个文件(文件名在 dentry 里面给出)是否在一个 inode 目录里面。

link: 在 inode 所给出的目录里面创建一个从第一个参数 dentry 文件到第三个参数 dentry 文件的硬链接(hard link)。

unlink: 从 inode 目录里面删除 dentry 所代表的文件。  
symlink: 用于在 inode 目录里面创建软链接(soft link)。  
mkdir: 用于在 inode 目录里面创建子目录。  
rmdir: 用于在 inode 目录里面删除子目录。  
mknod: 用于在 inode 目录里面创建设备文件。  
rename: 把第一个和第二个参数(inode, dentry)所定位的文件改名成第三个和第四个参数所定位的文件。  
readlink: 读取一个软链接所指向的文件名。  
follow\_link: VFS 调用这个函数跟踪一个软链接到它所指向的 inode。  
put\_link: VFS 调用这个函数释放 follow\_link 分配的一些资源。  
truncate: VFS 调用这个函数改变一个文件的大小。  
permission: VFS 调用这个函数得到对一个文件的访问权限。  
setattr: VFS 调用这个函数设置一个文件的属性。比如 chmod 系统调用就是调用这个函数。  
getattr: 查看一个文件的属性。比如 stat 系统调用就是调用这个函数。  
setattr: 设置一个文件的某项特殊属性。详细情况请查看 setattr 系统调用帮助。  
getattr: 查看一个文件的某项特殊属性。详细情况请查看 getattr 系统调用帮助。  
listxattr: 查看一个文件的所有特殊属性。详细情况请查看 listxattr 系统调用帮助。  
removexattr: 删除一个文件的特殊属性。详细情况请查看 removexattr 系统调用帮助。

### 目录项对象 dentry object

每个文件除了有一个索引节点 inode 数据结构外，还有一个目录项 dentry 数据结构。每个 dentry 代表路径中的一个组成部分。如：/，bin，vi 都属于目录项对象。目录项也可包括安装点，如：/mnt/cdrom/foo，/，mnt，cdrom，foo 都属于目录项对象。inode 结构代表的是物理意义上的文件，记录的是物理上的属性，对于一个具体的文件系统，其 inode 结构在磁盘上就有对应的映像。一个索引节点对象可能对应多个目录项对象。目录项对象作用是帮助实现文件的快速定位，还起到缓冲作用

```
struct dentry {
    atomic_t d_count; /* 目录项引用计数器 */
    unsigned int d_flags; /* 目录项标志 */
    struct inode *d_inode; /* 与文件名关联的索引节点 */
    struct inode *d_parent; /* 父目录的目录项 */
    struct list_head d_hash; /* 目录项形成的哈希表 */
    struct list_head d_lru; /* 未使用的 LRU 链表 */
    struct list_head d_child; /* 父目录的子目录项所形成的链表 */

    struct list_head d_subdirs; /* 该目录项的子目录所形成的链表 */
    struct list_head d_alias; /* 索引节点别名的链表 */
    int d_mounted; /* 目录项的安装点 */
    struct qstr d_name; /* 目录项名(可快速查找) */
    struct dentry_operations *d_op; /* 操作目录项的函数 */
    struct super_block *d_sb; /* 目录项树的根(即文件的超级块) */
    unsigned long d_vfs_flags; /* 具体文件系统的数据 */
    void *d_fsdata; /* 具体文件系统的数据 */
    unsigned char d_iname[DNAME_INLINE_LEN]; /* 短文件名 */
    .....
}
```

对目录项进行操作的一组函数，由 d\_op 指向 dentry\_operation 结构。  
struct dentry\_operations {  
d\_revalidate(): 判定目录项是否有效。  
d\_hash(): 生成一个哈希值。  
d\_compare(): 比较两个文件名。  
d\_delete(): 删除 d\_count 值为 0 的目录项对象。  
d\_release(): 释放一个目录项对象。  
d\_iput(): 调用该方法丢弃目录项对应的索引节点。

### VFS 的 dentry cache 与 inode cache

为了加速对经常使用的目录的访问，VFS 文件系统维护着一个目录项的缓存。  
为了加快文件的查找速度 VFS 文件系统维护一个 inode 节点的缓存以加速对所有装配的文件系统的访问。  
用 hash 表将缓存对象组织起来。

### File 对象

文件对象 file 表示进程已打开的文件，只有当文件被打开时才在内存中建立 file 对象的内容。  
该对象由相应的 open()系统调用创建，由 close()系统调用销毁。

```
struct file {
    struct list_head f_list; /* file 结构链表 */
    struct dentry *f_dentry; /* 指向与文件对象关联的 dentry 对象 */
    struct vfsmount *f_vfsmnt; /* 文件相应的 vfsmount 结构 */
    struct file_operations *f_op; /* 文件对象的操作集合 */
    atomic_t f_count; /* 文件打开的引用计数 */
    unsigned int f_flags; /* 使用 open() 时设置的标志 */
```

```
mode_t f_mode; /* 文件读写权限 */
loff_t f_pos; /* 对文件读写操作的前置位置 */
struct fown_struct f_owner;
.....
};
file_operations
lseek: 用于移动文件内部偏移量。
read: 读文件。
aio_read: 异步读，被 aio_submit 和其他的异步 IO 函数调用。
write: 写文件。
aio_write: 异步写，被 aio_submit 和其他的异步 IO 函数调用。
readerr: 当 VFS 需要读目录内容的时候调用这个函数。
poll: 当一个进程想检查一个文件是否有内容可读的时候，VFS 调用这个函数；一般来说，调用这个函数之后进程进入睡眠，直到文件中有内容读写就绪时被唤醒。详情请参考 select 和 poll 系统调用。
ioctl: 被系统调用 ioctl 调用。
unlocked_ioctl: 被系统调用 ioctl 调用；不需要 BKL(内核锁)的文件系统应该使用这个函数，而不是上面那个 ioctl。
compat_ioctl: 被系统调用 ioctl 调用；当在 64 位内核上使用 32 位系统调用时使用这个 ioctl 函数。
mmap: 被系统调用 mmap 调用。
open: 通过创建一个新的文件对象而打开一个文件，并把它链接到相应的索引节点对象。
flush: 被系统调用 close 调用，把一个文件内容写回磁盘。
release: 当对一个打开文件的最后引用关闭的时候，VFS 调用这个函数释放文件。
fsync: 被系统调用 fsync 调用。
fasync: 当对一个文件启用异步读写(非阻塞读写)的时候，被系统调用 fcntl 调用。
lock: fcntl 系统调用使用命令 F_GETLK, F_SETLK 和 F_SETLKW 的时候，调用这个函数。
```

```
struct file_system_type { /* 文件系统注册链表
const char *name; /* 文件系统名称
int fs_flags; /* 文件系统类型 flags
struct super_block * (*get_sb) (struct file_system_type *, int, const char *, void *) /* Method for reading a superblock
void (*kill_sb) (struct super_block *) /* Method for removing a superblock
struct module *owner; /* Pointer to the module implementing the filesystem
struct file_system_type * next; /* Pointer to the next element in the list of filesystem types
struct list_head fs_supers; /* Head of a list of superblock objects having the same filesystem type */
};
name: 文件系统类型名字，比如“ext2”，“vfat”等等。
fs_flags: mount 的文件系统的参数。
get_sb: 当这种类型的文件系统要被 mount 的时候，这个函数会被调用，用以得到相应文件系统的超级块。
kill_sb: 当这种类型的文件系统被 umount 的时候，这个函数被调用。
owner: VFS 内部使用，大多数情况下，你只需要初始化为 THIS_MODULE。
next: 文件系统类型链表的后继指针。VFS 内部使用，初始化为 NULL。
list_head fs_supers: 文件系统的超级块的双向链表。
```

### 文件的 open()操作

open()系统调用内核函数是 sys\_open()。  
第一个参数是打开文件的路径名。第二个参数是文件的访问标志。  
sys\_open():  
1. Invokes getname() to read the file pathname from the process address space.  
2. 用 get\_unused\_fd() 在 current->files->fd 所指向的文件对象指针数组中查找一个未使用的文件号，存储在局部变量 fd 中。

3. 调用 filp\_open() 函数，工作主要分两部分：  
第一步：调用 open\_name() 函数，找到目标节点(可以是文件、目录)所对应的 dentry 对象，与 dentry 对象相对应的 inode 对象此时也应该在物理内存中。  
第二步：调用 dentry\_open() 函数，该函数申请一个 file 对象的结构，然后初始化该对象，其中的第一步使 file 对象的 f\_dentry 指向已获得的 dentry 对象。  
4. 调用 fd\_install() 函数，将文件对象装入当前进程的打开文件表：current->files->fd[fd] = file; 然后返回文件号 fd。最重要的步骤是 open\_name() 完成的。函数的执行过程：  
a. 确定从哪一个 dentry 对象出发进行路径解析。根据指定文件路径名是相对路径还是绝对路径从 current 中得到相应的 dentry 对象。  
b. 调用 path\_walk() 函数进行路径解析，该函数执行一个循环，每一循环都是得到一个 dentry 对象，该对象对应文件路径的一个子路径。

### read()的实现

read() 在内核中的对应函数为 sys\_read()。

(1)根据文件描述符调用函数 fget()，找到相应的文件对象 file。  
(2)检查 file->f\_mode 是否允许读。  
(3)调用 file->f\_op->read() 函数执行读的操作。对于大部分文件系统实际是 generic\_file\_read() 函数。该函数根据文件位置和要读出的长度确定相应的页面，然后再检查该页面是否在 pagecache 中存在，如果不存在，就要调用 inode 节点的 l\_mapping->a\_ops->readpage() 方法，将其从磁盘读入。不同磁盘文件系统的 readpage 方法不同，ext2 文件系统相应的函数为 ext2\_readpage。  
为了提高性能，读操作采用了预读机制。

do\_page\_fault() 工作原理: compares the linear address that caused the Page Fault against the memory regions of the current process; it can thus determine the proper way to handle the exception  
Do\_page\_fault():

### 实验

#### 编译内核:

make clean 删除大多数的编译生成文件，但是会保留内核的配置文件.config，还有足够的编译支持来建立扩展模块  
make mrproper 删除所有的编译生成文件，还有内核配置文件，再加上各种备份文件  
make distclean mrproper 删除的文件，加上编辑备份文件和一些补丁文件。  
apt-get install kernel-package libncurses5-dev fakeroot wget bz2zip  
//安装工具包  
make config 是有问必答的方式，每个内核选项它都会问你，要、不要、模块，选错了一个就必须从头再来一遍；  
make menuconfig 提供了一个基于文本的图形界面，它依赖于 ncurses5 这个包，键鼠操作，可以修改选项，一般推荐用这个；  
make xconfig 需要你要有 x window system 支持，就是说你要在 KDE、GNOME 之类的 X 桌面环境下才可用，好处是支持鼠标，坏处是 X 本身占用系统周期，而且 X 环境容易引起编译器的不稳定  
make -j4 启动 4 个线程(双核)来编译内核文件生成 o 等中间文件  
内核文件 bzImage 的位置在 /usr/src/linux/arch/i386/boot 目录下。  
make modules\_install 安装模块  
make install 使用命令 make install 将 bzImage 和 System.map 拷贝到 boot 目录下。这样，Linux 在系统引导从 /boot 目录下读取内核映像到内存中

#### 添加系统调用:

system\_call()函数实现了系统调用中断处理程序：  
1.它首先把系统调用和该异常处理程序用到的所有 CPU 寄存器 保存到相应的栈中， SAVE\_ALL  
2.保存进进程 task\_struct (thread\_info) 结构的地址存放在 ebx 中  
3.对用户态进程传进来的系统调用号进行有效性检查。若调用号大于或等于 NR\_syscalls，系统调用处理程序返回。(sys\_call\_table)  
4.若系统调用号无效，函数就把 ENOSYS 值存放在栈中 eax 寄存器所在的单元，再跳到 ret\_from\_sys\_call()  
5.根据 eax 中所包含的系统调用号调用对应的特定服务例程

实验修改的主要有 3 处地方:

```
for (p = &init_task; (p = next_task(p)) != &init_task;)
//遍历进程
p->comm //comm 类型为 char[16],代表进程名
p->pid //当前进程号
p->state //当前进程的状态
-1 unrunnable, 0 runnable, >0 stopped
p->parent //指向父进程 task_struct 的地址
```

#### 添加文件系统:

```
/* Structure of a directory entry */
#define EXT2_NAME_LEN 255
/* The new version of the directory entry. Since EXT2 structures are stored in intel byte order, and the name_len field could never be bigger than 255 chars, it's safe to reclaim the extra byte for the file_type field. */
struct ext2_dir_entry_2 {
    u32 inode; /* Inode number */
    u16 rec_len; /* Directory entry length */
    u8 name_len; /* Name length */
    u8 file_type;
    char name[EXT2_NAME_LEN]; /* File name */
};
/* Ext2 directory file types. Only the low 3 bits are used. The other bits are reserved for now. */
enum {
    EXT2_FT_UNKNOWN,
```

```
EXT2_FT_REG_FILE,
EXT2_FT_DIR,
EXT2_FT_CHRDEV,
EXT2_FT_BLKDEV,
EXT2_FT_FIFO,
EXT2_FT_SOCK,
EXT2_FT_SYMLINK,
EXT2_FT_MAX
};
文件类型:
普通文件(文件名不超过 255)
目录文件
字符设备文件和块设备文件:
fd0 (for floppy drive)
hda (for harddisk a)
lp0 (for line printer 0)
tty (for teletype terminal)
管道(FIFO)文件
链接文件
socket 文件
文件系统分三大类:
基于磁盘的文件系统, 如 ext2/ext3/ext4、VFAT、NTFS 等。
网络文件系统, 如 NFS 等。
特殊文件系统, 如 proc 文件系统、devfs、sysfs (/sys) 等。
#define NR_OPEN 256
struct files_struct {
    int count; /* 该文件结构的计数值 */
    fd_set open_on_exec;
    fd_set open_fds;
    struct file *fd[NR_OPEN];
};
fd[] 每个元素是一个指向 file 结构体的指针，该数组组为进程打开文件表。
进程打开一个文件时，建立一个 file 结构体，并加入到系统打开文件表中，然后把该 file 结构体的首地址写入 fd[] 数组的第一个空闲元素中，一个进程所有打开的文件都记载在 fd[] 数组中。
```

dd: 用指定大小的块拷贝一个文件，并在拷贝的同时进行指定的转换  
命令语法: dd [选项]  
常用参数:  
if = 输入文件(或设备名称)  
of = 输出文件(或设备名称)  
bs = bytes 同时设置读/写缓冲区的字节数(等于设置 ibs 和 obs)  
count=blocks 只拷贝输入的 blocks 块  
conv = ucase 把字母由小写转换为大写  
conv = lcase 把字母由大写转换为小写。  
例: dd if=/dev/zero of=/myfs bs=1M count=1 /dev/zero; 零设备“0”  
/dev/loop: loopback device (回环设备、或虚拟设备) 是指用文件来模拟块设备

ext2 的操作都在结构 ext2\_dir\_inode\_operations 中定义  
fs/ext2/namei.c 下定义了具体实现方法

man -S2 open#选择第二个 section  
1 用户命令，2 系统调用，3 语言数据库调用，4 设备和网络界面，5 文件格式，6 游戏和示范，troff 的环境，7 表格和宏，8 关于系统维护的命令

### Crtl 快捷键

```
<Backspace>或<Ctrl-H>删除前一个字符
<Ctrl-U>删除当前行
<Ctrl-C>终止现在的命令，终止一个前台进程使用<Ctrl-Z>挂起一个前台进程
<Ctrl+D>退出当前的 shell，必须从登陆 shell 退出，必须关闭所有的 shell
<Ctrl-K>删除一行光标后字符
<Ctrl-P>上一次执行的命令，扫描过的不会再次出现
操作系统原理
线程 相对优先级
THREAD_PRIORITY_TIME_CRITICAL THREAD_PRIORITY_HIGHEST +2
THREAD_PRIORITY_ABOVE_NORMAL +1
THREAD_PRIORITY_NORMAL 0
THREAD_PRIORITY_BELOW_NORMAL -1
THREAD_PRIORITY_LOWEST -2
THREAD_PRIORITY_IDLE
基本优先级
REALTIME_PRIORITY_CLASS，基本优先权为 24。
HIGH_PRIORITY_CLASS，基本优先权为 13。
ABOVE_NORMAL_PRIORITY_CLASS，基本优先权为 10
//Windows NT 和 Windows Me/98/95: This value is not supported.
NORMAL_PRIORITY_CLASS，基本优先权为 8。
BELOW_NORMAL_PRIORITY_CLASS，基本优先权为 6
//Windows NT 和 Windows Me/98/95: This value is not supported.
IDLE_PRIORITY_CLASS，基本优先权为 4。
```

System Call 3 种传参方法，1. Register 2. Pass address to register 3. Parameter pushed to stack, popped by system  
ch1  
- The one program running at all times on the computer” is the kernel(内核)。  
- Efficient Fair Convenient  
What is the difference between kernel mode and user mode?  
Why is the difference important to an operating system?  
In kernel mode, the executing code has complete and unrestricted access to the underlying hardware. It can execute any CPU instruction and reference any memory address. In user mode, the executing code has no ability to directly access hardware or reference memory. Code running in user mode must delegate to system APIs to access hardware or memory.  
The difference rather protect the computer system resources, while preventing from errant users.  
bootstrap program is loaded at power-up or reboot  
- Typically stored in ROM or EEPROM, generally known as firmware  
- Initializes all aspects of system  
- Loads operating system kernel and starts execution  
批处理系统(batch System)  
注意同一时间间隔(并发)和同一时刻(并行)的区别。在多道程序环境下，一段时间内，宏观上有多道程序在同时执行，而在每一时刻，单道程序环境下实际只能有一道程序执行，故微观上这些程序还是在分时地交替执行。操作系统的非发性是通过分时得以实现的。  
- 多任务并发  
lack of interaction  
分时操作系统(Time-Sharing Systems ——Linux a multi-user time-sharing system)  
所谓分时技术就是把处理器的运行时间分成很短的时间片，按时间片轮流把处理器分配给各联机作业使用。若某个作业在分配给它的时间内不能完成其计算，则该作业暂时停止运行，把处理器让给其他作业使用，等待下一轮再继续运行。由于计算机速度很快，作业运行轮转很快，给每个用户的感受好像是自己独占一台计算机。  
实时操作系统  
实时操作系统的主要特点是及时性和可靠性。  
④  
运行机制  
- 时钟管理 中断机制  
- 系统控制的数据结构及处理  
- 系统中用来登记状态信息的数据结构很多，比如作业控制块、进程控制块(PCB)、设备控制块、各类链表、消息队列、缓冲区、空闲区登记表、内存分配表等。为了实现有效的管理，系统需要一些基本的操作，常见的操作有以下三种：  
- 进程管理：进程状态管理、进程调度和分派、创建与撤销进程控制块等。  
- 存储器管理：存储器的空间分配和回收、内存信息保护程序、代码对换程序等。  
- 设备管理：缓冲区管理、设备分配和回收等。

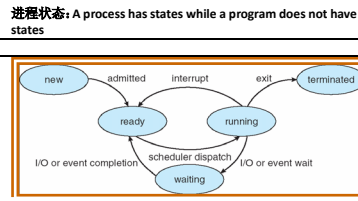
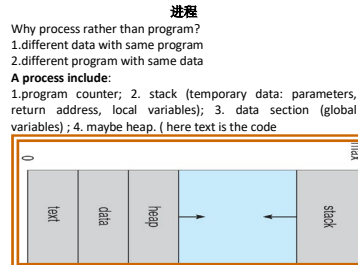
中断  
- 外中断(中断)：I/O 中断、时钟中断  
- 内中断(异常)：系统调用(陷入)、缺页异常、断点指令、其他程序性异常(如算数溢出)  
- 自愿中断(指令中断)、强迫中断(硬件故障、软件中断)  
中断处理: The operating system preserves the state of the CPU by storing registers and the program counter.  
从用户态转换为核心态的唯一途径是中断或异常  
由用户态进入核心态，不仅仅是状态需要切换。而且，所使用的堆栈也可能需要由用户堆栈切换为系统堆栈，但这个系统堆栈也是属于该进程的。  
I/O 方式: 程序 I/O (Programmed I/O) ④ 中断 I/O (Interrupt I/O) ④ 同步 I/O 和异步 I/O  
DMA 方式  
通道方式  
Synchronous I/O(同步 I/O): After I/O starts, control returns to user program only upon I/O completion.  
\* Wait for I/O completion, tow ways ④ wait instruction idles the CPU until the next interrupt wait loop (loop: jmp loop).  
\* At most one I/O request is outstanding at a time, no simultaneous I/O processing.  
Asynchronous I/O (异步 I/O): After I/O starts, control returns to user program without waiting for I/O completion.  
System call 系统调用 - request to the operating system to allow user to wait for I/O completion.  
Device-state table contains entry for each I/O device indicating its type, address, and state.  
Operating system indexes into I/O device table to determine device status and to modify table entry to include interrupt.  
DMA  
Used for high-speed I/O devices able to transmit information at close to memory speeds.  
\* Device controller transfers blocks of data from buffer storage

directly to main memory without CPU intervention.  
\* Only on interrupt is generated per block, rather than the one interrupt per byte. ❸  
System Calls :Programming interface to the services provided by the OS  
The function implementation is inside the OS, or we name it the OS kernel.  
Those functions are usually packed inside an object called the library file.  
Three general methods used to pass parameters to the OS  
- Pass the parameters in registers  
- Parameters stored in a block, or table, in memory, and address of block passed as a parameter in a register  
- This approach taken by Linux and Solaris  
- Parameters placed, or pushed, onto the stack by the program and popped off the stack by the operating system

Types of System Calls  
- Process control  
- File management  
- Device management  
- Information maintenance  
- Communication: Provide the mechanism for creating virtual connections among processes, users, and computer systems  
Simple Structure 简单结构  
Layered Approach 层次化结构  
Layered OS is more efficient than monolithic OS(X)错误, 巨型内核直接访问更加高效

Microkernel System Structure (Windows NT ... Windows 8、Windows 10、Mac OS)  
优点  
- Easier to extend a microkernel  
- Easier to port the operating system to new architectures  
- More reliable (less code is running in kernel mode)  
- More secure  
缺点  
- Performance overhead of user space to kernel space communication 费时间  
Monolithic Kernel Structure - packed into a single file (OS/360, VMS and Linux)  
优点  
- highly efficient because of direct communication between components  
缺点  
- Difficult to isolate source of bugs and other errors  
- Hard to modify and maintain  
- Kernel gets bigger as the OS develops.  
Module (Linux , Solaris )  
Most modern operating systems implement kernel modules  
- Uses object-oriented approach  
- Each core component is separate  
- Each talks to the others over known interfaces  
- Each is loadable as needed within the kernel  
Virtual Machines 虚拟机  
为了在通用操作系统管理下的计算机上运行一个程序, 需要经历几个步骤。但是, 不一定需要向操作系统预定运行时间  
需要, 用控制台监控程序执行过程, 将程序装入内存, 确定起始地址, 并从这个地址开始执行

进程  
Why process rather than program?  
1.different data with same program  
2.different program with same data  
A process include:  
1.program counter; 2. stack (temporary data: parameters, return address, local variables); 3. data section (global variables); 4. maybe heap. (here text is the code)



PCB 的内容:  
1.Process state  
2.Program counter  
3.CPU registers  
4.CPU scheduling information  
5.Memory-management information  
6.Accounting information  
7.I/O status information  
8.page table or relocation register and limit register  
9.file open table

调度队列的种类  
Processes migrate among the various queues  
1. Job queue – set of all processes in the system  
2. ready queue – set of all processes residing in main memory, ready and waiting to execute. Generally stored as a linked list  
3. device queues – set of processes waiting for a particular I/O device, eg. a tape driver, a disk

调度器 scheduler 分类  
long-term scheduler (or job scheduler) selects which processes should be brought into the ready queue  
short-term scheduler (or CPU scheduler) selects which process should be executed next and allocates CPU  
Attention: The long-term scheduler controls the degree of multiprogramming

Context-switch  
Context-switch time is overhead; the system does no useful work while switching.  
Process Termination  
Process executes last statement and asks the operating system to delete it (exit)  
1. Output data from child to parent (via wait)  
2. Process' resources are de-allocated by operating system  
Parent may terminate execution of children processes (abort)  
1. Child has exceeded allocated resources  
3. If parent is exiting  
Some operating system do not allow child to continue if its parent terminates--All children terminated - cascading termination  
Advantages of process cooperation  
Information sharing, Computation speed-up, Modularity and Convenience.  
Interprocess Communication  
1.shared memory and 2.message passing

进程同步种类: blocking (synchronous) vs. nonblocking (asynchronous)  
1.blocking send: the sending process is blocked until the message is received by the receiving process or by the mailbox  
2.nonblocking send: the sending process sends the message, and resumes operation  
3.blocking receive: the receiver blocks until a message is available  
4.nonblocking receive: the receiver retrieves either a valid message or a null, and resumes operation

生产者-消费者问题 (指利用了 n-1 个单元)  
item buffer[BUFFER\_SIZE];  
int in = 0;  
int out = 0;  
// producer  
item nextProduced; // local variable  
while (1) {  
/\* produce an item in nextProduced \*/  
while (((in + 1) % BUFFER\_SIZE) == out) ; // do nothing  
buffer[in] = nextProduced;  
in = (in + 1) % BUFFER\_SIZE;  
}  
// consumer  
item nextConsumed; // local variable  
while (1) {  
while (in == out) ; // buffer is empty  
nextConsumed = buffer[out];  
out = (out + 1) % BUFFER\_SIZE;  
/\* consume the item in nextConsumed \*/  
}

线程  
Benefits  
1.Responsiveness (e.g. thread as GUI engine)  
2.Resource Sharing (e.g. shared variable)  
3.Economy (e.g. save memory)  
4.Utilization of multiprocessor architectures  
Thread can be a basic scheduling unit, but not one of resource allocation.  
用户线程和内核线程  
- 用户级线程: 不依赖于 OS 核心 (内核不了解用户线程的存在), 应用进程利用 \*\*线程库\*\*提供创建、同步、调度和管理线程的函数来控制用户线程。

内核级线程: 依赖于 OS 核心, 由内核的内部需求\*\*进行创建和撤销, 用来执行一个指定的函数。一个线程发起系统调用而阻塞  
- 用户线程的维护由应用进程完成; \*\*内核不了解用户线程的存在; \*\*  
\*\*用户线程切换不需要内核特权; \*\*  
- 用户线程调度算法可针对应用优化:  
\*\*一个线程发起系统调用而阻塞, 则整个进程在等待。  
\*\*  
- 内核维护进程和线程的上下文信息: \*\*  
\*\*线程切换由内核完成; \*\*  
\*\*时间片分配给线程, ▲所以多线程的进程获得更多 CPU 时间; \*\*  
\*\*▲一个线程发起系统调用而阻塞, 不会影响其他线程的运行  
device, eg. a tape driver, a disk

synchronous signals - 1. an illegal memory access, a division by zero 2. delivered to the same process that cause the signal  
asynchronous signals  
1. a user keystroke (Ctrl-C), a timer expiration  
2. typically sent to another process  
Multithreading models  
Many to one:  
1.Many user-level threads mapped to single kernel thread  
2.thread management is done in user space, used on systems that do not support kernel threads  
3.drawbacks: whole process block if one thread blocks. unable to run parallel on multiprocessors  
One to one:  
1.Each user-level thread maps to a kernel thread: more concurrency than many-to-one. support multiprocessors  
2.drawback: overhead of creating kernel threads  
3.examples: Windows 95/98/NT/2000, OS/2

Many to many:  
1.multiplexes many user level threads to a smaller or equal number of kernel threads  
2. allows the programmer to create a sufficient number of user threads  
3. avoid bad blocking, support multiprocessors  
Examples: Solaris 2, Windows NT/2000 with the ThreadFiber package  
Two-level Model -- M:M + 1:1  
Similar to M:M, except that it allows a user thread to be bound to kernel thread.  
Examples: IRIX, HP-UX, Tru64 UNIX, Solaris 8 and earlier.

线程的终止: thread cancellation  
1.target thread: the thread to be cancelled  
2.asynchronous cancellation: one thread immediately terminates the target thread, may be cancelled in the middle of updating data shared with other threads, may not free a system-wide resource  
3.deferred cancellation: the target thread can periodically check if it should terminate (at so called cancellation points in pthread)

线程池  
Why the thread pools?  
1. avoid creation and termination overhead, so that it is faster to service a request  
2. put bound on number of threads, thus limit CPU and memory usage  
Considerations :  
1. number of CPUs  
2. amount of physical memory  
3. expected number of concurrent requests  
fork may cp all/current, exec replace whole.  
CPU 调度  
Maximum CPU utilization obtained with multiprogramming  
CPU scheduling  
Decisions may take place when a process, but not limited to:  
1.Switches from running to waiting state  
2.Switches from running to ready state  
3.Switches from waiting to ready  
4.Terminates  
Scheduling under 1 and 4 is nonpreemptive.  
All other scheduling is preemptive.

Dispatcher  
Dispatcher module gives control of the CPU to the process selected by the short-term scheduler; this involves:  
1.switching context  
2.switching to user mode  
3.jumping to the proper location in the user program to restart that program  
Dispatch latency – time it takes for the dispatcher to stop one process and start another running.  
CPU 调度的评判标准 Scheduling Criteria  
1.CPU utilization – keep CPU as busy as possible  
2.Throughput – number of processes that complete their execution per time unit

Turnaround time 周转时间– the interval from the time of submission of a process to the time of completion  
Waiting time – amount of time a process has been waiting in the ready queue  
Response time – amount of time it takes from when a request was submitted until the first response  
CPU 调度算法  
FCFS: 等待时间不稳定, 响应时间不稳定.  
特点: easy to understand, easy to implement (an FIFO queue), large variance of waiting/response time, convoy effect, nonpreemptive, bad for time-sharing system  
SJF: 分为 preemptive(SRTF) & nonpreemptive. 平均等待时间最短.  
Nonpreemptive 会把当前正在执行的推掉. 最优但不实际, 是其他算法的判定标准. 会造成 starvation(饥饿)  
Approximation: predict the next burst length.  
$$r_{n+1} = \alpha I_n + (1 - \alpha) r_n$$
  
Each successive term has less weight than its predecessor.  
Highest response Ratio Next, 最高响应比优先  
$$HRR = (W + T) / T$$
 W: Waiting Time T: Burst Cycle Time  
Priority: preemptive & nonpreemptive. 存在饥饿的问题, 可以通过 aging 的方法来弥补.  
Round Robin: q large  $\Rightarrow$  FIFO, q small  $\Rightarrow$  q must be large with respect to context switch, otherwise overhead is too high. No process waits more than (n-1)q time units. higher average turnaround than SJF, but better response. Typical time quanta: 10 to 100 ms. 80% of the CPU bursts should be shorter than the quantum.  
other  
Multilevel Queue: scheduling between the queues, commonly implemented as fixed-priority preemptive.  
Each queue has its own scheduling algorithm.  
foreground – RR background – FCFS  
Multilevel Feedback-Queue: a process that waits too long is moved to a higher-priority queue. Prevents starvation.  
Concurrent access to shared data may result in data inconsistency. Maintaining data consistency requires mechanisms to ensure the orderly execution of cooperating processes  
Atomic operation means an operation that completes in its entirety without interruption.  
临界区问题的解决  
1.Mutual Exclusion. If process P<sub>i</sub> is executing in its critical section, then no other processes can be executing in their critical sections.  
2.Progress. If no process is executing in its critical section and there exist some processes that wish to enter their critical section, then the selection of the process that will enter the critical section next cannot be postponed indefinitely.  
3.Bounded Waiting. A bound must exist on the number of times that other processes are allowed(被允许而不一定要进入) to enter their critical sections after a process, as made a request to enter its critical section and before that request is granted.

算法 3: Peterson's Solution  
do {  
flag[i] = true;  
turn = i+1;  
while (flag [1-i] and turn == 1-i);  
/\* Critical Section  
flag[i] = false;  
Remainder Section  
}  
while (1);  
分析: 满足了互斥、有空让进和有限等待。这个算法能够解决临界区问题。这里,turn 的含义是“当前的‘进入权’给了谁”。而 Flag[i]=ture 的含义为, P<sub>i</sub> 已经准备好进入自己的临界区了。  
Bakery Algorithm -- n processes  
boolean choosing[n] = {false};  
int number[n] = {0};  
// 某一个线程  
do  
choosing[i] = true; //表示进程 i 正在获取它的排队登记号, 保证计算 number[i] 的原子操作  
number[i] = max(number[0], number[1], ..., number[n - 1]) + 1; //是进程 i 的当前排队登记号。如果值为 0, 表示进程 i 未参加排队, 不想获得该资源。  
choosing[i] = false; // 获取排队登记号完成  
for(j = 0; j < n; j++)  
{  
while(choosing[j]); // 确保比较的原子性  
while((number[j] != 0) && (number[j],j) < number[i],j)); // a < c or (a = c and b < c) 理解 (a,b) = ticket#, process id# order  
}

}  
// critical section  
number[i] = 0; // 结束业务。退出排队  
// remainder section  
while(1);  
Special atomic hardware instructions  
基于 TestAndSet  
boolean TestAndSet(boolean \*target) {  
boolean rv = \*target;  
\*target = true;  
return rv;  
}  
boolean lock = false 对于第 i 个进程:  
do {  
while (TestAndSet(&lock));  
/\* Critical Section  
lock = false;  
Remainder Section  
}  
while(1);  
分析: 满足了互斥, 但不满足有限等待条件。这里某个号数的进程可以一直霸占进入临界区的机会, 使得其他号数的进程无限等待。  
基于 Swap 的算法:  
boolean lock = false 对于第 i 个进程  
do {  
key = true;  
while (key == true) Swap(lock,key);  
/\* Critical Section  
lock = false;  
Remainder Section  
}  
while(1);  
分析: 满足了互斥, 但不满足有限等待条件。这里某个号数的进程可以一直霸占进入临界区的机会, 使得其他号数的进程无限等待。  
Bounded-waiting mutual exclusion with TestAndSet()  
boolean waiting[n] = false;  
boolean lock = false;  
do {  
waiting[i] = true;  
key = true;  
while (waiting[i] & key) keys=TestAndSet(lock);  
waiting[i] = false;  
/\*critical section;  
j = (i+1) % n;  
while ((j == i) & !waiting[j]) j = (i+1) % n;  
if (j == i) lock = false;  
else waiting[j] = false;  
//remainder section;  
}  
while(1)

Semaphores 基于信号量的算法:  
把一个进程 P 屏蔽, 从就绪队列转移到 L 上去: wait(S)  
S.value--;  
if (S.value < 0){  
//add this process to S.List;  
block;  
}  
把一个进程 P 唤醒, 从 L 转移到就绪队列中:  
signal(S)  
S.value++;  
if (S.value <= 0){  
//remove a process P from S.List;  
wakeup(P);  
}  
S.value>0 表示有 S 个资源可用;  
S.value=0 表示无资源可用或表示不允许进程再进入临界区;  
S.value<0 则 |S.value| 表示在等待队列中进程的个数或表示等待进入临界区的进程个数。  
wait(S)表示申请一个资源; signal(S)表示释放一个资源。  
如果两个 wait 操作相邻, 那么它们的顺序至关重要, 而两个相邻的 signal 操作的顺序无关紧要。一个同步 wait 操作与一个互斥 wait 操作在一起时, 同步 wait 操作在互斥 wait 操作前。

Bounded-Buffer Problem  
有限缓冲区问题: 进程共享一块有界的缓冲区。如何实现安全同步?  
采用信号量 full, empty, mutex。  
初始值 full = 0, empty = n, mutex = 1。  
wait(full/EMPTY); CON/PRO  
wait(mutex);  
.....  
signal(mutex);  
signal(empty/FULL); PRO/CON  
.....  
while (1);

Reader's-Writers Problem  
mutex = 1, wrt = 1  
//readers  
readcount = 0  
wait(mutex);  
readcount++;  
if (readcount == 1) wait(r);  
signal(mutex);  
//reading is performed  
wait(mutex);  
readcount--;  
if (readcount == 0) signal(wrt);  
signal(mutex);  
//writers  
wait(wrt);  
//writing is performed  
signal(wrt);  
Dining-Philosophers Problem  
semaphore chopstick[5] = 1.  
对于第 i 哲学家:  
do {  
wait(chopstick[i]);  
wait(chopstick[(i+1) % 5]);  
eat;  
signal(chopstick[i]);  
signal(chopstick[(i+1) % 5]);  
think;  
} while (1);  
其他解决方法:  
1.allow at most four philosophers to be sitting simultaneously at the table--a spare resource.  
2.allow a philosopher to pick up chopsticks only if both are available, and pick up them simultaneously.  
3.use a asymmetric solution: an odd philosopher picks up first the left chopstick and then the right chopstick, whereas an even one picks up first the right and then the left.  
Deadlock and Starvation  
Deadlock – two or more processes are waiting indefinitely for an event that can be caused by only one of the waiting processes  
Starvation – indefinite blocking. A process may never be removed from the semaphore queue in which it is suspended  
死锁问题  
A set of blocked processes each holding a resource and waiting to acquire a resource held by another process in the set.  
1.Mutual exclusion: only one process at a time can use a resource.  
2.Hold and wait: a process holding at least one resource is waiting to acquire additional resources held by other processes.  
3. No preemption: a resource can be released only voluntarily by the process holding it, after that process has completed its task.  
4. Circular wait: there exists a set {P<sub>1</sub>, P<sub>2</sub>, ..., P<sub>n</sub>} of waiting processes such that P<sub>1</sub> is waiting for a resource that is held by P<sub>2</sub>, P<sub>2</sub> is waiting for a resource that is held by P<sub>3</sub>, ..., P<sub>n-1</sub> is waiting for a resource that is held by P<sub>n</sub>, and P<sub>n</sub> is waiting for a resource that is held by P<sub>1</sub>.

Handling Deadlocks  
1.Prevention 低设备使用率和系统吞吐率  
条件 2: 必须同时申请所有资源(之后不 wait); 或只有没有资源才可申请(之前不 hold)。资源利用率低 饥饿  
条件 3: 申请其它资源不能立即得到, 那么现有资源可被抢占  
适用 cpu 内存 不适用打印机 磁带驱动器  
条件 4: 资源排序, 递增申请  
2. Avoidance  
3. Detection  
4.Recovery  
5.Ignore the problem and pretend that deadlocks never occur in the system (Unix)  
Resource-Allocation Graph 该算法只适用单个实例资源  
If no cycle exists, then the allocation of the resource will leave the system in a safe state.  
Safe State  
System is in safe state if there exists a sequence <P<sub>1</sub>, P<sub>2</sub>, ..., P<sub>n</sub>> of all the processes is the systems such that for each P<sub>i</sub>, the resources that P<sub>i</sub> can still request can be satisfied by currently available resources + resources held by all the P<sub>j</sub>, with j < i.  
safe state  $\Rightarrow$  no deadlocks  
unsafe state  $\Rightarrow$  possibility of deadlock  
Banker's Algorithm(Deadlock Avoidance)  
Resource-Request Algorithm  
1. If Request  $\leq$  Needi go to step 2. Otherwise, raise error condition, since process has exceeded its maximum claim.

Reader's-Writers Problem  
mutex = 1, wrt = 1  
//readers  
readcount = 0  
wait(mutex);  
readcount++;  
if (readcount == 1) wait(r);  
signal(mutex);  
//reading is performed  
wait(mutex);  
readcount--;  
if (readcount == 0) signal(wrt);  
signal(mutex);  
//writers  
wait(wrt);  
//writing is performed  
signal(wrt);  
Dining-Philosophers Problem  
semaphore chopstick[5] = 1.  
对于第 i 哲学家:  
do {  
wait(chopstick[i]);  
wait(chopstick[(i+1) % 5]);  
eat;  
signal(chopstick[i]);  
signal(chopstick[(i+1) % 5]);  
think;  
} while (1);  
其他解决方法:  
1.allow at most four philosophers to be sitting simultaneously at the table--a spare resource.  
2.allow a philosopher to pick up chopsticks only if both are available, and pick up them simultaneously.  
3.use a asymmetric solution: an odd philosopher picks up first the left chopstick and then the right chopstick, whereas an even one picks up first the right and then the left.  
Deadlock and Starvation  
Deadlock – two or more processes are waiting indefinitely for an event that can be caused by only one of the waiting processes  
Starvation – indefinite blocking. A process may never be removed from the semaphore queue in which it is suspended  
死锁问题  
A set of blocked processes each holding a resource and waiting to acquire a resource held by another process in the set.  
1.Mutual exclusion: only one process at a time can use a resource.  
2.Hold and wait: a process holding at least one resource is waiting to acquire additional resources held by other processes.  
3. No preemption: a resource can be released only voluntarily by the process holding it, after that process has completed its task.  
4. Circular wait: there exists a set {P<sub>1</sub>, P<sub>2</sub>, ..., P<sub>n</sub>} of waiting processes such that P<sub>1</sub> is waiting for a resource that is held by P<sub>2</sub>, P<sub>2</sub> is waiting for a resource that is held by P<sub>3</sub>, ..., P<sub>n-1</sub> is waiting for a resource that is held by P<sub>n</sub>, and P<sub>n</sub> is waiting for a resource that is held by P<sub>1</sub>.

Handling Deadlocks  
1.Prevention 低设备使用率和系统吞吐率  
条件 2: 必须同时申请所有资源(之后不 wait); 或只有没有资源才可申请(之前不 hold)。资源利用率低 饥饿  
条件 3: 申请其它资源不能立即得到, 那么现有资源可被抢占  
适用 cpu 内存 不适用打印机 磁带驱动器  
条件 4: 资源排序, 递增申请  
2. Avoidance  
3. Detection  
4.Recovery  
5.Ignore the problem and pretend that deadlocks never occur in the system (Unix)  
Resource-Allocation Graph 该算法只适用单个实例资源  
If no cycle exists, then the allocation of the resource will leave the system in a safe state.  
Safe State  
System is in safe state if there exists a sequence <P<sub>1</sub>, P<sub>2</sub>, ..., P<sub>n</sub>> of all the processes is the systems such that for each P<sub>i</sub>, the resources that P<sub>i</sub> can still request can be satisfied by currently available resources + resources held by all the P<sub>j</sub>, with j < i.  
safe state  $\Rightarrow$  no deadlocks  
unsafe state  $\Rightarrow$  possibility of deadlock  
Banker's Algorithm(Deadlock Avoidance)  
Resource-Request Algorithm  
1. If Request  $\leq$  Needi go to step 2. Otherwise, raise error condition, since process has exceeded its maximum claim.

Bounded-Buffer Problem  
有限缓冲区问题: 进程共享一块有界的缓冲区。如何实现安全同步?  
采用信号量 full, empty, mutex。  
初始值 full = 0, empty = n, mutex = 1。  
wait(full/EMPTY); CON/PRO  
wait(mutex);  
.....  
signal(mutex);  
signal(empty/FULL); PRO/CON  
.....  
while (1);

Reader's-Writers Problem  
mutex = 1, wrt = 1  
//readers  
readcount = 0  
wait(mutex);  
readcount++;  
if (readcount == 1) wait(r);  
signal(mutex);  
//reading is performed  
wait(mutex);  
readcount--;  
if (readcount == 0) signal(wrt);  
signal(mutex);  
//writers  
wait(wrt);  
//writing is performed  
signal(wrt);  
Dining-Philosophers Problem  
semaphore chopstick[5] = 1.  
对于第 i 哲学家:  
do {  
wait(chopstick[i]);  
wait(chopstick[(i+1) % 5]);  
eat;  
signal(chopstick[i]);  
signal(chopstick[(i+1) % 5]);  
think;  
} while (1);  
其他解决方法:  
1.allow at most four philosophers to be sitting simultaneously at the table--a spare resource.  
2.allow a philosopher to pick up chopsticks only if both are available, and pick up them simultaneously.  
3.use a asymmetric solution: an odd philosopher picks up first the left chopstick and then the right chopstick, whereas an even one picks up first the right and then the left.  
Deadlock and Starvation  
Deadlock – two or more processes are waiting indefinitely for an event that can be caused by only one of the waiting processes  
Starvation – indefinite blocking. A process may never be removed from the semaphore queue in which it is suspended  
死锁问题  
A set of blocked processes each holding a resource and waiting to acquire a resource held by another process in the set.  
1.Mutual exclusion: only one process at a time can use a resource.  
2.Hold and wait: a process holding at least one resource is waiting to acquire additional resources held by other processes.  
3. No preemption: a resource can be released only voluntarily by the process holding it, after that process has completed its task.  
4. Circular wait: there exists a set {P<sub>1</sub>, P<sub>2</sub>, ..., P<sub>n</sub>} of waiting processes such that P<sub>1</sub> is waiting for a resource that is held by P<sub>2</sub>, P<sub>2</sub> is waiting for a resource that is held by P<sub>3</sub>, ..., P<sub>n-1</sub> is waiting for a resource that is held by P<sub>n</sub>, and P<sub>n</sub> is waiting for a resource that is held by P<sub>1</sub>.

Handling Deadlocks  
1.Prevention 低设备使用率和系统吞吐率  
条件 2: 必须同时申请所有资源(之后不 wait); 或只有没有资源才可申请(之前不 hold)。资源利用率低 饥饿  
条件 3: 申请其它资源不能立即得到, 那么现有资源可被抢占  
适用 cpu 内存 不适用打印机 磁带驱动器  
条件 4: 资源排序, 递增申请  
2. Avoidance  
3. Detection  
4.Recovery  
5.Ignore the problem and pretend that deadlocks never occur in the system (Unix)  
Resource-Allocation Graph 该算法只适用单个实例资源  
If no cycle exists, then the allocation of the resource will leave the system in a safe state.  
Safe State  
System is in safe state if there exists a sequence <P<sub>1</sub>, P<sub>2</sub>, ..., P<sub>n</sub>> of all the processes is the systems such that for each P<sub>i</sub>, the resources that P<sub>i</sub> can still request can be satisfied by currently available resources + resources held by all the P<sub>j</sub>, with j < i.  
safe state  $\Rightarrow$  no deadlocks  
unsafe state  $\Rightarrow$  possibility of deadlock  
Banker's Algorithm(Deadlock Avoidance)  
Resource-Request Algorithm  
1. If Request  $\leq$  Needi go to step 2. Otherwise, raise error condition, since process has exceeded its maximum claim.



2.If Requesti ≤ Available, go to step 3. Otherwise Pi must wait, since resources are not available.  
3. Pretend to allocate requested resources to Pi by modifying the state as follows:  
Available = Available - Requesti  
Allocationi = Allocationi + Requesti  
Needi = Needi - Requesti  
If safe ⇒ the resources are allocated to Pi  
If unsafe ⇒ Pi must wait, and the old resource-allocation state is restored

**Safety Algorithm**

1.Let Work and Finish be vectors of length m and n, respectively. Initialize:  
(a) Work = Available  
(b) Finish[i] = false for i = 0, 1, ..., n-1.  
2.Find an index i such that both:  
(a)Finish[i] = false  
(b)Requesti ≤ Work  
If no such i exists, go to step 4.  
3.Work = Work + Allocationi  
Finish[i] = true  
go to step 2.  
4.If Finish[i] = true for all i, then the system is in a safe state.

**Detection Algorithm(死锁检测算法)**

Almost the same with safety algorithm  
1.(b)For i = 1,2, ..., n, if Allocationi ≠ 0, then  
Finish[i] = false; otherwise, Finish[i] = true.  
4.If Finish[i] = false, for some i, 1 ≤ i ≤ n, then the system is in deadlock state.  
Moreover, if Finish[i] = false, then Pi is deadlocked.  
The algorithm requires an order of O(m × n<sup>2</sup>) operations to detect whether the system is in deadlock state)

Wait for 图算法: 简化了资源分配图, 只适用单个实例资源

**内存管理**

**Binding 地址绑定**  
Addresses in the program are generally symbolic.  
A compiler binds these symbolic addresses to relocatable addresses.  
A linkage editor (loader) binds these relocatable addresses to absolute addresses.  
1. Compile time: If memory location known a priori, absolute code can be generated; must recompile code if starting location changes.  
2. Load time: Must generate relocatable code if memory location is not known at compile time.  
3. Execution time: Binding delayed until run time if the process can be moved during its execution from one memory segment to another. Need hardware support for address maps (e.g., base and limit registers).

**Logical vs. Physical Address Space**

Central to proper memory management  
1.logical address – generated by the CPU; also referred to as virtual address  
2.physical address – address seen by the memory unit  
Differ in execution-time address-binding scheme and same in compile-time and load-time.

**Memory-Management Unit(内存管理单元)**

MMU is a hardware device that maps virtual to physical address.  
In MMU scheme, the value in the relocation register is added to every address generated by a user process at the time it is sent to memory.  
The user program deals with logical addresses; it never sees the real physical addresses.

**Dynamic Linking**

linking postponed until execution time. Particularly useful for (shared) libraries; better memory usage; easier library updates; operating system need to check if routine is in processes' memory address.

**Swapping**

Major part of swap time is transfer time  
Backing store – fast disk large enough

**Overlay**

An overlay is when a process replaces itself with the code of another program.  
法 1: 不能换出有待处理 I/O 的进程  
法 2: IO 只能操作系统缓冲区

**1.Contiguous Allocation**

Base register contains value of smallest physical address.

Limit register contains range of logical addresses.

Logical address must be less than limit register

MMU maps logical address dynamically

**Dynamic Storage-Allocation Problem**

**First-fit:** Allocate the first hole that is big enough  
**Best-fit:** Allocate the smallest hole that is big enough; must search entire list, unless ordered by size; Produces the smallest leftover hole

**Worst-fit:** Allocate the largest hole; must also search entire list. Produces the largest leftover hole.  
first fit and best fit are better than worst fit. first fir a little better than best fit.

**Fragmentation**

**External Fragmentation** – total memory space exists to satisfy a request, but it is not contiguous  
**Internal Fragmentation** – allocated memory may be slightly larger than requested memory; this size difference is memory internal to a partition, but not being used.  
Compaction is possible only if relocation is dynamic, and is done at execution time.

**2.Paging**

Divide physical memory into fixed-sized blocks called frames. Size is power of 2, 512B~ 8,192B).  
Divide logical memory into blocks of same size called pages.  
Page table: translate logical to physical addresses  
**Page number (p)** – used as an index into a page table which contains base address of each page in physical memory  
**Page offset (d)** – combined with base address to define the physical memory address that is sent to the memory unit

page number	page offset
$p$	$d$
$m - n$	$n$

**Implementation of Page Table**

Page table is kept in main memory.  
**Base register (PTBR)** points to the page table  
**Length register (PTLR)** indicates size of the table.  
Every data/instruction access requires two memory accesses. One for the page table and one for the data/instruction.  
Can be solved by the use of a special fast-lookup hardware cache called associative memory or translation look-aside buffers (TLBs).  
 $Base = (t \times m) + (t + r + d) \times (1 - m)$

**Memory Protection - Valid-invalid bit**

**Shared Pages**  
Shared code must appear in same location in the logical address space of all processes.

**Structure of the Page Table**

1.Hierarchical Paging (Two-Level Paging)  
Page-number(2-level) + Page-offset  
2.Hashed Page Tables  
The virtual page number is hashed into a page table. This page table contains a chain of elements hashing to the same location.  
3.Inverted Page Tables

**3.Segmentation**

**Segment table** – maps 2d physical addresses. 动态链接  
Logical address consists <segment-number, offset>. Each table entry has:

**base** – contains the starting physical address where the segments reside in memory

**limit** – specifies the length of the segment

**STBR and STLR**

Allocation: segments vary in length → dynamic storage-allocation problem ; first fit/best fit ; external fragmentation  
Sharing : same segment number

**Segmentation and paging**

Segmented paging is helpful when the page table becomes very large. A large contiguous section of the page table that is unused can be collapsed into a single segment table entry with a page-table. Paged segmentation handles the case of having very long segments that require a lot of time for allocation. By paging the segments, we reduce wasted memory due to external fragmentation and simplify the allocation.

**虚拟存储**

**Virtual memory** – separation of user logical memory from physical memory. 当物理内存较小时提供一个较大的外部虚拟存储  
1.Only part of the program needs to be in memory for execution. 2.logical address space can therefore be much larger than physical address space. 3.allows address spaces to be shared by several processes. 4.allows for more efficient process creation. 5.virtual memory allows other benefits during process creation: copy-on-write ; memory-mapped file.

**Demand Paging**

Bring a page into memory only when it is needed  
1.Less I/O needed 2.Less memory needed  
3.Faster response 4.More users

**Page Fault**

1.Operating system looks at another table to decide:(1)Invalid reference ⇒ abort (2)Just not in memory 2.Get empty frame 3.Swap page into frame 4.Reset tables 5.Set validation bit = v 6.Restart the instruction that caused the page fault

**Page Fault Rate**  $0 \leq p \leq 1.0$ .  $EAT = (1 - p) \times \text{memory access} + p \times \text{page fault overhead} + \text{swap page out} + \text{swap page in} + \text{restart overhead}$

**Copy-on-Write (COW)** allows both parent and child processes to initially share the same pages in memory until either process modifies. More efficient process creation.

**Page Replacement Algorithm 页置换算法**

Goal - lowest page-fault rate.

**Reference string:** page numbers only, with adjacent duplications eliminated.

1.**First-in-First-Out Algorithm** - use a FIFO queue to hold all pages in memory. When a page is brought into memory, insert it at the tail. When a free frame is needed, we replace the page at the queue. Belady's Anomaly  
2.**Optimal** Page Replacement OPT/MIN- replace the page that won't be used for the longest period of time  
3.**Least Recently Used** - associates with each page the time of that page's last use. When a page must be replaced, choose the page that has not been used for the longest period of time.

a. Counter implementation - A time-of-use field. Copy the clock into it when referenced. Look for smallest time-of-use field. Difficulties: requires a search of the page table; a write to memory for each memory access; page-table maintain in context switch; overflow of the clock.  
b. Stack implementation - Keep a stack of page numbers in a double linked list. Move a page to the top when referenced. Requires 6 pointers to be changed. Always replace the page at the bottom. No search for replacement.  
PS: both need special hardware support(TLB).

**4.LRU-Approximation Page Replacement**

Reference Bit - with each page associate a bit, initially = 0, when page is referenced, bit set to 1.  
a. Additional Reference Bits  
Use right-shift history byte for each page, current reference bit shift into the left-most bit. choose the page with lowest number.

**b. Second-Chance Algorithm**

FIFO + reference bit. Circular Queue. Inspect the current frame, if the reference bit is set, reset it, and skip to next frame; otherwise, replace it.  
c. Enhanced Second-Chance Algorithm  
Modify bit (dirty, bit). Replace the first page encountered in the lowest nonempty class. Drawback: may have to scan several times. (Mac)  
5.Counting-based Algorithms - keep a counter of number of references to each page.

a. Least-frequently-used (LFU) - page with the smallest count is replaced. counter shift right at regular interval , exponentially decaying average.  
b. Most-frequently-used (MFU) - based on the argument that page with the smallest count was probably just brought in and has yet to be used.

**Frame Allocation**

1.Fixed Allocation – Equal or Proportional(比例).  
2.Priority Allocation  
Global vs Local Replacement  
Global(more common used) - process selects from the set of all frames. High priority process can take a frame from a lower priority process. Results in greater system throughput.  
Problem: a process cannot control its own page-fault rate.  
Local - process selects from only its own set of allocated frames. the number of frames allocated to a process does not change.

**Thrashing & Working-Set Model**

If a process doesn't have "enough" pages, the page-fault rate is very high ⇒ 1.Low CPU utilization 2.os thinks it needs to increase the degree of multiprogramming 3.another process added to the system 4.worse the condition.  
Thrashing ⇒ A process is busy swapping pages in and out. Σ size of locality > total memory size, thrashing occur. Δ = working-set window = a fixed number of page references. Working-Set Size  
WSS = total number of pages referenced in the most recent Δ (varies in time). Δ too small will not encompass entire locality. A too large will encompass several localities. Δ = ∞ will encompass entire program.  
D = Σ WSSi = total demand frames. If D > m ⇒ thrashing. Policy: if D > m, then suspend one of the processes (and swap it out completely).  
Working-set strategy prevents thrashing while keeping the degree of multiprogramming as high as possible, it optimizes CPU utilization.

Implementation: approximate with interval timer + a reference bit.

Example: Δ = 10,000, timer interrupts every 5000 references, keep in memory 2 bits for each page. Whenever a timer interrupts, copy and sets the values of all reference bits to 0. If one of the bits in memory = 1 ⇒ page in working set

**Memory-Mapped Files**

Memory-mapped file I/O allows file I/O to be treated as routine memory access by mapping a disk block to a page in memory. 1.Simplifies file access by tring file I/O through memory rather than read() write() system calls. 2.Allows several processes to

map the same file - pages in memory to be shared.

**Allocating Kernel Memory**

1.**Buddy System** - Allocates memory from fixed-size segment consisting of physically contiguous pages. (用高位地址判断是否 Buddy)  
2.**Slab Allocator - Slab** is one or more physically contiguous pages. Cache consists of one or more slabs. Single cache for each unique kernel data structure. When cache created, filled with objects marked as free. When structures stored, objects marked as used. If slab is full of used objects, next object allocated from empty slab. If no empty slabs, new slab allocated. Benefits include no fragmentation, fast memory request satisfaction.

**文件系统接口**

**File Attributes**

1.Name 2.Identifier – unique tag identifies file within file system 3.Type 4.Location 5.Size 6.Protection 7.Time, date, and user identification  
8.Information about files kept in the directory structure, which is maintained on the disk.  
**File Operations** - 1.Create 2.Write 3.Read 4.Reposition within file 5.Delete 6.Truncate 截短  
**Open-file table**  
a. file can be opened simultaneously by more than one processes in name of several users  
b. per-process table tracks all files that a process has open. May have a pointer to an entry in the system-wide table.  
c. a system-wide table contains process-independent information. Have an open count tracks the number of processes.  
info associated with an open file: 1.File pointer 2.File-open count 3.Disk location of the file 4.Access rights

**File Types – Name, Extension**

**Access Methods - Sequential & Direct Access**

**File Directory**

A directory is a collection of nodes containing information about all files. Both the directory structure and the files reside on disk.  
Information in dir: name, type, address , current length, maximum length, date last accessed (for archival), date last updated (for dump), owner ID (who pays), protection information.  
Operations: 1.Search 2.Create 3.Delete 4.List a directory  
5.Rename 6.Traverse the file system

**Directory Structure**

Goal: 1.Efficiency 2.Naming 3.Grouping  
1.Single-Level (2/3)X limited name length  
2.Two-Level (3)X  
3.Tree-Structured Absolute or relative path name  
4.Acyclic-Graph have shared subdirectories and files. A file may have multiple absolute path names. Distinct file names may refer to the same

Delete file ⇒ dangling pointer. To solve, using (1)Backpointers using a daisy chain (2) Every-entry-count. 保证无圈方法:a.只允许指向文件的链接;b.每次增加link时检查有无圈.  
5.General Graph – Need garbage collection

**File System Mounting**

A file system must be mounted before it can be accessed. An unmounted file system is mounted at a mount point.A directory structure can be built out of multiple partitions.

**File Sharing**

Desirable on multi-user systems. May be done through a protection scheme.  
1)User IDs and Group IDs  
2)Remote File Systems: NFS, CIFS

**File Protection**

Mode of access: read, write, execute  
three classes of users  
a) owner access 7 ⇒ 1 1 1  
b) group access 6 ⇒ 1 1 0  
c) public access 1 ⇒ 0 0 1

**文件系统实现**

**File-System Structure**

File structure goal:  
1.Logical storage unit - data blocks  
2.Collection of related information - FCB  
File system organized into layers

devices I/O control basic file system file-organization module logical file system application programs

File system resides on secondary storage

Device driver controls the physical device

**File-System Implementation**

On disk:

1.Boot control block contains info needed by system to boot OS from that volume.  
2.Boot block(UFS) & partition boot sector(NTFS)  
3.Volume control block contains volume details  
# of blocks, free block count and pointers, free FCB count and pointers)  
superblock(UFS) & master file table(NTFS)  
3.Directory structure organizes the files.  
node numbers(UFS) & master file table(NTFS)  
4.Per-file File Control Block (FCB) contains many details about the file.  
n memory

1.in-memory mount table  
2. in-memory directory structure

3. system-wide open-file table - contains a copy of the FCB of each open files, open count, etc.  
4.per-process open-file table - contains a pointer to the entry in the system-wide open-file table, read/write position, etc.

**File Control Block**

1.file permissions 2.file dates(create, access, write) 3.file owner, group, ACL 4. file size 5. file data blocks or pointers to file data blocks.

**File Open**

1.the directory structure is searched for the given file name 2.if found, the FCB is copied into a system-wide open-file table in memory 3.an entry is made in the per-process open-file table 4.returns a pointer to the entry in the per-process open-file table – file descriptor in UNIX, file handle in Windows

**File Close**

1.the per-process table entry is removed. 2.the system-wide table entry's open count is decremented. If the count hits zero, all updated info is copied back to the disk, and the system-wide table entry is removed.

**Virtual File Systems**

Provides an object-oriented way of implementing file systems. Allows the same system call interface(API) to be used for different types of file systems. The API is to the VFS interface, rather than any specific type of file system.  
Three Layers: 1. file-system interface – based on open, close, read, write and on file descriptors.  
2. virtual file system.

**Directory Implementation**

1.Linear list of file names with pointer to the data blocks 2.Hash table – linear list with hash data structure

**Allocation Methods 磁盘分配方法**

**Contiguous allocation**

Benefits: simple, only starting location and length (number of blocks) are required; good support of random access.  
Problem: dynamic storage-allocation - first fit, best fit; external fragmentation – compaction.  
Difficulty: files grow, difficult to determine how large a file.  
Extent-Based Systems - a modified contiguous allocation scheme. It allocate disk blocks in extents; an extent is a contiguous block on the disk; extents are allocated for file allocation; a file consists of one or more extents.  
LA/512 Block X + start address Displacement Y

**Linked allocation**

Benefits: simple, need only starting address in directory entry; no external fragmentation.  
Drawbacks: no random access, can be used effectively only for sequential-access files; space required for the pointers in every block - use clusters(簇) 按簇(几个块一起)分配,减少指针消耗,但会增加内部碎片  
Reliability: pointers scatter all over the disk, and may lost or damage - double FAT!  
LA/511 Block X Displacement Y+1

**Indexed allocation**

Brings all pointers into the index block; directory entry contains the address of the index block.  
Benefits: support random access; no external fragmentation  
Drawbacks: waste space(whole block may contains one or two pointers)

Index table of variant size files  
1.Linked scheme 2.Multilevel index 3.combined scheme(12+256+256\*2+256\*3)  
LA/512 X - displacement in index table  
Y=displacement in block

**File Allocate Table(FAT)**

Disk-space allocation used by MS-DOS and OS/2  
A section at the beginning of each partition is set aside to contain the table. A FAT is used as a linked list. 0 means free block; a special value (-1) means end-of-file; double FAT to increase the reliability. FAT cause a lot of disk head seeks.  
FAT is usually cached in memory; random access time is improved, we can find the location of any block by reading the

FAT (provided it is cached in memory). 类似于游标的一种实现。

**Free-Space Management 空闲空间管理**

1.bit map, or bit vector (n blocks)  
Each block is represented by 1 bit(1 free 0 occupied). Bit map requires extra space.  
First free block number calculation=(number of bits per word) \* (number of 0-value words) + offset of first 1 bit  
Ex: block size = 512B disk size = 1GB n = 1GB / 512B = 2M = 256K clustering blocks in group (say, of four blocks) reduces the number to 64KB.  
Easy to get contiguous files  
2. Linked Free-space List  
Cannot get contiguous space easily; no waste of space; an FAT implementation may help to improve performance; FIFO or stack?

**磁盘系统**

Transfer rate is rate at which data flow between drive and computer.  
Positioning time (random-access time) is time to move disk arm to desired cylinder (seek time) and time for desired sector to rotate under the disk head (rotational latency)  
Drive attached to computer via I/O bus

1.Busses vary, including EIDE, ATA, SATA, USB, Fibre Channel, SCSI.  
2.Host controller in computer uses bus to talk to disk controller built into drive or storage array.

**Moving-head Disk Mechanism**

A read-write head "flies" just above each surface of every platter. The heads are attached to a disk arm that moves all the heads as a unit. The surface of a platter is logically divided into circular tracks, which are subdivided into sectors. The set of tracks that are at one arm position makes up a cylinder.

**Disk Structure**

Disk drives are addressed as large 1-dimensional arrays of logical blocks, where the logical block is the smallest unit of transfer - typical called a sector, 512B or 1KB  
The 1-dimensional array of logical blocks is mapped into the sectors of the disk sequentially.  
1.sector 0 is the first sector of the first track on the outermost cylinder.  
2.mapping proceeds in order through that track, then the rest of the tracks in that cylinder, and then through the rest of the cylinders from outermost to innermost.

**Disk Scheduling**

The operating system is responsible for using hardware efficiently – for the disk drives, this means having a fast access time and disk bandwidth.  
Access time has two major components  
1.Seek time is the time for the disk are to move the heads to the cylinder containing the desired sector. 2.Rotational latency is the additional time waiting for the disk to rotate the desired sector to the disk head  
Minimize seek time  
Seek time = seek distance  
Disk bandwidth is the total number of bytes transferred, divided by the total time between the first request for service and the completion of the last transfer

**Disk Scheduling Algorithm 磁盘调度算法**

1.FCFS – First Come First Service  
2.SSTF – 选寻道时间最短的。May cause starvation  
3.SCAN(elevator algorithm) The disk arm starts at one end of the disk, and moves toward the other end, servicing requests until it gets to the other end of the disk, where the head movement is reversed and servicing continues. 扫到头  
4.C-SCAN – 基于 SCAN. 磁盘是圆形的, 从一头倒另一头很快且不执行 JOB. 扫到头, 快速扫回  
5.LOOK 扫过头的 C 最近的位置, 不用到头  
6.C-LOOK 快速扫回, 注意不是环形, 只是快速扫回, 中途不折活, 比如设备只能单向存取。

**Disk Management**

1.Low-level formatting, or physical formatting – Dividing a disk into sectors that the disk controller can read and write.  
2.To use a disk to hold files, the operating system still needs to record its own data structures on the disk. a. Partition the disk into one or more groups of cylinders. b. Logical formatting or "making a file system" c. To increase efficiency most file systems group blocks into clusters. (1)Disk I/O done in blocks. (2)File I/O done in clusters

**Master Boot Record**

**Raid Structure**

RAID – multiple disk drives provides reliability via redundancy. increases the mean time to failure  
Frequently combined with NVRAM(Non-volatile RAM) to improve write performance. RAID is arranged into six different levels.  
RAID 0: non-redundant striping.

RAID 1: mirrored disks.  
RAID 2: memory-style error correcting codes.  
RAID 3: bit-interleaved parity.  
RAID 4:block-interleaved parity.  
RAID 5: block-interleaved distributed parity.  
RAID 6: P + Q redundancy.

I/O 系统

An I/O channel is a special processor.  
Virtual devices are implemented with **spoofing!**

Basic I/O Concepts

device drivers—present a uniform device access interface to the I/O subsystem  
I/O subsystem – separates the rest of the kernel from the complexity of managing I/O devices  
**Port** – a device communicates with the machine via a connection point, e.g. a serial port  
**Bus** – a set of wires and a **protocol**  
**Controller**– a collection of electronics that operate a port, a bus, or a device: the processor communicates with the controller by reading and writing bit patterns in data and control registers; **Special I/O instruction or Memory-mapped I/O.**

Polling 轮询

Determines state of device. Command-ready; busy; error.  
Inefficient when speed differ much.

Interrupt 中断

1.CPU interrupt request line triggered by I/O device.  
2.**Interrupt handler** receives interrupts.  
3. Maskable to ignore or delay some interrupts  
Two lines: Nonmaskable and maskable interrupt.  
4.Interrupt vector to dispatch interrupt to correct handler: based on priority; Some nonmaskable.  
5.Interrupt mechanism also used for exceptions.

DMA(Direct Memory Access)

1.Used to avoid programmed I/O for large data movement  
2.requires DMA controller- a special-purpose processor.  
3.Bypasses CPU to transfer data directly between I/O device and memory. Handshake between the main processor and the DMA controller: two wires - DMA-request wire and DMA-acknowledge wire

Application I/O Interface I/O 应用接口

1.I/O system calls encapsulate device behaviors in generic classes. 2.Device-driver layer hides differences among I/O controllers from kernel.  
Block and Character Devices    Network Devices  
Clocks and Timers Blocking and Nonblocking I/O

Kernel I/O Subsystem

1.I/O scheduling    Device-status Table

2.Buffering

3.Caching

4.Spooling and Device Reservation

5.error handling

6.I/O Protection

I/O port registers

1.status register    2.control register

3.data-in register(s) 4.data-out register(s)

Kernel I/O subsystem summary

1. management of name of files and devices  
2. access control to files and devices  
3. file system space allocation  
4. device allocation (including reservation)  
5. buffering, caching, spooling  
6. I/O scheduling  
7. device status monitoring, error handling  
8. device driver configuration

logical address

CPU

pid p d

search

pid p

page table

physical address

physical memory

logical address

CPU

s d

segment table

trap: addressing error

physical memory

operating system

reference

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92

93

94

95

96

97

98

99

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

280

281

282

283

284

285

286

287

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

330

331

332

333

334

335

336

337

338

339

340

341

342

343

344

345

346

347

348

349

350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

382

383

384

385

386

387

388

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

777

778

779

780

781

782

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

999

1000

logical address

CPU

pid p d

search

pid p

page table

physical address

physical memory

logical address

CPU

s d

segment table

trap: addressing error

physical memory

operating system

reference

1

2

3

4

5

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44

45

46

47

48

49

50

51

52

53

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92

93

94

95

96

97

98

99

100

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

280

281

282

283

284

285

286

287

288

289

290

291

292

293

294

295

296

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

330

331

332

333

334

335

336

337

338

339

340

341

342

343

344

345

346

347

348

349

350

351

352

353

354

355

356

357

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

382

383

384

385

386

387

388

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443</